

# Distributionally robust risk evaluation with an isotonic constraint

Yu Gui, Rina Foygel Barber, and Cong Ma

Department of Statistics, University of Chicago

September 4, 2024

## Abstract

Statistical learning under distribution shift is challenging when neither prior knowledge nor fully accessible data from the target distribution is available. Distributionally robust learning (DRL) aims to control the worst-case statistical performance within an uncertainty set of candidate distributions, but how to properly specify the set remains challenging. To enable distributional robustness without being overly conservative, in this paper we propose a shape-constrained approach to DRL, which incorporates prior information about the way in which the unknown target distribution differs from its estimate—specifically, we assume the unknown density ratio between the target distribution and its estimate is isotonic with respect to some partial order. At the population level, we provide a solution to the shape-constrained optimization problem that does not involve the isotonic constraint. At the sample level, we provide consistency results for an empirical estimator of the target in a range of different settings. Empirical studies on both synthetic and real data examples demonstrate the improved efficiency of the proposed shape-constrained approach.

## 1 Introduction

Evaluating the performance of an estimator is of significant importance in statistics. To give several motivating examples, we first consider supervised learning settings, where our observations consist of features  $X \in \mathcal{X} \subseteq \mathbb{R}^d$  and a response  $Y \in \mathcal{Y} \subseteq \mathbb{R}$ :

- Given a fitted model  $\hat{\mu} : \mathcal{X} \rightarrow \mathbb{R}$ , we may want to estimate the expected value of the squared error  $(Y - \hat{\mu}(X))^2$  with respect to target distribution on  $(X, Y)$ .
- Or, in predictive inference, suppose we have constructed a prediction band  $\hat{C}_{1-\alpha}$ , where  $\hat{C}_{1-\alpha}(X) \subseteq \mathbb{R}$  is a confidence region for the response  $Y$  given features  $X$ , and  $1 - \alpha$  denotes the target coverage level. Then to determine whether  $\hat{C}_{1-\alpha}$  does in fact achieve coverage at level  $1 - \alpha$  for data points drawn from some target distribution, we would like to estimate the the expected value of  $\mathbb{1}\{Y \notin \hat{C}_{1-\alpha}(X)\}$  with respect to this target distribution. This is exactly the probability that our interval *fails* to cover the response.

We can also consider unsupervised learning settings, with observations  $X \in \mathcal{X} \subseteq \mathbb{R}^d$  only:

- In principal component analysis (PCA), suppose we have obtained a set of pre-fitted principal components  $\hat{V}_K = \{\hat{v}_1, \dots, \hat{v}_K\}$  which forms an orthonormal basis for a  $K$ -dimensional subspace of  $\mathbb{R}^d$ . To evaluate how well the variance in  $X$  is explained by the principal components, it would be of interest to analyze the expected value of the reconstruction error  $\|X - \sum_{k=1}^K (X^\top \hat{v}_k) \hat{v}_k\|^2$  with respect to the distribution of  $X$ .

- Another example is density estimation. In this case, given a density estimate  $P_\theta$  learned from data, we may want to evaluate its performance using the expected log-likelihood  $-\log dP_\theta(X)$  over a target distribution  $P_{\text{target}}$ . In fact,  $\mathbb{E}_{P_{\text{target}}}[-\log dP_\theta(X)]$  is the cross-entropy of  $P_\theta$  relative to  $P_{\text{target}}$ .

A key challenge for any of these problems is that the target distribution (say, the distribution of the general population) may be unknown, and our available data (say, individuals who participate in our study) may be drawn from a different distribution than the general population.

## 1.1 Problem formulation

To make the problem more concrete, and unify the examples mentioned above, here we introduce some notation to formulate the question at hand.

**The unsupervised setting.** Let  $R : \mathcal{X} \rightarrow \mathbb{R}_+$  denote a *risk function*, where our goal is to evaluate the expected value  $\mathbb{E}_{P_{\text{target}}}[R(X)]$  with respect to some target distribution  $P_{\text{target}}$  on  $X$ . However, the available data only provides information about  $P$ , a potentially different distribution.

For instance, in density estimation, after obtaining the density estimate  $P_\theta$ , we can then estimate  $\mathbb{E}_P[R(X)]$  using a calibration data set, which consists of samples  $X_1, \dots, X_n$  drawn from  $P$ . Instead, our aim is to provide a bound on the risk  $\mathbb{E}_{P_{\text{target}}}[R(X)]$ , or in other words, to bound the difference in risks (often called the *excess risk*),  $\mathbb{E}_{P_{\text{target}}}[R(X)] - \mathbb{E}_P[R(X)]$ . If we assume that the unknown distribution  $P_{\text{target}}$  lies in some class  $\mathcal{Q}$  (to be specified later on), then defining the *worst-case excess risk*

$$\Delta(R; \mathcal{Q}) = \sup_{Q \in \mathcal{Q}} \mathbb{E}_Q[R(X)] - \mathbb{E}_P[R(X)], \quad (1)$$

we can then bound

$$\mathbb{E}_{P_{\text{target}}}[R(X)] \leq \mathbb{E}_P[R(X)] + \Delta(R; \mathcal{Q}).$$

The right hand side provides an upper bound on the risk of our estimator under the target distribution  $P_{\text{target}}$ .

**The supervised setting: covariate shift assumption.** In the supervised learning setting, the data contains both features  $X$  and a response  $Y$ , so the setup is somewhat different. Here we will consider a loss function  $r : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$ , for instance,  $r(x, y) = (y - \hat{\mu}(x))^2$  for the squared error in a regression, or  $r(x, y) = \mathbb{1}\{y \notin \hat{C}_{1-\alpha}(x)\}$  for characterizing the (mis)coverage of a prediction interval in predictive inference.

Throughout this paper, for the supervised learning setting, we will assume the *covariate shift* setting, where the distribution of the available data and the target distribution may differ in the marginal distribution of the covariates  $X$ , but share the same conditional distribution  $Y | X$ . To make this concrete, if our calibration data consists of  $n$  data points  $(X_1, Y_1), \dots, (X_n, Y_n)$  drawn from  $\tilde{P}$ , while our goal is to control the expected loss with respect to the target distribution  $\tilde{P}_{\text{target}}$  on  $(X, Y)$ , we will assume that we can write

$$\begin{aligned} \text{data distribution: } \tilde{P} &= P \times P_{Y|X}, \\ \text{target distribution: } \tilde{P}_{\text{target}} &= P_{\text{target}} \times P_{Y|X}, \end{aligned}$$

so that  $\tilde{P}$  and  $\tilde{P}_{\text{target}}$  share the same conditional distribution  $P_{Y|X}$  for  $Y | X$ .

In fact, under covariate shift, this supervised setting can be unified with the unsupervised one by defining the risk

$$R(X) = \mathbb{E}[r(X, Y) | X],$$

which is the conditional expectation of  $r(X, Y)$  under *either*  $\tilde{P}$  or  $\tilde{P}_{\text{target}}$ , since we have assumed that both of these distributions share the same conditional distribution,  $P_{Y|X}$ . The quantity of interest is then given by  $\mathbb{E}_{P_{\text{target}}}[R(X)] = \mathbb{E}_{\tilde{P}_{\text{target}}}[r(X, Y)]$ , but our calibration data, which is sampled from  $P$ , instead enables us

to estimate  $\mathbb{E}_P[R(X)] = \mathbb{E}_{\tilde{P}}[r(X, Y)]$ . If we again assume that  $P_{\text{target}} \in \mathcal{Q}$ , then  $\Delta(R; \mathcal{Q})$  again allows us to bound the risk of our estimator under the target distribution, which is now given by  $\tilde{P}_{\text{target}}$ :

$$\mathbb{E}_{\tilde{P}_{\text{target}}} [r(X, Y)] \leq \mathbb{E}_{\tilde{P}} [r(X, Y)] + \Delta(R; \mathcal{Q}).$$

**Estimating the error or tuning the model?** In this paper, we consider the setting where our estimator—say, a prediction band  $\hat{C}_{1-\alpha}$ —is *pretrained*, meaning that we have available calibration data sampled from  $P$  (in the unsupervised setting) or  $\tilde{P}$  (in the supervised setting) that is independent of the fitted estimator. Consequently, our available calibration data provides us with an unbiased estimate of  $\mathbb{E}_P[R(X)]$  (or, equivalently in the supervised setting,  $\mathbb{E}_{\tilde{P}}[r(X, Y)]$ ).

In some settings, the goal may be to estimate the risk of each estimator within a family of (pretrained) options, in order to tune the choice of estimator. Returning again to the example of a prediction band, suppose, with any confidence level  $1 - a \in [0, 1]$ , we actually are given a nested family of prediction bands,  $\{\hat{C}_{1-a} : a \in [0, 1]\}$ . Choosing  $R_a(X) = \mathbb{P}_{P_{Y|X}}(Y \notin \hat{C}_{1-a}(X))$  or accordingly,  $r_a(X, Y) = \mathbb{1}\{Y \notin \hat{C}_{1-a}(X)\}$ , then, if we can compute a bound on the miscoverage rate  $\mathbb{E}_{P_{\text{target}}}[R_a(X)]$  of  $\hat{C}_{1-a}$  relative to the target distribution for each  $a$ , then we can choose a value of  $a$  that achieves some desired level of coverage. More generally, we may do the same in other settings as well—that is, given a family of candidate estimators, bounding the risk of each one under the target distribution  $P_{\text{target}}$  provides an intermediate step towards choosing the tuning parameter.

## 1.2 Prior work: distributionally robust learning

Our work builds upon the distributionally robust learning (DRL) literature (Ben-Tal and Nemirovski, 1998; El Ghaoui et al., 1998; Lam, 2016; Duchi and Namkoong, 2018), which is a well established framework for risk evaluation under distribution shift. In this framework, the target distribution  $P_{\text{target}}$  is assumed to lie in some neighborhood around the distribution  $P$  of the available data—for instance, we might assume that  $D_{\text{KL}}(P_{\text{target}} \| P) \leq \rho$ , where  $D_{\text{KL}}$  denotes the Kullback–Leibler (KL) divergence. DRL takes a conservative approach and evaluate the performance on  $P_{\text{target}}$  via its upper bound, i.e., the worst-case performance over all distributions within the specified neighborhood of  $P$ ,

$$\mathbb{E}_{P_{\text{target}}}[R(X)] \leq \sup \{ \mathbb{E}_Q[R(X)] : D_{\text{KL}}(Q \| P) \leq \rho \}. \quad (2)$$

Equivalently, we can write this upper bound as

$$\mathbb{E}_{P_{\text{target}}}[R(X)] \leq \mathbb{E}_P[R(X)] + \Delta(R; \mathcal{Q}_{\text{KL}}(\rho)),$$

where  $\Delta(R; \mathcal{Q}_{\text{KL}}(\rho))$  is defined as in (1) above by defining the constraint set as  $\mathcal{Q} = \mathcal{Q}_{\text{KL}}(\rho) = \{Q : D_{\text{KL}}(Q \| P) \leq \rho\}$ . More generally, we can consider divergence measures beyond the KL distance, as we will describe in more detail below.

## 1.3 Our proposal: iso-DRL

If the assumption  $D_{\text{KL}}(P_{\text{target}} \| P) \leq \rho$  is correct, then the upper bound (2) is valid. However, since this bound uses only the KL divergence to define the constraint  $P_{\text{target}} \in \mathcal{Q}$  on the target distribution, it could be quite conservative. In many practical settings, additional side information or prior knowledge on the structure of the distribution shift may allow for a tighter bound, which would be less conservative than the worst-case excess risk of DRL (2). This raises the following key question:

*Can we use side information on the distribution shift between the data distribution  $P$  and the target distribution  $P_{\text{target}}$ , to improve the worst-case excess risk of DRL in risk evaluation?*

In this paper, we study one specific example of this type of setting: we assume that the density ratio  $\frac{dP_{\text{target}}}{dP}(x)$  between the target distribution and the data distribution is isotonic (i.e., monotone) with respect to some order or partial order on  $\mathcal{X}$ .

**Motivation: recalibration of an estimated density ratio.** To motivate the use of such side information, consider a practical supervised setting where we have an initial estimate  $w_0$  for the density ratio:

$$w_0(x) \approx \frac{dP_{\text{target}}}{dP}(x).$$

This is possible in addition to labeled data (i.e.,  $(X, Y)$  pairs) sampled from the data distribution  $P \times P_{Y|X}$ , we also have access to unlabeled (i.e.,  $X$  only) data from the target population  $P_{\text{target}}$ . We may use these two data sets to train  $w_0$ . Although there is no guarantee that the estimate  $w_0$  is accurate, the shape or relative magnitude of  $w_0$  may provide us with useful side information: large values of  $w_0$  can identify portions of the target population that are *underrepresented* under the data distribution  $P$ . This motivates us to recalibrate  $w_0$  within the set of density ratios that are isotonic in  $w_0$ .

To express this scenario in the notation of the problem formulation above, we assume that the target distribution  $P_{\text{target}}$  satisfies an isotonicity constraint,  $P_{\text{target}} \in \mathcal{Q}_{\text{iso}}(w_0)$ , where

$$\mathcal{Q}_{\text{iso}}(w_0) = \left\{ Q : \frac{dQ}{dP}(x) \text{ is a monotonically nondecreasing function of } w_0(x) \right\}.$$

If we assume as before that the target distribution  $P_{\text{target}}$  satisfies  $D_{\text{KL}}(P_{\text{target}} \| P) \leq \rho$ , then we can bound

$$\mathbb{E}_{P_{\text{target}}}[R(X)] \leq \mathbb{E}_P[R(X)] + \Delta(R; \mathcal{Q}_{\text{KL}}(\rho) \cap \mathcal{Q}_{\text{iso}}(w_0)). \quad (3)$$

**The benefits of iso-DRL.** What are the benefits of iso-DRL, as compared to the existing DRL framework? Of course, thus far the idea is quite straightforward—if we have stronger constraints on  $P_{\text{target}}$ , then we can place a tighter bound on the excess risk  $\mathbb{E}_{P_{\text{target}}}[R(X)] - \mathbb{E}_P[R(X)]$ . But as we will see below, adding the isotonic constraint plays a crucial role in enabling DRL to provide bounds that are useful in practical scenarios. Specifically, consider a practical setting where the bound  $\rho$  on the distribution shift is a positive constant. As we will see below, the existing worst-case excess risk  $\Delta(R; \mathcal{Q}_{\text{KL}}(\rho))$  of DRL is often quite large, leading to extremely conservative statistical conclusions; in contrast, the worst-case excess risk  $\Delta(R; \mathcal{Q}_{\text{KL}}(\rho) \cap \mathcal{Q}_{\text{iso}}(w_0))$  given by iso-DRL is often vanishingly small, leading to much more informative conclusions. Moreover, surprisingly, this improvement in the bound does not incur any additional computational challenges—even though the constraint set  $\mathcal{Q}_{\text{KL}}(\rho) \cap \mathcal{Q}_{\text{iso}}(w_0)$  appears more complex than the original set  $\mathcal{Q}_{\text{KL}}(\rho)$ , we will see that  $\Delta(R; \mathcal{Q}_{\text{KL}}(\rho) \cap \mathcal{Q}_{\text{iso}}(w_0))$  can be computed as easily as the original  $\Delta(R; \mathcal{Q}_{\text{KL}}(\rho))$ . In addition, we further show in the appendix that the worst-case excess risk of iso-DRL can be consistently estimated with noisy observations of  $R(X)$ , while the estimation of the worst-case excess risk of DRL can be challenging even with bounded risks.

**Example: predictive inference for the wine quality dataset.** To illustrate the advantage of the proposed approach, Figure 1 presents a numerical example for a predictive inference problem on the **wine quality dataset**.<sup>1</sup> (See Section 5.2 for full details of this experiment.)

We are given a pretrained family of prediction bands  $\widehat{C}_{1-a}$ , indexed by the target coverage level  $1 - a$ . At each value  $a \in [0, 1]$ , we define  $R_a(X) = \mathbb{P}(Y \notin \widehat{C}_{1-a}(X) \mid X)$ , the probability of the prediction band failing to cover the true response value  $Y$  given features  $X$ . Our goal is to return a prediction band with 90% coverage—that is, we would like to choose a value of  $a$  such that the expected risk

$$\mathbb{E}_{P_{\text{target}}}[R_a(X)] = \mathbb{P}_{\widetilde{P}_{\text{target}}}(Y \notin \widehat{C}_{1-a}(X))$$

is bounded by  $0.1 = 1 - 90\%$ . In our experiment, the available data is given by all samples that are white wines (with distribution  $\widetilde{P}$ ), while the target population is comprised of the samples that are red wines (with a different distribution  $\widetilde{P}_{\text{target}}$ ).

In Figure 1 below, we compare four methods (see Section 5.2 for details):

<sup>1</sup><https://archive.ics.uci.edu/dataset/186/wine+quality>

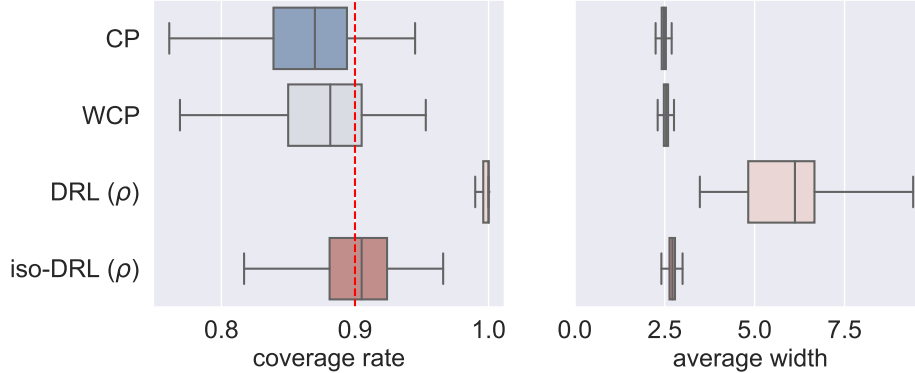


Figure 1: Coverage rate and average width of intervals for the `wine quality` dataset.

- An uncorrected interval—using conformal prediction (CP) (Vovk et al., 2005): the value  $a$  is chosen by tuning on the calibration data set (i.e., we choose  $a$  to satisfy  $\mathbb{E}_P[R_a(X)] \leq 0.1$ ), without correcting for the distribution shift.
- A corrected interval—using weighted conformal prediction (WCP) (Tibshirani et al., 2019): the value  $a$  is chosen by tuning on the calibration data set using an estimated density ratio  $w_0$  to correct for the covariate shift between distributions  $\tilde{P}$  and  $\tilde{P}_{\text{target}}$ . Since  $w_0$  is estimated from data, this correction is imperfect.
- The DRL interval: we choose  $a$  to satisfy  $\mathbb{E}_P[R_a(X)] + \Delta(R_a; \mathcal{Q}_{\text{KL}}(\rho)) \leq 0.1$ , where  $\mathbb{E}_P[R_a(X)]$  and  $\Delta(R_a; \mathcal{Q}_{\text{KL}}(\rho))$  are estimated using the calibration data.
- The iso-DRL interval: we choose  $a$  to satisfy  $\mathbb{E}_P[R_a(X)] + \Delta(R_a; \mathcal{Q}_{\text{KL}}(\rho) \cap \mathcal{Q}_{\text{iso}}(w_0)) \leq 0.1$ , where  $\mathbb{E}_P[R_a(X)]$  and  $\Delta(R_a; \mathcal{Q}_{\text{KL}}(\rho) \cap \mathcal{Q}_{\text{iso}}(w_0))$  are estimated using the calibration data.<sup>2</sup>

As we can see in Figure 1 below, the CP and WCP intervals both undercover—for CP, this is because the method does not correct for distribution shift, while for WCP, this is because the ratio  $w_0$  that corrects for distribution shift is imperfectly estimated. In contrast, DRL shows substantial overcoverage with extremely wide prediction intervals due to the worst-case nature of the bound  $\Delta(R_a; \mathcal{Q}_{\text{KL}}(\rho))$ . In contrast, our proposed method, iso-DRL, achieves the target coverage rate 90% without excessive increase in the size of the prediction interval, showing the benefit of adding the isotonic constraint to the DRL framework.

The motivating example demonstrates that, when we have access to meaningful—but imperfect—side information (e.g., in the form of the density ratio  $w_0$ ), adding the isotonic constraint to iso-DRL can provide an estimate of the risk that is *more reliable* than a non-distributionally-robust approach, but *less conservative* than the original DRL approach.

## 1.4 Organization of paper

Section 2 introduces a general class of uncertainty sets for candidate distributions and further studies the property of the worst-case excess risk defined in (1) for generic DRL. For the worst-case excess risk with the isotonic constraint, we prove that it is equivalent to the worst-case excess risk for a projected risk function without the isotonic constraint in Section 3. In Section 4, we propose an estimator of the worst-case excess risk with the isotonic constraint and establish the estimation error bounds. Numerical results for both synthetic and real data are shown in Section 5 and additional related work is summarized in Section 6. We defer technical proofs and additional simulations to the appendix.

<sup>2</sup>For both the DRL and iso-DRL methods, the parameter  $\rho$  is an estimate of the actual KL distance  $D_{\text{KL}}(P_{\text{target}}\|P)$ .

**Notation.** Before proceeding, we introduce useful notation for theoretical developments later on. To begin with, we denote by  $L_p(P)$  ( $1 \leq p \leq \infty$ ) the  $L_p$  function space under the probability measure  $P$ , i.e., when  $p \neq \infty$ ,

$$L_p(P) = \left\{ f : \|f\|_p = \left( \int_{\mathcal{X}} f^p(x) dP(x) \right)^{1/p} < \infty \right\}.$$

When  $p = \infty$ , the set  $L_\infty(P)$  consists of measurable functions that are bounded almost surely under  $P$ . In addition, for a measurable function  $w : \mathcal{X} \rightarrow \mathcal{W}$  and a measure  $P$  on  $\mathcal{X}$ , the pushforward measure  $w_\#P$  denotes the measure satisfying that  $(w_\#P)(B) = P(w^{-1}(B))$  for any  $B \in \mathcal{Z}$ , where  $w^{-1}(B) = \{x \in \mathcal{X} : w(x) \in B\}$  denotes the preimage of  $B$  under  $w$ . In other words, if  $X \sim P$ , then  $w(X)$  follows the distribution  $w_\#P$ . We say a function  $h$  is  $A_h$ -bounded if  $\sup_x |h(x)| \leq A_h$ .

Fix a partial (pre)order  $\preceq$  on  $\mathcal{X} \subseteq \mathbb{R}^d$ . A function  $g$  is isotonic in  $x$  if  $g(x_1) \leq g(x_2)$  for any  $x_1 \preceq x_2$ . Correspondingly, we define the cone of isotonic functions by

$$\mathcal{C}_{\preceq}^{\text{iso}} = \{w : w \text{ is isotonic w.r.t. partial order } \preceq\}.$$

Lastly, to compare two probability distributions  $Q$  and  $P$ , we define the convex ordering  $\preceq^{cvx}$  by

$$Q' \preceq^{cvx} Q \quad \text{if and only if} \quad \mathbb{E}_{Q'}[\psi(X)] \leq \mathbb{E}_Q[\psi(X)] \quad \text{for all convex functions } \psi.$$

## 2 The distributional robustness framework

As we have explained in Section 1.1, both the unsupervised setting and supervised setting under covariate shift can be unified. Therefore, from now on, to develop our theoretical results we will use the notation of the unsupervised setting with the risk function  $R(X)$ , with the understanding that this also covers the supervised setting under covariate shift.

Recall that  $\mathcal{X}$  is the feature domain. In this paper we consider a bounded risk function  $R : \mathcal{X} \rightarrow [0, B_R]$  with  $0 < B_R < \infty$ , and the goal is to evaluate (or bound) the target risk  $\mathbb{E}_{P_{\text{target}}}[R(X)]$  using samples from  $P$ , by assuming that the target distribution  $P_{\text{target}}$  is in some sense similar to the available distribution  $P$ —more concretely, by assuming that the target distribution  $P_{\text{target}}$  lies in some neighborhood  $\mathcal{Q}$  around the distribution  $P$  of the available data.

**Reformulating the neighborhood.** To unify the different examples of constraints described in Section 1, we will start by considering settings where we can express the constraint  $Q \in \mathcal{Q}$  using conditions on the density ratio  $w = \frac{dQ}{dP}$ . This type of framework includes the sensitivity analysis setting via the bounds constraint on  $w$  (Cornfield et al., 1959; Rosenbaum, 1987; Tan, 2006; Ding and VanderWeele, 2016; Zhao et al., 2019b; Yadlowsky et al., 2018; Jin et al., 2022, 2023; Sahoo et al., 2022), and  $f$ -divergence constraints (e.g., bounding  $D_f(Q||P) = \mathbb{E}_P[f(dQ/dP(X))]$ ) (Duchi et al., 2021; Namkoong and Duchi, 2017; Duchi and Namkoong, 2018; Cauchois et al., 2020).<sup>3</sup>

Concretely, we can reparameterize the distribution  $Q$  using the density ratio  $w(x) = \frac{dQ}{dP}(x)$ . Then we can reformulate the constraint  $Q \in \mathcal{Q}$  into a constraint on this density ratio, i.e.,

$$Q \in \mathcal{Q} \iff w_\#P \in \mathcal{B},$$

where  $\mathcal{B}$  is a set of distributions, and where  $w_\#P$  denotes the pushforward measure (as defined in Section 1.4). To facilitate understanding, let us consider several examples.

<sup>3</sup>We note that DRL with optimal transport divergences is not covered by this framework (Shafieezadeh Abadeh et al., 2015; Blanchet and Murthy, 2019; Blanchet et al., 2019; Esfahani and Kuhn, 2015) We discuss this in Section 7.

