

Learning from a Biased Sample

Roshni Sahoo
rsahoo@stanford.edu

Lihua Lei
lihuallei@stanford.edu

Stefan Wager
swager@stanford.edu

Stanford University

Abstract

The empirical risk minimization approach to data-driven decision making assumes that we can learn a decision rule from training data drawn under the same conditions as the ones we want to deploy it under. However, in a number of settings, we may be concerned that our training sample is biased, and that some groups (characterized by either observable or unobservable attributes) may be under- or over-represented relative to the general population; and in this setting empirical risk minimization over the training set may fail to yield rules that perform well at deployment. Building on concepts from distributionally robust optimization and sensitivity analysis, we propose a method for learning a decision rule that minimizes the worst-case risk incurred under a family of test distributions whose conditional distributions of outcomes Y given covariates X differ from the conditional training distribution by at most a constant factor, and whose covariate distributions are absolutely continuous with respect to the covariate distribution of the training data. We apply a result of Rockafellar and Uryasev to show that this problem is equivalent to an augmented convex risk minimization problem. We give statistical guarantees for learning a robust model using the method of sieves and propose a deep learning algorithm whose loss function captures our robustness target. We empirically validate our proposed method in simulations and a case study with the MIMIC-III dataset.

1 Introduction

When learning a data-driven decision rule, sampling bias in the data collection process may prevent practitioners from accessing training data from the distribution that they intend to deploy the rule on. The performance of a learned decision rule may suffer when deployed on populations that differ from the population its training data was drawn from. For example, suppose a practitioner plans to deploy a decision rule across the country but only has data from a handful of states; they may need to reason about how findings from one state generalize to others. Meanwhile, in randomized trials for estimating treatment effects, participants often volunteer or apply to be a part of the study: Attanasio et al. [2011] measures the effect of a vocational training program on labor market outcomes in a randomized trial where participants needed to apply to be a part of the study; and the effectiveness of antidepressants is typically assessed in randomized trials involving volunteers [Wang et al., 2018]. In such studies, participants may differ from non-participants in fundamental ways, and a decision rule based on trial data may perform poorly when deployed on non-participants.

When learning a data-driven decision rule, practitioners typically select the rule that minimizes the risk incurred on the training data. Formally, we suppose that the practitioner observes i.i.d. samples (X, Y) from the training distribution P , where $X \in \mathcal{X}$ is the observed covariate vector and $Y \in \mathcal{Y}$ is the outcome. Given decision rules h and a loss function $L(h(X), Y)$, the practitioner aims to compute

$$\operatorname{argmin}_h \mathbb{E}_P [L(h(X), Y)]. \quad (1)$$

However, a decision rule that satisfies (1) is not guaranteed to perform well on test distributions Q that differ from P , which may arise under sampling bias. The test distribution may differ from the training distribution

Draft version: September 2022. We are grateful for advice from Alyssa Chen regarding the experimental setup of our MIMIC-III case study. Code available at https://github.com/roshni714/ru_regression.

along the observable covariates X . Furthermore, if there exists an unobserved covariate vector $U \in \mathcal{U}$ that affects the outcome Y and the distribution of U changes from train-time to test-time, then the conditional distribution of Y given X of Q will differ from that of P . In this work, we propose a method for learning decision rules that are robust to these shifts.

Many previous works on learning models that are robust to unknown distribution shift apply distributionally robust optimization (DRO) [Ben-Tal et al., 2013]. The goal of DRO is to minimize the worst-case risk over a family of plausible test distributions \mathcal{S} (the robustness set), i.e.

$$\operatorname{argmin}_h \sup_{Q \in \mathcal{S}} \mathbb{E}_Q [L(h(X), Y)]. \quad (2)$$

In this work, we also adopt the DRO framework. To construct our robustness sets, we take inspiration from previous works in the sensitivity analysis literature [Andrews and Oster, 2019, Dorn et al., 2021, Nie et al., 2021, Yadlowsky et al., 2018], which studies the robustness of causal conclusions drawn from observational data to selection bias. In sensitivity analysis, the statistician assumes a sensitivity model, which places assumptions on the selection bias, and analyzes how the causal conclusions are affected by the assumed selection bias. Motivated by the widely-adopted Γ -marginal sensitivity model of Tan [2006], which assumes uniform upper and lower bounds on the amount of selection bias, we define our robustness sets with upper and lower bounds on the likelihood ratio between the conditional test and conditional train distributions. In particular, we consider robustness sets $\mathcal{S}_\Gamma(P, Q_X)$, where $Q \in \mathcal{S}_\Gamma(P, Q_X)$ if Q has conditional distribution of Y given X that differs from the conditional distribution of P by at most a factor, so for some $\Gamma > 1$,

$$\Gamma^{-1} \leq \frac{dQ_{Y|X=x}(y)}{dP_{Y|X=x}(y)} \leq \Gamma, \quad \forall x \in \mathcal{X}, y \in \mathcal{Y}, \quad (3)$$

and marginal distribution equal to Q_X .

The main contribution of this work is a method for solving

$$\operatorname{argmin}_h \sup_{Q \in \mathcal{S}_\Gamma(P, Q_X)} \mathbb{E}_Q [L(h(X), Y)] \quad (4)$$

for any distribution Q_X that is absolutely continuous to P_X . We show that the solution to the following risk minimization problem

$$(h_\Gamma^*, \alpha_\Gamma^*) = \operatorname{argmin}_{h, \alpha} \mathbb{E}_P [L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y)], \quad (5)$$

featuring data drawn from the training distribution P and a particular loss function L_{RU}^Γ , also solves (4) for any distribution Q_X that is absolutely continuous with respect to P_X . We call the minimization problem in (5) Rockafellar-Uryasev (RU) Regression and L_{RU}^Γ the RU loss because we derive them from the results of Rockafellar and Uryasev [2000]. A notable aspect of our proposed method is that it does not require any knowledge of Q_X because it relies on the fact that the minimization of the worst-case risk over a sufficiently flexible class of functions is equivalent to minimization of the conditional worst-case risk for every $x \in \mathcal{X}$.

The remainder of the paper investigates RU Regression theoretically and empirically. In Section 3.1, we demonstrate useful properties of the population RU risk, including convexity, differentiability, existence and uniqueness of the minimizer, and strong convexity around the minimizer. In Section 3.2, these properties enable us to derive estimation guarantees using the method of sieves [Geman and Hwang, 1982]. Furthermore, the useful properties of the population RU risk also suggest that for practical implementation, the problem in (5) can be solved via stochastic gradient descent. As a result, we propose to perform the optimization in (5) by joint-training of neural networks, one for each of h and α , with the RU loss as the objective. In Section 4, we validate our approach in simulations and a case study with the MIMIC-III dataset [Johnson et al., 2016a].

1.1 Related Work

The DRO framework is widely used for learning models that are robust to unknown distribution shift [Duchi and Namkoong, 2021, Duchi et al., 2020, Hu et al., 2018, Michel et al., 2022, Mohajerin Esfahani and Kuhn, 2018, Oberst et al., 2021, Oren et al., 2019, Sagawa et al., 2019, Thams et al., 2022]. Previous works that

apply DRO for learning robust models typically specify a robustness set of interest and provide a method for either evaluating the worst-case risk over the set, learning the solution that minimizes the worst-case risk over the set, or both. These works vary in how they define the robustness set and whether they consider robustness sets over the conditional distribution of Y given X , the marginal distribution over X , or the joint distribution over (X, Y) .

Most of the previous works on learning models that are robust to unknown distribution shift limit the extent to which the joint distribution over (X, Y) can shift [Duchi and Namkoong, 2021, Duchi et al., 2020, Hu et al., 2018, Michel et al., 2022, Mohajerin Esfahani and Kuhn, 2018, Oren et al., 2019, Sagawa et al., 2019]. Duchi and Namkoong [2021] consider f -divergence balls about the joint training distribution as the robustness sets. Duchi et al. [2020] propose *marginal*-DRO, predominantly focusing on the setting where the conditional distribution of Y given X is fixed and the robustness set bounds the amount of covariate shift. Similar to other works that directly place restrictions on the deviations from the joint training distribution, marginal-DRO implicitly limits the amount of shift in the joint distribution by holding the conditional distribution fixed. Hu et al. [2018], Oren et al. [2019], Sagawa et al. [2019] consider *group*-DRO, where the joint training distribution is a mixture of m groups and each group $g \in [m]$ has a corresponding joint distribution P_g and the robustness set consists of all mixtures of these distributions. Michel et al. [2022] model the robustness set by parametrizing the likelihood ratio between the joint training distribution and the worst-case distribution. Mohajerin Esfahani and Kuhn [2018] consider Wasserstein balls about the joint training distribution as the robustness sets. In contrast, the robustness sets that we propose restrict the amount of shift in the conditional distribution but potentially allow for large deviations from the joint training distribution. By directly placing constraints on the conditional distribution, the worst-case risk minimization problem in (4) can be solved conditionally for every $x \in \mathcal{X}$, resulting in a solution that is agnostic to almost arbitrary covariate shift. Our target problem (4) and results are closely related to those of Duchi and Namkoong [2021] and Duchi et al. [2020], and we discuss these connections in Section 2.1.