**Example 1: bound-constrained distribution shift.** In sensitivity analysis, it is common to assume that the likelihood ratio  $\frac{dP_{\text{target}}}{dP}$  is bounded from above and below. This corresponds to a constraint set of the form

$$\mathcal{Q} = \left\{ Q : a \leq \frac{dQ}{dP}(X) \leq b \text{ } P\text{-almost surely} \right\},$$

for some constants  $0 \leq a < 1 < b < +\infty$ . In particular, when  $a = \Gamma^{-1}$  and  $b = \Gamma$  for some  $\Gamma > 1$ , this constraint set represents the marginal  $\Gamma$ -selection model for the density ratio in sensitivity analysis (Rosenbaum, 1987; Tan, 2006). By defining

$$\mathcal{B} = \mathcal{B}_{a,b} = \left\{ \tilde{Q} : \mathbb{E}_{Z \sim \tilde{Q}}[Z] = 1, \mathbb{P}_{Z \sim \tilde{Q}}(a \leq Z \leq b) = 1 \right\},$$

we can verify that

$$Q \in \mathcal{Q} \iff w_{\#}P \in \mathcal{B}_{a,b} \text{ with } w(x) = \frac{dQ}{dP}(x).$$

**Example 2:  $f$ -constrained distribution shift.** For  $f$ -constrained distribution shift, we consider the constraint set

$$\mathcal{Q} = \left\{ Q : \mathbb{E}_P \left[ f \left( \frac{dQ}{dP}(X) \right) \right] \leq \rho \right\}.$$

For instance, if we take  $\mathcal{Q} = \mathcal{Q}_{\text{KL}}(\rho) = \{Q : D_{\text{KL}}(Q||P) \leq \rho\}$ , this corresponds to choosing  $f(x) = x \log(x)$  in  $f$ -divergence above (Rényi, 1961). Choosing

$$\mathcal{B} = \mathcal{B}_{f,\rho} = \left\{ \tilde{Q} : \mathbb{E}_{Z \sim \tilde{Q}}[Z] = 1, \mathbb{E}_{Z \sim \tilde{Q}}[f(Z)] \leq \rho, \mathbb{P}_{Z \sim \tilde{Q}}(Z \geq 0) = 1 \right\},$$

we can verify that

$$Q \in \mathcal{Q} \iff w_{\#}P \in \mathcal{B}_{f,\rho} \text{ with } w(x) = \frac{dQ}{dP}(x).$$

## 2.1 Worst-case excess risk with DRL

In this section, we explore some properties of the generic DRL, without the isotonic constraint. Building this framework will help us to introduce the isotonic constraint in the next section.

Based on the equivalence of  $\mathcal{Q}$  and  $\mathcal{B}$  in representing the uncertainty set, we focus on the following equivalent representation of  $\Delta(R; \mathcal{Q})$ :

$$\begin{aligned} \Delta(R; \mathcal{B}) &= \sup_{w \geq 0} \mathbb{E}_P [w(X)R(X)] - \mathbb{E}_P [R(X)] \\ &\text{subject to } w_{\#}P \in \mathcal{B}, \end{aligned} \tag{4}$$

where abusing notation we now write  $\Delta(\cdot; \mathcal{B})$  to express that  $\mathcal{B}$  is a constraint on the distribution of the density ratio  $w(X) = \frac{dQ}{dP}(X)$ , where previously we instead wrote  $\Delta(\cdot; \mathcal{Q})$ . We will say that  $\Delta(R; \mathcal{B})$  is *attainable* if this supremum is attained by some  $w^*$  in the constraint set.

Throughout the paper, we assume that the set  $\mathcal{B}$  satisfies the following condition.

**Condition 2.1.** The set  $\mathcal{B}$  is closed under convex ordering, that is, if  $Q \in \mathcal{B}$ , then for any  $Q' \stackrel{cvx}{\preceq} Q$ , it holds that  $Q' \in \mathcal{B}$ .

This condition enables the following reformulation of the quantity of interest,  $\Delta(R; \mathcal{B})$ :

**Proposition 2.2.** *Assume Condition 2.1 holds. Then  $\Delta(R; \mathcal{B})$  can be equivalently written as*

$$\begin{aligned} \Delta(R; \mathcal{B}) &= \sup_{\phi: \mathbb{R} \rightarrow \mathbb{R}_+} \mathbb{E}_P[(\phi \circ R)(X)R(X)] - \mathbb{E}_P[R(X)] \\ &\text{subject to } (\phi \circ R)_\# P \in \mathcal{B}, \quad \phi \text{ is nondecreasing.} \end{aligned}$$

Moreover, if  $\Delta(R; \mathcal{B})$  is attainable (i.e., the supremum is attained by some  $w^*$  satisfying the constraints), then the equivalent formulation is attainable as well (i.e., the supremum is attained by some  $w^* = \phi^* \circ R$ , where  $\phi^*$  is nondecreasing).

See Section A.1 for the proof. In words, this proposition shows that the excess risk is maximized by considering functions  $w(x)$  that are monotonically nondecreasing with respect to  $R(x)$  (i.e.,  $w = \phi \circ R$  for some nondecreasing  $\phi$ ). This is intuitive, since maximizing the expected value of  $w(X)R(X)$  implies that we should choose a function  $w$  that is large when  $R$  is large.

Most importantly, Proposition 2.2 implies that for a class of constrained sets  $\mathcal{B}$ , the optimal value in the constrained optimization problem (4) only depends on covariates  $x$  through the risk function  $R(x)$ , or equivalently, only depends on the distribution of  $X$  through the distribution of  $R(X)$ . As a corollary of Proposition 2.2, the worst-case excess risk  $\Delta(R; \mathcal{B})$  is also monotonically nondecreasing in  $R$ . This property of  $\Delta(R; \mathcal{B})$  is commonly known as the (strict) monotonicity of the functional  $\Delta(R; \mathcal{B})$  in  $R(X)$  under the usual stochastic order, which is treated as a condition on the functional  $\Delta(R; \mathcal{B})$  in Shapiro and Pichler (2023); Shapiro (2017). We note that, in the special case when  $\mathcal{B}$  is specified in term of an  $f$ -divergence (as in Example 2 above), the conclusion of Proposition 2.2 is established by Donsker and Varadhan (1976); Lam (2016); Namkoong et al. (2022).

Next, we return to the two earlier examples of the constraint set  $\mathcal{B}$  to verify that this result holds in those settings.

**Returning to Example 1: bound-constrained distribution shift.** Recall that in this example, we take the constraint set  $\mathcal{B}$  to be  $\mathcal{B} = \mathcal{B}_{a,b} = \left\{ \tilde{Q} : \mathbb{E}_{Z \sim \tilde{Q}}[Z] = 1, \mathbb{P}_{Z \sim \tilde{Q}}(a \leq Z \leq b) = 1 \right\}$ , for some  $0 \leq a < 1 < b < +\infty$ . It is straightforward to verify that  $\mathcal{B}_{a,b}$  satisfies Condition 2.1, implying that Proposition 2.2 can be applied.

Moreover, in this specific example, we can actually calculate the maximizing density ratio  $w^*(x)$  explicitly. If the distribution of  $R(X)$  is continuous, the worst-case density ratio that attains the worst-case excess risk takes the form

$$w^*(x) = a \cdot \mathbb{1} \left\{ R(x) \leq q_R \left( \frac{b-1}{b-a} \right) \right\} + b \cdot \mathbb{1} \left\{ R(x) > q_R \left( \frac{b-1}{b-a} \right) \right\},$$

where  $q_R(t) = \inf\{r \in \mathbb{R} \mid F_R(r) \geq t\}$  and  $F_R$  is the cumulative distribution function of  $R_\# P$ —that is,  $q_R(t)$  is the  $t$ -quantile of the distribution of  $R(X)$  under  $X \sim P$ . For general  $R(X)$ , we have

$$w^*(x) = a \cdot \mathbb{1} \left\{ R(x) < q_R \left( \frac{b-1}{b-a} \right) \right\} + b \cdot \mathbb{1} \left\{ R(x) > q_R \left( \frac{b-1}{b-a} \right) \right\} + c \cdot \mathbb{1} \left\{ R(x) = q_R \left( \frac{b-1}{b-a} \right) \right\},$$

where

$$c = a + \frac{(b-a)t^* - (b-1)}{\mathbb{P} \left\{ R(X) = q_R \left( \frac{b-1}{b-a} \right) \right\}} \quad \text{with} \quad t^* = \inf \left\{ t \in \text{range}(F_R) \mid t \geq \frac{b-1}{b-a} \right\}.$$

In particular, we can see that  $w^*(x)$  is nondecreasing in  $R(x)$ , i.e., we can write  $w = \phi^* \circ R$  for some nondecreasing  $\phi^*$ , thus validating that the conclusion of Proposition 2.2 holds in this example.

**Returning to Example 2:  $f$ -constrained distribution shift.** Recall that for an  $f$ -divergence constraint, we define  $\mathcal{B} = \mathcal{B}_{f,\rho} = \left\{ \tilde{Q} : \mathbb{E}_{Z \sim \tilde{Q}}[Z] = 1, \mathbb{E}_{Z \sim \tilde{Q}}[f(Z)] \leq \rho, Z \geq 0 \right\}$ . Since  $f$  must be convex (for



$f$ -divergence to be well-defined), this immediately implies that  $\mathcal{B}_{f,\rho}$  satisfies Condition 2.1. By the results of Donsker and Varadhan (1976); Lam (2016), the worst-case excess risk  $\Delta_\rho(R; \mathcal{B}_{f,\rho})$  is attained at

$$w^*(x) = w(x; \lambda^*, \nu^*) = \left\{ (f')^{-1} \left( \frac{R(x) - \nu^*}{\lambda^*} \right) \right\}_+,$$

where  $a_+$  is the positive part of  $a \in \mathbb{R}$  and  $\lambda^*, \nu^*$  are the solutions to the dual problem

$$\inf_{\lambda \geq 0, \nu} \left\{ \lambda \rho + \nu + \mathbb{E}_P \left[ w(X; \lambda, \nu) (R(X) - \nu) - \lambda f(w(X; \lambda, \nu)) \right] \right\}. \quad (5)$$

Since  $f$  is convex, the inverse of derivative  $(f')^{-1}$  is then nondecreasing, meaning that  $w^*(x)$  is nondecreasing in  $R(x)$ , which again validates the result in Proposition 2.2.

### 3 Worst-case excess risk with an isotonic constraint

In this section, we will now formally introduce our iso-DRL method, adding an isotonic constraint to the DRL framework developed in Section 2 above.

Recall the cone of isotonic functions

$$\mathcal{C}_{\preceq}^{\text{iso}} = \{w : \mathcal{X} \rightarrow \mathbb{R} : w \text{ is isotonic w.r.t. partial order } \preceq\}.$$

In this paper, we actually allow  $\preceq$  to be a partial *preorder* rather than a partial order, meaning that it may be the case that both  $x \preceq x'$  and  $x' \preceq x$ , even when  $x \neq x'$ . As an example, we denote  $\mathcal{C}_{w_0}^{\text{iso}} = \{w : w(x) \text{ is a monotonically nondecreasing function of } w_0(x)\}$ —this is obtained by the (pre)order given by  $x \preceq x'$  whenever  $w_0(x) \leq w_0(x')$ .

Our focus is the worst-case excess risk with the isotonic constraint:

$$\begin{aligned} \Delta^{\text{iso}}(R; \mathcal{B}) &= \sup_{w \geq 0} \mathbb{E}_P [w(X)R(X)] - \mathbb{E}_P[R(X)] \\ &\text{subject to } w_{\#}P \in \mathcal{B}, \quad w \in \mathcal{C}_{\preceq}^{\text{iso}}. \end{aligned} \quad (6)$$

To make this more concrete with a specific example, in the bound (3), this example corresponds to choosing  $\mathcal{B} = \mathcal{B}_{f,\rho}$  for the  $f$ -divergence  $f(x) = x \log x$  that is related to the KL distance. In particular, the bound (3) assumed two constraints on the distribution  $P_{\text{target}}$ —first,  $D_{\text{KL}}(P_{\text{target}} \| P) \leq \rho$  (which corresponds to assuming  $(dP_{\text{target}}/dP)_{\#}P \in \mathcal{B}_{f,\rho}$ , in our new notation), and second,  $P_{\text{target}} \in \mathcal{Q}_{\text{iso}}(w_0)$  (which is expressed by assuming  $w \in \mathcal{C}_{w_0}^{\text{iso}}$ , in our new notation, when we take the partial (pre)order defined as  $x \preceq x'$  whenever  $w_0(x) \leq w_0(x')$ ).

#### 3.1 Equivalent formulation

Optimization problems with isotonic constraints may be difficult to tackle both theoretically and computationally, since the isotonic cone, despite being convex, may be challenging to optimize over when working with an infinite-dimensional object such as the density ratio. In this section, we will show that the maximization problem (6) can equivalently be reformulated as an optimization problem *without* an isotonic constraint, by drawing a connection to the original (not isotonic) DRL maximization problem (4).

Given the probability measure  $P$ , we will define  $\pi$  as the projection to the isotonic cone  $\mathcal{C}_{\preceq}^{\text{iso}}$  with respect to  $L_2(P)$ :

$$\pi(a) = \operatorname{argmin}_{b \in \mathcal{C}_{\preceq}^{\text{iso}}} \int (a(x) - b(x))^2 dP(x).$$

As  $L_2(P)$  is reflexive and strictly convex, the projection  $\pi(a)$  exists and is unique (up to sets of measure zero) for all  $a \in L_2(P)$  (Megginson, 2012).

With the projection  $\pi$  in place, we are ready to state our main equivalence result.

**Theorem 3.1.** For any  $\mathcal{B}$  and any partial (pre)order  $\preceq$  on  $\mathcal{X}$ , it holds that

$$\Delta^{\text{iso}}(R; \mathcal{B}) \leq \Delta(\pi(R); \mathcal{B}).$$

If in addition Condition 2.1 holds, then we have

$$\Delta^{\text{iso}}(R; \mathcal{B}) = \Delta(\pi(R); \mathcal{B}),$$

and moreover,  $\Delta^{\text{iso}}(R; \mathcal{B})$  is attainable if and only if  $\Delta(\pi(R); \mathcal{B})$  is attainable.

See Section B.1 for the proof.

To interpret this theorem, recall from the definition (4) that we have

$$\begin{aligned} \Delta(\pi(R); \mathcal{B}) &= \sup_{w \geq 0} \mathbb{E}_P[w(X)[\pi(R)](X)] - \mathbb{E}_P[[\pi(R)](X)] \\ &\text{subject to } w_{\#}P \in \mathcal{B}. \end{aligned} \tag{7}$$

Compared with the formulation (6) that defines the isotonic worst-case risk  $\Delta^{\text{iso}}(R; \mathcal{B})$ , we see that this equivalent formulation removes the constraint  $w \in \mathcal{C}_{\preceq}^{\text{iso}}$  by replacing  $R$  with its isotonic projection  $\pi(R)$ . This brings computational benefits. The equivalent formulation (7) separates two constraints  $w_{\#}P \in \mathcal{B}$  and  $w \in \mathcal{C}_{\preceq}^{\text{iso}}$ , allowing us to first project the risk function onto  $\mathcal{C}_{\preceq}^{\text{iso}}$  and solve a problem that is as simple as the problem stated earlier in (4). More concretely, as seen in Examples 1 and 2, for many common choices of  $\mathcal{B}$ , we have closed-form solutions to (7) in terms of the projected risk  $\pi(R)$ .

### 3.2 Setting: iso-DRL with estimated density ratio

We now return to the scenario described in (3) in Section 1.3, where we would like to recalibrate a pretrained density ratio  $w_0$  that estimates the distribution shift  $\frac{dP_{\text{target}}}{dP}$ . As the shape or relative magnitude of  $w_0$  could contain useful information about the true density ratio, we assume that the true density ratio is an isotonic function of  $w_0$ —that is, we assume

$$\frac{dP_{\text{target}}}{dP}(x) = \phi(w_0(x))$$

for some nondecreasing function  $\phi$ , for  $P$ -almost every  $x$ . Equivalently, defining the partial (pre)order

$$x \preceq x' \iff w_0(x) \leq w_0(x'), \tag{8}$$

we are essentially assuming that  $\frac{dP_{\text{target}}}{dP} \in \mathcal{C}_{\preceq}^{\text{iso}}$  for this particular partial order. We will denote this specific cone as  $\mathcal{C}_{w_0}^{\text{iso}}$  and its isotonic projection as  $\pi_{w_0}$ , and abusing notation, we write  $\Delta^{\text{iso}}(R; \mathcal{B}, w_0)$  to denote the excess risk for this particular setting, to emphasize the role of  $w_0$ .

By Theorem 3.1, if we assume  $\mathcal{B}$  satisfies Condition 2.1 then we have the equivalence

$$\Delta^{\text{iso}}(R; \mathcal{B}, w_0) = \Delta(\pi_{w_0}(R); \mathcal{B}). \tag{9}$$

To understand the projection onto the cone  $\mathcal{C}_{w_0}^{\text{iso}}$  more straightforwardly, we can derive a further simplification, with a few more definitions. First, write  $\pi_1$  to denote the isotonic projection of functions  $\mathbb{R} \rightarrow \mathbb{R}$  under the measure  $(w_0)_{\#}P$ , and define a function  $\tilde{R} : \mathbb{R} \rightarrow \mathbb{R}$  to satisfy

$$\tilde{R}(w_0(X)) = \mathbb{E}_P[R(X) \mid w_0(X)]$$

$P$ -almost surely. We then have the following simplified equivalence:

**Proposition 3.2.** *Assume Condition 2.1 holds. We have the equivalence*

$$\Delta^{\text{iso}}(R; \mathcal{B}, w_0) = \Delta(\pi_1(\tilde{R}) \circ w_0; \mathcal{B}, w_0),$$

where we recall

$$\begin{aligned} \Delta(R; \mathcal{B}, w_0) &= \sup_{h: h \circ w_0 \geq 0} \mathbb{E}_P [(h \circ w_0)(X)R(X)] - \mathbb{E}_P [R(X)] \\ &\text{subject to } (h \circ w_0)_{\#} P \in \mathcal{B}. \end{aligned}$$

Compared to the equivalence (9), the new equivalence in the proposition relies on an isotonic projection with respect to the canonical order on the real line (i.e., the projection  $\pi_1$ ), as opposed to projecting to the cone  $\mathcal{C}_{w_0}^{\text{iso}}$ , which uses the more complicated partial preorder defined in (8).

### 3.3 A misspecified isotonic constraint

When the true distribution shift does not obey the isotonic constraint exactly, we can nonetheless provide a bound on the worst-case excess risk, which is tighter than the (non-iso) DRL bound whenever the isotonic constraint provides a reasonable approximation.

Denote  $\tilde{w}^*$  as the underlying density ratio  $dP_{\text{target}}/dP$  and  $\Delta^*(R) = \mathbb{E}_P[\tilde{w}^*(X)R(X)] - \mathbb{E}_P[R(X)]$  as the true excess risk. Then, we have the following connections between  $\Delta^*(R)$  and  $\Delta^{\text{iso}}(R; \mathcal{B})$ .

**Proposition 3.3.** *Assume Condition 2.1 holds. If  $\tilde{w}_{\#}^* P \in \mathcal{B}$ , then we have*

$$\Delta^*(R) \leq \Delta^{\text{iso}}(R; \mathcal{B}) + \mathbb{E}_P \left[ [\tilde{w}^* - \pi(\tilde{w}^*)](X) \cdot [R - \pi(R)](X) \right].$$

In particular, if either  $\tilde{w}^* \in \mathcal{C}_{\geq}^{\text{iso}}$  or  $R \in \mathcal{C}_{\geq}^{\text{iso}}$ , then one has  $\Delta^*(R) \leq \Delta^{\text{iso}}(R; \mathcal{B})$ .

The result states that when the isotonic constraint is violated, the worst-case excess risk of iso-DRL will be no worse than the true excess risk minus a gap which can be controlled by the correlation between  $[\tilde{w}^* - \pi(\tilde{w}^*)](X)$  and  $[R - \pi(R)](X)$ . In particular, if *either* the risk or the true density ratio is itself isotonic, the excess risk calculation  $\Delta^{\text{iso}}(R; \mathcal{B})$ , which is tighter than the (non-iso) DRL bound  $\Delta(R; \mathcal{B})$ , will never underestimate the true risk  $\Delta^*(R)$ .

## 4 Estimation of worst-case excess risk with isotonic constraint

So far, our focus is on the population level, namely we assume full access to the data distribution  $P$  and the risk function  $R$ . In practice, however, we may only access the data distribution  $P$  via samples drawn from  $P$ , and we may only be able to learn about the risk function  $R$  via noisy evaluations of  $R(X)$  on each sampled point  $X$ .

In this section, focusing on the supervised setting, we propose a fully data dependent estimator for the worst-case excess risk  $\Delta^{\text{iso}}(R; \mathcal{B})$ . Moreover, we characterize the estimation error for different choices of  $\mathcal{B}$ , including the bounds constraint and the  $f$ -divergence constraint for the distribution shift.

In this case, we observe  $\{(X_i, Y_i)\}_{i \leq n}$  drawn i.i.d. from  $P \times P_{Y|X}$ , and only observe the risk  $r : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}_+$  that approximates the true underlying risk function  $R$  in the sense that

$$R(X) = \mathbb{E}[r(X, Y) | X], \quad \text{almost surely.}$$

To make this concrete, as an example we can recall the regression setting with squared loss from Section 1, given by  $R(X) = \mathbb{E}[r(X, Y) | X]$  for  $r(x, y) = (y - \hat{\mu}(x))^2$ . Then our information about the (expected) risk

is often limited to evaluating  $r(X_i, Y_i) = (Y_i - \hat{\mu}(X_i))^2$  on the samples  $(X_i, Y_i)$  drawn from  $\tilde{P}$ —that is, we only observe *noisy* estimates of  $R(X)$ , at only *finitely many* randomly sampled values of  $X$ .

Given a truncation level  $1 \leq \Omega < +\infty$ , we propose to estimate the worst-case excess risk via the following optimization problem:

$$\begin{aligned} \hat{\Delta}_{\Omega}^{\text{iso}}(r; \mathcal{B}) &:= \max_{w \geq 0} \frac{1}{n} \sum_{i \leq n} w(X_i) r(X_i, Y_i) - \frac{1}{n} \sum_{i \leq n} r(X_i, Y_i) \\ \text{subject to} & \quad w_{\#} \hat{P}_n \in \mathcal{B}, \quad w \in \mathcal{C}_{\preceq}^{\text{iso}}, \quad \|w\|_{\infty} \leq \Omega. \end{aligned} \quad (10)$$

Here,  $\hat{P}_n$  denotes the empirical distribution of i.i.d. observations  $\{X_i\}_{i \leq n}$  drawn from  $P$ . Since the feasible set is compact, the maximum is attainable, and we denote it by  $\hat{w}_r^{\text{iso}}$ . When  $\Omega = +\infty$ , we write  $\hat{\Delta}^{\text{iso}}(r; \mathcal{B}) = \hat{\Delta}_{\infty}^{\text{iso}}(r; \mathcal{B})$ .<sup>4</sup>

From now on, when  $d \geq 2$ , we consider a bounded domain  $\mathcal{X}$  equipped with the componentwise order (Han et al., 2019; Deng and Zhang, 2020; Gao and Wellner, 2007), i.e.  $x \preceq z$  if and only if  $x_j \leq z_j$  for all  $j \in [d]$ . Similar to Han et al. (2019), in the multivariate case, we assume  $0 < m_0 \leq \inf_{x \in \mathcal{X}} dP(x) \leq \sup_{x \in \mathcal{X}} dP(x) \leq M_0 < \infty$  with constants  $m_0$  and  $M_0$ .

## 4.1 Computation: estimation after projection

In view of Theorem 3.1, we may accelerate the computation of (10) via an equivalent optimization problem without the isotonic constraint.

To be more specific, denote  $r^{\text{iso}} = (r_i^{\text{iso}})_{i \leq n} \in \mathbb{R}^n$  as the isotonic projection of  $(r(X_i, Y_i))_{i \leq n}$  with respect to the empirical distribution  $\hat{P}_n$  under the partial order  $\preceq$ . Then, consider the optimization problem

$$\begin{aligned} \hat{\Delta}_{\Omega}(r^{\text{iso}}; \mathcal{B}) &:= \max_{w \geq 0} \frac{1}{n} \sum_{i \leq n} w(X_i) r_i^{\text{iso}} - \frac{1}{n} \sum_{i \leq n} r_i^{\text{iso}} \\ \text{subject to} & \quad w_{\#} \hat{P}_n \in \mathcal{B}, \quad \|w\|_{\infty} \leq \Omega. \end{aligned} \quad (11)$$

By Theorem 3.1, we have  $\hat{\Delta}_{\Omega}(r^{\text{iso}}; \mathcal{B}) = \hat{\Delta}_{\Omega}^{\text{iso}}(r; \mathcal{B})$ .

Note that in iso-DRL with estimated density ratio in Section 3.2, we can simply apply the isotonic regression for  $(r(X_i, Y_i))_{i \leq n}$  on  $(w_0(X_i))_{i \leq n}$  to obtain the projected risk.

## 4.2 Reduction to noiseless risk

To control the estimation error, we will consider an oracle estimator with perfect knowledge of the noiseless risk  $R$ . That is, we consider the following optimization problem:

$$\begin{aligned} \hat{\Delta}_{\Omega}^{\text{iso}}(R; \mathcal{B}) &:= \max_{w \geq 0} \frac{1}{n} \sum_{i \leq n} w(X_i) R(X_i) - \frac{1}{n} \sum_{i \leq n} r(X_i, Y_i) \\ \text{subject to} & \quad w_{\#} \hat{P}_n \in \mathcal{B}, \quad w \in \mathcal{C}_{\preceq}^{\text{iso}}, \quad \|w\|_{\infty} \leq \Omega. \end{aligned} \quad (12)$$

In comparison to (10), in the first sum in the maximization, the noisy risk  $r(X_i, Y_i)$  is replaced by the noiseless counterpart  $R(X_i)$ .

Recall our population-level target  $\Delta^{\text{iso}}(R; \mathcal{B})$  defined in (6). It is clear via the triangle inequality that the estimation error can be decomposed into two parts:

$$|\hat{\Delta}_{\Omega}^{\text{iso}}(r; \mathcal{B}) - \Delta^{\text{iso}}(R; \mathcal{B})| \leq |\hat{\Delta}_{\Omega}^{\text{iso}}(r; \mathcal{B}) - \hat{\Delta}_{\Omega}^{\text{iso}}(R; \mathcal{B})| + |\hat{\Delta}_{\Omega}^{\text{iso}}(R; \mathcal{B}) - \Delta^{\text{iso}}(R; \mathcal{B})|.$$

We first show in the following theorem that the convergence of the first term does not depend on the specific choice of  $\mathcal{B}$ .

<sup>4</sup>We note that  $\hat{\Delta}_{\Omega}^{\text{iso}}(r; \mathcal{B})$  can also be written as  $\hat{\Delta}^{\text{iso}}(r; \mathcal{B}_{\Omega})$ , where  $\mathcal{B}_{\Omega} = \{Q \in \mathcal{B} : \mathbb{P}_{Z \sim Q}(Z \leq \Omega) = 1\}$ .

**Theorem 4.1.** *Assume both  $R$  and  $r$  are  $B_R$ -bounded, there exist constants  $C > 0$  and  $\gamma_d$  that only depends on  $d$  such that, with probability at least  $1 - 2n^{-1}$ ,*

$$\left| \widehat{\Delta}_{\Omega}^{\text{iso}}(r; \mathcal{B}) - \widehat{\Delta}_{\Omega}^{\text{iso}}(R; \mathcal{B}) \right| \leq C \frac{\log^{\gamma_d/2} n}{n^{1/\max\{d,2\}}}.$$

See Section C.1 for the proof. As a result, bounding  $|\widehat{\Delta}_{\Omega}^{\text{iso}}(r; \mathcal{B}) - \Delta^{\text{iso}}(R; \mathcal{B})|$  is reduced to bounding the estimation error with the noiseless risk  $R$ , i.e., controlling the second term  $|\widehat{\Delta}_{\Omega}^{\text{iso}}(R; \mathcal{B}) - \Delta^{\text{iso}}(R; \mathcal{B})|$ .

### 4.3 Estimation error with the noiseless risk

In this section, we bound the error  $|\widehat{\Delta}_{\Omega}^{\text{iso}}(R; \mathcal{B}) - \Delta^{\text{iso}}(R; \mathcal{B})|$  induced by the sampling of  $X$ . To simplify the presentation, we focus on two canonical examples with bounds constraints and  $f$ -divergence constraints. For the population level worst-case excess risk, we also denote  $\Delta_{\Omega}^{\text{iso}}(R; \mathcal{B})$  as a modification of  $\Delta^{\text{iso}}(R; \mathcal{B})$  with an additional constraint  $\|w\|_{\infty} \leq \Omega$ , i.e.

$$\Delta_{\Omega}^{\text{iso}}(R; \mathcal{B}) = \Delta^{\text{iso}}(R; \mathcal{B}_{\Omega}), \quad \text{where } \mathcal{B}_{\Omega} = \{Q \in \mathcal{B} : \mathbb{P}_{Z \sim Q}(Z \leq \Omega) = 1\}.$$

We will start with the analysis of the error  $|\widehat{\Delta}_{\Omega}^{\text{iso}}(R; \mathcal{B}) - \Delta_{\Omega}^{\text{iso}}(R; \mathcal{B})|$ , and then validate that the remaining bias term  $|\Delta_{\Omega}^{\text{iso}}(R; \mathcal{B}) - \Delta^{\text{iso}}(R; \mathcal{B})|$  is zero when  $\Omega$  is sufficiently large.

Recall the definitions of  $\mathcal{B}_{a,b}$  and  $\mathcal{B}_{f,\rho}$  in Section 2. We further denote the function class  $\mathcal{G}_{\mathcal{B}_{a,b}} = \{w : a \leq w \leq b, w \in \mathcal{C}_{\geq}^{\text{iso}}\}$  and  $\mathcal{G}_{\mathcal{B}_{f,\rho}} = \{w : w \in \mathcal{C}_{\geq}^{\text{iso}}, w \in [-\Omega \vee B_f, \Omega \vee B_f]\}$ , where we assume  $\sup_{t \in [0, \Omega]} f(t) = B_f < +\infty$ . Define the empirical Rademacher complexity of the function class  $\mathcal{G}$  by

$$\mathcal{R}_n(\mathcal{G}) := \mathbb{E}_{\sigma} \left[ \sup_{g \in \mathcal{G}} \left| \frac{1}{n} \sum_{i \leq n} \sigma_i g(Z_i) \right| \right],$$

where  $\{Z_i\}_{i \leq n}$  is a sample of size  $n$  from  $P$  and  $\{\sigma_i\}_{i \leq n}$  are independent random variables drawn from the Rademacher distribution. Then, we have the following theorem for the estimation error bound of  $|\widehat{\Delta}_{\Omega}^{\text{iso}}(R; \mathcal{B}) - \Delta_{\Omega}^{\text{iso}}(R; \mathcal{B})|$ .

**Theorem 4.2.** *Assume  $R$  is  $B_R$ -bounded and  $\sup_{t \in [0, \Omega]} f(t) = B_f < +\infty$ . For any  $\Omega > 0$ , and for either choice of constraint set  $\mathcal{B} = \mathcal{B}_{a,b}$  or  $\mathcal{B} = \mathcal{B}_{f,\rho}$ , there exists a constant  $C$  that does not depend on  $n$  such that with probability at least  $1 - 2n^{-1}$ , it holds that*

$$\left| \widehat{\Delta}_{\Omega}^{\text{iso}}(R; \mathcal{B}) - \Delta_{\Omega}^{\text{iso}}(R; \mathcal{B}) \right| \leq C \left( \sqrt{\frac{\log n}{n}} + \mathcal{R}_n(\mathcal{G}_{\mathcal{B}}) \right). \quad (13)$$

We provide two concrete examples to make apparent the dependence of the empirical Rademacher complexity on the sample size.

- (1) When  $d = 1$ , e.g., in the setting of density ratio recalibration in Section 3.2, both  $\mathcal{B}_{a,b}$  and  $\mathcal{B}_{f,\rho}$  are contained in the set of uniformly bounded unimodal functions, then, similar to the results of Chatterjee and Lafferty (2019), one can show by Dudley's theorem (Dudley, 1967) that  $\mathcal{R}_n(\mathcal{G}_{\mathcal{B}}) \lesssim n^{-1/2}$ , ignoring logarithmic factors.
- (2) For  $\mathbb{R}^d$  with a fixed dimension  $d \geq 2$  and a bounded domain  $\mathcal{X}$  equipped with the componentwise order, by Han et al. (2019), if  $0 < m_0 \leq \inf_{x \in \mathcal{X}} dP(x) \leq \sup_{x \in \mathcal{X}} dP(x) \leq M_0 < \infty$ , we have  $\mathcal{R}_n(\mathcal{G}_{\mathcal{B}}) \lesssim n^{-1/d}$ , ignoring logarithmic factors, for which we provide more details in Appendix C.2 and C.3.

To conclude, we note that the estimation error in Theorem 4.2 is for the truncated population worst-case excess risk  $\Delta_{\Omega}^{\text{iso}}(R; \mathcal{B})$ . We will show that the bias term  $|\Delta_{\Omega}^{\text{iso}}(R; \mathcal{B}) - \Delta^{\text{iso}}(R; \mathcal{B})|$  is zero for  $\mathcal{B}_{a,b}$  and  $\mathcal{B}_{f,\rho}$  as follows:

- For  $\mathcal{B}_{a,b}$ , the bound (13) still holds for  $|\widehat{\Delta}^{\text{iso}}(R; \mathcal{B}_{a,b}) - \Delta^{\text{iso}}(R; \mathcal{B}_{a,b})|$  when  $\Omega \geq b$ .
- For the  $f$ -constrained problem, with  $\sup_{t \in [0, \Omega]} f(t) = B_f < +\infty$ , as we will show in Appendix C.3 that the worst-case excess risk  $\Delta^{\text{iso}}(R; \mathcal{B}_{f,\rho})$  is attained at  $w_{f,\rho}^{\text{iso}} \in \mathcal{C}_{\geq}^{\text{iso}}$  with  $\|w_{f,\rho}^{\text{iso}}\|_{\infty} < \infty$  almost surely, which implies that  $\Delta_{\Omega}^{\text{iso}}(R; \mathcal{B}_{f,\rho}) = \Delta^{\text{iso}}(R; \mathcal{B}_{f,\rho})$  whenever  $\Omega \geq \|w_{f,\rho}^{\text{iso}}\|_{\infty}$ .

Combining the results in Theorem 4.1 and 4.2, it holds that with  $\mathcal{B} = \mathcal{B}_{a,b}$  or  $\mathcal{B}_{f,\rho}$  and adequately large  $\Omega$ , the estimation error  $|\widehat{\Delta}_{\Omega}^{\text{iso}}(r; \mathcal{B}) - \Delta^{\text{iso}}(R; \mathcal{B})| \lesssim n^{-1/\max\{d,2\}}$ , up to logarithmic factors. As we will discuss in Appendix C.4, the isotonic constraint  $w \in \mathcal{C}_{\geq}^{\text{iso}}$  plays an important role in the convergence result. Even in the simple case with  $\mathcal{B} = \mathcal{B}_{a,b}$ , we will present an example in Appendix C.4 where the estimation error of the (non-iso) DRL risk does not converge to zero.

## 5 Numerical experiments

In this section, we demonstrate the benefits of iso-DRL in calibrating prediction sets under covariate shift with empirical examples, as previewed in Section 1.3. Throughout all experiments, we have data  $\mathcal{D}_{\text{train}} = \{(X_i, Y_i)\}_{i \leq N}$  drawn from the data distribution and the test set  $\mathcal{D}_{\text{test}} = \{(\tilde{X}_i, \tilde{Y}_i)\}_{i \leq M}$  drawn from the target distribution with  $M \ll N$ . We consider both synthetic and real datasets. Code to reproduce all experiments is available at <https://github.com/yugjerry/iso-DRL>.

**Background.** When covariate shift is present, Tibshirani et al. (2019) proposes the weighted conformal prediction (WCP) method, which produces a prediction set  $C_{1-\alpha}^{w_0}(X)$  with an estimated density ratio  $w_0$ , which is valid for the covariate distribution  $\widehat{P}$  defined by  $d\widehat{P} \propto w_0 \cdot dP$ . The validity for the target distribution  $\tilde{P}_{\text{target}}$  is only guaranteed up to a coverage gap due to the estimation error or potential misspecification in  $w_0$  (Lei and Candès, 2020; Candès et al., 2023; Gui et al., 2023, 2024)—that is, if  $w_0$  is a reasonably accurate estimate of the true density ratio  $\frac{dP_{\text{target}}}{dP}$  of the covariate shift, then WCP will lead to coverage at approximately  $(1 - \alpha)$  level relative to  $\tilde{P}_{\text{target}}$ . In comparison to our approach, Cauchois et al. (2020); Ai and Ren (2024) share similar idea with the generic DRL to adjust the target level  $\alpha$ , but focuses on a different setting with distribution shift on the joint distribution of  $(X, Y)$ . More related work on conformal prediction is discussed in Section 6.

**Dataset partition.** The datasets  $\mathcal{D}_{\text{train}}$  and  $\mathcal{D}_{\text{test}}$  are partitioned as follows:

- First, we use a subset  $\mathcal{D}_1 \subset \mathcal{D}_{\text{train}}$  of the training data of size  $|\mathcal{D}_1| = n_{\text{pre}}$ , and a subset  $\mathcal{D}_{\text{test},1} \subseteq \mathcal{D}_{\text{test}}$  of the test data of size  $|\mathcal{D}_{\text{test},1}| = n_{\text{pre}}$ , to train the estimator of the covariate shift, i.e., the function  $w_0$ .
- Next, we use a subset  $\mathcal{D}_2 \subset \mathcal{D}_{\text{train}} \setminus \mathcal{D}_1$  of the training data of size  $|\mathcal{D}_2| = n_{\text{train}}$  to train CP or WCP prediction intervals.
- Then,  $\mathcal{D}_3 = \mathcal{D}_{\text{train}} \setminus (\mathcal{D}_1 \cup \mathcal{D}_2)$  is used to for estimating upper bounds on the excess risk for the DRL and iso-DRL methods. We further define  $n = |\mathcal{D}_3|$  to ease notations.
- Finally,  $\mathcal{D}_{\text{test},0} = \mathcal{D}_{\text{test}} \setminus \mathcal{D}_{\text{test},1}$  with  $|\mathcal{D}_{\text{test},0}| = n_{\text{test}}$  is used for estimating the actual performance of each method relative to the target distribution. We will measure the coverage rate on  $\mathcal{D}_{\text{test},0}$  to assess each method's performance:

$$\text{Coverage rate}(C, \alpha) = \frac{1}{n_{\text{test}}} \sum_{i \in \mathcal{D}_{\text{test},0}} \mathbb{1} \left\{ \tilde{Y}_i \in C(\tilde{X}_i) \right\}.$$

We next turn to the details of how each of these steps are carried out.

**Initial density ratio estimation.** Using data from  $\mathcal{D}_1$  and  $\mathcal{D}_{\text{test},1}$ , we construct a data set comprised of the covariate  $X$  (from either the training data points  $\mathcal{D}_1$  or the test data points  $\mathcal{D}_{\text{test},1}$ ) and a binary label  $L \in \{0, 1\}$  (0 for the training points, 1 for the test points). We then fit a logistic regression model and obtain the estimated probability  $\hat{p}(x)$  for  $\mathbb{P}(L = 1 \mid X = x)$ , with which we define

$$w_0(x) = \frac{\hat{p}(x)}{1 - \hat{p}(x)}$$

**(Weighted) Split conformal prediction.** With data from  $\mathcal{D}_2$ , we use Ordinary Least Squares (OLS) as the base algorithm, where we denote  $\hat{\mu}$  as the fitted regression model, and, following [Tibshirani et al. \(2019\)](#); [Lei et al. \(2018\)](#), apply split conformal prediction with the nonconformity score  $V(x, y) = |y - \hat{\mu}(x)|$  to obtain the following prediction intervals for comparison:

- CP: conformal prediction interval  $C_{1-\alpha}$  without adjusting for covariate shift;
- WCP-oracle: weighted conformal prediction interval  $C_{1-\alpha}^{w^*}$  with true density ratio  $w^* = dP_{\text{target}}/dP$ ;
- WCP: weighted conformal prediction interval  $C_{1-\alpha}^{w_0}$  with estimated  $w_0$ ;

**DRL methods: estimation of worst-case excess risks.** Next we give details on how we implement the two distributionally robust methods, which we denote by DRL (i.e., without an isotonic constraint) and iso-DRL- $w_0$  (i.e., our proposed method, with the constraint that the distribution shift is monotone with respect to the estimated covariate shift function  $w_0$ ).

Using the subset  $\mathcal{D}_3$  of the training data, the observed risks can be calculated by

$$r_i = \mathbb{1}\{Y_i \notin C_{1-\alpha}(X_i)\}, \quad i \in \mathcal{D}_3.$$

We adopt the KL-constraint  $D_{\text{KL}}(Q\|P) \leq \rho$  to measure the magnitude of distribution shift, with which we can obtain the following estimated worst-case excess risk

$$\begin{aligned} \hat{\Delta}(\alpha) = \max \quad & \frac{1}{n} \sum_{i \in \mathcal{D}_3} w_i r_i - \frac{1}{n} \sum_{i \in \mathcal{D}_3} r_i \\ \text{subject to} \quad & \frac{1}{n} \sum_{i \in \mathcal{D}_3} w_i = 1, \quad \frac{1}{n} \sum_{i \in \mathcal{D}_3} w_i \log w_i \leq \rho, \quad 0 \leq w_i \leq \Omega, \end{aligned} \quad (14)$$

with the upper bound set as  $\Omega = 100$  throughout the experiments. Given the estimated density ratio  $w_0$ , we run isotonic regression for  $(r_i)_{i \leq n}$  on  $(w_0(X_i^{(3)}))_{i \leq n}$  to obtain the projected risk  $(r_i^{\text{iso}})_{i \in \mathcal{D}_3}$ , with which we can calculate the worst-case excess risk

$$\begin{aligned} \hat{\Delta}^{\text{iso}}(\alpha) = \max \quad & \frac{1}{n} \sum_{i \in \mathcal{D}_3} w_i r_i^{\text{iso}} - \frac{1}{n} \sum_{i \in \mathcal{D}_3} r_i^{\text{iso}} \\ \text{subject to} \quad & \frac{1}{n} \sum_{i \in \mathcal{D}_3} w_i = 1, \quad \frac{1}{n} \sum_{i \in \mathcal{D}_3} w_i \log w_i \leq \rho, \quad 0 \leq w_i \leq \Omega. \end{aligned} \quad (15)$$

Given these estimates of the worst-case excess risks, we compare the following methods:

- DRL: CP interval  $C_{1-\tilde{\alpha}}$ , where  $\tilde{\alpha} = \max\{0, \alpha - \hat{\Delta}(\alpha)\}$ .<sup>5</sup>

<sup>5</sup>To explain this construction, recall from Section 1 that we can use the excess risk estimate to choose a tuning parameter that achieves a desired bound on risk. Specifically, for any value of  $\tilde{\alpha}$ , we can bound the risk (i.e., the miscoverage) for the CP interval  $C_{1-\tilde{\alpha}}$  as  $\mathbb{E}_{\tilde{P}_{\text{target}}}[Y \notin C_{1-\tilde{\alpha}}(X)] \leq \mathbb{E}_{\tilde{P}}[Y \notin C_{1-\tilde{\alpha}}(X)] + \Delta(R_{\tilde{\alpha}}; \mathcal{B}_{f,\rho}) \leq \tilde{\alpha} + \Delta(R_{\tilde{\alpha}}; \mathcal{B}_{f,\rho})$  (where  $R_{\tilde{\alpha}}$  is the risk defined by the CP interval  $C_{1-\tilde{\alpha}}$ , for any value of  $\tilde{\alpha}$ ). Since  $a \mapsto R_a$  is nondecreasing, this also implies that  $a \mapsto \Delta(R_a; \mathcal{B}_{f,\rho})$  is nondecreasing (recall from Section 2.1 that  $\Delta(R; \mathcal{B})$  is monotone in  $R$ , as a corollary of Proposition 2.2). Thus, for  $\tilde{\alpha} \leq \alpha$  we have  $\mathbb{E}_{\tilde{P}_{\text{target}}}[Y \notin C_{1-\tilde{\alpha}}(X)] \leq \tilde{\alpha} + \Delta(R_{\tilde{\alpha}}; \mathcal{B}_{f,\rho}) \approx \tilde{\alpha} + \hat{\Delta}(\alpha)$ . Consequently, the above choice of  $\tilde{\alpha}$  ensures that miscoverage will be (approximately) bounded by  $\alpha$ . A similar argument also holds for iso-DRL- $w_0$ .



- iso-DRL- $w_0$ : CP interval  $C_{1-\alpha_{\text{iso}}}$ , where  $\alpha_{\text{iso}} = \max\{0, \alpha - \widehat{\Delta}^{\text{iso}}(\alpha)\}$ .

## 5.1 Synthetic dataset

We start with a synthetic example, in which we fix  $n_{\text{train}} = n = n_{\text{test}} = 500$  and will vary  $n_{\text{pre}}$  to see how will the initial density ratio estimation  $w_0$  affect the result. We will consider two settings—the “well-specified” and “misspecified” settings, where model class within which  $w_0$  is estimated does, or does not, contain the true density ratio  $\frac{dP_{\text{target}}}{dP}(x)$ . Specifically, for the marginal distributions of  $X$ , we set

$$\text{Well-specified setting: } \begin{cases} \text{data distribution} & P : X \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d), \\ \text{target distribution} & P_{\text{target}} : X \sim \mathcal{N}(\mu, \mathbf{I}_d), \end{cases}$$

or

$$\text{Misspecified setting: } \begin{cases} \text{data distribution} & P : X \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d), \\ \text{target distribution} & P_{\text{target}} : X \sim \mathcal{N}(\mu, \mathbf{I}_d + \frac{\zeta}{d} \mathbf{1}_d \mathbf{1}_d^\top), \end{cases}$$

where  $d = 20$ ,  $\mu = (2/\sqrt{d}) \cdot (1, \dots, 1)^\top$ , and  $\zeta = 6$ . Since the estimate  $w_0$  for the density ratio will be fitted via logistic regression as described above, the first setting is indeed well-specified since, due to the fact that  $P$  and  $P_{\text{target}}$  have the same covariance, the logistic model is correct for the distribution shift from  $P$  to  $P_{\text{target}}$ . In contrast, the second setting is misspecified since, due to the change in covariance matrix, the underlying log-density ratio is no longer a linear function of  $\mu^\top X$ , which cannot be characterized by logistic regression.

Finally, for the conditional distribution of  $Y | X$ , we set

$$Y | X \sim 0.2 \cdot \mathcal{N}(X^\top \beta + \sin(X_1) + 0.4X_3^3 + 0.2X_4^2, 1)$$

for both training and target distributions, where  $\beta \sim \mathcal{N}(\mathbf{0}_d, \mathbf{I}_d)$ .

### 5.1.1 Results with varying sample size $n_{\text{pre}}$ for estimating $w_0$

We first consider the scenario with an estimated density ratio  $w_0$ . Recall that we use the subsets  $\mathcal{D}_1 \subset \mathcal{D}_{\text{train}}$  and  $\mathcal{D}_{\text{test},1} \subset \mathcal{D}_{\text{test}}$  with  $|\mathcal{D}_1| = |\mathcal{D}_{\text{test},1}| = n_{\text{pre}}$  for estimating  $w_0$ ; consequently, for larger values of  $n_{\text{pre}}$ , we will expect a more accurate  $w_0$  (but will then have a lower sample size for the remaining steps of the workflow). By varying  $n_{\text{pre}}$ , we aim to investigate the robustness of WCP and iso-DRL w.r.t. the accuracy in  $w_0$ . The sample size  $n_{\text{pre}}$  varies in  $\{20, 40, 60, 80, 100\}$  and we fix  $\rho = \rho^* := D_{\text{KL}}(P_{\text{target}} \| P)$ .

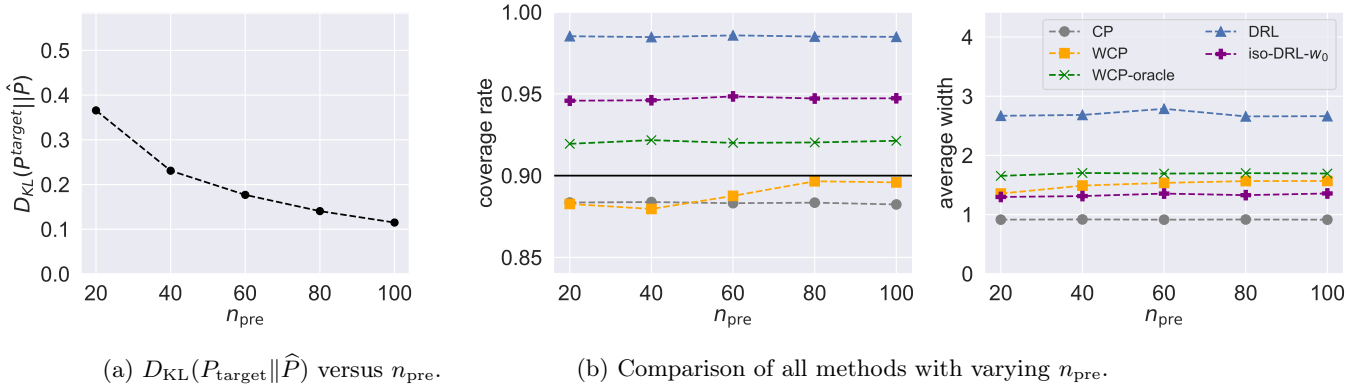


Figure 2: Results in the well-specified setting.



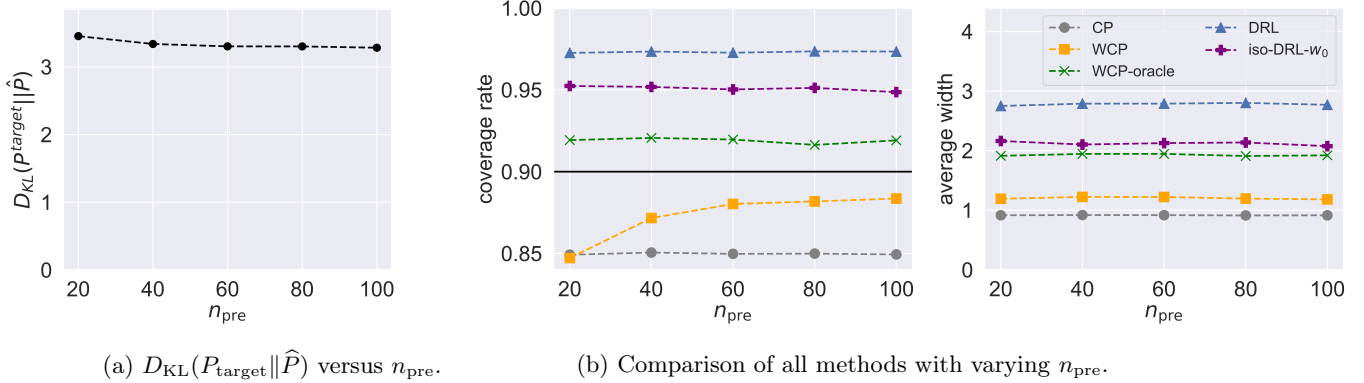


Figure 3: Results in the misspecified setting.

**Well-specified setting.** In Figure 2b, we consider the well-specified setting for generating the data. We can see that the uncorrected CP exhibits undercoverage due to the mismatch between  $P_{\text{target}}$  and  $P$  while the coverage of WCP using  $w_0$  increases to 90% as  $n_{\text{pre}}$  increases, since  $w_0$  becomes more accurate with larger  $n_{\text{pre}}$  (cf. Figure 2a). The generic DRL, even with  $\rho = \rho^*$ , tends to be conservative and has the widest interval. In comparison, iso-DRL- $w_0$  has coverage very close to the target level—indeed, the width is even shorter than WCP-oracle, due to the limited effective sample size of WCP-oracle.

**Misspecified  $w_0$ .** In Figure 3b, we show results for the misspecified setting. Since  $w_0$  is estimated from a model class that does not contain the true density ratio, consequently  $D_{\text{KL}}(P_{\text{target}} \|\hat{P})$  does not converge to zero as  $n_{\text{pre}}$  increases (cf. Figure 3a). As a result, both uncorrected CP and WCP (which is weighted with the misspecified  $w_0$ ) exhibit undercoverage. Proposed iso-DRL- $w_0$  has coverage slightly above 90% but has interval width close to that of WCP-oracle (which uses the correct weight function), while DRL is overly conservative.