Also related to our work, Oberst et al. [2021], Thams et al. [2022] consider distribution shifts that arise from changes in the conditional distribution. Oberst et al. [2021] focuses on learning linear models that are robust to changes in the distribution of unobserved variables of bounded magnitude using noisy proxies for the unobserved variables. Thams et al. [2022] proposes a method for evaluation of the worst-case loss under a parametric robustness set, which consists of interpretable shifts in the distribution of observed variables. Both of these works make more fine-grained assumptions about the distribution shift, such as access to proxy variables or parametric shifts, while our shift model is nonparametric and may be viewed as more pessimistic.

Nevertheless, our shift model is based on the Γ -marginal sensitivity model [Tan, 2006], which is used for modeling selection bias due to unmeasured confounding in causal inference. The Γ -marginal sensitivity model and its extensions are used by many previous works that aim to study the sensitivity of causal inference when treatment assignments may depend on unobserved confounders [Andrews and Oster, 2019, Dorn et al., 2021, Jin et al., 2022, Nie et al., 2021, Yadlowsky et al., 2018]. These previous works focus on obtaining partial identification bounds for treatment effects (i.e., the bound of $\operatorname{argmin}_h \mathbb{E}_Q [L(h(X), Y)]$) under the Γ -sensitivity model, whereas we aim to develop decision rules with performance guarantees that are as good as possible across deployment environments that may differ from the training environment across unobservables.

The broader literature on data-driven decision making has been active in recent years, including contributions from Athey and Wager [2021], Bertsimas and Kallus [2020], Elmachoub and Grigas [2022], Foster and Syrgkanis [2019], Kallus and Zhou [2021], Kitagawa and Tetenov [2018], Manski [2004], Nie and Wager [2021], Stoye [2009], Swaminathan and Joachims [2015] and Zhao et al. [2012]. A recurring theme of this line of work is in choosing loss functions $L(\cdot)$ that capture relevant aspects of various decision tasks. Depending on the setting, these loss functions may directly reflect the decision maker’s incurred loss [Bertsimas and Kallus, 2020], leverage implicit representations via importance weighting [Kitagawa and Tetenov, 2018, Swaminathan and Joachims, 2015], or even rely on pre-computed estimates of nuisance components [Athey and Wager, 2021, Foster and Syrgkanis, 2019]. Our results pair naturally with this line of work, in that our approach can be applied with generic loss functions to learn decision rules that are robust to potential test/train bias. We also draw attention to Kallus and Zhou [2021], who consider learning optimal treatment rules from confounded data, i.e., where the “treated” and “control” samples available for training may be biased according to unobservable attributes. Our work is related to that of Kallus and Zhou [2021] in that we both consider using robust optimization techniques to learn from data potentially corrupted via biased

sampling; however, the type of bias we consider (test/train vs. treatment/control), and resulting algorithmic and conceptual remedies, are different.

2 Robustness to Unknown Distribution Shifts

In this section, we formally specify the distribution shifts we consider and our goal of learning a decision rule that is robust to these shifts. Next, we show that learning the robust decision rule is equivalent to an augmented risk minimization problem, motivating the use of modern machine learning techniques for finding the robust decision rule. Finally, we conclude the section with how our setup relates to previously studied DRO problems.

We consider the general loss minimization setting. We have access to $i = 1, \dots, n$ covariate-outcome pairs $(X_i, Y_i) \in \mathcal{X} \times \mathcal{Y}$, where $\mathcal{X} \times \mathcal{Y} \subset \mathbb{R}^d \times \mathbb{R}$, drawn independently from a distribution P , i.e., $(X_i, Y_i) \sim P$. Given a loss function $L(\hat{y}, y)$, we seek to learn a decision rule $h : \mathcal{X} \rightarrow \mathcal{Y}$ such that it achieves low loss $L(h(X), Y)$. At test-time, the data (X, Y) is distributed following a shifted distribution Q , which may differ from P . We do not know Q . However, following the Γ -marginal sensitivity model [Tan, 2006], we assume that there exists $\Gamma > 1$ such that

$$\Gamma^{-1} \leq \frac{dQ_{Y|X=x}(y)}{dP_{Y|X=x}(y)} \leq \Gamma, \quad \forall x \in \mathcal{X}, y \in \mathcal{Y}. \quad (6)$$

In other words, we assume that $Q_{Y|X}$ may be shifted by up to a factor of Γ from $P_{Y|X}$. Then, for any marginal covariate distribution Q_X , let $S_\Gamma(P, Q_X)$ denote the set of all distributions Q that both satisfy the constraint (6) and have marginal distribution over X equal to Q_X . Given any choice of Q_X , the Γ -sensitivity model motivates targeting a robust decision rule

$$h_{Q_X, \Gamma}^* \in \operatorname{argmin}_{h \in L^2(P_X, \mathcal{X})} \sup \left\{ \mathbb{E}_Q [L(h(X), Y)] : Q \in S_\Gamma(P, Q_X) \right\}, \quad (7)$$

where $L^2(P_X, \mathcal{X})$ denotes the space of square-integrable measurable functions with respect to P_X .

The formulation (7) may look challenging to use as the basis for a practical approach to learning. First, it is formulated in terms of the marginal distribution Q_X which may sometimes be known [e.g., Nie et al., 2021], but often is not known. Second, the optimization problem (7) has a min-max form that is not obviously amenable to statistical learning. However, as shown below, both of these concerns can be resolved: There exists a single function h_Γ^* that solves the problem (7) simultaneously for any Q_X that is absolutely continuous with respect to P_X , and furthermore this h_Γ^* can be characterized as the minimizer of a convex loss defined in terms of the observed data distribution P . Specifically, we show the following.

Theorem 1. *Suppose that $(X, Y) \in \mathcal{X} \times \mathcal{Y}$ are drawn i.i.d. with respect to a distribution P for some $\mathcal{X} \subseteq \mathbb{R}^d$ and $\mathcal{Y} \subseteq \mathbb{R}$. Let $L(z, y)$ be a loss function that is convex in z for any $y \in \mathcal{Y}$, and let $\Gamma > 1$. Then the following augmented loss function,*

$$L_{RU}^\Gamma(z, a, y) = \Gamma^{-1}L(z, y) + (1 - \Gamma^{-1})a + (\Gamma - \Gamma^{-1})(L(z, y) - a)_+, \quad (8)$$

is convex in (z, a) for any $y \in \mathcal{Y}$. Furthermore, any solution

$$\{h_\Gamma^*, \alpha_\Gamma^*\} \in \inf_{(h, \alpha) \in L^2(P_X, \mathcal{X}) \times L^2(P_X, \mathcal{X})} \mathbb{E}_P [L_{RU}^\Gamma(h(X), \alpha(X), Y)] \quad (9)$$

is also a solution to (7) for any Q_X that is absolutely continuous with respect to P_X , i.e., $Q_X \ll P_X$. Proof in Appendix C.3.

We name the loss function in (8) the Rockafellar-Uryasev (RU) loss and the minimization problem in (9) RU Regression because they are derived from the results in Rockafellar and Uryasev [2000]. We sketch the proof of Theorem 1 here. First, we define notation that is used in the proof sketch, as well as the remainder of the paper. Let $F_{x;h(x)}(z)$ be the c.d.f. of $L(h(x), Y)$, where Y is distributed according to $P_{Y|X=x}$. So, $F_{x;h(x)}(z)$ is the distribution over the conditional losses when $X = x$. Define the function $q_\eta^L(x; h(x))$ to be the η -th quantile of distribution over the conditional losses when $X = x$, i.e.

$$q_\eta^L(x; h(x)) = F_{x;h(x)}^{-1}(\eta). \quad (10)$$

Also, define

$$\eta(\Gamma) = \frac{\Gamma}{\Gamma + 1}. \quad (11)$$

Now, we proceed with the main proof. Convexity of L_{RU}^Γ follows immediately using the standard rules for composing convex function. We focus on the second claim of Theorem 1. Crucially, we realize that minimizing the worst-case loss in (7) is equivalent to minimizing the worst-case loss conditionally for each $x \in \mathcal{X}$,

$$\min_{h(x) \in \mathbb{R}} \sup \left\{ \mathbb{E}_{Q_{Y|X}} [L(h(x), Y) \mid X = x] : Q \in S_\Gamma(P, Q_X) \right\}. \quad (12)$$

Given our sensitivity model (6), the Neyman-Pearson Lemma yields that

$$\begin{aligned} & \sup \{ \mathbb{E}_{Q_{Y|X}} [L(h(X), Y) \mid X = x] : Q \in S_\Gamma(P, Q_X) \} \\ &= \mathbb{E}_{P_{Y|X}} \left[L(h(X), Y) \left(\Gamma^{-1} + (\Gamma - \Gamma^{-1}) \mathbb{I}(L(h(X), Y) \geq q_{\eta(\Gamma)}^L(X; h(X))) \right) \mid X = x \right]. \end{aligned} \quad (13)$$