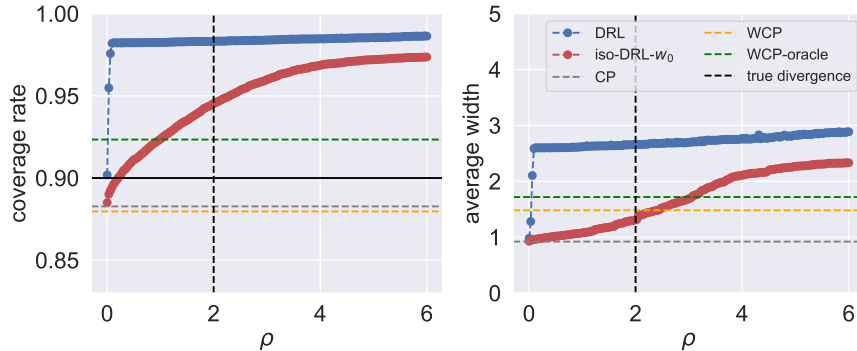


Figure 4: Results with varying  $\rho$  in the well-specified setting.

### 5.1.2 Results with varying $\rho$

In the previous section, the parameter  $\rho$ , which is used to measure the size of the distribution shift, was assumed to be known. In practice, of course, we can only estimate it. In this section, we investigate the sensitivity of each approach (DRL and iso-DRL- $w_0$ ) to the choice of  $\rho$ . Of course, the other methods considered previously (CP, WCP, and WCP-oracle) do not have  $\rho$  as an input; for comparison, we will display these methods' outputs as constant over  $\rho$ .

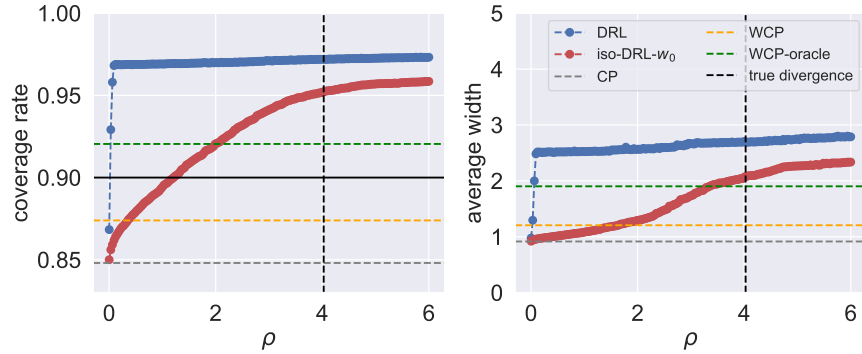


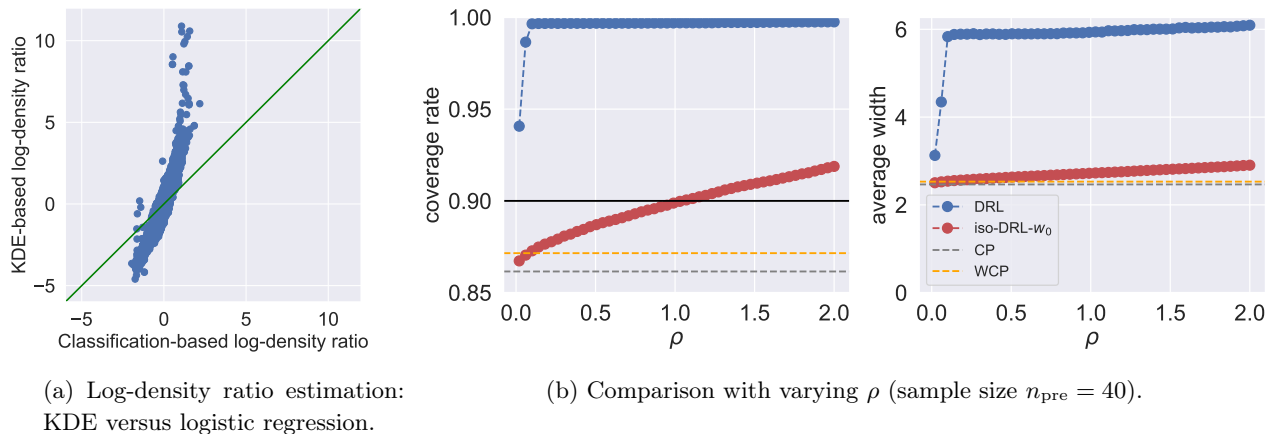
Figure 5: Results with varying  $\rho$  in the misspecified setting.

By fixing  $n_{\text{pre}} = 50$ , we vary  $\rho$  in  $[0.002, 6]$  and the underlying true KL-divergence  $\rho^* = D_{\text{KL}}(P_{\text{target}} \| P)$  is marked by the vertical dashed line. The uncorrected CP, WCP with the true density ratio and estimated density ratio  $w_0$  behave in the same way as shown in the previous section.

We can see from both plots that the prediction intervals produced by DRL is quite conservative (is much wider than the oracle interval) across nearly the entire range of  $\rho$ , even values  $\rho$  much smaller than the true distribution shift magnitude  $\rho^* = D_{\text{KL}}(P_{\text{target}} \| P)$ . In comparison, for iso-DRL- $w_0$ , when  $\rho = \rho^*$ , the width of intervals is comparable to the oracle interval in both cases, and the coverage and width vary slowly as we change the value of  $\rho$ . From this we can see that the isotonic constraint offers a significant gain in accuracy if we have a reasonable estimate of  $\rho^*$ .

## 5.2 Real data: wine quality dataset

We also consider a real dataset: the wine quality dataset (<https://archive.ics.uci.edu/dataset/186/wine+quality>). The dataset includes 12 variables that measure the physicochemical properties of wine and we treat the variable `quality` as the response of interest. The entire dataset consists of two groups: the white and red variants of the Portuguese “Vinho Verde” wine, which are unbalanced (1599 data points for the red wine and 4898 data points for the white wine). The subset of red wine is treated as the test dataset and that of white wine is viewed as the training set. All variables are nonnegative and we scale each variable by its largest value such that the entries are bounded by 1. Similar to the dataset partition in synthetic simulation, we fix  $n_{\text{pre}} = 40$ ,  $n_{\text{train}} = n = 1900$ , and  $n_{\text{test}} = 1000$ .



(a) Log-density ratio estimation: KDE versus logistic regression.

(b) Comparison with varying  $\rho$  (sample size  $n_{\text{pre}} = 40$ ).

Figure 6: Results for wine quality dataset.

We first fit a kernel density estimator (Gaussian kernel with a bandwidth suggested by cross-validation) using the entire dataset as a proxy of the oracle density ratio. Figure 6a plots this against the log-density ratio obtained from logistic regression fitted on  $n_{\text{pre}}$  samples from each group. It can be seen that the two density ratios exhibit an approximately isotonic trend. This motivates us to consider the isotonic constraint with respect to the initial density ratio estimate  $w_0$ .

To assess the performance of the proposed approach, we estimate  $w_0$  using the same procedure as for the simulated data, with sample size  $n_{\text{pre}} = 40$  for estimating the initial density ratio  $w_0$ . We consider the uncertainty set of distribution shifts defined by KL-divergence and choose  $\rho$  from 50 uniformly located grid points in  $[0.02, 2]$ . In Figure 6b, similar to the performance in Section 5.1 for simulated data, DRL tends to be conservative: the coverage rate quickly approaches 1 while  $\rho$  is still below 0.1. Moreover, the widths of intervals are nearly three times those produced by iso-DRL- $w_0$ . In the meantime, iso-DRL- $w_0$  captures the approximate isotonic trend in Figure 6a and achieves valid coverage by recalibrating the weighted approach. The key message is that in the real data case, even when there is no oracle information for selecting  $\rho$  and the isotonic trend is not exact, the proposed iso-DRL- $w_0$  with the isotonic constraint with respect to the pre-fitted density ratio is robust to the selection of  $\rho$ . Additional details for this experiment are given in Appendix D.2.

## 6 Additional related work

In this section, we discuss some additional literature in several related areas, including transfer learning, DRL, sensitivity analysis, shape-constrained learning, and conformal prediction.

**Transfer learning.** Our investigation in this paper falls into the area of transfer learning (Hu et al., 2019; Hu and Lei, 2020; Mei et al., 2010; Sun and Hu, 2016; Turki et al., 2017; Weng et al., 2020), in which data from one distribution is used to improve performance on a related but different distribution. Transfer learning is mostly studied in the supervised learning setting where we have  $(X, Y) \in \mathcal{X} \times \mathcal{Y}$ , and it is categorized into domain adaptation and inductive transfer learning (Redko et al., 2020).

Domain adaptation focuses on the scenario with covariate shift, where the conditional distribution of  $Y | X$  is assumed to be unchanged. From the theoretical side, the performance of machine learning models including hardness results is analyzed in Ben-David et al. (2010); Ben-David and Uner (2012, 2013); Johansson et al. (2019); Zhao et al. (2019a); Pathak et al. (2022); Pathak and Ma (2024), etc. The covariate shift assumption is further relaxed in Hanneke and Kpotufe (2019) to study the value of target data in adaptation. To implement efficient predictions, weighted methods are adopted as the first trial to draw  $P$  closer to  $Q$  after re-weighting the labeled samples (Cortes et al., 2008; Gretton et al., 2009; Ma et al., 2023; Ge et al., 2023). Another attempt is to require a small number of labeled target samples, which can be feasible in reality and related works include Chen et al. (2011); Chattopadhyay et al. (2013); Yang et al. (2012), etc.

For inductive transfer learning, the marginal distribution of  $X$  is assumed to be the same for training and target distributions. In the regression setting, the performance of the least square estimator with side information from the target domain is studied by Bastani (2021). The minimax theorem is further presented for nonparametric classification by Cai and Wei (2019). In the high-dimensional case, Li et al. (2021) consider transfer learning with Lasso, and Tian and Feng (2021) extend transfer learning with generalized linear models.

**Distributionally robust learning (DRL).** Our work is directly related to DRL (Ben-Tal and Nemirovski, 1998; El Ghaoui and Lebret, 1997; El Ghaoui et al., 1998), which is a popular technique in transfer learning that aims to control certain statistical risks uniformly over a set of candidate distributions for the target distribution. Different choices of the uncertainty set are studied in the literature: in one line of research, the distribution shift is measured in terms of the optimal transport discrepancy (Shafieezadeh Abadeh

et al., 2015; Blanchet and Murthy, 2019; Blanchet et al., 2019; Esfahani and Kuhn, 2015); another line of research adopts the uncertainty set defined by  $f$ -divergence (Duchi et al., 2021; Namkoong and Duchi, 2017; Duchi and Namkoong, 2018; Cauchois et al., 2020; Ai and Ren, 2024; Weiss et al., 2023).

Further constraints on the uncertainty set as the improvement of DRL are explored by Duchi et al. (2019); Setlur et al. (2023); Esteban-Pérez and Morales (2022); Liu et al. (2023), while in an earlier work, Popescu (2007) considers certain families of uncertainty sets in which distributions preserve similar structural properties. The recent work of Wang et al. (2023) considers the constraint that the unseen target distribution is a weighted average of data distribution from multiple sources. Additionally, Shapiro and Pichler (2023) propose the conditional distributional robust optimization to incorporate side information. Some recent works based on the DRL framework also focus on the quantification of stability against distribution shift, among which Gupta and Rothenhäusler (2021); Namkoong et al. (2022); Rothenhäusler and Bühlmann (2023) quantify the smallest possible divergence from  $P$  (e.g.,  $D_{\text{KL}}(Q\|P)$ ) with a fixed lower bound of the worst-case risk, which can be viewed as a dual formulation of (1).

It is also worth mentioning related works addressing distribution shift based on multicalibration (Hébert-Johnson et al., 2018; Deng et al., 2023; Kim et al., 2022), which guarantees the performance (e.g., coverage in uncertainty quantification) within certain function classes.

**Sensitivity analysis.** Sensitivity analysis is closely related to DRL but is particularly widely studied in the field of causal inference (Cornfield et al., 1959; Rosenbaum, 1987; Tan, 2006; Ding and VanderWeele, 2016; Zhao et al., 2019b; De Bartolomeis et al., 2023) with the goal of evaluating the effect of unmeasured confounders and relaxing untestable assumptions. Sensitivity models can be viewed as a specific example of constraints on distribution shift. For example, if we consider a treatment  $T \in \{0, 1\}$ , the marginal  $\Gamma$ -selection model (Tan, 2006) implies that

$$\frac{1}{\Gamma} \leq \frac{dP_{Y(1)|X,T=0}}{dP_{Y(1)|X,T=1}} \leq \Gamma,$$

which imposes the bounds constraint on the distribution shift from the data distribution  $P_{Y(1)|X,T=1}$  to the counterfactual  $P_{Y(1)|X,T=0}$ . Recent works investigate the performance of estimation, prediction, and inference under the sensitivity model from the perspective of DRL, such as the works of Yadlowsky et al. (2018); Jin et al. (2022, 2023); Sahoo et al. (2022).

**Statistical learning with shape constraints.** Our work also borrows ideas from shape-constrained learning. Shape constraints, including monotonicity, convexity and log-concavity constraints, have been used for many decades across various applications (Grenander, 1956; Schell and Singh, 1997; Matzkin, 1991). The monotonicity (or isotonic) constraint is the most common one among these. The nonstandard asymptotic behavior of estimator with the isotonic constraint is identified by Rao (1969), since which the properties of isotonic regression are well studied in the literature (Brunk et al., 1957, 1972; Zhang, 2002; Han et al., 2019; Yang and Barber, 2019; Durot and Lopuhaä, 2018; Bogdan et al., 2015; Su and Candes, 2016). Moreover, the isotonic constraint is also widely applied to calibration for distributions in regression and classification settings (Zadrozny and Elkan, 2002; Niculescu-Mizil and Caruana, 2012; van der Laan et al., 2023; Henzi et al., 2021; Berta et al., 2024).

**Conformal prediction.** One important application of our distributionally robust risk evaluation with an isotonic constraint is to recalibrate prediction intervals from conformal prediction. Conformal prediction, proposed by Vovk et al. (2005); Shafer and Vovk (2008), provides a framework for distribution-free uncertainty quantification, which constructs confidence intervals that are valid with exchangeable data from any underlying distribution and with any “black-box” algorithm. When covariate shift is present between training and target distributions, Tibshirani et al. (2019) firstly introduce the notion of weighted exchangeability and the weighted conformal prediction approach to maintain validity with the oracle information of the density

ratio. However, with the estimated density ratio, the validity of WCP only holds up to a coverage gap (Lei and Candès, 2020; Candès et al., 2023; Gui et al., 2024); building on this, Jin et al. (2023) further establish a robust guarantee via sensitivity analysis. Besides the weighted approaches, there are other solutions in the literature: Cauchois et al. (2020); Ai and Ren (2024) address the issue of joint distribution shift via the DRL; Qiu et al. (2023); Yang et al. (2024); Chen and Lei (2024) formulate the covariate shift problem within the semiparametric/nonparametric framework and utilize the doubly-robust theory to correct the distributional bias.

## 7 Discussion

In this paper, we focus on distributionally robust risk evaluation with the isotonic constraint on the density ratio. We provide an efficient approach to solve the shape-constrained optimization problem via an equivalent reformulation. Estimation error bounds for the worst-case excess risk are also provided when only noisy observations of the risk function can be accessed.

To conclude, we provide further discussions on the proposed iso-DRL framework and single out several open questions.

**Isotonic constraint as regularization on distribution shift.** The isotonic constraint on the density ratio, which is the key difference between DRL and iso-DRL, is related to regularization on distribution shifts. The worst-case density ratio for the generic DRL will always align with the risk function even when the risk is highly non-smooth, which results in over-conservativeness. By adding an isotonic constraint, we aim to avoid over-pessimistic choices of the density ratio. This is similar in flavor to many tools in high-dimensional statistical learning, where regularization/inductive bias is introduced to improve generalization. More broadly, how to explicitly quantify the validity-accuracy tradeoff under distribution shift is an important problem.

**Stability against distribution shift.** Excess risk can also be interpreted from the perspective of stability against distribution shift (Lam, 2016; Namkoong et al., 2022). Suppose we have a fixed budget  $\varepsilon \ll 1$  for the excess risk, it is of interest to characterize the largest tolerance of distribution shift such that the excess risk is under control. Taking the  $f$ -constrained problem as an example, if we aim at the budget  $\Delta_\rho(R; \mathcal{B}_{f,\rho}) \leq \varepsilon \ll 1$ , then only an infinitesimal  $\rho$  that is quadratic in  $\varepsilon$  will be allowed (Lam, 2016; Duchi and Namkoong, 2018; Blanchet and Shapiro, 2023), i.e.,  $\rho$  needs to obey

$$\rho \leq \frac{f''(1)}{2\text{Var}(R(X))} \cdot \varepsilon^2 + o(\varepsilon^2).$$

However, with the additional isotonic constraint on the density ratio and the same budget  $\varepsilon$ , we can tolerate larger distribution shift:

$$\rho \leq \frac{f''(1)}{2\text{Var}([\pi(R)](X))} \cdot \varepsilon^2 + o(\varepsilon^2).$$

(Note that the variance in the denominator may be substantially smaller here—for instance, if  $R(X)$  is uncorrelated with  $X$ , we might have  $\text{Var}(R(X)) \asymp 1$  but  $\text{Var}([\pi(R)](X)) \approx 0$ .) This improvement drives the following findings:

1. When side information of the underlying distribution shift is provided, e.g., the shape constraints of the density ratio, risk evaluation will be less sensitive to the hyperparameters describing the uncertainty set (e.g.,  $\rho$ ), thus is more robust with the presence of distribution shift.
2. Moreover, the denominator  $\text{Var}([\pi(R)](X))$  also implies that when the shape of the uncertainty set is well-designed such that the projected risk  $[\pi(R)](X)$  has small variance, then the out-of-sample risk

within the uncertainty set will be more distributionally robust. Thus, it remains an open question on how to construct the variance-reduction projection and design the uncertainty set based on noisy observations of the risk function.

**From risk evaluation to distributionally robust optimization.** Different from risk evaluation, distributionally robust optimization (DRO) focuses on the optimization problem with a loss function  $\ell_\theta(x)$ :

$$\hat{\theta} \in \operatorname{argmin}_{\theta \in \Theta} \sup_{Q \in \mathcal{Q}} \mathbb{E}_Q \ell_\theta(X).$$

Under smoothness conditions on  $\ell_\theta$ , asymptotic normality for  $\hat{\theta}$  is established in the literature (Duchi and Namkoong, 2018). The DRO framework is shown to regularize  $\hat{\theta}$  in terms of variance penalization (Lam, 2016; Duchi and Namkoong, 2018) or explicit norm regularization (Blanchet and Murthy, 2019). It is interesting to incorporate the isotonic constraint into DRO and to understand the effect of the isotonic constraint in the asymptotics of  $\hat{\theta}^{\text{iso}}$ .

**Extension to the optimal transport discrepancy.** Finally, we should note that there is a rich literature on DRL with the optimal transport discrepancy, in which case the distribution shift cannot be simply represented by density ratios (Shafieezadeh Abadeh et al., 2015; Blanchet and Murthy, 2019; Blanchet et al., 2019; Esfahani and Kuhn, 2015). Suppose we have side information about the functional  $\sigma(P, P_{\text{target}})$  of two distributions, of which  $w_0 = dP_{\text{target}}/dP$  is an example, it will be an open question regarding how to utilize  $\sigma(P, P_{\text{target}})$  in guiding the constraint on the candidate distributions or the choice of the cost function in the optimal transport discrepancy.

## Acknowledgements

R.F.B. was supported by the Office of Naval Research via grant N00014-20-1-2337, and by the National Science Foundation via grant DMS-2023109. C.M. was partially supported by the National Science Foundation via grant DMS-2311127.

# Contents

<b>A Proofs of results in Section 2</b>	<b>23</b>
A.1 Proof of Proposition 2.2	23
<b>B Proofs of results in Section 3</b>	<b>24</b>
B.1 Proof of Theorem 3.1	25
B.2 Proof of Proposition 3.3	27
B.3 Proof of Proposition B.1	28
<b>C Proofs of results in Section 4</b>	<b>29</b>
C.1 Proof of Theorem 4.1	29
C.2 Proof of Theorem 4.2 with $\mathcal{B} = \mathcal{B}_{a,b}$	31
C.2.1 Explicit rate of convergence	33
C.2.2 Proof of Lemma C.1	34
C.2.3 Proof of Lemma C.2	34
C.3 Proof of Theorem 4.2 with $\mathcal{B} = \mathcal{B}_{f,\rho}$	34
C.3.1 Explicit rate of convergence	37
C.3.2 Proof of Lemma C.3	37
C.3.3 Proof of Lemma C.4	38
C.3.4 Proof of Lemma C.5	39
C.4 Hardness of consistent estimation without isotonic constraints	40
C.4.1 Proof of Proposition C.6	40
<b>D Additional simulation results</b>	<b>41</b>
D.1 iso-DRL under componentwise order	41
D.2 Wine quality simulation: a proxy of the oracle KL-divergence	42

## A Proofs of results in Section 2

### A.1 Proof of Proposition 2.2

It is straightforward to check that  $\Delta(R; \mathcal{B})$  is always an upper bound of the new formulation stated in Proposition 2.2, simply by taking  $w = \phi \circ R$ . Therefore, it remains to show the converse:  $\Delta(R; \mathcal{B})$  is also a *lower* bound of the new formulation stated in Proposition 2.2.

To this end, it suffices to prove that for any  $w_{\#}P \in \mathcal{B}$ , there exists a nondecreasing function  $\phi$  such that  $(\phi \circ R)_{\#}P \in \mathcal{B}$ , and

$$\mathbb{E}_P [w(X)R(X)] \leq \mathbb{E}_P [\phi(R(X))R(X)].$$

We construct such a function  $\phi$  in two steps.

**Step 1: conditioning.** For any  $w$  such that  $w_{\#}P \in \mathcal{B}$ , we define  $g$  as a measurable function satisfying

$$g(R(X)) = \mathbb{E}[w(X) \mid R(X)], \quad P\text{-almost surely.}$$

(Note that  $g$  is not necessarily a monotone function.) As a result, by the tower law, we have

$$\mathbb{E}_P [w(X)R(X)] = \mathbb{E}_P [g(R(X))R(X)]. \tag{16}$$

Since  $w_{\#}P \in \mathcal{B}$ , by Jensen's inequality, for any convex function  $\psi$ , we have

$$\mathbb{E}_P[\psi(g(R(X)))] = \mathbb{E}[\psi(\mathbb{E}[w(X) | R(X)])] \leq \mathbb{E}_P[\psi(w(X))],$$

which implies  $(g \circ R)_{\#}P \in \mathcal{B}$  by Condition 2.1.

**Step 2: rearrangement.** Denote  $F_1$  and  $F_2$  as the cumulative distribution functions of  $g(R(X))$  and  $R(X)$ , respectively. Let  $U \sim \text{Unif}([0, 1])$ . Then, we have  $F_1^{-1}(U) \stackrel{d}{=} g(R(X))$  and  $F_2^{-1}(U) \stackrel{d}{=} R(X)$ , where  $F_k^{-1}$  is the generalized inverse of  $F_k$ ,  $k = 1, 2$ . Moreover,  $F_1^{-1}$  is nondecreasing and

$$g(F_2^{-1}(U)) \stackrel{d}{=} g(R(X)) \stackrel{d}{=} F_1^{-1}(U),$$

which implies that  $F_1^{-1}$  is the monotone rearrangement of  $g \circ F_2^{-1}$ . By inequality (378) in Hardy et al. (1952), we have

$$\mathbb{E}_P[g(R(X))R(X)] = \mathbb{E}[g(F_2^{-1}(U))F_2^{-1}(U)] \leq \mathbb{E}[F_1^{-1}(U)F_2^{-1}(U)]. \quad (17)$$

Next, let  $\phi$  be a measurable function satisfying

$$\phi(F_2^{-1}(U)) = \mathbb{E}[F_1^{-1}(U) | F_2^{-1}(U)],$$

almost surely with respect to the distribution  $U \sim \text{Unif}([0, 1])$ . Since  $F_k^{-1}$  is the generalized inverse of a CDF  $F_k$ , for each  $k = 1, 2$ , it is therefore monotone nondecreasing. Therefore, we can choose  $\phi$  to be a monotone nondecreasing function. Moreover, to verify that  $(\phi \circ R)_{\#}P \in \mathcal{B}$ , we will check that  $\phi(R(X)) \stackrel{cux}{\preceq} g(R(X))$  (and use Condition 2.1, along with the fact that  $(g \circ R)_{\#}P \in \mathcal{B}$  as established above). For any convex function  $\psi$ , we have

$$\mathbb{E}_P[\psi(\phi(R(X)))] \stackrel{d}{=} \mathbb{E}[\psi(\phi(F_2^{-1}(U)))] = \mathbb{E}[\psi(\mathbb{E}[F_1^{-1}(U) | F_2^{-1}(U)])] \leq \mathbb{E}[\psi(F_1^{-1}(U))] = \mathbb{E}_P[\psi(g(R(X)))],$$

where the inequality holds by Jensen's inequality.

We then have

$$\begin{aligned} \mathbb{E}[F_1^{-1}(U)F_2^{-1}(U)] &= \mathbb{E}[\mathbb{E}[F_1^{-1}(U) | F_2^{-1}(U)] F_2^{-1}(U)] \\ &= \mathbb{E}[\phi(F_2^{-1}(U))F_2^{-1}(U)] = \mathbb{E}_P[\phi(R(X))R(X)]. \end{aligned}$$

This equality, combined with (16) and (17), yields the desired outcome:  $\mathbb{E}_P[w(X)R(X)] \leq \mathbb{E}_P[\phi(R(X))R(X)]$ . We hence complete the proof.

## B Proofs of results in Section 3

To ease notation, we denote  $\langle a, b \rangle_P = \int_{\mathcal{X}} a(x)b(x)dP(x)$  for any functions  $a, b$ .

Before proceeding to the proof of Theorem 3.1, denote

$$\begin{aligned} \Delta_2(R; \mathcal{B}) &= \sup_{w \geq 0, w \in L_2(P)} \mathbb{E}_P[w(X)R(X)] - \mathbb{E}_P[R(X)] \\ &\text{subject to } w_{\#}P \in \mathcal{B}. \end{aligned} \quad (18)$$

The only difference between this new definition, and the quantity  $\Delta(R; \mathcal{B})$  defined in (4), is that we now add the constraint  $w \in L_2(P)$ —the  $L_2$  constraint will be useful for some of our theoretical results in this section. In fact, the following result shows that adding the  $L_2$  constraint does not change the optimal value.



**Proposition B.1.** *Under the notation and definitions above, it holds that  $\Delta(R; \mathcal{B}) = \Delta_2(R; \mathcal{B})$ .*

We defer the proof of this proposition to Section B.3.