So, for each $x \in \mathcal{X}$, our conditional risk minimization problem in (12) can be written as

$$\min_{h(x) \in \mathbb{R}} \mathbb{E}_{P_{Y|X}} \left[L(h(x), Y) \left(\Gamma^{-1} + (\Gamma - \Gamma^{-1}) \mathbb{I}(L(h(x), Y) \geq q_{\eta(\Gamma)}^L(X; h(x))) \right) \mid X = x \right]. \quad (14)$$

We realize that the objective in (14) is closely related to the conditional value-at-risk (CVaR) [Rockafellar and Uryasev, 2000], which is widely considered in the finance literature. For a continuous random variable W with quantile function (inverse c.d.f.) q_W and $\eta \in (0, 1)$, the η -CVaR of W is given by

$$\text{CVaR}_\eta(W) = \mathbb{E} [W \mid W \geq q_W(\eta)].$$

As a result, we can show that (14) can be re-expressed as

$$\min_{h(x) \in \mathbb{R}} \Gamma^{-1} \mathbb{E}_{P_{Y|X}} [L(h(x), Y) \mid X = x] + (1 - \Gamma^{-1}) \cdot \text{CVaR}_{\eta(\Gamma)}(L(h(x), Y)). \quad (15)$$

From Rockafellar and Uryasev [2000], computing the CVaR itself can be formulated as an optimization problem. In particular,

$$\text{CVaR}_{\eta(\Gamma)}(L(h(x), Y)) = \min_{\alpha(x) \in \mathbb{R}} \alpha(x) + \frac{1}{1 - \eta(\Gamma)} \mathbb{E}_{P_{Y|X}} [(L(h(x), Y) - \alpha(x))_+ \mid X = x], \quad (16)$$

and the solution to the joint optimization problem

$$\begin{aligned} & \min_{h(x), \alpha(x) \in \mathbb{R}} \Gamma^{-1} \cdot \mathbb{E}_{P_{Y|X}} [L(h(x), Y) \mid X = x] \\ & + (1 - \Gamma^{-1}) \cdot \left(\alpha(x) + \frac{1}{1 - \eta(\Gamma)} \mathbb{E}_{P_{Y|X}} [(L(h(x), Y) - \alpha(x))_+ \mid X = x] \right) \end{aligned} \quad (17)$$

also solves (15). By the definition of L_{RU}^Γ in (8), the joint optimization problem in (17) can be written as

$$\min_{h(x), \alpha(x) \in \mathbb{R}} \mathbb{E}_{P_{Y|X}} [L_{\text{RU}}^\Gamma(h(x), \alpha(x), Y) \mid X = x]. \quad (18)$$

Lastly, functions $h_\Gamma^*, \alpha_\Gamma^*$ that for every $x \in \text{supp}(P_X)$ solve (18) also solve (9).

Now, we can show that h_Γ^* solves (7) for any Q_X that is absolutely continuous to P_X . Let \mathcal{T} be any set with nonzero measure with respect to Q_X . Then \mathcal{T} must also have nonzero measure with respect to P_X because $Q_X \ll P_X$. So, for any $h \in L^2(Q_X, \mathcal{X})$

$$\sup_{Q_{Y|X}: Q \in S_\Gamma(P, Q_X)} \mathbb{E}_{Q_{Y|X}} [L(h_\Gamma^*(X), Y) \mid X \in \mathcal{T}] \leq \sup_{Q_{Y|X}: Q \in S_\Gamma(P, Q_X)} \mathbb{E}_{Q_{Y|X}} [L(h(X), Y) \mid X \in \mathcal{T}]$$

for any set \mathcal{T} with nonzero measure with respect to Q_X . This is sufficient to show that h_Γ^* is a solution to (7) for any $Q_X \ll P_X$.

Thus, we can show that our robust optimization problem in (7) can be formulated as (9), a risk minimization problem under the training distribution P that involves learning an auxiliary function α along with the decision rule h .

A key aspect of RU Regression is that the optimal decision rule is agnostic to the test covariate distribution Q_X as long as it is absolutely continuous with respect to the training covariate distribution P_X . This is because we propose to learn the minimizer of the worst-case loss *conditionally for every* $x \in \mathcal{X}$. So, the minimizer is a conditional quantity. We can simply study (9) to learn a decision rule that is robust to conditional shifts of the form in (6) and almost arbitrary covariate shifts.

In order for conditional risk minimization to be equivalent to the population risk minimization, we require the decision rule h and auxiliary function α to come from a flexible class, such as $L^2(P_X, \mathcal{X})$. For practical implementation, in Section 4, we propose to use joint optimization of deep neural networks to learn the solution of (9). We will use one neural network to represent h and another neural network to represent α and train the networks with the RU loss using a standard optimization algorithm, such as stochastic gradient descent or its variants.

2.1 Connections to Other DRO Frameworks

At a high level, our approach to learning decision rules under unknown conditional shifts is an instance of DRO [Ben-Tal et al., 2013]. We draw connections between our work and related works that also use the DRO framework to address distribution shift and point out subtleties of our shift model that require new results.

Duchi and Namkoong [2021] consider worst-case shifts in the *joint* distribution over (X, Y) and robustness sets that are f -divergence balls about the training distribution P ; they propose to solve

$$\operatorname{argmin}_h \sup \left\{ \mathbb{E}_Q [L(h(X), Y)] : D_f(Q|P) \leq \rho \right\}, \quad D_f(Q|P) = \int f\left(\frac{dQ}{dP}\right) dP, \quad (19)$$

where D_f is an f -divergence. Clearly, this problem is similar to (7), but we outline some key differences.

One difference between our DRO problem (7) and (19) is the robustness sets that are considered in each problem. To cast (6) as a constraint of the form $D_f(Q|P) \leq \rho$, we would need to consider an “improper” f -divergence, i.e. with

$$f(z) = \begin{cases} 0 & \Gamma^{-1} \leq z \leq \Gamma \\ \infty & \text{else} \end{cases}. \quad (20)$$

The fact that this function is discontinuous and unbounded means that the formal results (and proof strategies) of Duchi and Namkoong [2021] cannot be applied in our setting.

A second difference between our DRO problem (7) and the problem in (19) is that (7) involves constraints on the distribution shift that hold conditionally on x , as arises naturally from our motivating problem (and also in the work of Dorn et al. [2021], Jin et al. [2022], Nie et al. [2021], Yadlowsky et al. [2018]). Constraining conditional shifts simultaneously for every x results in a substantially more complicated optimization problem, requiring more delicate methods and analysis. For example, Levy et al. [2020] proposes a mini-batch gradient-descent algorithm for learning the solution to (19); however, this algorithm cannot be used with conditional constraints (unless one can gather multiple observations for every x , which is impossible for continuous-valued x).

Lastly, our result in Theorem 1 resembles the dual formulation of the problem (19) from Duchi and Namkoong [2021]:

$$\operatorname{argmin}_h \inf_{\lambda, \eta \geq 0} \left\{ \mathbb{E}_P \left[\lambda f^* \left(\frac{L(h(X), Y) - \eta}{\lambda} \right) \right] + \lambda \rho + \eta \right\}, \quad (21)$$

where f^* is the Fenchel conjugate of f . Similar to (9), (21) is an augmented risk minimization problem. Nevertheless, our method, which entails jointly optimizing over the arguments of the augmented risk minimization problem, to our knowledge, has not been considered in the literature. In fact, comments in Namkoong and Duchi [2016] suggest that analogous algorithms would be ill-conditioned with general f -divergences due to the dependence on λ^{-1} in the first term. Nevertheless, for the improper function f defined in (20), $f^*(u) = \Gamma(u)_+ - \Gamma^{-1}(u)_-$. It is not hard to see that λ can be removed from the optimization

problem (21). So, our approach exploits special structure in our distribution shift model that is not present in the problems studied in Duchi and Namkoong [2021].

Duchi et al. [2020] also applies the DRO framework and is very related to our work. They use DRO to learn a model that is robust to worst-case shifts in the *marginal* distribution over X while the conditional distribution of Y given X is held fixed. They consider a set of plausible subpopulations

$$\mathcal{P}_{\alpha_0, X} := \{Q_{X,0} : P_X = \alpha Q_{X,0} + (1 - \alpha)Q_{X,1} \text{ for some } \alpha \geq \alpha_0, \text{ distribution } Q_{X,1} \text{ on } \mathcal{X}\} \quad (22)$$

and aim to minimize the worst-case subpopulation risk

$$\min_h \sup_{Q_{X,0} \in \mathcal{P}_{\alpha_0, X}} \mathbb{E}_{Q_{X,0}} [\mathbb{E}_{P_{Y|X}} [L(h(X), Y) | X]]. \quad (23)$$

The constraint on the allowable covariate shifts given in (22) is equivalent to requiring that $0 \leq \frac{dQ_{X,0}(x)}{dP_X(x)} \leq \frac{1}{\alpha_0}$ for all $x \in \mathcal{X}$, which is similar to the bounds on the likelihood ratio between the conditional test and train distributions we impose in (6). Nevertheless, our target problem (7) differs from (23) because (6) places restrictions on shifts in the conditional distribution, which causes our minimization problem to be agnostic to shifts in the covariate distribution. Notably, this relationship is not symmetric—placing restrictions on the covariate distribution does not yield a minimization problem that is agnostic to shifts in the conditional distribution.