In addition, by the property of the projection onto a closed and convex set (Theorem 3.14 in [Bauschke and Combettes \(2019\)](#), Proposition 1.12.4 in [Edwards \(2012\)](#)) and the fact that constant functions  $\pm 1 \in \mathcal{C}_{\geq}^{\text{iso}}$ , we have

$$\langle R - \pi(R), 1 \rangle_P \leq 0 \quad \text{and} \quad \langle R - \pi(R), -1 \rangle_P \leq 0,$$

which implies that  $\langle R - \pi(R), 1 \rangle_P = \mathbb{E}_P[[\pi(R)](X)] - \mathbb{E}_P[R(X)] = 0$ .

## B.1 Proof of Theorem 3.1

We split the proof into three steps:

1. prove that  $\Delta^{\text{iso}}(R; \mathcal{B}) \leq \Delta(\pi(R); \mathcal{B})$ ;
2. prove that  $\Delta^{\text{iso}}(R; \mathcal{B}) = \Delta(\pi(R); \mathcal{B})$  provided that Condition 2.1 holds;
3. prove the claim on attainability of minimizers.

**Step 1: Prove  $\Delta^{\text{iso}}(R; \mathcal{B}) \leq \Delta(\pi(R); \mathcal{B})$ .** By the definition of  $\Delta^{\text{iso}}(R; \mathcal{B})$  as a supremum, for any  $\varepsilon > 0$ , there exists  $w_\varepsilon \in \mathcal{C}_{\geq}^{\text{iso}}$  such that

$$\Delta^{\text{iso}}(R; \mathcal{B}) - \varepsilon \leq \langle w_\varepsilon, R \rangle_P - \mathbb{E}_P[R(X)] \leq \Delta^{\text{iso}}(R; \mathcal{B}). \quad (19)$$

By the properties of projection onto a convex and closed cone  $\mathcal{C}_{\geq}^{\text{iso}}$ , we have  $\mathbb{E}_P[[\pi(R)](X)] = \mathbb{E}_P[R(X)]$  and

$$\langle w_\varepsilon - \pi(R), R - \pi(R) \rangle_P = \langle w_\varepsilon, R - \pi(R) \rangle_P \leq 0. \quad (20)$$

Clearly we also have  $\langle w_\varepsilon, \pi(R) \rangle_P \leq \Delta(\pi(R); \mathcal{B})$ . This combined with (19), (20), and the property  $\mathbb{E}_P[[\pi(R)](X)] = \mathbb{E}_P[R(X)]$  yields

$$\Delta^{\text{iso}}(R; \mathcal{B}) - \varepsilon \leq \langle w_\varepsilon, R \rangle_P - \mathbb{E}_P[R(X)] \leq \langle w_\varepsilon, \pi(R) \rangle_P - \mathbb{E}_P[R(X)] \leq \Delta(\pi(R); \mathcal{B}).$$

Since  $\varepsilon > 0$  is arbitrary, we obtain the desired result  $\Delta^{\text{iso}}(R; \mathcal{B}) \leq \Delta(\pi(R); \mathcal{B})$ .

**Step 2: Prove  $\Delta^{\text{iso}}(R; \mathcal{B}) = \Delta(\pi(R); \mathcal{B})$ .** With a similar proof with that of Proposition 2.2 (see Section A.1), we have the following equivalent formulation of  $\Delta_2(R; \mathcal{B})$ :

$$\begin{aligned} \Delta_2(R; \mathcal{B}) &= \sup_{\phi: \mathbb{R} \rightarrow \mathbb{R}_+} \mathbb{E}_P[(\phi \circ R)(X)R(X)] - \mathbb{E}_P[R(X)] \\ &\text{subject to} \quad (\phi \circ R)_{\#} P \in \mathcal{B}, \quad \phi \circ R \in L_2(P), \quad \phi \text{ is nondecreasing.} \end{aligned} \quad (21)$$

Accordingly, we also have  $\Delta(\pi(R); \mathcal{B}) = \Delta_2(\pi(R); \mathcal{B})$ . Then, by (21) and the definition of supremum, for any  $\varepsilon > 0$ , there exists  $\tilde{w}_\varepsilon = g_\varepsilon \circ [\pi(R)] \in L_2(P)$  where  $g_\varepsilon(t)$  is nondecreasing in  $t$  such that

$$\Delta(\pi(R); \mathcal{B}) - \varepsilon = \Delta_2(\pi(R); \mathcal{B}) - \varepsilon \leq \langle \tilde{w}_\varepsilon, \pi(R) \rangle_P - \mathbb{E}_P[R(X)] \leq \Delta_2(\pi(R); \mathcal{B}) = \Delta^{\text{iso}}(R; \mathcal{B}). \quad (22)$$

More importantly, since  $\pi(R) \in \mathcal{C}_{\geq}^{\text{iso}}$  and  $g_\varepsilon \in \mathcal{C}_{\geq}^{\text{iso}}$ , we know that  $\tilde{w}_\varepsilon \in \mathcal{C}_{\geq}^{\text{iso}}$ , which implies that

$$\langle \tilde{w}_\varepsilon, R \rangle_P - \mathbb{E}_P[R(X)] \leq \Delta^{\text{iso}}(R; \mathcal{B}). \quad (23)$$

In addition, assume for the moment that we can represent the isotonic cone  $\mathcal{C}_{\geq}^{\text{iso}}$  using a  $\sigma$ -lattice—a claim we will prove in the end this section. We can then invoke Theorem 1 in [Brunk \(1963\)](#) (or Corollary 3.1 in [Brunk \(1965\)](#)) to conclude that

$$\langle h(\pi(R)), R - \pi(R) \rangle_P = 0, \quad \text{for all measurable functions } h \text{ such that } h(\pi(R)) \in L_2(P).$$

In particular, since  $\tilde{w}_\varepsilon = g_\varepsilon \circ [\pi(R)] \in L_2(P)$ , we have

$$\langle \tilde{w}_\varepsilon, R - \pi(R) \rangle_P = 0. \quad (24)$$

Combining [\(22\)](#), [\(23\)](#), and [\(24\)](#), we arrive at the conclusion

$$\Delta(\pi(R); \mathcal{B}) - \varepsilon \leq \langle \tilde{w}_\varepsilon, \pi(R) \rangle_P - \mathbb{E}_P[R(X)] = \langle \tilde{w}_\varepsilon, R \rangle_P - \mathbb{E}_P[R(X)] \leq \Delta^{\text{iso}}(R; \mathcal{B}). \quad (25)$$

Since  $\varepsilon > 0$  is arbitrary, we have  $\Delta(\pi(R); \mathcal{B}) \leq \Delta^{\text{iso}}(R; \mathcal{B})$ . In all, we combine Steps 1 and 2 to finish the proof.

**Step 3: attainability of minimizers.** Suppose  $\Delta(\pi(R); \mathcal{B})$  is attained at  $\tilde{w}$ . Then, it holds that

$$\langle \tilde{w}, \pi(R) \rangle_P - \mathbb{E}_P[R(X)] = \Delta(\pi(R); \mathcal{B}) = \Delta^{\text{iso}}(R; \mathcal{B}).$$

Moreover, by [Proposition 2.2](#), without loss of generality, we assume that  $\tilde{w} \in \mathcal{C}_{\geq}^{\text{iso}}$ .

By [Proposition B.1](#), for any  $\eta > 0$ , there exists  $\tilde{w}_\eta \in L_2(P)$  such that

$$\max \left\{ \left| \langle \tilde{w}, \pi(R) \rangle_P - \langle \tilde{w}_\eta, \pi(R) \rangle_P \right|, \left| \langle \tilde{w}, R \rangle_P - \langle \tilde{w}_\eta, R \rangle_P \right| \right\} \leq \eta.$$

By the formulation [\(21\)](#), we assume that  $\tilde{w}_\eta \in \mathcal{C}_{\geq}^{\text{iso}}$ . In addition, by Theorem 1 in [Brunk \(1963\)](#) (or Corollary 3.1 in [Brunk \(1965\)](#)), we further have  $\langle \tilde{w}_\eta, R - \pi(R) \rangle_P = 0$ , thus

$$\left| \langle \tilde{w}, R \rangle_P - \mathbb{E}_P[R(X)] - \Delta^{\text{iso}}(R; \mathcal{B}) \right| \leq \eta.$$

Since  $\eta$  is arbitrary,  $\Delta^{\text{iso}}(R; \mathcal{B})$  is also attained at  $\tilde{w}$ .

On the other hand, suppose  $\Delta^{\text{iso}}(R; \mathcal{B})$  is attained at  $w^{*\text{iso}}$ . It holds that

$$\langle w^{*\text{iso}}, R \rangle_P - \mathbb{E}_P[R(X)] = \Delta(\pi(R); \mathcal{B}) = \Delta^{\text{iso}}(R; \mathcal{B}).$$

Suppose there exists  $\varepsilon > 0$  such that  $\langle w^{*\text{iso}}, \pi(R) \rangle_P - \mathbb{E}_P[R(X)] \leq \Delta(\pi(R); \mathcal{B}) - \varepsilon$ . By the definition of supremum together with [Proposition B.1](#) and the equivalent formulation in [\(21\)](#), there exists  $\tilde{w}_\varepsilon = g_\varepsilon \circ R \in L_2(P)$ , where  $g_\varepsilon(t)$  is nondecreasing in  $t$ , such that

$$\langle \tilde{w}_\varepsilon, \pi(R) \rangle_P > \Delta(\pi(R); \mathcal{B}) - \varepsilon.$$

By the property of projection and the optimality of  $w^{*\text{iso}}$ , we obtain

$$\langle \tilde{w}_\varepsilon, R \rangle_P \leq \langle w^{*\text{iso}}, R \rangle_P \leq \langle w^{*\text{iso}}, \pi(R) \rangle_P < \langle \tilde{w}_\varepsilon, \pi(R) \rangle_P.$$

[Brunk \(1963, 1965\)](#) further gives us  $\langle g_\varepsilon \circ R, R - \pi(R) \rangle_P = 0$ , which yields  $\langle w^{*\text{iso}}, R \rangle_P = \langle \tilde{w}_\varepsilon, \pi(R) \rangle_P > \Delta(\pi(R); \mathcal{B}) - \varepsilon$  and draws the contradiction.

**Representation of the isotonic cone with  $\sigma$ -lattice.** For  $\mathcal{X} \subseteq \mathbb{R}^d$  with  $d \geq 1$ , consider any partial order  $\preceq$  on  $\mathcal{X}$ . Recall the isotonic cone

$$\mathcal{C}_{\preceq}^{\text{iso}} = \{g \mid g(x) \leq g(x') \text{ for all } x \preceq x'\}.$$

Consider a class of subsets  $\mathcal{A} \subseteq \mathbb{R}^d$  such that for each  $A \in \mathcal{A}$ , we have  $\{x \mid a \preceq x\} \subseteq A$  for any  $a \in A$ . Denote the set

$$\mathcal{R}(\mathcal{A}) = \{g \mid \{x \mid g(x) > \tau\} \in \mathcal{A} \text{ for all } \tau \in \mathbb{R}\}.$$

One can verify that  $\mathcal{A}$  is closed under countable union and countable intersection, thus  $\mathcal{R}(\mathcal{A})$  forms a sigma lattice. For any  $g \in \mathcal{C}_{\preceq}^{\text{iso}}$  and any  $\tau \in \mathbb{R}$ , we have

$$\{x \mid g(x) > \tau\} \in \mathcal{A} \quad \implies \quad \mathcal{C}_{\preceq}^{\text{iso}} \subseteq \mathcal{R}(\mathcal{A}).$$

On the other hand, for any  $g \in \mathcal{R}(\mathcal{A})$ . Suppose  $g \notin \mathcal{C}_{\preceq}^{\text{iso}}$ , i.e. there exists  $x_1 \preceq x_2$  such that  $g(x_1) > g(x_2)$ . By choosing  $\tau = (g(x_1) + g(x_2))/2$ , we have  $A_\tau := \{x \mid g(x) > \tau\} \in \mathcal{A}$  and  $x_1 \in A_\tau$ ,  $x_2 \notin A_\tau$ . However, since  $x_1 \preceq x_2$ , we have  $\{x \mid x_1 \preceq x\} \not\subseteq A_\tau$  which draws the contradiction to  $A_\tau \in \mathcal{A}$ . Hence, we have verified that  $\mathcal{C}_{\preceq}^{\text{iso}} = \mathcal{R}(\mathcal{A})$ .

## B.2 Proof of Proposition 3.3

The proof relies on the following property of the set  $\mathcal{B}$ , which we establish in the end of this section.

**Lemma B.2.** *Assume Condition 2.1 holds. The set  $\mathcal{B}$  is closed under the isotonic projection, that is, for any  $w_{\#}P \in \mathcal{B}$ , it holds  $\pi(w)_{\#}P \in \mathcal{B}$ .*

Recall that  $\tilde{w}^*$  is the underlying density ratio  $dP_{\text{target}}/dP$ . Since  $\mathcal{B}$  is closed under  $\pi$  by Lemma B.2, we have  $\pi(\tilde{w}^*) \in \mathcal{C}_{\preceq}^{\text{iso}}$  and  $\pi(\tilde{w}^*)_{\#}P \in \mathcal{B}$ . By optimality, we have

$$\Delta^{\text{iso}}(R; \mathcal{B}) \geq \mathbb{E}_P \left[ [\pi(\tilde{w}^*)](X)R(X) \right] - \mathbb{E}_P[R(X)].$$

**Case 1.** We first assume that  $\Delta^{\text{iso}}(R; \mathcal{B}) > \mathbb{E}_P [[\pi(\tilde{w}^*)](X)R(X)] - \mathbb{E}_P[R(X)]$ . For any  $0 < \varepsilon < \Delta^{\text{iso}}(R; \mathcal{B}) - \mathbb{E}_P [[\pi(\tilde{w}^*)](X)R(X)] + \mathbb{E}_P[R(X)]$ , by the definition of  $\Delta^{\text{iso}}(R; \mathcal{B})$  and the definition of supremum, there exists  $w_\varepsilon^{*\text{iso}} \in \mathcal{C}_{\preceq}^{\text{iso}}$  such that

$$\Delta^{\text{iso}}(R; \mathcal{B}) - \varepsilon \leq \mathbb{E}_P[w_\varepsilon^{*\text{iso}}(X)R(X)] - \mathbb{E}_P[R(X)] \leq \Delta^{\text{iso}}(R; \mathcal{B}), \quad (w_\varepsilon^{*\text{iso}})_{\#}P \in \mathcal{B}.$$

Then, by the choice of  $\varepsilon$ , we have

$$\mathbb{E}_P[w_\varepsilon^{*\text{iso}}(X)R(X)] - \mathbb{E}_P[R(X)] \geq \Delta^{\text{iso}}(R; \mathcal{B}) - \varepsilon \geq \mathbb{E}_P \left[ [\pi(\tilde{w}^*)](X)R(X) \right] - \mathbb{E}_P[R(X)],$$

which yields

$$\begin{aligned} \Delta^{\text{iso}}(R; \mathcal{B}) - \varepsilon &\geq \mathbb{E}_P[\tilde{w}^*(X)R(X)] - \mathbb{E}_P[R(X)] + \left\{ \mathbb{E}_P \left[ [\pi(\tilde{w}^*)](X)R(X) \right] - \mathbb{E}_P[\tilde{w}^*(X)R(X)] \right\} \\ &= \Delta^*(R) - \mathbb{E}_P \left[ [\tilde{w}^* - \pi(\tilde{w}^*)](X)R(X) \right]. \end{aligned}$$

In addition, by the property of the projection onto a convex and closed cone, we have

$$\mathbb{E}_P \left[ [\tilde{w}^* - \pi(\tilde{w}^*)](X)[\pi(R)](X) \right] \leq 0,$$

which further yields

$$\begin{aligned}\Delta^*(R) + \varepsilon &\leq \Delta^{\text{iso}}(R; \mathcal{B}) + \mathbb{E}_P [[\tilde{w}^* - \pi(\tilde{w}^*)](X)R(X)] - \mathbb{E}_P [[\tilde{w}^* - \pi(\tilde{w}^*)](X)[\pi(R)](X)] \\ &= \Delta^{\text{iso}}(R; \mathcal{B}) + \mathbb{E}_P [[\tilde{w}^* - \pi(\tilde{w}^*)](X) \cdot [R - \pi(R)](X)].\end{aligned}$$

Since  $\varepsilon > 0$  is arbitrary, we obtain

$$\Delta^*(R) \leq \Delta^{\text{iso}}(R; \mathcal{B}) + \mathbb{E}_P [[\tilde{w}^* - \pi(\tilde{w}^*)](X) \cdot [R - \pi(R)](X)]. \quad (26)$$

**Case 2.** If the equality holds such that

$$\Delta^{\text{iso}}(R; \mathcal{B}) = \mathbb{E}_P [[\pi(\tilde{w}^*)](X)R(X)] - \mathbb{E}_P [R(X)],$$

it implies that worst-case excess risk  $\Delta^{\text{iso}}(R; \mathcal{B})$  can be attained by some  $w^{*\text{iso}}$ . Then, we can follow the same proof as above to show (26).

In particular, if either  $\tilde{w}^* \in \mathcal{C}_{\leq}^{\text{iso}}$  or  $R \in \mathcal{C}_{\leq}^{\text{iso}}$ , it holds that  $\Delta^*(R) \leq \Delta^{\text{iso}}(R; \mathcal{B})$ .

**Proof of Lemma B.2.** Since  $\mathcal{B}$  is closed under convex ordering, it suffices to prove that  $\pi(w)_{\#}P \stackrel{cvx}{\preceq} w_{\#}P$ . To this end, fix any convex function  $\psi$ . We know the expectation of the Bregman divergence is nonnegative, i.e.,

$$\langle \psi(w) - \psi(\pi(w)), 1 \rangle_P - \langle \psi'(\pi(w)), w - \pi(w) \rangle_P \geq 0.$$

Moreover, due to the monotonicity of  $\psi'$  and the orthogonality property mentioned before in (Brunk, 1963, 1965), it holds that

$$\langle \psi'(\pi(w)), w - \pi(w) \rangle_P = 0.$$

Consequently, we obtain

$$\langle \psi(\pi(w)), 1 \rangle_P \leq \langle \psi(w), 1 \rangle_P.$$

This implies  $\pi(w)_{\#}P \in \mathcal{B}$ .

### B.3 Proof of Proposition B.1

For any  $w \geq 0$ , we define the sequence of truncated functions  $\{w_n\}_{n \in \mathbb{N}}$  via

$$w_n(x) = w(x) \cdot \mathbb{1}\{w(x) \leq n\} + L_n \cdot \mathbb{1}\{w(x) > n\},$$

where  $L_n = \mathbb{E}[w(X) \mid w(X) > n]$ , which implies that  $\mathbb{E}_P[w_n(X)] = 1$  and, since  $\max\{n, L_n\} = L_n < \infty$ ,  $w_n \in L_2(P)$  for each  $n \geq 1$ .

**Step 1: feasibility of  $w_n$ .** We first prove the feasibility of  $w_n$ . To see this, as  $\mathbb{E}_P[w_n(X)] = 1$  by construction, we need to show that  $(w_n)_{\#}P \in \mathcal{B}$ . By Condition 2.1, since  $\mathcal{B}$  is closed under the convex ordering, it suffices to show that

$$\mathbb{E}_P [\psi(w_n(X))] \leq \mathbb{E}_P [\psi(w(X))] \quad \text{for any convex function } \psi.$$

This is true by Jensen's inequality, since, by construction,  $\mathbb{E}_P[w(X) \mid w_n(X)] = w_n(X)$ .

**Step 2: convergence of  $\mathbb{E}_P[w_n(X)R(X)]$ .** To verify the convergence of  $\mathbb{E}_P[w_n(X)R(X)]$ , consider

$$\begin{aligned}
& \left| \mathbb{E}_P[w_n(X)R(X)] - \mathbb{E}_P[w(X)R(X)] \right| \\
&= \left| \int_{w(x) > n} (L_n - w(x))R(x) dP(x) \right| \\
&\leq B_R \int_{w(x) > n} |L_n - w(x)| dP(x) \\
&\leq B_R \left( \int_{w(x) > n} w(x) dP(x) + L_n \mathbb{P}(w(X) > n) \right) \\
&= 2\mathbb{E}_P[w(X) \cdot \mathbb{1}\{w(X) > n\}].
\end{aligned}$$

Finally, since  $\mathbb{E}_P[w(X)] = 1$  (i.e., we know that  $w \in L_1(P)$ ), this means that  $\lim_{n \rightarrow \infty} \mathbb{E}_P[w(X) \cdot \mathbb{1}\{w(X) > n\}] = 0$ .

**Summary.** For any  $\varepsilon > 0$ , there exists  $w \geq 0$  such that  $\mathbb{E}_P[w(X)] = 1$ ,  $w \# P \in \mathcal{B}$ , and

$$\mathbb{E}_P[w(X)R(X)] - \mathbb{E}_P[R(X)] \geq \Delta(R; \mathcal{B}) - \varepsilon/2$$

Then, based on the previous steps, there exist  $\tilde{w} \in L_2(P)$  such that  $\tilde{w}$  is feasible for (4) and  $\mathbb{E}_P[\tilde{w}(X)R(X)] \geq \mathbb{E}_P[w(X)R(X)] - \varepsilon/2$ , which implies that

$$\mathbb{E}_P[\tilde{w}(X)R(X)] \geq \Delta(R; \mathcal{B}) - \varepsilon \quad \implies \quad \Delta_2(R; \mathcal{B}) \geq \Delta(R; \mathcal{B}) - \varepsilon.$$

Since  $\varepsilon$  is arbitrary and  $\Delta_2(R; \mathcal{B}) \leq \Delta(R; \mathcal{B})$ , we have shown that  $\Delta_2(R; \mathcal{B}) = \Delta(R; \mathcal{B})$ , which completes the proof.

## C Proofs of results in Section 4

### C.1 Proof of Theorem 4.1

To ease the notations, denote  $r_i = r(X_i, Y_i)$  and  $R_i = R(X_i)$ . We abuse the notation and denote  $\hat{w}^R = (\hat{w}_i^R)_{i \leq n} = (\hat{w}_{R, \Omega}^{\text{iso}}(X_i))_{i \leq n}$  as the maximizer to (12) and  $\hat{w}^r = (\hat{w}_i^r)_{i \leq n} = (\hat{w}_{r, \Omega}^{\text{iso}}(X_i))_{i \leq n}$  as the maximizer to (10), which are both attainable due to the compactness of a truncated feasible set. Then we have

$$\begin{aligned}
\hat{\Delta}_{\Omega}^{\text{iso}}(R; \mathcal{B}) - \hat{\Delta}_{\Omega}^{\text{iso}}(r; \mathcal{B}) &= \frac{1}{n} \sum_{i \leq n} \hat{w}_i^R R_i - \frac{1}{n} \sum_{i \leq n} \hat{w}_i^r r_i \\
&= \frac{1}{n} \sum_{i \leq n} \hat{w}_i^R (R_i - r_i) + \frac{1}{n} \sum_{i \leq n} (\hat{w}_i^R - \hat{w}_i^r) r_i.
\end{aligned}$$

We can similarly rewrite

$$\hat{\Delta}_{\Omega}^{\text{iso}}(R; \mathcal{B}) - \hat{\Delta}_{\Omega}^{\text{iso}}(r; \mathcal{B}) = \frac{1}{n} \sum_{i \leq n} \hat{w}_i^r (R_i - r_i) + \frac{1}{n} \sum_{i \leq n} (\hat{w}_i^R - \hat{w}_i^r) R_i.$$

By the optimality of  $\hat{w}^R$  for (12) and of  $\hat{w}^r$  for (10) (since the constraints are identical for the two optimization problems), we have

$$\begin{aligned}
\frac{1}{n} \sum_{i \leq n} (\hat{w}_i^R - \hat{w}_i^r) R_i &= \frac{1}{n} \sum_{i \leq n} \hat{w}_i^R R_i - \frac{1}{n} \sum_{i \leq n} \hat{w}_i^r R_i \geq 0, \\
\frac{1}{n} \sum_{i \leq n} (\hat{w}_i^R - \hat{w}_i^r) r_i &= \frac{1}{n} \sum_{i \leq n} \hat{w}_i^R r_i - \frac{1}{n} \sum_{i \leq n} \hat{w}_i^r r_i \leq 0.
\end{aligned}$$

Consequently, we have the lower and upper bounds

$$\frac{1}{n} \sum_{i \leq n} \widehat{w}_i^r (R_i - r_i) \leq \widehat{\Delta}_\Omega^{\text{iso}}(R; \mathcal{B}) - \widehat{\Delta}_\Omega^{\text{iso}}(r; \mathcal{B}) \leq \frac{1}{n} \sum_{i \leq n} \widehat{w}_i^R (R_i - r_i). \quad (27)$$

Denote  $\varepsilon_i = r_i - R_i$  and  $\mathcal{G}_1 = \{w \geq 0 : w \in \mathcal{C}_{\geq}^{\text{iso}}, \|w\|_\infty \leq \Omega\}$ . Then, to control both the lower and upper bounds in (27), as both  $\widehat{w}^R$  and  $\widehat{w}^r$  are isotonic in  $X_i$ 's with respect to the componentwise order in  $\mathbb{R}^d$ , it suffices to control the following quantify

$$\sup_{w \in \mathcal{G}_1} \left| \frac{1}{n} \sum_{i \leq n} \varepsilon_i w(X_i) \right|. \quad (28)$$

We first control the expectation of (28). For any  $w \in \mathcal{G}_1$ , it holds that  $\mathbb{E}[w(X_i)\varepsilon_i] = 0$ , thus by symmetrization (Wellner et al. (2013) Theorem 2.3.1), we have

$$\mathbb{E} \left[ \sup_{w \in \mathcal{G}_1} \left| \frac{1}{n} \sum_{i \leq n} \varepsilon_i w(X_i) \right| \right] \leq 2\mathbb{E} \left[ \sup_{w \in \mathcal{G}_1} \left| \frac{1}{n} \sum_{i \leq n} \sigma_i \cdot \varepsilon_i w(X_i) \right| \right], \quad (29)$$

where  $\sigma_i$ 's are independent Rademacher random variables. Since  $\{\varepsilon_i\}_{i \leq n}$  are bounded within  $[-2B_R, 2B_R]$ , by Ledoux-Talagrand contraction lemma (Ledoux and Talagrand (2013) Theorem 4.12) with  $\phi_i(t) = \varepsilon_i t$ , we further have

$$\mathbb{E} \left[ \sup_{w \in \mathcal{G}_1} \left| \frac{1}{n} \sum_{i \leq n} \varepsilon_i w(X_i) \right| \right] \leq 4B_R \mathbb{E} \left[ \sup_{w \in \mathcal{G}_1} \left| \frac{1}{n} \sum_{i \leq n} \sigma_i \cdot w(X_i) \right| \right] = 4B_R \mathbb{E} [\mathcal{R}_n(\mathcal{G}_1)].$$

**Case 1:**  $d = 1$ . In the univariate case, by Birgé (1987); Chatterjee and Lafferty (2019), for the empirical Rademacher complexity of uniformly bounded isotonic function class, there exists a constant  $C_0$  such that with probability at least  $1 - n^{-1}$ ,

$$\mathcal{R}_n(\mathcal{G}_1) \leq C_0 B_R \sqrt{\frac{\log n}{n}}.$$

This bound is obtained via Dudley's theorem (Dudley, 1967).

**Case 2:**  $d \geq 2$ . In the multivariate case, when  $\mathcal{X}$  is bounded and  $0 < m_0 \leq \inf_{x \in \mathcal{X}} dP(x) \leq \sup_{x \in \mathcal{X}} dP(x) \leq M_0 < \infty$ , by Han et al. (2019) (Proposition 9), there exists a constant  $\tilde{C}_0$  such that

$$\mathbb{E} \left[ \sup_{w \in \mathcal{G}_1} \left| \frac{1}{n} \sum_{i \leq n} \sigma_i \cdot w(X_i) \right| \right] \leq \tilde{C}_0 B_R^2 \frac{\sqrt{\log^{\gamma_d} n}}{n^{1/d}},$$

where  $\gamma_1 = 1$ ,  $\gamma_2 = 8$ , and  $\gamma_d = d(d+1)$  for  $d \geq 3$ .

Combining the cases above, there exists a constant  $C$  such that

$$\mathbb{E} \left[ \sup_{w \in \mathcal{G}_1} \left| \frac{1}{n} \sum_{i \leq n} \varepsilon_i w(X_i) \right| \right] \leq 4B_R \mathbb{E} \left[ \sup_{w \in \mathcal{G}_1} \left| \frac{1}{n} \sum_{i \leq n} \sigma_i \cdot w(X_i) \right| \right] \leq C \frac{\sqrt{\log^{\gamma_d} n}}{n^{1/\max\{d, 2\}}}.$$

Since  $\mathcal{G}_1$  is a uniformly bounded function class and (28) is a function of  $n$  random vectors  $\{(X_i, \varepsilon_i)\}_{i \leq n}$ , if we substitute  $(X_i, \varepsilon_i)$  by  $(X'_i, \varepsilon'_i)$ , by the inequality  $|a - b| \geq ||a| - |b||$ , it holds that

$$\begin{aligned} & \left| \sup_{w \in \mathcal{G}_1} \left| \sum_{j \neq i} \varepsilon_j w(X_j) + \varepsilon_i w(X_i) \right| - \sup_{w \in \mathcal{G}_1} \left| \sum_{j \neq i} \varepsilon_j w(X_j) + \varepsilon'_i w(\tilde{X}_i) \right| \right| \\ & \leq \sup_{w \in \mathcal{G}_1} \left| \left| \sum_{j \neq i} \varepsilon_j w(X_j) + \varepsilon_i w(X_i) \right| - \left| \sum_{j \neq i} \varepsilon_j w(X_j) + \varepsilon'_i w(\tilde{X}_i) \right| \right| \\ & \leq \sup_{w \in \mathcal{G}_1} \left| \varepsilon_i w(X_i) - \varepsilon'_i w(X'_i) \right| \leq 2\Omega B_R < \infty. \end{aligned}$$

Then, by McDiarmid's inequality (McDiarmid et al., 1989), there exist constants  $\tilde{C}$  and  $\tilde{c}$  such that with probability at least  $1 - n^{-1}$ ,

$$\sup_{w \in \mathcal{G}_1} \left| \frac{1}{n} \sum_{i \leq n} \varepsilon_i w(X_i) \right| \leq \mathbb{E} \left[ \sup_{w \in \mathcal{G}_1} \left| \frac{1}{n} \sum_{i \leq n} \varepsilon_i w(X_i) \right| \right] + \tilde{c} \sqrt{\frac{\log n}{n}} \leq \tilde{C} \frac{\log^{\gamma_{d/2}} n}{n^{1/\max\{d,2\}}},$$

which implies that with probability at least  $1 - n^{-1}$ ,

$$|\widehat{\Delta}_{\Omega}^{\text{iso}}(R; \mathcal{B}) - \widehat{\Delta}_{\Omega}^{\text{iso}}(r; \mathcal{B})| \leq \max \left\{ \left| \frac{1}{n} \sum_{i \leq n} \widehat{w}_i^r(R_i - r_i) \right|, \left| \frac{1}{n} \sum_{i \leq n} \widehat{w}_i^R(R_i - r_i) \right| \right\} \leq \tilde{C} \frac{\log^{\gamma_{d/2}} n}{n^{1/\max\{d,2\}}}.$$

## C.2 Proof of Theorem 4.2 with $\mathcal{B} = \mathcal{B}_{a,b}$

For any  $\gamma > -1$ , we consider a relaxed optimization problem:

$$\begin{aligned} \widetilde{\Delta}(\gamma) &:= \max_{w \geq 0} \quad \frac{1}{n} \sum_{i \leq n} w(X_i) R(X_i) - \frac{1}{n} \sum_{i \leq n} r(X_i, Y_i) \\ &\text{subject to} \quad \frac{1}{n} \sum_{i \leq n} w(X_i) \leq 1 + \gamma, \quad a \leq w \leq b, \\ &\quad w \in \mathcal{C}_{\leq}^{\text{iso}}. \end{aligned} \tag{30}$$

In the constraints, when the inequality  $\frac{1}{n} \sum_{i \leq n} w(X_i) \leq 1 + \gamma$  is replaced with equality, we denote the resulting optimal value by  $\widetilde{\Delta}^*(\gamma)$ .

Similarly, we also consider a relaxed optimization problem at the population level:

$$\begin{aligned} \Delta(\gamma) &:= \sup_{w \geq 0} \quad \mathbb{E}_P[w(X)R(X)] - \mathbb{E}_P[R(X)] \\ &\text{subject to} \quad \mathbb{E}_P[w(X)] \leq 1 + \gamma, \quad a \leq w \leq b, \\ &\quad w \in \mathcal{C}_{\leq}^{\text{iso}}. \end{aligned} \tag{31}$$

The quantity  $\Delta^*(\gamma)$  is defined similarly. With these definitions in place, we have the following equivalence results.

**Lemma C.1.** *For any  $\gamma > -1$ , we have*

$$\Delta^*(\gamma) = \Delta(\gamma), \quad \widetilde{\Delta}^*(\gamma) = \widetilde{\Delta}(\gamma).$$

*In particular, we have*

$$\Delta^{\text{iso}}(R; \mathcal{B}_{a,b}) = \Delta(0), \quad \widehat{\Delta}^{\text{iso}}(R; \mathcal{B}_{a,b}) = \widetilde{\Delta}(0).$$

We are now ready to prove the theorem.<sup>6</sup>

**Step 1: Bounding the deviation.** Denote the maximizer of (30) by  $\widehat{w}(\cdot; \gamma)$  and the maximizer to (31) by  $w^*(\cdot; \gamma)$ . By Hoeffding's inequality, there exists a universal constant  $C_1$  such that with probability at

<sup>6</sup>We note that as  $\mathbb{E}[r(X_i, Y_i)] = \mathbb{E}_P[R(X)]$  and  $\{r(X_i, Y_i)\}_{i \leq n}$  are independent and bounded, the sample average  $n^{-1} \sum_{i \leq n} r(X_i, Y_i)$  converges to its mean  $\mathbb{E}_P[R(X)]$  at the parametric rate by Hoeffding's inequality. Since the parametric rate will not exceed the dominating rate in analysis, we focus on the convergence of the first term in  $\widehat{\Delta}_{\Omega}^{\text{iso}}(r; \mathcal{B})$  from now on.

least  $1 - n^{-1}$ ,

$$\left| \frac{1}{n} \sum_{i \leq n} w^*(X_i; 0) - \mathbb{E}_P[w^*(X; 0)] \right| \leq C_1 \sqrt{\frac{\log n}{n}} =: \gamma_n,$$

$$\text{and } \left| \frac{1}{n} \sum_{i \leq n} w^*(X_i; 0) R(X_i) - \mathbb{E}_P[w^*(X; 0) R(X)] \right| \leq C_1 \sqrt{\frac{\log n}{n}}.$$

As a result of the first bound, with probability at least  $1 - n^{-1}$ , the density ratio  $w^*(\cdot; -\gamma_n)$  is feasible for the sample problem (30) with  $\gamma = 0$ . In addition, the second bound implies a lower bound for  $\tilde{\Delta}(0)$ , namely,

$$\widehat{\Delta}^{\text{iso}}(R; \mathcal{B}_{a,b}) = \tilde{\Delta}(0) \geq \Delta(-\gamma_n) - \gamma_n. \quad (32)$$

On the other hand, denote  $\mathcal{G}_{\mathcal{B}_{a,b}} = \{w \geq 0 : a \leq w \leq b, w \in \mathcal{C}_{\mathcal{B}_{a,b}}^{\text{iso}}\}$ . For the maximizer  $\widehat{w}(\cdot; 0)$ , we have

$$\left| \frac{1}{n} \sum_{i \leq n} \widehat{w}(X_i, 0) - \mathbb{E}_P[\widehat{w}(X, 0)] \right| \leq \sup_{w \in \mathcal{G}_{\mathcal{B}_{a,b}}} \left| (\mathbb{E}_n - \mathbb{E}_P)[w(X)] \right|,$$

$$\left| \frac{1}{n} \sum_{i \leq n} \widehat{w}(X_i, 0) R(X_i) - \mathbb{E}_P[\widehat{w}(X, 0) R(X)] \right| \leq \sup_{w \in \mathcal{G}_{\mathcal{B}_{a,b}}} \left| (\mathbb{E}_n - \mathbb{E}_P)[w(X) R(X)] \right|.$$

By the uniform law of large numbers (e.g., Theorem 4.10 in [Wainwright \(2019\)](#)), with probability at least  $1 - n^{-1}$ ,

$$\sup_{w \in \mathcal{G}_{\mathcal{B}_{a,b}}} \left| (\mathbb{E}_n - \mathbb{E}_P)[w(X)] \right| \leq C_2 \sqrt{\frac{\log n}{n}} + C_3 \mathcal{R}_n(\mathcal{G}_{\mathcal{B}_{a,b}}) =: \varepsilon_n^{(1)},$$

and similarly,

$$\sup_{w \in \mathcal{G}_{\mathcal{B}_{a,b}}} \left| (\mathbb{E}_n - \mathbb{E}_P)[w(X) R(X)] \right| \leq C_2 \sqrt{\frac{\log n}{n}} + C_3 \mathcal{R}_n(\{w \cdot R \mid w \in \mathcal{G}_{\mathcal{B}_{a,b}}\}),$$

where  $\mathcal{R}_n(\mathcal{G})$  is the empirical Rademacher complexity of the function class  $\mathcal{G}$ , and  $C_2, C_3$  are universal constants. Moreover, as  $R$  is non-negative and  $B_R$ -bounded, by Ledoux-Talagrand contraction lemma ([Ledoux and Talagrand \(2013\)](#) Theorem 4.12) with  $\phi_i(t) = R(X_i)t$ , we have

$$\mathcal{R}_n(\{w \cdot R : w \in \mathcal{G}\}) \leq B_R \cdot \mathcal{R}_n(\mathcal{G}).$$

Hence, we have

$$\sup_{w \in \mathcal{G}_{\mathcal{B}_{a,b}}} \left| (\mathbb{E}_n - \mathbb{E}_P)[w(X) R(X)] \right| \leq C_2 \sqrt{\frac{\log n}{n}} + C_3 B_R \mathcal{R}_n(\mathcal{G}_{\mathcal{B}_{a,b}}) =: \varepsilon_n^{(2)}.$$

We denote  $\varepsilon_n := \max\{\varepsilon_n^{(1)}, \varepsilon_n^{(2)}\}$ .

Therefore, the sample maximizer  $\widehat{w}(\cdot; 0)$  is feasible for the population problem (31) with  $\gamma = \varepsilon_n$ . In addition, we have

$$\Delta(\varepsilon_n) \geq \tilde{\Delta}(0) - \varepsilon_n.$$

Combining the pieces above, we conclude that

$$- \{\Delta(0) - \Delta(-\gamma_n)\} - \gamma_n \leq \tilde{\Delta}(0) - \Delta(0) \leq \{\Delta(\varepsilon_n) - \Delta(0)\} + \varepsilon_n \quad (33)$$

holds with probability at least  $1 - 2n^{-1}$ .



**Step 2: Perturbation analysis.** In view of (33), our goal boils down to controlling the terms  $\Delta(0) - \Delta(-\gamma_n)$  and  $\Delta(\varepsilon_n) - \Delta(0)$ . This calls for the stability analysis of the population problem (31), which is supplied in the following lemma.

**Lemma C.2.** *For any  $\varepsilon > 0$  and  $\gamma \leq (b-1)/2$ , we have*

$$\begin{aligned}\Delta(\varepsilon) - \Delta(0) &\leq q_0^* \varepsilon, \\ \Delta(0) - \Delta(-\gamma) &\leq q_{-(b-1)/2}^* \gamma,\end{aligned}$$

where  $q_\gamma^* = q_R\left(\frac{b-1-\gamma}{b-a}\right)$  and  $q_R$  is the quantile function of  $R_{\#}P$ .

**Step 3: Combining pieces.** When  $n$  is sufficiently large, we have  $\gamma_n \leq (b-1)/2$ . Hence we can combine Lemma C.2 with (33) to obtain

$$\left| \widehat{\Delta}^{\text{iso}}(R; \mathcal{B}_{a,b}) - \Delta^{\text{iso}}(R; \mathcal{B}_{a,b}) \right| = \left| \widetilde{\Delta}(0) - \Delta(0) \right| \leq \max \left\{ \left(1 + q_{-(b-1)/2}^*\right) \gamma_n, (1 + q_0^*) \varepsilon_n \right\}.$$

Recalling the definitions of  $\varepsilon_n$  and  $\gamma_n$ , we complete the proof.

### C.2.1 Explicit rate of convergence

The proof above establishes that there exists a constant  $\widetilde{C} > 0$  such that with probability at least  $1 - 2n^{-1}$ ,

$$\left| \widehat{\Delta}^{\text{iso}}(R; \mathcal{B}_{a,b}) - \Delta^{\text{iso}}(R; \mathcal{B}_{a,b}) \right| \leq \widetilde{C} \left( \sqrt{\frac{\log n}{n}} + \mathcal{R}_n(\mathcal{G}_{\mathcal{B}_{a,b}}) \right).$$

However, this result is only meaningful if we can verify that  $\mathcal{R}_n(\mathcal{G}_{\mathcal{B}_{a,b}})$  is indeed small in settings of interest. To explicitly show a rate of convergence for this result, we consider the following cases.

- When  $d = 1$ , by Dudley's theorem (Dudley, 1967), we can obtain a tighter bound  $\mathcal{R}_n(\mathcal{G}_{\mathcal{B}_{a,b}}) \lesssim n^{-1/2}$  (Birgé, 1987; Chatterjee and Lafferty, 2019).
- When  $d \geq 2$ , with a bounded domain  $\mathcal{X}$  equipped with the componentwise order, by Han et al. (2019), if  $0 < m_0 \leq \inf_{x \in \mathcal{X}} \mathbf{d}P(x) \leq \sup_{x \in \mathcal{X}} \mathbf{d}P(x) \leq M_0 < \infty$ , they have shown that

$$\mathbb{E}[\mathcal{R}_n(\mathcal{G}_{\mathcal{B}_{a,b}})] \lesssim n^{-1/d}.$$

Since  $\mathcal{G}_{\mathcal{B}_{a,b}}$  consists of uniformly bounded functions, by McDiarmid's inequality (McDiarmid et al., 1989), the empirical Rademacher complexity concentrates to its mean such that

$$\mathcal{R}_n(\mathcal{G}_{\mathcal{B}_{a,b}}) \lesssim n^{-1/2} + n^{-1/d} \lesssim n^{-1/d}.$$

To sum up, for any fixed  $d \geq 1$ , there exists constant  $\widetilde{C}$  that does not depend on  $n$  such that with probability at least  $1 - 2n^{-1}$ ,

$$\left| \widehat{\Delta}^{\text{iso}}(R; \mathcal{B}_{a,b}) - \Delta^{\text{iso}}(R; \mathcal{B}_{a,b}) \right| \leq \widetilde{C} \left( \frac{1}{n} \right)^{1/d}.$$

Combining the upper bound above with the result in Theorem 4.1, we obtain

$$\left| \widehat{\Delta}^{\text{iso}}(r; \mathcal{B}_{a,b}) - \Delta^{\text{iso}}(R; \mathcal{B}_{a,b}) \right| \lesssim n^{-1/\max\{d,2\}}.$$

### C.2.2 Proof of Lemma C.1

Suppose  $\Delta_{[a,b]}(\gamma)$  is attained at  $\tilde{w} \in \mathcal{C}_{\geq}^{\text{iso}}$  in the interior of the feasible set, i.e.

$$\mathbb{E}_P[\tilde{w}(X)] < 1 + \gamma, \quad a \leq \tilde{w} \leq b.$$

Then, let  $w_\eta = (\tilde{w} + \eta) \wedge b$ , which satisfies that  $\mathbb{E}_P[w_\eta(X)] = 1 + \gamma$  and

$$\mathbb{E}_P[w_\eta(X)R(X)] - \mathbb{E}_P[R(X)] - \Delta(\gamma) = \eta \cdot \mathbb{E}_P[R(X) \mathbb{1}\{\tilde{w}(X) < b\}] > 0.$$

This draws the contradiction. Thus  $\Delta_{[a,b]}(\gamma)$  is attained on the boundary of the feasible set and particularly, when  $\gamma = 0$ , we have  $\Delta^{\text{iso}}(R; \mathcal{B}_{a,b}) = \Delta_{[a,b]}(0)$ .

### C.2.3 Proof of Lemma C.2

Recall that  $F_{\pi(R)}$  and  $q_{\pi(R)}$  denote the cumulative density function and the quantile function of  $(\pi(R))_{\#}P$ , respectively. The maximizer of (31) takes the form

$$w^*(x, \gamma) = (a - \eta_\gamma^*) \cdot \mathbb{1}\{[\pi(R)](x) < q_\gamma\} + (b - \eta_\gamma^*) \cdot \mathbb{1}\{[\pi(R)](x) > q_\gamma\} + \eta_\gamma^*,$$

where we define  $q_\gamma = q_{\pi(R)}(t_\gamma^*)$ ,

$$t_\gamma^* = \inf \left\{ t \in \text{range}(F_{\pi(R)}) : t \geq \frac{b-1-\gamma}{b-a} \right\}, \quad \text{and} \quad \eta_\gamma^* = a + \frac{(b-a)t_\gamma^* - (b-1-\gamma)}{\mathbb{P}(R(X) = q_\gamma)}.$$

Note that we have  $t_0^* \geq t_\gamma^*$  and  $q_0 \geq q_\gamma$  by definition. Then, we can show that

$$\Delta(\gamma) = a\mathbb{E}_P\left[[\pi(R)](X)\right] + (b-a)\mathbb{E}\left[[\pi(R)](X) \cdot \mathbb{1}\{[\pi(R)](X) > q_\gamma\}\right] + q_\gamma(b-a)t_\gamma^* - q_\gamma(b-1-\gamma).$$

Thus, for any  $a-1 < \gamma < b-1$ , it holds that

$$\begin{aligned} \Delta(\gamma) - \Delta(0) &= (b-a)\mathbb{E}\left[[\pi(R)](X) \cdot \mathbb{1}\{q_\gamma < [\pi(R)](X) \leq q_0\}\right] \\ &\quad + (b-a)(q_\gamma t_\gamma^* - q_0 t_0^*) + (q_0 - q_\gamma)(b-1) + q_\gamma \gamma \\ &\leq (b-a)q_0(t^* - t_\gamma^*) + (b-a)(q_\gamma t_\gamma^* - q_0 t_0^*) + (q_0 - q_\gamma)(b-1) + q_\gamma \gamma \\ &= (q_0 - q_\gamma) [b-1 - (b-a)t_\gamma^*] + q_\gamma \gamma \\ &\leq q_0 \gamma. \end{aligned}$$

The last inequality results from the fact that  $t_\gamma^* \geq (b-1-\gamma)/(b-a)$ . Similarly, we have the bound

$$\begin{aligned} \Delta(0) - \Delta(-\gamma) &= (b-a)\mathbb{E}\left[[\pi(R)](X) \cdot \mathbb{1}\{q_0 < [\pi(R)](X) \leq q_{-\gamma}\}\right] \\ &\quad + (b-a)(q_0 t_0^* - q_{-\gamma} t_{-\gamma}^*) + (q_{-\gamma} - q_0)(b-1) + q_{-\gamma} \gamma \\ &\leq q_{-\gamma} \gamma. \end{aligned}$$

When  $\gamma < (b-1)/2$ , we have  $q_{-\gamma} \leq q_{-(b-1)/2}$ , which completes the proof.

### C.3 Proof of Theorem 4.2 with $\mathcal{B} = \mathcal{B}_{f,\rho}$

Our proof in the case with  $\mathcal{B} = \mathcal{B}_{f,\rho}$  follows a similar route as that with  $\mathcal{B} = \mathcal{B}_{a,b}$ . To begin with, we define

$$\begin{aligned} \Delta(\rho) := \Delta^{\text{iso}}(R; \mathcal{B}_{f,\rho}) &= \sup_{w \geq 0} \mathbb{E}_P[w(X)R(X)] - \mathbb{E}_P[R(X)] \\ \text{subject to} &\quad \mathbb{E}_P[w(X)] = 1, \quad 0 \leq w \leq \Omega, \\ &\quad \mathbb{E}_P[f(w(X))] \leq \rho, \quad w \in \mathcal{C}_{\geq}^{\text{iso}}, \end{aligned} \tag{34}$$

where we write  $\rho$  as an argument of  $\Delta(\cdot)$  to help stability analysis later on.

In addition, we define the relaxed optimization problem in the sample space:

$$\begin{aligned} \tilde{\Delta}(\gamma, \rho) &:= \max_{w \geq 0} \quad \frac{1}{n} \sum_{i \leq n} w(X_i)R(X_i) - \frac{1}{n} \sum_{i \leq n} r(X_i, Y_i) \\ &\text{subject to} \quad \frac{1}{n} \sum_{i \leq n} w(X_i) \leq 1 + \gamma, \quad 0 \leq w \leq \Omega, \\ &\quad \quad \quad \frac{1}{n} \sum_{i \leq n} f(w(X_i)) \leq \rho, \quad w \in \mathcal{C}_{\geq}^{\text{iso}}. \end{aligned} \quad (35)$$

When we replace the inequality  $\frac{1}{n} \sum_{i \leq n} w(X_i) \leq 1 + \gamma$  in the constraints with equality, we denote the resulting optimal value  $\tilde{\Delta}^*(\gamma, \rho)$ . Similarly, we consider a relaxed optimization problem in the population level:

$$\begin{aligned} \Delta(\gamma, \rho) &:= \sup_{w \geq 0} \quad \mathbb{E}_P[w(X)R(X)] - \mathbb{E}_P[R(X)] \\ &\text{subject to} \quad \mathbb{E}_P[w(X)] \leq 1 + \gamma, \quad 0 \leq w \leq \Omega, \\ &\quad \quad \quad \mathbb{E}_P[f(w(X))] \leq \rho, \quad w \in \mathcal{C}_{\geq}^{\text{iso}}. \end{aligned} \quad (36)$$

Also, we denote the optimal value to be  $\Delta^*(\gamma, \rho)$  when we replace the inequality  $\mathbb{E}_P[w(X)] \leq 1 + \gamma$  with equality.

**Lemma C.3.** *For  $-1 < \gamma \leq 0$ , we have*

$$\Delta^*(\gamma, \rho) = \Delta(\gamma, \rho), \quad \tilde{\Delta}^*(\gamma, \rho) = \tilde{\Delta}(\gamma, \rho).$$

*In particular, we have*

$$\Delta^{\text{iso}}(R; \mathcal{B}_{f, \rho}) = \Delta(0, \rho) \quad \text{and} \quad \widehat{\Delta}^{\text{iso}}(R; \mathcal{B}_{f, \rho}) = \tilde{\Delta}(0, \rho).$$

**Step 1: Bounding the deviation.** Denote the maximizer of (36) by  $w^*(\cdot; \gamma, \rho)$  and maximizer of (35) by  $\widehat{w}(\cdot; \gamma, \rho)$ . By Hoeffding's inequality, there exists a universal constant  $\tilde{C}_1$  such that with probability at least  $1 - n^{-1}$ ,

$$\begin{aligned} \left| \frac{1}{n} \sum_{i \leq n} w^*(X_i; 0, \rho) - \mathbb{E}_P[w^*(X; 0, \rho)] \right| &\leq C_1 \sqrt{\frac{\log n}{n}} =: \gamma_n, \\ \left| \frac{1}{n} \sum_{i \leq n} f(w^*(X_i; 0, \rho)) - \mathbb{E}_P[f(w^*(X; 0, \rho))] \right| &\leq C_1 \sqrt{\frac{\log n}{n}}, \\ \left| \frac{1}{n} \sum_{i \leq n} w^*(X_i; 0, \rho)R(X_i) - \mathbb{E}_P[w^*(X; 0, \rho)R(X)] \right| &\leq C_1 \sqrt{\frac{\log n}{n}}. \end{aligned}$$

As a result, with probability at least  $1 - n^{-1}$ , the maximizer  $w^*(\cdot; -\gamma_n, \rho - \gamma_n)$  is feasible for the sample problem (35) with  $\gamma = 0$ , and hence

$$\tilde{\Delta}(0, \rho) \geq \Delta(-\gamma_n, \rho - \gamma_n) - \gamma_n.$$

On the other hand, define  $\mathcal{G}_1 = \{w \geq 0 : w \in \mathcal{C}_{\geq}^{\text{iso}}, w \in [0, \Omega]\}$  and  $\mathcal{G}_{f,2} = \{f \circ w \mid w \in \mathcal{G}_1\}$ . Since we assume  $f$  is bounded on  $[0, \Omega]$ , we further have  $\mathcal{G}_{f,2} \subset \{f \circ w : w \in \mathcal{C}_{\geq}^{\text{iso}}, f \circ w \in [0, \Omega]\}$ . For the maximizer

$\widehat{w}(X; 0, \rho)$ , we have

$$\begin{aligned} \left| \mathbb{E}_P[\widehat{w}(X, 0)] - \frac{1}{n} \sum_{i \leq n} \widehat{w}(X_i, 0) \right| &\leq \sup_{w \in \mathcal{G}_1} \left| (\mathbb{E}_n - \mathbb{E}_P)[w(X)] \right|, \\ \left| \frac{1}{n} \sum_{i \leq n} f(\widehat{w}(X_i, 0)) - \mathbb{E}_P[f(\widehat{w}(X, 0))] \right| &\leq \sup_{w \in \mathcal{G}_1} \left| (\mathbb{E}_n - \mathbb{E}_P)[f(w(X))] \right|, \\ \left| \frac{1}{n} \sum_{i \leq n} \widehat{w}(X_i, 0)R(X_i) - \mathbb{E}_P[\widehat{w}(X, 0)R(X)] \right| &\leq B_R \sup_{w \in \mathcal{G}_1} \left| (\mathbb{E}_n - \mathbb{E}_P)[w(X)] \right|, \end{aligned}$$

where the last inequality is based on Ledoux-Talagrand contraction lemma ([Ledoux and Talagrand \(2013\)](#) Theorem 4.12).