3 Theoretical Guarantees

In this section, we first demonstrate useful properties of the population RU risk. We demonstrate that the population RU risk has a unique minimizer and is strongly convex and smooth about the minimizer. These properties enable us to obtain nonparametric estimation guarantees by applying the method of sieves in Section 3.2.

3.1 Properties of Population RU Risk

First, we consider the problem of minimizing the population RU risk with respect to (h, α) over $L^2(P_X, \mathcal{X}) \times L^2(P_X, \mathcal{X})$. We consider the following norm on this product space

$$\|(h, \alpha)\|_{L^2(P_X, \mathcal{X})} = \sqrt{\|h\|_{L^2(P_X, \mathcal{X})}^2 + \|\alpha\|_{L^2(P_X, \mathcal{X})}^2}.$$

Under the following two assumptions, we can show that any minimizer of the population RU risk lies in a bounded subset of $L^2(P_X, \mathcal{X}) \times L^2(P_X, \mathcal{X})$.

Assumption 1. $\mathcal{X} \times \mathcal{Y}$ is compact.

Assumption 2. The loss function $L(\hat{y}, y) = \ell(y - \hat{y})$ for some function $\ell(z)$ that is $C_{L,\ell}$ -strongly convex, twice-differentiable and is minimized at $\ell(0) = 0$.

Since \mathcal{Y} is bounded, $\mathcal{Y} \subset [-B, B]$. So, we can define a bounded class of decision rules

$$\mathcal{H} = \{h \in L^2(P_X, \mathcal{X}) \mid \|h\|_\infty \leq 2B\}.$$

We define a constant M_u such that

$$\sup_{h \in \mathcal{H}, x \in \mathcal{X}} q_{\eta(\Gamma)}^L(x; h(x)) < M_u, \quad (24)$$

and note that $M_u < \infty$ because \mathcal{H} is bounded and $\mathcal{X} \times \mathcal{Y}$ is compact. We define the bounded class \mathcal{A} for the auxiliary functions

$$\mathcal{A} = \{\alpha \in L^2(P_X, \mathcal{X}) \mid 0 \leq \alpha(x) \leq M_u \quad \forall x \in \mathcal{X}\}.$$

Let $\Theta = \mathcal{H} \times \mathcal{A}$. In the following result, we show that minimizing the population RU risk over $L^2(P_X, \mathcal{X}) \times L^2(P_X, \mathcal{X})$ is equivalent to minimizing the population RU risk over Θ .

Lemma 2. Under Assumption 1, 2, if any minimizer of $(h, \alpha) \mapsto \mathbb{E}_P [L_{RU}^\Gamma(h(X), \alpha(X), Y)]$ exists over $L^2(P_X, \mathcal{X}) \times L^2(P_X, \mathcal{X})$, then it must lie in Θ . *Proof in Appendix C.4.*

From now on, we will only consider minimization of the population RU risk over Θ . We can show that the population RU risk has at least one minimizer on Θ .

Lemma 3. Under Assumption 1, 2, $\mathbb{E}_P [L_{RU}^\Gamma(h(X), \alpha(X), Y)]$ has at least one minimizer on Θ .

To show that the population RU risk is strictly convex on Θ , we make the following assumption on the conditional distribution $P_{Y|X=x}$.

Assumption 3. For every $x \in \mathcal{X}$, we assume that $P_{Y|X=x}(y)$ is differentiable and strictly increasing in its argument and has positive density on \mathcal{Y} . We assume that $\sup_{x \in \mathcal{X}, y \in \mathbb{R}} p_{Y|X=x}(y) \leq C_{p,u}$, where $0 < C_{p,u} < \infty$.

Lemma 4. Under Assumptions 1, 2, 3, $\mathbb{E}_P [L_{RU}^\Gamma(h(X), \alpha(X), Y)]$ is strictly convex in (h, α) on Θ . *Proof in Appendix C.6.*

As a consequence of strict convexity on Θ , the population RU risk must have at most one minimizer over Θ . Meanwhile, Lemma 3 gives that it has at least one minimizer over Θ , as well. Combining these results gives that the population RU risk has a unique minimizer over Θ . Because of Lemma 2, this means that the population RU risk also has a unique minimizer over all of $L^2(P_X, \mathcal{X}) \times L^2(P_X, \mathcal{X})$.

Theorem 5. Under Assumptions 1, 2, 3, $\mathbb{E}_P [L_{RU}^\Gamma(h(X), \alpha(X), Y)]$ has a unique minimizer $(h_\Gamma^*, \alpha_\Gamma^*)$ over Θ . *Proof in Appendix C.7.*

In addition, we can develop an interpretation of α_Γ^* that minimizes the population RU risk.

Lemma 6. Under Assumptions 1, 2, 3,

$$\alpha_\Gamma^*(x) = q_{\eta(\Gamma)}^L(x; h_\Gamma^*(x)),$$

and there exists $M_l > 0$ such that

$$\alpha_\Gamma^*(x) > M_l \quad \forall x \in \mathcal{X}.$$

Proof in Appendix C.8.

Using Lemma 6, we can show that the population RU risk is strongly convex near the minimizer. We define constants that will be used in the proof of strong convexity. Recall that under Assumption 2, we can rewrite $L(\hat{y}, y) = \ell(y - \hat{y})$. Let ℓ_1^{-1} be the inverse of $\ell(z)$ where $z > 0$. Let ℓ_2^{-1} be the inverse of $\ell(z)$ where $z \leq 0$. Define

$$C_{a,u} := \sup_{i \in \{1,2\}} |\ell'(\ell_i^{-1}(M_u))|, \quad (25)$$

$$C_{a,l} := \inf_{i \in \{1,2\}} |\ell'(\ell_i^{-1}(M_l))|. \quad (26)$$

To define the next set of constants, we define $q_c^Y(x)$ to be the c -th quantile of Y where Y is distributed according to $P_{Y|X=x}$.

$$C_{p,l} := \inf_{c \in [1 - \frac{\eta(\Gamma)}{2}, 1 + \frac{\eta(\Gamma)}{2}], x \in \mathcal{X}} p_{Y|X=x}(q_c^Y(x)), \quad (27)$$

$$\kappa_1 := (1 - \Gamma^{-1}) \cdot \frac{C_{L,l} \cdot C_{p,l}}{2C_{p,u} \cdot (C_{a,u} \cdot C_{a,l} + 1) + C_{L,l} \cdot C_{a,l}} \cdot \frac{C_{a,l}}{C_{a,u}}. \quad (28)$$

Additionally, let

$$C_{a,l,\delta} := \inf_{i \in \{1,2\}} |\ell'(\ell_i^{-1}(M_l - \delta))|, \quad (29)$$

$$C_{p,l,\epsilon} := \inf_{c \in [1 - \frac{\eta(\Gamma)}{2}, 1 + \frac{\eta(\Gamma)}{2}], b \in [-\epsilon, \epsilon], x \in \mathcal{X}} p_{Y|X=x}(q_c^Y(x) + b). \quad (30)$$

We can show that in a $\|\cdot\|_\infty$ -ball about the minimizer, the population RU loss is strongly convex, where the constant of strong convexity approaches κ_1 as the ball's radius shrinks.

Theorem 7. Suppose Assumptions 1, 2, 3, hold. Let $\mathcal{C}_\delta = \{(h, \alpha) \in \Theta \mid \|(h, \alpha) - (h_\Gamma^*, \alpha_\Gamma^*)\|_\infty < \delta\}$, and let $0 < \epsilon < \frac{1-\eta(\Gamma)}{2C_{p,u}}$. There exists $0 < \delta(\epsilon) < M_l$ such that $\mathbb{E}_P [L_{RU}^\Gamma(h(X), \alpha(X), Y)]$ is $\kappa_{1,\epsilon}$ -strongly convex on $\mathcal{C}_{\delta(\epsilon)}$, where

$$\kappa_{1,\epsilon} := (\Gamma - \Gamma^{-1}) \frac{C_{L,l} \cdot (1 - \eta(\Gamma) - 2C_{p,u} \cdot \epsilon) \cdot C_{p,l,\epsilon}}{2C_{p,u} \cdot (C_{a,u} \cdot C_{a,l,\delta(\epsilon)} + 1) + C_{L,l} \cdot C_{a,l,\delta(\epsilon)}} \cdot \frac{C_{a,l,\delta(\epsilon)}}{C_{a,u}}. \quad (31)$$

As $\epsilon \rightarrow 0$,

$$\kappa_{1,\epsilon} \rightarrow \kappa_1.$$

Proof in Appendix C.9.