Moreover, with probability at least  $1 - n^{-1}$ , we have

$$\begin{aligned} \sup_{w \in \mathcal{G}_1} \left| (\mathbb{E}_n - \mathbb{E}_P)[w(X)] \right| &\leq C_2 \left\{ \sqrt{\frac{\log n}{n}} + \mathcal{R}_n(\mathcal{G}_1) \right\} =: \varepsilon_n, \\ \sup_{w \in \mathcal{G}_1} \left| (\mathbb{E}_n - \mathbb{E}_P)[f(w(X))] \right| &\leq C_3 \left\{ \sqrt{\frac{\log n}{n}} + \mathcal{R}_n(\mathcal{G}_{f,2}) \right\} =: \tilde{\varepsilon}_n. \end{aligned}$$

where  $C_2, C_3$  are universal constants and  $\mathcal{R}_n(\mathcal{G}_1), \mathcal{R}_n(\mathcal{G}_{f,2})$  are the empirical Rademacher complexities of  $\mathcal{G}_1$  and  $\mathcal{G}_{f,2}$ , respectively.

Denote  $t_f^* = \operatorname{argmin}_{t \in [0, \Omega]} f(t)$ . We have the decomposition

$$f(t) = f(t) \cdot \mathbb{1}\{f(t) \geq t_f^*\} + f(t) \cdot \mathbb{1}\{f(t) < t_f^*\} =: f_1 + f_2,$$

where both  $f_1$  and  $-f_2$  are nondecreasing. Then, for any  $g = f \circ w \in \mathcal{G}_{f,2}$ , we have the decomposition  $g = f_1 \circ w + f_2 \circ w$ , where  $f_1 \circ w \in \mathcal{C}_{\geq}^{\text{iso}}$ ,  $-f_2 \circ w \in \mathcal{C}_{\geq}^{\text{iso}}$ , and both functions are uniformly bounded within  $[-B_f, B_f]$ . Recall the notation  $\mathcal{G}_{\mathcal{B}_{f,\rho}} = \{w : w \in \mathcal{C}_{\geq}^{\text{iso}}, w \in [-\Omega \vee B_f, \Omega \vee B_f]\}$ . Hence, we have  $\mathcal{G}_{f,2} \subseteq \mathcal{G}_{\mathcal{B}_{f,\rho}} - \mathcal{G}_{\mathcal{B}_{f,\rho}}$  and by the property of the Rademacher complexity, we obtain  $\mathcal{R}_n(\mathcal{G}_1) \leq \mathcal{R}_n(\mathcal{G}_{\mathcal{B}_{f,\rho}})$  and  $\mathcal{R}_n(\mathcal{G}_{f,2}) \leq 2\mathcal{R}_n(\mathcal{G}_{\mathcal{B}_{f,\rho}})$ , which implies that there exists a constant  $C_4$  such that

$$\max\{\varepsilon_n, \tilde{\varepsilon}_n\} \leq C_4 \left\{ \sqrt{\frac{\log n}{n}} + \mathcal{R}_n(\mathcal{G}_{\mathcal{B}_{f,\rho}}) \right\}.$$

Therefore, the maximizer  $\widehat{w}(X; 0, \rho)$  is feasible for the population problem (36) with  $\gamma = \varepsilon_n$  and  $\rho + \tilde{\varepsilon}_n$ , and we have the upper bound for  $\tilde{\Delta}(0, \rho)$ :

$$\Delta(\varepsilon_n, \rho + \tilde{\varepsilon}_n) \geq \tilde{\Delta}(0, \rho) - \varepsilon_n.$$

Combining the pieces above, we obtain

$$\begin{aligned} & - \{ \Delta(0, \rho) - \Delta(-\gamma_n, \rho - \gamma_n) \} - \gamma_n \\ & \leq \tilde{\Delta}(0, \rho) - \Delta(0, \rho) \leq \{ \Delta(\varepsilon_n, \rho + \tilde{\varepsilon}_n) - \Delta(0, \rho) \} + \varepsilon_n. \end{aligned} \tag{37}$$

It remains to characterize the stability of  $\Delta(\gamma, \rho)$  w.r.t. both  $\gamma$  and  $\rho$ .

**Step 2: Perturbation analysis.** According to (37), bounding the estimation error bound boils down to bounding  $\Delta(0, \rho) - \Delta(-\gamma_n, \rho - \gamma_n)$  and  $\Delta(\varepsilon_n, \rho + \tilde{\varepsilon}_n) - \Delta(0, \rho)$  respectively, for which we have the following lemma on the stability analysis of the population problem (36).

**Lemma C.4.** Assume  $\sup_{t \in [0, \Omega]} |f(t)| < \infty$ . For any  $\varepsilon > 0$ , there exists  $0 < \bar{\lambda} < \infty$ ,  $0 < L_f < \infty$  and  $\delta > 0$  such that when  $|\gamma| < \delta$ , it holds that

$$\begin{aligned}\Delta(\varepsilon, \rho + \varepsilon) - \Delta(0, \rho) &\leq \bar{\lambda}\varepsilon + (B_R + \bar{\lambda}|f'(1)|)\varepsilon, \\ \Delta(0, \rho) - \Delta(-\gamma, \rho - \gamma) &\leq \bar{\lambda}\gamma + (B_R + \bar{\lambda}L_f)\gamma.\end{aligned}$$

**Step 3: Combining pieces.** Based on Lemma C.4, bounds in (37) can be written as

$$\left| \tilde{\Delta}(0, \rho) - \Delta(0, \rho) \right| \leq \max \left\{ (1 + B_R + \bar{\lambda}L_f + \bar{\lambda})\gamma_n, (1 + B_R + \bar{\lambda}|f'(1)|)\varepsilon_n + \bar{\lambda}\tilde{\varepsilon}_n \right\}.$$

Recall the definitions of  $\varepsilon_n$  and  $\gamma_n$  to finish the proof.

### C.3.1 Explicit rate of convergence

We first present the following lemma to ensure the boundedness of the maximizer of  $\Delta^{\text{iso}}(R; \mathcal{B}_{f, \rho})$ .

**Lemma C.5.** Assume  $R$  is  $B_R$ -bounded and convex function  $f$  is differentiable with  $f(1) = 0$ . The excess risk  $\Delta^{\text{iso}}(R; \mathcal{B}_{f, \rho})$  is attained at  $w_{f, \rho}^{\text{iso}} \in \mathcal{B}_{f, \rho} \cap \mathcal{C}_{\geq}^{\text{iso}}$  with  $\|w_{f, \rho}^{\text{iso}}(X)\|_{\infty} < \infty$  almost surely.

Then, when  $0 < \Omega < \|w_{f, \rho}^{\text{iso}}\|$  and  $n$  is sufficiently large, there exist a constant  $\tilde{C} = \tilde{C}(\rho) > 0$  such that with probability at least  $1 - 2n^{-1}$ ,

$$\left| \hat{\Delta}^{\text{iso}}(R; \mathcal{B}_{f, \rho}) - \Delta^{\text{iso}}(R; \mathcal{B}_{f, \rho}) \right| \leq \tilde{C} \left( \sqrt{\frac{\log n}{n}} + \mathcal{R}_n(\mathcal{G}_{\mathcal{B}_{f, \rho}}) \right).$$

To derive the explicit rate of convergence, it suffices to bound the empirical Rademacher complexity. Since  $\mathcal{G}_{\mathcal{B}_{f, \rho}}$  is contained in the bounded isotonic function class equipped with the componentwise order, then when  $\mathcal{X}$  is bounded, similar to the proof of the previous theorem in Section C.2.1, we have

$$\mathcal{R}_n(\mathcal{G}_{\mathcal{B}_{f, \rho}}) \lesssim \left( \frac{1}{n} \right)^{1/d}.$$

Moreover, when  $d = 1$ , by Dudley's theorem (Dudley, 1967), we can obtain a tighter bound  $\mathcal{R}_n(\mathcal{G}_{\mathcal{B}_{f, \rho}}) \lesssim n^{-1/2}$  (Birgé, 1987; Chatterjee and Lafferty, 2019). To sum up, for any fixed  $d \geq 1$ , there exists a constant  $\tilde{C}$  that does not depend on  $n$  such that with probability at least  $1 - 2n^{-1}$ ,

$$\left| \hat{\Delta}^{\text{iso}}(R; \mathcal{B}_{f, \rho}) - \Delta^{\text{iso}}(R; \mathcal{B}_{f, \rho}) \right| \leq \tilde{C} \left( \frac{1}{n} \right)^{1/d}.$$

Combining the upper bound above with the result in Theorem 4.1, we obtain

$$\left| \hat{\Delta}^{\text{iso}}(r; \mathcal{B}_{f, \rho}) - \Delta^{\text{iso}}(R; \mathcal{B}_{f, \rho}) \right| \lesssim n^{-1/\max\{d, 2\}}.$$

### C.3.2 Proof of Lemma C.3

Without loss of generality, we assume  $t = 1$  is the minimizer of  $f^7$ . With  $-1 < \gamma \leq 0$ , recall the KKT condition of problem (36):

$$\begin{aligned}-[\pi(R)](x) + \lambda f'(w^*(x, \gamma, \rho)) + \nu &= 0, \\ \lambda (\mathbb{E}_P[f(w^*(X, \gamma, \rho))] - \rho) &= 0, \\ \nu (\mathbb{E}_P[w^*(X, \gamma, \rho)] - 1 - \gamma) &= 0.\end{aligned}$$

<sup>7</sup>To see this, for any  $c \in \mathbb{R}$  and  $\tilde{f}(t) = f(t) - c(t-1)$ , we have  $\Delta^{\text{iso}}(R; \mathcal{B}_{f, \rho}) = \Delta^{\text{iso}}(R; \mathcal{B}_{\tilde{f}, \rho})$ . Specifically, the equality holds for  $\tilde{f}(t) = f(t) - f'(1)(t-1)$

Suppose  $\mathbb{E}_P[w^*(X, \gamma, \rho)] < 1 + \gamma$ , then  $\nu = 0$ , which implies that

$$\lambda f'(w(x)) = [\pi(R)](x).$$

Denote  $\mathcal{G} = \{w \geq 0 : \mathbb{E}_P[f(w(X))] \leq \rho\}$ , which is a convex set.

**Case 1.** If  $\lambda = 0$ , the optimal density ratio  $w^*(x, \gamma, \rho)$  is an interior point of  $\mathcal{G}$  and there exists  $\eta > 0$  such that  $\tilde{w}(x, \gamma, \rho) = w^*(x, \gamma, \rho) + \eta \in \mathcal{G}$ . In the meantime,

$$\mathbb{E}_P \left[ \tilde{w}(X, \gamma, \rho) [\pi(R)](X) \right] = \mathbb{E}_P \left[ w^*(X, \gamma, \rho) [\pi(R)](X) \right] + \eta \mathbb{E}_P \left[ [\pi(R)](X) \right] > \mathbb{E}_P \left[ w^*(X, \gamma, \rho) [\pi(R)](X) \right],$$

which draws the contradiction.

**Case 2.** If  $\lambda > 0$ , we have  $f'(w^*(x, \gamma, \rho)) \propto [\pi(R)](x)$ . Since  $\pi(R) \geq 0$  and  $f'(t) \geq 0$  if and only if  $t \geq 1$ , we obtain  $w^*(x, \gamma, \rho) \geq 1$ . However, when  $\gamma \leq 0$ , as  $\mathbb{E}_P[w^*(X, \gamma, \rho)] < 1 + \gamma \leq 1$ , it draws the contradiction.

### C.3.3 Proof of Lemma C.4

We prove Lemma C.4 via analyzing the dual formulation of the worst-case excess risk.

**Analysis of the dual formulation.** Recall that the optimization problem (36) has the dual formulation:

$$\Delta(\gamma, \rho) = \inf_{\lambda \geq 0, \nu} \left\{ \lambda \rho + \nu(1 + \gamma) + \mathbb{E}_P \left[ w(X) ([\pi(R)](X) - \nu) - \lambda f(w(X)) \right] \right\}, \quad (38)$$

where

$$w(x) = \mathcal{P}_{[0, \Omega]} \left\{ (f')^{-1} \left( \frac{[\pi(R)](x) - \nu}{\lambda} \right) \right\}.$$

Denote  $\lambda^*(\gamma, \rho)$  and  $\nu^*(\gamma, \rho)$  as the minimizers and

$$H(\lambda, \nu; \rho, \gamma) = \lambda \rho + \nu(1 + \gamma) + \mathbb{E}_P \left[ w(X) ([\pi(R)](X) - \nu) - \lambda f(w(X)) \right].$$

Then, with any  $\varepsilon > 0$  and  $\tilde{\varepsilon} > 0$ , we have

$$\Delta(\varepsilon, \rho + \tilde{\varepsilon}) - \Delta(0, \rho) \leq H(\lambda^*(0, \rho), \nu^*(0, \rho); \rho + \tilde{\varepsilon}, \varepsilon) - H(\lambda^*(0, \rho), \nu^*(0, \rho); \rho, 0) \quad (39)$$

$$\leq \lambda^*(0, \rho) \tilde{\varepsilon} + \nu^*(0, \rho) \varepsilon. \quad (40)$$

Similarly, for any  $\gamma > 0$ , we have

$$\Delta(0, \rho) - \Delta(-\gamma, \rho - \gamma) \leq \lambda^*(-\gamma, \rho - \gamma) \gamma + \nu^*(-\gamma, \rho - \gamma) \gamma.$$

Our goal is to derive upper bounds for  $\lambda^*(\gamma, \rho)$  and  $\nu^*(\gamma, \rho)$ . Note that  $\nu^*(\gamma, \rho)$  is the parameter for standardization, thus to guarantee  $\mathbb{E}_P[w^*(X; \gamma, \rho)] = 1 + \gamma$ , we have

$$(f')^{-1} \left( \frac{B_R - \nu^*(\gamma, \rho)}{\lambda^*(\gamma, \rho)} \right) \geq \sup_{x \in \mathcal{X}} w^*(X; \gamma, \rho) \geq 1 + \gamma.$$

This implies that

$$\nu^*(\gamma, \rho) \leq B_R - \lambda^*(\gamma, \rho) f'(1 + \gamma), \quad (41)$$

then it suffices to show that  $\lambda^*(\gamma, \rho)$  is finite. Moreover, as  $(f')^{-1}(-\nu^*(\gamma, \rho)/\lambda^*(\gamma, \rho)) \leq \min_{x \in \mathcal{X}} w^*(X; \gamma, \rho) \leq 1$ , we have  $\nu^*(\gamma, \rho) \geq -\lambda^*(\gamma, \rho) f'(1)$ .

**Bounding the dual minimizers.** Denote  $\phi(\xi) = \mathcal{P}_{[0, \Omega]} \{(f')^{-1}(\xi)\}$  and  $\psi(\xi) = \xi\phi(\xi) - f(\phi(\xi))$ . Recall the dual formulation

$$\Delta(\gamma, \rho) = \inf_{\lambda \geq 0, \nu} \left\{ \lambda\rho + \nu(1 + \gamma) + \lambda \cdot \mathbb{E}_P \left[ \psi \left( \frac{[\pi(R)](X) - \nu}{\lambda} \right) \right] \right\}.$$

Denote  $\nu = \nu(\lambda)$  as the solution to  $\mathbb{E}_P[\phi(\xi)] = 1 + \gamma$  and

$$L(\lambda, \rho, \gamma) = \lambda\rho + \nu(\lambda)(1 + \gamma) + \lambda \cdot \mathbb{E}_P \left[ \psi \left( \frac{[\pi(R)](X) - \nu(\lambda)}{\lambda} \right) \right].$$

Suppose the optimal dual variable  $\lambda^*(\gamma, \rho)$  is unbounded. Due to the fact that  $R$  is bounded and

$$\mathbb{E}_P \left[ \phi \left( \frac{[\pi(R)](X) - \nu(\lambda)}{\lambda} \right) \right] = 1 + \gamma,$$

we have  $\nu(\lambda) = -c \cdot \lambda$  as  $\lambda \rightarrow \infty$ , which satisfies that  $(f')^{-1}(c) = 1 + \gamma$ , which yields  $c = f'(1 + \gamma)$ . Then, we can verify that

$$\mathbb{E}_P \left\{ f \left( \phi \left( \frac{[\pi(R)](X) - \nu(\lambda)}{\lambda} \right) \right) \right\} \rightarrow f(1 + \gamma), \quad \text{as } \lambda \rightarrow +\infty,$$

which implies that

$$\frac{\partial}{\partial \lambda} L(\lambda, \rho, \gamma) = \rho - \mathbb{E}_P \left\{ f \left( \phi \left( \frac{[\pi(R)](X) - \nu(\lambda)}{\lambda} \right) \right) \right\} \rightarrow \rho - f(1 + \gamma), \quad \text{as } \lambda \rightarrow +\infty.$$

For any  $\delta \in (0, \min\{\Omega - 1, 1\}/2)$ , since  $f$  is convex and  $\sup_{t \in [0, \Omega]} |f(t)| < \infty$ , we define  $L_f = \sup_{|t-1| < \delta} |f'(t)| < \infty$ . For any  $\delta' < \min\{\delta, \rho/(2L_f)\}$ , as  $f(1 + \gamma) < L_f \delta' < \rho/2$  when  $|\gamma| < \delta'$ , there exists  $\bar{\lambda} > 0$  such that when  $\lambda > \bar{\lambda}$ , we have  $\partial_\lambda L(\lambda, \rho, \gamma) \geq \rho/2$ , which draws the contradiction to the assumption that the dual problem is minimized at  $\lambda^*(\gamma, \rho) = \infty$ . This implies that both  $\lambda^*(0, \rho)$  and  $\lambda^*(-\gamma, \rho - \gamma)$  can be bounded from above by  $\bar{\lambda}$  when  $0 < \gamma < \delta'$ .

**Combining pieces.** Consequently, with  $L_f = \sup_{|t-1| \leq \delta} |f'(t)| < \infty$  and  $0 < \gamma < \delta'$  with  $\delta' < \min\{\delta, \rho/(2L_f)\}$ , by (39) and (41), it holds that

$$\begin{aligned} \Delta(\varepsilon, \rho + \tilde{\varepsilon}) - \Delta(0, \rho) &\leq \bar{\lambda} \tilde{\varepsilon} + (B_R + \bar{\lambda} |f'(1)|) \varepsilon, \\ \Delta(0, \rho) - \Delta(-\gamma, \rho - \gamma) &\leq \bar{\lambda} \gamma + (B_R + \bar{\lambda} L_f) \gamma, \end{aligned}$$

which completes the proof.

### C.3.4 Proof of Lemma C.5

Recall the dual formulation. There exists a pair  $(\lambda^*, \nu^*)$  such that

$$w_{f, \rho}^{* \text{iso}}(x) = \mathcal{P}_{[0, +\infty)} \left\{ (f')^{-1} \left( \frac{[\pi(R)](x) - \nu^*}{\lambda^*} \right) \right\}.$$

According to the proof of Lemma C.4, we have already shown that  $\lambda^* \leq \bar{\lambda} < \infty$ . Since  $-\lambda^* f'(1) \leq \nu^* \leq B_R - f'(1)\lambda^*$  as is shown in the proof of Lemma C.4, if  $\underline{\lambda} < \lambda^* < \bar{\lambda}$ , there exists  $\bar{\nu} < \infty$  such that  $|\nu^*| \leq \bar{\nu}$ , thus it holds that

$$\|w_{f, \rho}^{* \text{iso}}\|_\infty \leq (f')^{-1} \left( \frac{B_R + \bar{\nu}}{\underline{\lambda}} \right) < \infty.$$

Then, it remains to prove that  $\lambda^* \neq 0$ . To see this, consider the KKT condition:

$$\begin{aligned} -[\pi(R)](x) + \lambda^* f'(w_{f,\rho}^{\text{iso}}(x)) + \nu^* &= 0, \\ \lambda^* (\mathbb{E}_P[f(w_{f,\rho}^{\text{iso}}(X))] - \rho) &= 0, \\ \nu^* (\mathbb{E}_P[w_{f,\rho}^{\text{iso}}(X)] - 1) &= 0. \end{aligned}$$

If  $\lambda^* = 0$ , we have  $[\pi(R)](X) = \nu^*$   $P$ -almost surely, which implies that  $w_{f,\rho}^{\text{iso}}(X) = 1$   $P$ -almost surely, in which case  $w_{f,\rho}^{\text{iso}}$  is also bounded. Combining pieces above, we have shown that  $\|w_{f,\rho}^{\text{iso}}\|_\infty < \infty$ .

## C.4 Hardness of consistent estimation without isotonic constraints

We should note that the concentration bound in Theorem 4.1 with the presence of the noisy risk is not possible without the isotonic constraint. To see this, consider a counterexample with the risk function  $R(x) = 1/2$ , which leads to  $\Delta(R; \mathcal{B}) = 0$ , and the noisy risk  $r_i = \text{Bern}(R(X_i)) = \text{Bern}(1/2)$  independently. We consider the bound-constraint  $\mathcal{B}_{a,b}$  with  $0 \leq a \leq 1 \leq b$ . To estimate the excess risk without the isotonic constraint, consider the optimization problem

$$\widehat{\Delta}(r; \mathcal{B}_{a,b}) = \max_{w \geq 0} \frac{1}{n} \sum_{i \leq n} w(X_i) r_i - \frac{1}{n} \sum_{i \leq n} r_i \quad \text{subject to} \quad \mathbb{E}_{\widehat{P}_n}[w(X)] = 1, \quad w_{\#} \widehat{P}_n \in \mathcal{B}_{a,b},$$

for which we have the following proposition.

**Proposition C.6.** *Assume  $R(X) = 1/2$   $P$ -almost surely and  $\min\{1 - a, b - 1\} > 0$ . Then, there exists a strictly positive constant  $\underline{C} > 0$  depending on  $a$  and  $b$  such that*

$$\widehat{\Delta}(r; \mathcal{B}_{a,b}) - 0 \geq \underline{C} > 0 \quad \text{with probability at least } 1 - n^{-1}.$$

We defer the proof to Section C.4.1.

### C.4.1 Proof of Proposition C.6

We follow the setup with  $r_i \sim \text{Bern}(R(X_i)) = \text{Bern}(1/2)$  independently. Consider the bounds constraint  $\mathcal{B} = \mathcal{B}_{a,b}$ . The estimated excess takes the form

$$\widehat{\Delta}(r; \mathcal{B}_{a,b}) = \max_{w \geq 0} \frac{1}{n} \sum_{i \leq n} w(X_i) r_i - \frac{1}{n} \sum_{i \leq n} R(X_i) \quad \text{subject to} \quad \mathbb{E}_{\widehat{P}_n}[w(X)] = 1, \quad w_{\#} \widehat{P}_n \in \mathcal{B}_{a,b}.$$

According to Section 2, the worst-case weights take the form  $w_i = w(X_i) = c_1 \cdot \mathbb{1}\{r_i = 0\} + c_2 \cdot \mathbb{1}\{r_i = 1\}$ , where  $a \leq c_1 \leq 1 \leq c_2 \leq b$ . Moreover, by the KKT condition, at least one of  $c_1 = a$  and  $c_2 = b$  holds. Then, the estimated excess risk can be expressed as

$$\widehat{\Delta}(r; \mathcal{B}_{a,b}) = \frac{c_2}{n} \sum_{i \leq n} r_i - \frac{1}{n} \sum_{i \leq n} R(X_i).$$

Since  $n^{-1} \sum_{i \leq n} w_i = 1$ , we have

$$\frac{1}{n} \sum_{i \leq n} (1 - r_i) = \frac{c_2 - 1}{c_2 - c_1}.$$

In the meantime, by Hoeffding's inequality, there exists a universal constant  $\tilde{c} > 0$  such that with probability at least  $1 - n^{-1}$ ,

$$\left| \frac{1}{n} \sum_{i \leq n} r_i - \frac{1}{n} \sum_{i \leq n} R(X_i) \right| \leq \frac{\tilde{c}}{2} \sqrt{\frac{\log n}{n}}, \quad \left| \frac{1}{n} \sum_{i \leq n} R(X_i) - \mathbb{E}_P[R(X)] \right| \leq \frac{\tilde{c}}{2} \sqrt{\frac{\log n}{n}},$$



which implies

$$c_2 - 1 \geq \left( \mathbb{E}_P[1 - R(X)] - \tilde{c} \sqrt{\frac{\log n}{n}} \right) \cdot \min\{1 - a, b - 1\}.$$

If  $n$  is large enough such that  $\min\{(\mathbb{E}_P[1 - R(X)])^2, (\mathbb{E}_P[R(X)])^2\}n \geq (1 - \kappa)^{-2} \tilde{c}^2 \log n$ , where  $1/2 \leq \kappa < 1$  satisfies that

$$\kappa + \kappa^2 \mathbb{E}_P[1 - R(X)] \cdot \min\{1 - a, b - 1\} > 1 + \underline{c}, \quad (42)$$

for some  $0 < \underline{c} < \mathbb{E}_P[1 - R(X)] \cdot \min\{1 - a, b - 1\}$ . Then, with probability at least  $1 - n^{-1}$ , we have

$$c_2 \geq 1 + \kappa \mathbb{E}_P[1 - R(X)] \cdot \min\{1 - a, b - 1\}.$$

Consequently, for the excess risk, with probability at least  $1 - n^{-1}$ , we have

$$\begin{aligned} \widehat{\Delta}(r; \mathcal{B}_{a,b}) &= \frac{c_2 - 1}{n} \sum_{i \leq n} r_i + \frac{1}{n} \sum_{i \leq n} (r_i - R(X_i)) \\ &= (c_2 - 1) \mathbb{E}_P[R(X)] + \frac{c_2 - 1}{n} \sum_{i \leq n} (r_i - \mathbb{E}_P[R(X)]) + \frac{1}{n} \sum_{i \leq n} (r_i - R(X_i)) \\ &\geq (c_2 - 1) \mathbb{E}_P[R(X)] - c_2 \cdot \tilde{c} \sqrt{\frac{\log n}{n}} \\ &\geq (\kappa c_2 - 1) \mathbb{E}_P[R(X)] \geq \underline{c} \cdot \mathbb{E}_P[R(X)] = \frac{\underline{c}}{2}, \end{aligned}$$

where the last inequality holds according to (42). Choosing  $\underline{C} = \underline{c}/2$  completes the proof.