To show that the population RU risk is smooth in an $\|\cdot\|_\infty$ -ball around the minimizer, we require an additional assumption on the loss function L . Essentially, we need L to be $C_{L,u}$ -smooth for some constant $0 < C_{L,u} < \infty$.

Assumption 4. The second derivative of $\ell(z)$ as defined in Assumption 2 is upper bounded by $C_{L,u}$, where $0 < C_{L,u} < \infty$.

The constant for smoothness depends on the constant $C_{p,u}$ from Assumption 3, $C_{a,u}$ from (25), $C_{a,l,\delta}$ from (29), and $C_{L,u}$ from Assumption 4. Let

$$\kappa_2 := (\Gamma - \Gamma^{-1}) \cdot \left(2C_{p,u} \left(C_{a,u} + \frac{1}{C_{a,l}} \right) \right) + \Gamma \cdot C_{L,u}. \quad (32)$$

We can show that in an $\|\cdot\|_\infty$ ball about the minimizer, the population RU risk is smooth, where the constant for smoothness of the population RU risk approaches κ_2 as the radius of the ball decreases.

Theorem 8. Suppose Assumptions 1, 2, 3, 4 hold. Let $\mathcal{C}_\delta = \{(h, \alpha) \in \Theta \mid \|(h, \alpha) - (h_\Gamma^*, \alpha_\Gamma^*)\|_\infty < \delta\}$. For every $0 < \epsilon < \frac{1-\eta(\Gamma)}{2C_{p,u}}$, there is $0 < \delta(\epsilon) < M_l$ such that $\mathbb{E}_P [L_{RU}^\Gamma(h(X), \alpha(X), Y)]$ is $\kappa_{2,\epsilon}$ -smooth in (h, α) on $\mathcal{C}_{\delta(\epsilon)}$ where

$$\kappa_{2,\epsilon} := (\Gamma - \Gamma^{-1}) \cdot \left(2C_{p,u} \left(C_{a,u} + \frac{1}{C_{a,l,\delta(\epsilon)}} \right) \right) + \Gamma \cdot C_{L,u}. \quad (33)$$

As $\epsilon \rightarrow 0$,

$$\kappa_{2,\epsilon} \rightarrow \kappa_2.$$

Proof in Appendix C.10.

3.2 Estimation Guarantees via Method of Sieves

To simplify notation, we denote $\theta := (h, \alpha)$ and rewrite the population RU risk as $\mathbb{E}_P [L_{RU}^\Gamma(\theta(X), Y)]$. The empirical risk is accordingly

$$\widehat{\mathbb{E}}_P [L_{RU}^\Gamma(\theta(X), Y)] = \frac{1}{n} \sum_{i=1}^n L_{RU}^\Gamma(\theta(X_i), Y_i). \quad (34)$$

In addition, we will denote the minimizer of the population RU risk as simply $\theta^* := (h_\Gamma^*, \alpha_\Gamma^*)$, omitting the dependence on Γ .

Thus far, we have demonstrated that θ^* is the minimizer of the population RU risk over the infinite-dimensional space Θ . In practice, we aim to minimize the empirical RU loss (34). However, due to the computational difficulties of estimating infinite-dimensional models using finite-samples, we do not minimize the empirical risk over Θ directly. Instead, we apply the method of sieves [Geman and Hwang, 1982]; we consider optimizing the empirical risk over an increasing sequence of sieves $\Theta_1 \subseteq \Theta_2 \subseteq \dots \subseteq \Theta$, which are finite-dimensional parameter spaces. The sieves we consider have the property that $\inf_{\theta \in \Theta_m} \|\theta - \theta^*\|_\infty \rightarrow 0$ as $m \rightarrow \infty$. To ensure consistency, we increase the complexity of the sieves with the sample size. We let

$$\hat{\theta}_n = \underset{\theta \in \Theta_n}{\operatorname{argmin}} \widehat{\mathbb{E}}_P [L_{RU}^\Gamma(\theta(X), Y)].$$

To estimate h , it is sufficient to consider sieves that consist functions bounded between $-2B$ and $2B$. To estimate α , it is sufficient to consider sieves that consist of nonnegative, bounded functions because any minimizer of the population RU risk has $0 \leq \alpha^*(x) \leq M_u$ for all $x \in \mathcal{X}$. In order to make our sieve-based estimates $h(x), \alpha(x)$ be bounded, we use the same strategy as in Jin et al. [2022]; we truncate standard sieve space to bounded functions. The following two natural examples of truncated sieve spaces were also discussed by Jin et al. [2022]:

Example 1 (Polynomials). Let $\text{Pol}(J_n)$ be the space of polynomials on $[0, 1]$ of degree J_n or less; that is

$$\text{Pol}(J_n) = \left\{ x \mapsto \sum_{k=0}^{J_n} a_k x^k, x \in [0, 1] : a_k \in \mathbb{R} \right\}.$$

Let $\text{Pol}(J_n, a, b)$ be the space of polynomials on $[0, 1]$ of degree J_n or less that are bounded between a and b ; that is

$$\text{Pol}(J_n, a, b) = \left\{ x \mapsto \min(\max(f(x), a), b), x \in [0, 1] : f \in \text{Pol}(J_n) \right\}.$$

Then, we define the sieve *with truncation* as $\Theta_n = \mathcal{H}_n \times \mathcal{A}_n$, where $\mathcal{H}_n = \{x \mapsto \prod_{k=1}^d f_k(x_k) : f_k \in \text{Pol}(J_n, -2B, 2B), k = 1, \dots, d\}$ and $\mathcal{A}_n = \{x \mapsto \prod_{k=1}^d f_k(x_k) : f_k \in \text{Pol}(J_n, 0, M_u)\}$ for $J_n \rightarrow \infty$. We can also define the sieve without truncation as $\tilde{\Theta}_n = \tilde{\mathcal{H}}_n \times \tilde{\mathcal{A}}_n$, where $\tilde{\mathcal{H}}_n = \{x \mapsto \prod_{k=1}^d f_k(x_k) : f_k \in \text{Pol}(J_n)\}$ for $J_n \rightarrow \infty$ and $\tilde{\mathcal{A}}_n = \{x \mapsto \prod_{k=1}^d f_k(x_k) : f_k \in \text{Pol}(J_n)\}$ for $J_n \rightarrow \infty$.

Example 2 (Univariate Splines). Let J_n be a positive number, and let $t_0, t_1, \dots, t_{J_n}, t_{J_n+1}$ be real numbers with $0 = t_0 < t_1 < \dots < t_{J_n} < t_{J_n+1} = 1$. Partition $[0, 1]$ into $J_n + 1$ subintervals $I_j = [t_j, t_{j+1}]$, $j = 0, \dots, J_n - 1$ and $I_{J_n} = [t_{J_n}, t_{J_n+1}]$. We assume that the knots t_1, t_2, \dots, t_{J_n} have bounded mesh ratio:

$$\frac{\max_{0 \leq j \leq J_n} (t_{j+1} - t_j)}{\min_{0 \leq j \leq J_n} (t_{j+1} - t_j)} \leq c \text{ for some constant } c > 0.$$

Let $r \geq 1$ be an integer. A spline of order r with knots $t_1 \dots t_{J_n}$ is given by

$$\text{Spl}(r, J_n) = \left\{ \sum_{k=0}^{r-1} a_k x^k + \sum_{j=1}^{J_n} b_j [\max\{x - t_j, 0\}]^{r-1}, x \in [0, 1] : a_k, b_j \in \mathbb{R} \right\}.$$

Let $\text{Spl}(r, J_n, a, b)$ be the space of splines that are bounded between a and b ; that is

$$\text{Spl}(r, J_n, a, b) = \left\{ x \mapsto \min(\max(f(x), a), b), x \in [0, 1] : f \in \text{Spl}(r, J_n) \right\}.$$

Then, we define the sieve *with truncation* as $\Theta_n = \mathcal{H}_n \times \mathcal{A}_n$, where $\mathcal{H}_n = \{x \mapsto \prod_{k=1}^d f_k(x_k) : f_k \in \text{Spl}(r, J_n, -2B, 2B), k = 1, \dots, d\}$ and $\mathcal{A}_n = \{x \mapsto \prod_{k=1}^d f_k(x_k) : f_k \in \text{Spl}(r, J_n, 0, M_u)\}$ for $J_n \rightarrow \infty$. We can also define the sieve without truncation as $\tilde{\Theta}_n = \tilde{\mathcal{H}}_n \times \tilde{\mathcal{A}}_n$, where $\tilde{\mathcal{H}}_n = \{x \mapsto \prod_{k=1}^d f_k(x_k) : f_k \in \text{Spl}(r, J_n)\}$ for $J_n \rightarrow \infty$ and $\tilde{\mathcal{A}}_n = \{x \mapsto \prod_{k=1}^d f_k(x_k) : f_k \in \text{Spl}(r, J_n)\}$ for $J_n \rightarrow \infty$.

We prove results that demonstrate the consistency of the sieve estimation procedure. Let

$$\theta_m^* = \underset{\theta \in \Theta_m}{\text{argmin}} \mathbb{E}_P [L_{\text{RU}}^\Gamma(\theta(X), Y)]. \quad (35)$$

First, we show that θ_m^* is the unique minimizer of the population RU risk over the sieve space Θ_m . Then, we prove that the sieve approximation error, the bias that results from minimizing the population RU risk over a finite-dimensional sieve space, converges to zero as the dimension of the sieve spaces goes to infinity. Then, we consider $\hat{\theta}_{m,n}$, the minimizer of the empirical risk over Θ_m , i.e.

$$\hat{\theta}_{m,n} = \underset{\theta \in \Theta_m}{\text{argmin}} \widehat{\mathbb{E}}_P [L_{\text{RU}}^\Gamma(\theta(X), Y)]$$

for a sufficiently large integer m . We can show that the estimation error, the error that results from estimating the minimizer of the empirical risk (in finite samples) in a fixed sieve space, converges to zero in probability.

Lemma 9. Under Assumptions 1, 2, 3, $\mathbb{E}_P [L_{RU}^\Gamma(\theta(X), Y)]$ has a unique minimizer over Θ_m called θ_m^* . *Proof in Appendix C.11*

Theorem 10. Under Assumptions 1, 2, 3, as $m \rightarrow \infty$,

$$\|\theta_m^* - \theta^*\|_{L^2(P_{X, \mathcal{X}})} \rightarrow 0.$$

Proof in Appendix C.12

Lemma 11. Under Assumptions 1, 2, 3, $\hat{\theta}_{m,n}$ exists with probability approaching 1 and

$$\hat{\theta}_{m,n} \xrightarrow{P} \theta_m^*$$

as $n \rightarrow \infty$ and m sufficiently large. *Proof in Appendix C.13.*

Combining Theorem 10 and Lemma 11 implies the consistency of the sieve estimation procedure: as $m, n \rightarrow \infty$,

$$\|\hat{\theta}_{m,n} - \theta^*\|_{L^2(P_{X, \mathcal{X}})} \leq \|\hat{\theta}_{m,n} - \theta_{m,n}^*\|_{L^2(P_{X, \mathcal{X}})} + \|\theta_{m,n}^* - \theta^*\|_{L^2(P_{X, \mathcal{X}})} \xrightarrow{P} 0.$$

To obtain a rate of convergence, we consider the classes of sufficiently smooth functions. Given a d -tuple $\beta = (\beta_1, \dots, \beta_d)$ of nonnegative integers, set $[\beta] = \beta_1 + \beta_2 + \dots + \beta_d$ and let D^β denote the differential operator defined by $D^\beta = \frac{\partial^{|\beta|}}{\partial x_1^{\beta_1} \dots \partial x_d^{\beta_d}}$. A real-valued function h on \mathcal{X} is p -smooth if it is m times continuously differentiable on \mathcal{X} and $D^\beta h$ satisfies a Hölder condition (Definition 3) with exponent γ for all d -tuples β of nonnegative integers with $[\beta] = m$. Denote the Hölder class, or the class of all p -smooth real-valued functions on \mathcal{X} , by $\Lambda^p(\mathcal{X})$, and the space of all m -times differentiable real-valued functions on \mathcal{X} by $C^m(\mathcal{X})$. Define a Hölder ball with smoothness $p = m + \gamma$ as

$$\Lambda_c^p(\mathcal{X}) = \left\{ h \in C^m(\mathcal{X}) : \sup_{[\beta] \leq m} \sup_{x \in \mathcal{X}} |D^\beta h(x)| \leq c, \sup_{\substack{[\beta] = m \\ x, y \in \mathcal{X}, \\ x \neq y}} \frac{|D^\beta h(x) - D^\beta h(y)|}{|x - y|_2^\gamma} \leq c \right\}.$$

To ensure that h, α are bounded, we define the truncated function class

$$\Lambda_c^p(\mathcal{X}, a, b) := \{x \mapsto \min(\max(f(x), a), b), f \in \Lambda_c(\mathcal{X})\}.$$

To obtain a rate of convergence for the estimators, we impose the following assumption on the true optimizer.

Assumption 5. Assume that $\theta^* \in \Lambda_c^p(\mathcal{X}, -2B, 2B) \times \Lambda_c^p(\mathcal{X}, 0, M_u)$ for some $c > 0$. We redefine $\Theta := \Lambda_c^p(\mathcal{X}, -2B, 2B) \times \Lambda_c^p(\mathcal{X}, 0, M_u)$.

We also require that the second moment of Y , where Y is distributed following $P_{Y|X=x}$, is bounded for all $x \in \mathcal{X}$.

Assumption 6. We assume that $\sup_{x \in \mathcal{X}} \mathbb{E}_{P_{Y|X}} [Y^2 | X = x] < \infty$.

In addition, we require the following condition on the density of P_X .

Assumption 7. P_X has a density that is bounded away from 0 and ∞ , i.e. $0 < \inf_{x \in \mathcal{X}} p_X(x) < \sup_{x \in \mathcal{X}} p_X(x) < \infty$ for all $x \in \mathcal{X}$.

Under this last assumption, $\|\cdot\|_{L^2(P_{X, \mathcal{X}})} \asymp \|\cdot\|_{L^2(\lambda, \mathcal{X})}$, where λ is the Lebesgue measure. Finally, with these assumptions, we can apply a result from Chen [2007] to show the following rate of convergence. The proof of the result requires balancing the sieve approximation error and estimation error. To get a handle on the sieve approximation error, we use the result from Timan [2014] that for the sieves $\tilde{\Theta}_{J_n}$ in Example 1 and 2 and $\theta^* \in \Lambda_c^p(\mathcal{X}) \times \Lambda_c^p(\mathcal{X})$ for \mathcal{X} compact,

$$\inf_{\theta \in \tilde{\Theta}_{J_n}} \|\theta - \theta^*\|_\infty = O(J_n^{-p}).$$

Theorem 12. Let $J_n = \left(\frac{n}{\log n}\right)^{\frac{1}{2p+d}}$. Under Assumptions 1, 2, 3, 4, 5, 6, 7,

$$\|\hat{\theta}_n - \theta^*\|_{L^2(P_{X, \mathcal{X}})} = O_P\left(\left(\frac{\log n}{n}\right)^{\frac{p}{2p+d}}\right).$$

Proof in Appendix C.14.

4 Experiments

We evaluate the empirical performance of RU Regression when neural networks are used to learn h, α . First, we demonstrate that RU Regression enables us to learn models that are robust to conditional shifts in simulation experiments with synthetic data. Second, we apply RU Regression in a semi-synthetic distribution shift experiment with patient length-of-stay data from the MIMIC-III dataset [Johnson et al. \[2016a\]](#). The code for our experiments is available in https://github.com/roshni714/ru_regression.

4.1 Simulations with Synthetic Data

We perform two simulations with synthetic data. We first consider a one-dimensional toy example because it permits visualization of the data distributions and the learned models. Next, we show that similar trends hold in a high-dimensional simulation. Implementation details for these experiments can be found in [Appendix A](#).

4.1.1 Methods

We compare the following three baselines.

1. Standard ERM - We fit a neural network model with the squared loss function

$$L(\hat{y}, y) = (y - \hat{y})^2 \tag{36}$$

on the training data.

2. Oracle ERM - We fit a neural network model with the squared loss function (36) on data sampled from the test distribution.
3. Rockafellar-Uryasev Regression (RU Regression) - We fit two neural networks with the RU loss function, where one network learns h and the other network learns α , on the training data. We set L to the squared loss function (36). A visualization for the model architecture is provided in [Figure 3](#) in [Appendix A](#).

The Standard ERM and RU Regression methods are not necessarily trained on data from the same distribution as the test distribution (see further details in [Section 4.1.2](#) and [Section 4.1.3](#)). The Oracle ERM method is used to provide an upper bound on the performance of Standard ERM and RU Regression because it is trained on the same distribution as the test distribution.

4.1.2 One-Dimensional Toy Example

Data Generation. We generate a synthetic dataset of samples of the form (X_i, Y_i, U_i) , where $X_i \in \mathbb{R}$ represents observed covariates, $Y_i \in \mathbb{R}$ represents the outcome, and $U_i \in \{0, 1\}$ represents an unobserved variable that influences the outcome Y_i . We suppose the data is distributed as follows

$$X_i \sim \text{Uniform}[0, 10], \quad U_i \sim \text{Bernoulli}(p), \quad Y_i | X_i \sim N(\sqrt{X_i} + U_i(3\sqrt{X_i} + 1), 1). \tag{37}$$

The outcomes Y_i can be clustered into two bands corresponding to $U_i = 1$ and $U_i = 0$, respectively. In this simulation, we consider distribution shift that results from varying p , the probability that $U_i = 1$. We suppose that in the training distribution $p = 0.2$, so we are less likely to observe examples with $U_i = 1$. At test-time, we evaluate the learned models on data distributions where $p \in [0.1, 0.2, 0.5, 0.7, 0.9]$. These data distributions are visualized in [Figure 1](#).