## D Additional simulation results

### D.1 iso-DRL under componentwise order

In Section 5, we mainly focused on the partial order with respect to  $w_0(x)$ . In this section, to demonstrate the effect of various choices of the partial (pre)order, we further consider an alternative choice of the partial (pre)order: the componentwise order where

$$x \preceq x' \quad \text{if and only if} \quad x_j \leq x'_j, \text{ for all } j \in [m],$$

where we set  $m = 5 < d = 20$ . Let iso-DRL-comp denote the CP interval with calibrated target level  $\alpha'_{\text{iso}} = \max\{0, \alpha - \tilde{\Delta}^{\text{iso}}\}$ , where

$$\begin{aligned} \tilde{\Delta}^{\text{iso}} &= \max \quad \frac{1}{n} \sum_{i \in \mathcal{D}_3} w_i \tilde{r}_i^{\text{iso}} - \frac{1}{n} \sum_{i \in \mathcal{D}_3} r_i \\ \text{subject to} \quad &\frac{1}{n} \sum_{i \in \mathcal{D}_3} w_i = 1, \quad \frac{1}{n} \sum_{i \in \mathcal{D}_3} w_i \log w_i \leq \rho, \quad 0 \leq w_i \leq \Omega, \end{aligned} \quad (43)$$

and  $(\tilde{r}_i)_{i \in \mathcal{D}_3}$  is the isotonic projection of  $(r_i)_{i \in \mathcal{D}_3}$  with respect to the componentwise order.

We follow the exactly same settings with Section 5.1 with  $n_{\text{pre}} = 50$  and vary  $\rho$  in  $[0.002, 6]$ . From Figure 7 and 8, each of the coverage rate and average interval width of iso-DRL-comp lies between that of DRL and iso-DRL- $w_0$ , which indicates that additional constraints will relieve the conservativeness of DRL, but only a proper choice of the partial (pre)order will lead to desired performance close to the oracle weighted CP.

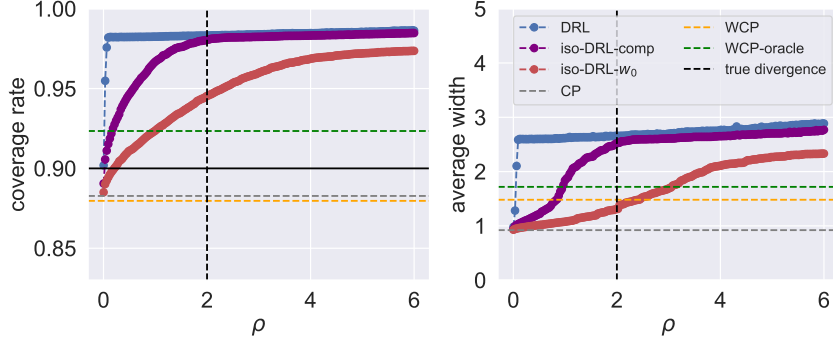


Figure 7: Results with varying  $\rho$  in the well-specified setting. (See Appendix D.1 for details.)

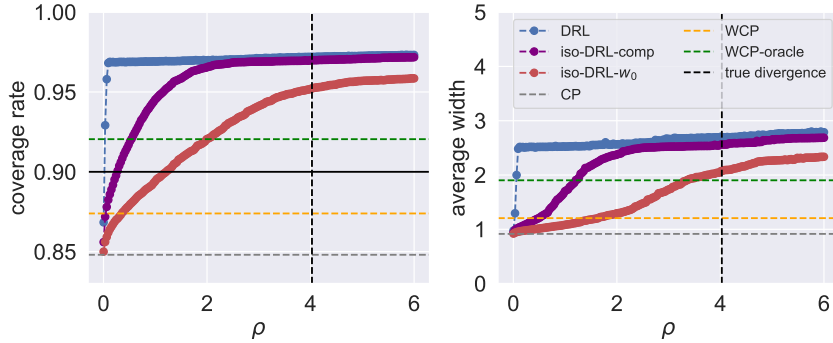


Figure 8: Results with varying  $\rho$  in the misspecified setting. (See Appendix D.1 for details.)

## D.2 Wine quality simulation: a proxy of the oracle KL-divergence

In this section, we examine the choice of  $\rho$  in the wine quality data experiment from Section 5.2. In a real data setting, the true KL divergence,  $D_{\text{KL}}(P_{\text{target}} \| P)$ , is of course unknown, so we need to use a data-driven choice of  $\rho$  in order to implement a DRL procedure (with or without an isotonic constraint).

As is shown in Section 5.2, we denote  $\hat{w}_{\text{kde}}$  as the density ratio obtained by kernel density estimation (Gaussian kernel with bandwidth 0.125). Accordingly, let  $d\hat{Q}_{\text{kde}} = \hat{w}_{\text{kde}} \cdot dP$  be an estimate of  $P_{\text{target}}$ . With a subsample  $\{X_i\}_{i \leq K}$  drawn from the group of white wine (data distribution  $P$ ), a reasonable value for  $\hat{\rho}$  (i.e., an estimate of the true divergence  $\rho$  between the distributions  $P$  and  $P_{\text{target}}$ ) can be calculated by

$$\begin{aligned} \hat{\rho} &= \frac{1}{K} \sum_{i \leq K} \hat{w}_{\text{kde}}(X_i) \log(\hat{w}_{\text{kde}}(X_i)) \\ &\approx \mathbb{E}_P \left\{ \frac{d\hat{Q}_{\text{kde}}}{dP} \log \left( \frac{d\hat{Q}_{\text{kde}}}{dP} \right) \right\} = D_{\text{KL}}(\hat{Q}_{\text{kde}} \| P). \end{aligned}$$

To show the range for values of  $\hat{\rho}$ , we repeatedly fit KDE on the 50% samples from each group (white and red wine groups respectively). Figure 9 shows the histogram of  $\hat{\rho}$  with 1000 repetitions, of which the median is approximately 0.859.

## References

Ai, J. and Ren, Z. (2024). Not all distributional shifts are equal: Fine-grained robust conformal inference. *arXiv preprint arXiv:2402.13042*.

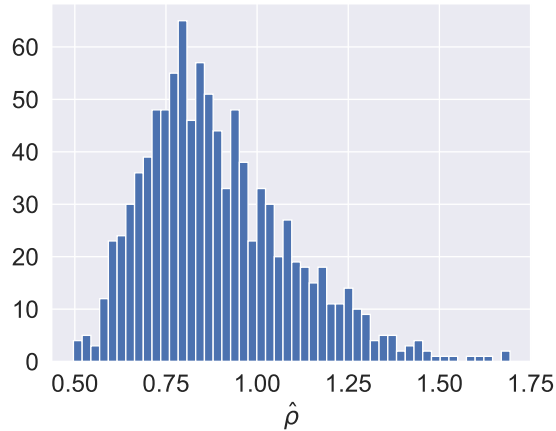


Figure 9: Histogram of  $\hat{\rho}$ . (See Appendix D.2 for details.)

- Bastani, H. (2021). Predicting with proxies: Transfer learning in high dimension. *Manag. Sci.*, 67:2964–2984.
- Bauschke, H. and Combettes, P. (2019). Convex analysis and monotone operator theory in hilbert spaces, corrected printing.
- Ben-David, S., Lu, T., Luu, T., and Pál, D. (2010). Impossibility theorems for domain adaptation. In *AISTATS*.
- Ben-David, S. and Urner, R. (2012). On the hardness of domain adaptation and the utility of unlabeled target samples. In *ALT*.
- Ben-David, S. and Urner, R. (2013). Domain adaptation—can quantity compensate for quality? *Annals of Mathematics and Artificial Intelligence*, 70:185–202.
- Ben-Tal, A. and Nemirovski, A. (1998). Robust convex optimization. *Mathematics of operations research*, 23(4):769–805.
- Berta, E., Bach, F., and Jordan, M. (2024). Classifier calibration with roc-regularized isotonic regression. In *International Conference on Artificial Intelligence and Statistics*, pages 1972–1980. PMLR.
- Birgé, L. (1987). Estimating a density under order restrictions: Nonasymptotic minimax risk. *The Annals of Statistics*, pages 995–1012.
- Blanchet, J., Kang, Y., and Murthy, K. (2019). Robust wasserstein profile inference and applications to machine learning. *Journal of Applied Probability*, 56(3):830–857.
- Blanchet, J. and Murthy, K. (2019). Quantifying distributional model risk via optimal transport. *Mathematics of Operations Research*, 44(2):565–600.
- Blanchet, J. and Shapiro, A. (2023). Statistical limit theorems in distributionally robust optimization. *arXiv preprint arXiv:2303.14867*.
- Bogdan, M., Van Den Berg, E., Sabatti, C., Su, W., and Candès, E. J. (2015). Slope—adaptive variable selection via convex optimization. *The annals of applied statistics*, 9(3):1103.
- Brunk, H. (1963). On an extension of the concept conditional expectation. *Proceedings of the American Mathematical Society*, 14(2):298–304.

- Brunk, H. (1965). Conditional expectation given a  $\sigma$ -lattice and applications. *The Annals of Mathematical Statistics*, 36(5):1339–1350.
- Brunk, H., Barlow, R. E., Bartholomew, D. J., and Bremner, J. M. (1972). Statistical inference under order restrictions.(the theory and application of isotonic regression). *International Statistical Review*, 41:395.
- Brunk, H., Ewing, G., and Utz, W. (1957). Minimizing integrals in certain classes of monotone functions. *Pacific Journal of Mathematics*.
- Cai, T. T. and Wei, H. (2019). Transfer learning for nonparametric classification: Minimax rate and adaptive classifier. *ArXiv*, abs/1906.02903.
- Candès, E., Lei, L., and Ren, Z. (2023). Conformalized survival analysis. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 85(1):24–45.
- Cauchois, M., Gupta, S., Ali, A., and Duchi, J. C. (2020). Robust validation: Confident predictions even when distributions shift. *arXiv preprint arXiv:2008.04267*.
- Chatterjee, S. and Lafferty, J. (2019). Adaptive risk bounds in unimodal regression. *Bernoulli*.
- Chattopadhyay, R., Fan, W., Davidson, I., Panchanathan, S., and Ye, J. (2013). Joint transfer and batch-mode active learning. In *ICML*.
- Chen, M., Weinberger, K. Q., and Blitzer, J. (2011). Co-training for domain adaptation. In *NIPS*.
- Chen, Y. and Lei, J. (2024). De-biased two-sample u-statistics with application to conditional distribution testing. *arXiv preprint arXiv:2402.00164*.
- Cornfield, J., Haenszel, W., Hammond, E. C., Lilienfeld, A. M., Shimkin, M. B., and Wynder, E. L. (1959). Smoking and lung cancer: recent evidence and a discussion of some questions. *Journal of the National Cancer institute*, 22(1):173–203.
- Cortes, C., Mohri, M., Riley, M., and Rostamizadeh, A. (2008). Sample selection bias correction theory. *ArXiv*, abs/0805.2775.
- De Bartolomeis, P., Abad, J., Donhauser, K., and Yang, F. (2023). Hidden yet quantifiable: A lower bound for confounding strength using randomized trials. *arXiv preprint arXiv:2312.03871*.
- Deng, H. and Zhang, C.-H. (2020). Isotonic regression in multi-dimensional spaces and graphs.
- Deng, Z., Dwork, C., and Zhang, L. (2023). Happymap: A generalized multi-calibration method. *arXiv preprint arXiv:2303.04379*.
- Ding, P. and VanderWeele, T. J. (2016). Sensitivity analysis without assumptions. *Epidemiology (Cambridge, Mass.)*, 27(3):368.
- Donsker, M. D. and Varadhan, S. S. (1976). Asymptotic evaluation of certain markov process expectations for large time—iii. *Communications on pure and applied Mathematics*, 29(4):389–461.
- Duchi, J. and Namkoong, H. (2018). Learning models with uniform performance via distributionally robust optimization. *arXiv preprint arXiv:1810.08750*.
- Duchi, J. C., Glynn, P. W., and Namkoong, H. (2021). Statistics of robust optimization: A generalized empirical likelihood approach. *Mathematics of Operations Research*, 46(3):946–969.
- Duchi, J. C., Hashimoto, T., and Namkoong, H. (2019). Distributionally robust losses against mixture covariate shifts. *Under review*, 2(1).

- Dudley, R. M. (1967). The sizes of compact subsets of hilbert space and continuity of gaussian processes. *Journal of Functional Analysis*, 1(3):290–330.
- Durot, C. and Lopuhaä, H. P. (2018). Limit Theory in Monotone Function Estimation. *Statistical Science*, 33(4):547 – 567.
- Edwards, R. E. (2012). *Functional analysis: theory and applications*. Courier Corporation.
- El Ghaoui, L. and Lebret, H. (1997). Robust solutions to least-squares problems with uncertain data. *SIAM Journal on matrix analysis and applications*, 18(4):1035–1064.
- El Ghaoui, L., Oustry, F., and Lebret, H. (1998). Robust solutions to uncertain semidefinite programs. *SIAM Journal on Optimization*, 9(1):33–52.
- Esfahani, P. M. and Kuhn, D. (2015). Data-driven distributionally robust optimization using the wasserstein metric: Performance guarantees and tractable reformulations. *arXiv preprint arXiv:1505.05116*.
- Esteban-Pérez, A. and Morales, J. M. (2022). Partition-based distributionally robust optimization via optimal transport with order cone constraints. *4OR*, 20(3):465–497.
- Gao, F. and Wellner, J. A. (2007). Entropy estimate for high-dimensional monotonic functions. *Journal of Multivariate Analysis*, 98(9):1751–1764.
- Ge, J., Tang, S., Fan, J., Ma, C., and Jin, C. (2023). Maximum likelihood estimation is all you need for well-specified covariate shift. *arXiv preprint arXiv:2311.15961*.
- Grenander, U. (1956). On the theory of mortality measurements. *Skandinavisk Aktuarietidskrift*, 39:1–55.
- Gretton, A., Smola, A., Huang, J., Schmittfull, M., Borgwardt, K. M., Schölkopf, B., Candela, Q., Sugiyama, M., Schwaighofer, A., and Lawrence, N. D. (2009). Covariate shift by kernel mean matching. In *NIPS 2009*.
- Gui, Y., Barber, R., and Ma, C. (2024). Conformalized matrix completion. *Advances in Neural Information Processing Systems*, 36.
- Gui, Y., Hore, R., Ren, Z., and Barber, R. F. (2023). Conformalized survival analysis with adaptive cut-offs. *Biometrika*, page asad076.
- Gupta, S. and Rothenhäusler, D. (2021). The  $s$ -value: evaluating stability with respect to distributional shifts. *arXiv preprint arXiv:2105.03067*.
- Han, Q., Wang, T., Chatterjee, S., and Samworth, R. J. (2019). Isotonic regression in general dimensions. *The Annals of Statistics*.
- Hanneke, S. and Kpotufe, S. (2019). On the value of target data in transfer learning. In *NeurIPS*.
- Hardy, G. H., Littlewood, J. E., and Pólya, G. (1952). *Inequalities*. Cambridge university press.
- Hébert-Johnson, U., Kim, M., Reingold, O., and Rothblum, G. (2018). Multicalibration: Calibration for the (computationally-identifiable) masses. In *International Conference on Machine Learning*, pages 1939–1948. PMLR.
- Henzi, A., Ziegel, J. F., and Gneiting, T. (2021). Isotonic distributional regression. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 83(5):963–993.
- Hu, X. and Lei, J. (2020). A distribution-free test of covariate shift using conformal prediction. *arXiv: Methodology*.

- Hu, Y., Li, M., Lu, Q., Weng, H., Wang, J., Zekavat, S. M., Yu, Z., Li, B., Gu, J., Muchnik, S. K., Shi, Y., Kunkle, B. W., Mukherjee, S., Natarajan, P., Naj, A. C., Kuzma, A., Zhao, Y., Crane, P. K., Lu, H., and Zhao, H. (2019). A statistical framework for cross-tissue transcriptome-wide association analysis. *Nature Genetics*, 51:568–576.
- Jin, Y., Ren, Z., and Candès, E. J. (2023). Sensitivity analysis of individual treatment effects: A robust conformal inference approach. *Proceedings of the National Academy of Sciences*, 120(6):e2214889120.
- Jin, Y., Ren, Z., and Zhou, Z. (2022). Sensitivity analysis under the  $f$ -sensitivity models: a distributional robustness perspective. *arXiv preprint arXiv:2203.04373*.
- Johansson, F. D., Sontag, D. A., and Ranganath, R. (2019). Support and invertibility in domain-invariant representations. *ArXiv*, abs/1903.03448.
- Kim, M. P., Kern, C., Goldwasser, S., Kreuter, F., and Reingold, O. (2022). Universal adaptability: Target-independent inference that competes with propensity scoring. *Proceedings of the National Academy of Sciences*, 119(4):e2108097119.
- Lam, H. (2016). Robust sensitivity analysis for stochastic systems. *Mathematics of Operations Research*, 41(4):1248–1275.
- Ledoux, M. and Talagrand, M. (2013). *Probability in Banach Spaces: isoperimetry and processes*. Springer Science & Business Media.
- Lei, J., G’Sell, M., Rinaldo, A., Tibshirani, R. J., and Wasserman, L. (2018). Distribution-free predictive inference for regression. *Journal of the American Statistical Association*, 113(523):1094–1111.
- Lei, L. and Candès, E. J. (2020). Conformal inference of counterfactuals and individual treatment effects. *arXiv preprint arXiv:2006.06138*.
- Li, S., Cai, T. T., and Li, H. (2021). Transfer learning for high-dimensional linear regression: Prediction, estimation and minimax optimality. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*.
- Liu, J., Wu, J., Wang, T., Zou, H., Li, B., and Cui, P. (2023). Geometry-calibrated dro: Combating over-pessimism with free energy implications. *arXiv preprint arXiv:2311.05054*.
- Ma, C., Pathak, R., and Wainwright, M. J. (2023). Optimally tackling covariate shift in rkhs-based non-parametric regression. *The Annals of Statistics*, 51(2):738–761.
- Matzkin, R. L. (1991). Semiparametric estimation of monotone and concave utility functions for polychotomous choice models. *Econometrica: Journal of the Econometric Society*, pages 1315–1327.
- McDiarmid, C. et al. (1989). On the method of bounded differences. *Surveys in combinatorics*, 141(1):148–188.
- Meggison, R. E. (2012). *An introduction to Banach space theory*, volume 183. Springer Science & Business Media.
- Mei, S., Fei, W., and Zhou, S. (2010). Gene ontology based transfer learning for protein subcellular localization. *BMC Bioinformatics*, 12:44 – 44.
- Namkoong, H. and Duchi, J. C. (2017). Variance-based regularization with convex objectives. *Advances in neural information processing systems*, 30.

- Namkoong, H., Ma, Y., and Glynn, P. W. (2022). Minimax optimal estimation of stability under distribution shift. *arXiv preprint arXiv:2212.06338*.
- Niculescu-Mizil, A. and Caruana, R. A. (2012). Obtaining calibrated probabilities from boosting. *arXiv preprint arXiv:1207.1403*.
- Pathak, R. and Ma, C. (2024). On the design-dependent suboptimality of the lasso. *arXiv preprint arXiv:2402.00382*.
- Pathak, R., Ma, C., and Wainwright, M. (2022). A new similarity measure for covariate shift with applications to nonparametric regression. In *International Conference on Machine Learning*, pages 17517–17530. PMLR.
- Popescu, I. (2007). Robust mean-covariance solutions for stochastic optimization. *Operations Research*, 55(1):98–112.
- Qiu, H., Dobriban, E., and Tchetgen Tchetgen, E. (2023). Prediction sets adaptive to unknown covariate shift. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, page qkad069.
- Rao, B. P. (1969). Estimation of a unimodal density. *Sankhyā: The Indian Journal of Statistics, Series A*, pages 23–36.
- Redko, I., Morvant, E., Habrard, A., Sebban, M., and Bennani, Y. (2020). A survey on domain adaptation theory: learning bounds and theoretical guarantees. *arXiv: Learning*.
- Rényi, A. (1961). On measures of entropy and information. In *Proceedings of the fourth Berkeley symposium on mathematical statistics and probability, volume 1: contributions to the theory of statistics*, volume 4, pages 547–562. University of California Press.
- Rosenbaum, P. R. (1987). Sensitivity analysis for certain permutation inferences in matched observational studies. *Biometrika*, 74(1):13–26.
- Rothenhäusler, D. and Bühlmann, P. (2023). Distributionally robust and generalizable inference. *Statistical Science*, 38(4):527–542.
- Sahoo, R., Lei, L., and Wager, S. (2022). Learning from a biased sample. *arXiv preprint arXiv:2209.01754*.
- Schell, M. J. and Singh, B. (1997). The reduced monotonic regression method. *Journal of the American Statistical Association*, 92(437):128–135.
- Setlur, A., Dennis, D., Eysenbach, B., Raghunathan, A., Finn, C., Smith, V., and Levine, S. (2023). Bitrate-constrained dro: Beyond worst case robustness to unknown group shifts. *arXiv preprint arXiv:2302.02931*.
- Shafer, G. and Vovk, V. (2008). A tutorial on conformal prediction. *Journal of Machine Learning Research*, 9(3).
- Shafieezadeh Abadeh, S., Mohajerin Esfahani, P. M., and Kuhn, D. (2015). Distributionally robust logistic regression. *Advances in Neural Information Processing Systems*, 28.
- Shapiro, A. (2017). Interchangeability principle and dynamic equations in risk averse stochastic programming. *Operations Research Letters*, 45(4):377–381.
- Shapiro, A. and Pichler, A. (2023). Conditional distributionally robust functionals. *Operations Research*.
- Su, W. and Candes, E. (2016). Slope is adaptive to unknown sparsity and asymptotically minimax.

- Sun, Y. and Hu, Y.-J. (2016). Integrative analysis of multi-omics data for discovery and functional studies of complex human diseases. *Advances in genetics*, 93:147–90.
- Tan, Z. (2006). A distributional approach for causal inference using propensity scores. *Journal of the American Statistical Association*, 101(476):1619–1637.
- Tian, Y. and Feng, Y. (2021). Transfer learning under high-dimensional generalized linear models. *ArXiv*, abs/2105.14328.
- Tibshirani, R. J., Barber, R. F., Candès, E. J., and Ramdas, A. (2019). Conformal prediction under covariate shift. In *NeurIPS*.
- Turki, T., Wei, Z., and Wang, J. T.-L. (2017). Transfer learning approaches to improve drug sensitivity prediction in multiple myeloma patients. *IEEE Access*, 5:7381–7393.
- van der Laan, L., Ulloa-Pérez, E., Carone, M., and Luedtke, A. (2023). Causal isotonic calibration for heterogeneous treatment effects. *arXiv preprint arXiv:2302.14011*.
- Vovk, V., Gammerman, A., and Shafer, G. (2005). *Algorithmic learning in a random world*, volume 29. Springer.
- Wainwright, M. J. (2019). *High-dimensional statistics: A non-asymptotic viewpoint*, volume 48. Cambridge university press.
- Wang, Z., Bühlmann, P., and Guo, Z. (2023). Distributionally robust machine learning with multi-source data. *arXiv preprint arXiv:2309.02211*.
- Weiss, A., Lancho, A., Bu, Y., and Wornell, G. W. (2023). A bilateral bound on the mean-square error for estimation in model mismatch. In *2023 IEEE International Symposium on Information Theory (ISIT)*, pages 2655–2660. IEEE.
- Wellner, J. et al. (2013). *Weak convergence and empirical processes: with applications to statistics*. Springer Science & Business Media.
- Weng, C., Shah, N. H., and Hripcsak, G. (2020). Deep phenotyping: Embracing complexity and temporality—towards scalability, portability, and interoperability. *Journal of Biomedical Informatics*, 105:103433 – 103433.
- Yadlowsky, S., Namkoong, H., Basu, S., Duchi, J., and Tian, L. (2018). Bounds on the conditional and average treatment effect with unobserved confounding factors. *arXiv preprint arXiv:1808.09521*.
- Yang, F. and Barber, R. F. (2019). Contraction and uniform convergence of isotonic regression. *Electronic Journal of Statistics*.
- Yang, L., Hanneke, S., and Carbonell, J. G. (2012). A theory of transfer learning with applications to active learning. *Machine Learning*, 90:161–189.
- Yang, Y., Kuchibhotla, A. K., and Tchetgen Tchetgen, E. (2024). Doubly robust calibration of prediction sets under covariate shift. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, page qkae009.
- Zadrozny, B. and Elkan, C. (2002). Transforming classifier scores into accurate multiclass probability estimates. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 694–699.
- Zhang, C.-H. (2002). Risk bounds in isotonic regression. *The Annals of Statistics*, 30(2):528–555.



Zhao, H., des Combes, R. T., Zhang, K., and Gordon, G. J. (2019a). On learning invariant representations for domain adaptation. In *ICML*.

Zhao, Q., Small, D. S., and Bhattacharya, B. B. (2019b). Sensitivity analysis for inverse probability weighting estimators via the percentile bootstrap. *Journal of the Royal Statistical Society Series B: Statistical Methodology*, 81(4):735–761.