Results. From [Table 1](#), Standard ERM achieves low test MSE on test distributions that are similar to the training distribution $p \in [0.1, 0.2]$. However, the test MSE of Standard ERM increases as p increases. The RU Regression methods achieve higher test MSE on the original training distribution than Standard ERM but are more robust than Standard ERM as p increases. Note that Oracle ERM outperforms both the Standard ERM and RU Regression methods; this is expected because the Oracle ERM model is trained on data from the same distribution as the test distribution. We note that RU Regression matches the performance of the Oracle ERM model when $p = 0.5$.

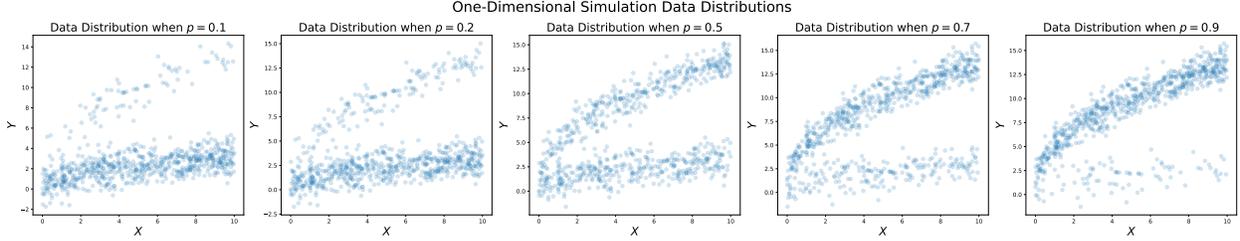


Figure 1: From left to right, we visualize the distribution over (X_i, Y_i) as p varies in $[0.1, 0.2, 0.5, 0.7, 0.9]$. We note that as p increases the proportion of samples where $U_i = 1$ increases.

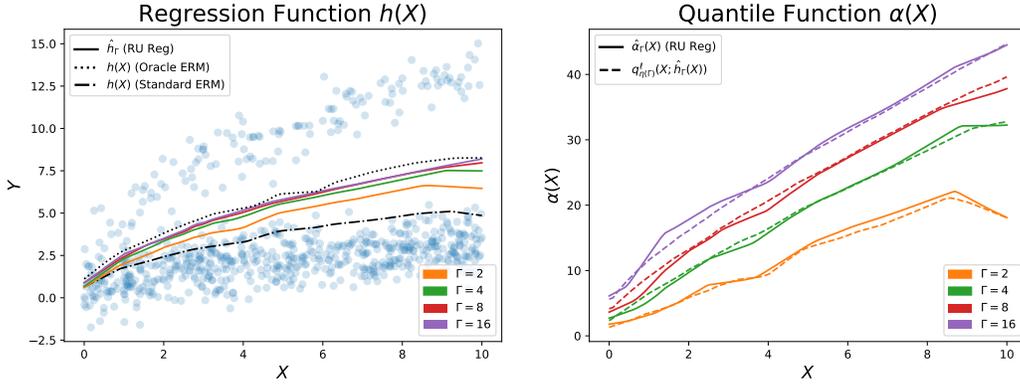


Figure 2: **Left:** We visualize the decision rules \hat{h} that are learned in our one-dimensional toy example. Standard ERM incurs low error on samples with $U_i = 0$ but high error on samples with $U_i = 1$. Models learned via RU Regression incur lower error on samples with $U_i = 1$. The Oracle ERM model visualized here is the model that is trained on the data distribution when $p = 0.5$. **Right:** We visualize the auxiliary function $\hat{\alpha}_\Gamma$ that is learned in the RU Regression methods for our one-dimensional toy example. We realize that the learned $\hat{\alpha}_\Gamma$ closely tracks $q_{\eta(\Gamma)}^L(x; \hat{h}_\Gamma(x))$ as expected.

In addition, we visualize the regression functions learned from each of the methods. From the left plot of Figure 2, it is clear that the regression model learned via Standard ERM incurs high error on samples with $U_i = 1$ and low error on samples with $U_i = 0$, which explains why the method performs poorly on distributions with higher p (higher proportion of samples with $U_i = 1$). Furthermore, we observe that increasing Γ yields regression functions that incur lower error on samples with $U_i = 1$, relative to the Standard ERM model. The Oracle ERM model visualized in Figure 2 is the model that is trained on data generated when $p = 0.5$. We see that this model makes similar predictions as the RU Regression models, which explains why the RU Regression models perform similarly to the Oracle ERM model on the $p = 0.5$ test distribution.

Furthermore, we verify that the solution learned by the neural network is consistent with Theorem 6, which states that

$$\alpha_\Gamma^*(x) = q_{\eta(\Gamma)}^L(x; h_\Gamma^*(x)) \quad \forall x \in \mathcal{X}.$$

For each RU Regression method, we plot the function $\hat{\alpha}_\Gamma(x)$ learned by the neural network. In addition, with access to the data generating process, we can explicitly compute the function $q_{\eta(\Gamma)}^L(X; \hat{h}_\Gamma(X))$. In the right plot of Figure 2, we observe that $\hat{\alpha}_\Gamma$ closely matches $q_{\eta(\Gamma)}^L(X; \hat{h}_\Gamma(X))$ across the possible values of X .

4.1.3 High-Dimensional Experiment

Data Generation. We generate a synthetic dataset of samples of the form (X_i, Y_i, U_i) , where $X_i \in \mathbb{R}^d$ represents observed covariates, $Y_i \in \mathbb{R}$ represents the outcome, and $U_i \in \{0, 1\}$ represents an unobserved variable that influences the outcome Y_i . Since we aim to consider a high-dimensional example, we set $d = 18$.

Method	Test MSE				
	$p = 0.1$	$p = 0.2$	$p = 0.5$	$p = 0.7$	$p = 0.9$
Standard ERM	6.939 ± 0.174	10.480 ± 0.126	20.866 ± 0.304	27.880 ± 0.484	34.913 ± 0.668
RU Regression ($\Gamma = 2$)	10.074 ± 0.247	11.846 ± 0.138	17.029 ± 0.236	20.522 ± 0.308	24.046 ± 0.540
RU Regression ($\Gamma = 4$)	12.456 ± 0.431	13.388 ± 0.309	16.093 ± 0.179	17.898 ± 0.300	19.750 ± 0.584
RU Regression ($\Gamma = 8$)	13.419 ± 0.306	14.057 ± 0.255	15.895 ± 0.133	17.119 ± 0.143	18.388 ± 0.304
RU Regression ($\Gamma = 16$)	13.613 ± 0.400	14.197 ± 0.308	15.873 ± 0.140	16.983 ± 0.262	18.142 ± 0.480
Oracle ERM	6.306 ± 0.187	10.480 ± 0.126	15.743 ± 0.152	13.341 ± 0.123	6.274 ± 0.176

Table 1: Results from the one-dimensional simulation experiment. We report the mean and standard deviation of the test MSE from 6 random trials, where the randomness is over the dataset generation. Standard ERM incurs high test MSE for high values of p . RU Regression is more robust to deviations from the training distribution than Standard ERM. RU Regression matches the performance of Oracle ERM at $p = 0.5$.

Method	Test MSE				
	$p = 0.1$	$p = 0.2$	$p = 0.5$	$p = 0.7$	$p = 0.9$
Standard ERM	0.028 ± 0.000	0.043 ± 0.000	0.088 ± 0.002	0.118 ± 0.002	0.148 ± 0.003
RU Regression ($\Gamma = 2$)	0.041 ± 0.002	0.049 ± 0.001	0.071 ± 0.001	0.086 ± 0.003	0.100 ± 0.004
RU Regression ($\Gamma = 4$)	0.054 ± 0.008	0.057 ± 0.006	0.067 ± 0.002	0.073 ± 0.006	0.080 ± 0.011
RU Regression ($\Gamma = 8$)	0.056 ± 0.003	0.058 ± 0.002	0.066 ± 0.001	0.071 ± 0.002	0.076 ± 0.004
RU Regression ($\Gamma = 16$)	0.057 ± 0.003	0.059 ± 0.002	0.066 ± 0.000	0.070 ± 0.002	0.074 ± 0.003
Oracle ERM	0.025 ± 0.000	0.043 ± 0.000	0.066 ± 0.000	0.056 ± 0.000	0.025 ± 0.000

Table 2: Results from the high-dimensional ($d = 16$) simulation experiment. We report the mean and standard deviation of the test MSE from 6 random trials, where the randomness is over the dataset generation. Standard ERM incurs high test MSE for high values of p . RU Regression is more robust to deviations from the training distribution than Standard ERM. RU Regression matches the performance of Oracle ERM at $p = 0.5$.

We suppose the data is distributed as follows

$$X_i \sim \text{Uniform}[0, 1]^d, \quad U_i \sim \text{Bernoulli}(p), \quad Y_i | X_i \sim N(\mathbf{a}^T X_i + 0.5 \cdot U_i, 0.1), \quad (38)$$

where $\mathbf{a} \in \mathbb{R}^d$ is a constant vector. Similar to the one-dimensional example, the outcomes Y_i can be clustered into two hyperplanes $Y_i = \mathbf{a}^T X_i + 0.5$ for samples with $U_i = 1$ and $Y = \mathbf{a}^T X_i$ for samples with $U_i = 0$. As in the one-dimensional example, we consider distribution shifts which result from varying p , the probability that $U_i = 1$. For the training distribution, we set $p = 0.2$, so examples with $U_i = 1$ occur with lower frequency in the training set. At test-time, we evaluate the learned models on data distributions where $p \in [0.1, 0.2, 0.5, 0.7, 0.9]$.

Results. The results from the high-dimensional simulation are consistent with those from the one-dimensional simulation. From Table 2, Standard ERM achieves low test MSE on test distributions that are similar to the training distribution $p \in [0.1, 0.2]$. However, the test MSE of Standard ERM increases as p increases. The RU Regression methods achieve higher test MSE on the original training distribution than Standard ERM but are more robust than Standard ERM as p increases. We note that RU Regression matches the performance of the Oracle ERM model when $p = 0.5$.

4.2 MIMIC-III Data

Accurate patient length-of-stay predictions are useful for scheduling and hospital resource management [Harutyunyan et al., 2019]. Many recent works study the problem of predicting patient length-of-stay from patient covariates [Daghistani et al., 2019, Morton et al., 2014, Sotoodeh and Ho, 2019]. In this experiment, we evaluate our approach on electronic health record data drawn from the publicly available MIMIC-III

Method	Weighted Test MSE	
	No Shift	Shift
Standard ERM	3.230 ± 0.079	6.926 ± 0.265
RU Regression ($\Gamma = 1.50$)	3.227 ± 0.074	6.663 ± 0.253
RU Regression ($\Gamma = 2.00$)	3.274 ± 0.072	6.607 ± 0.247
RU Regression ($\Gamma = 2.50$)	3.349 ± 0.070	6.313 ± 0.237
RU Regression ($\Gamma = 3.00$)	3.441 ± 0.068	6.060 ± 0.224

Table 3: Results from MIMIC-III Experiment. We report the weighted test MSE and bootstrap standard error with 5000 bootstrap samples.

dataset [Johnson et al., 2016a]. We study the robustness of regression models when the distribution of patients observed at test time differs from the distribution of patients observed at train-time.

Data. In this experiment, the observed covariates X_i consist of 17 different medical measurements of a patient recorded within the first 24 hours of hospital stay (see Appendix A.3 for details on the particular covariates). The outcome Y_i is the patient length-of-stay in the ICU in days. To simulate the biased test distributions, we compute a weight

$$w_i = \frac{Y_i}{\sum_i Y_i}$$

for each test example i , and we report the weighted test MSE.

$$\text{Weighted Test MSE} = \sum_{i=1}^{n_{\text{test}}} w_i L(h(X_i), Y_i).$$

Reporting the weighted test MSE allows us to simulate the test MSE on a biased test distribution where patients with length of stay equal to Y_i are sampled with probability w_i .

We compare the following methods.

1. Standard ERM - We fit a neural network model with the squared loss function (Equation 36) on the training data.
2. Rockafellar-Uryasev Regression (RU Regression) - We fit two neural networks with the L_{RU}^Γ loss function, where one network learns h and the other network learns α , on the training data. A figure for the training pipeline is provided in Figure 3 in Appendix A.

Results. As seen in Table 3, RU Regression trades performance in the environment with no shift for improved performance under the shifted distribution. The RU Regression performs worse than standard ERM in the environment with no shift but is more accurate than the standard ERM in the shifted environment. Thus, at a modest cost in shift-free accuracy, our method achieved considerable improvements in a new shifted environment. We emphasize that RU Regression was not given any information on how the test set might differ from the training set; we simply posited that the shift is some re-weighting of the type (6), and asked RU Regression to be robust to any such shift (up to a factor $\Gamma = 3$). We report the bootstrap standard error obtained with 5000 bootstrap samples.

References

- Isaiah Andrews and Emily Oster. A simple approximation for evaluating external validity bias. *Economics Letters*, 178:58–62, 2019.
- Susan Athey and Stefan Wager. Policy learning with observational data. *Econometrica*, 89(1):133–161, 2021.
- Orazio Attanasio, Adriana Kugler, and Costas Meghir. Subsidizing vocational training for disadvantaged youth in colombia: Evidence from a randomized trial. *American Economic Journal: Applied Economics*, 3(3):188–220, 2011.

- Aharon Ben-Tal, Dick Den Hertog, Anja De Waegenare, Bertrand Melenberg, and Gijs Rennen. Robust solutions of optimization problems affected by uncertain probabilities. *Management Science*, 59(2):341–357, 2013.
- Dimitris Bertsimas and Nathan Kallus. From predictive to prescriptive analytics. *Management Science*, 66(3):1025–1044, 2020.
- Xiaohong Chen. Large sample sieve estimation of semi-nonparametric models. *Handbook of econometrics*, 6:5549–5632, 2007.
- Xiaohong Chen and Xiaotong Shen. Sieve extremum estimates for weakly dependent data. *Econometrica*, pages 289–314, 1998.
- Tahani A Daghistani, Radwa Elshawi, Sherif Sakr, Amjad M Ahmed, Abdullah Al-Thwayee, and Mouaz H Al-Mallah. Predictors of in-hospital length of stay among cardiac patients: a machine learning approach. *International journal of cardiology*, 288:140–147, 2019.
- Jacob Dorn, Kevin Guo, and Nathan Kallus. Doubly-valid/doubly-sharp sensitivity analysis for causal inference with unmeasured confounding. *arXiv preprint arXiv:2112.11449*, 2021.
- John Duchi, Tatsunori Hashimoto, and Hongseok Namkoong. Distributionally robust losses for latent covariate mixtures. *arXiv preprint arXiv:2007.13982*, 2020.
- John C Duchi and Hongseok Namkoong. Learning models with uniform performance via distributionally robust optimization. *The Annals of Statistics*, 49(3):1378–1406, 2021.
- Adam N Elmachtoub and Paul Grigas. Smart “predict, then optimize”. *Management Science*, 68(1):9–26, 2022.
- Dylan J Foster and Vasilis Syrgkanis. Orthogonal statistical learning. *arXiv preprint arXiv:1901.09036*, 2019.
- Stuart Geman and Chii-Ruey Hwang. Nonparametric maximum likelihood estimation by the method of sieves. *The annals of Statistics*, pages 401–414, 1982.
- Ary L Goldberger, Luis AN Amaral, Leon Glass, Jeffrey M Hausdorff, Plamen Ch Ivanov, Roger G Mark, Joseph E Mietus, George B Moody, Chung-Kang Peng, and H Eugene Stanley. Physiobank, physiotoolkit, and physionet: components of a new research resource for complex physiologic signals. *Circulation*, 101(23):e215–e220, 2000.
- Hrayr Harutyunyan, Hrant Khachatrian, David C Kale, Greg Ver Steeg, and Aram Galstyan. Multitask learning and benchmarking with clinical time series data. *Scientific data*, 6(1):1–18, 2019.
- Weihua Hu, Gang Niu, Issei Sato, and Masashi Sugiyama. Does distributionally robust supervised learning give robust classifiers? In *International Conference on Machine Learning*, pages 2029–2037. PMLR, 2018.
- Ying Jin, Zhimei Ren, and Zhengyuan Zhou. Sensitivity analysis under the f -sensitivity models: Definition, estimation and inference. *arXiv preprint arXiv:2203.04373*, 2022.
- Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. Mimic-iii, a freely accessible critical care database. *Scientific data*, 3(1):1–9, 2016a.
- Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. Mimic-iii, a freely accessible critical care database. *Scientific data*, 3:160035, 2016b.
- Nathan Kallus and Angela Zhou. Minimax-optimal policy learning under unobserved confounding. *Management Science*, 67(5):2870–2890, 2021.

- Toru Kitagawa and Aleksey Tetenov. Who should be treated? empirical welfare maximization methods for treatment choice. *Econometrica*, 86(2):591–616, 2018.
- Daniel Levy, Yair Carmon, John C Duchi, and Aaron Sidford. Large-scale methods for distributionally robust optimization. *Advances in Neural Information Processing Systems*, 33:8847–8860, 2020.
- Charles F Manski. Statistical treatment rules for heterogeneous populations. *Econometrica*, 72(4):1221–1246, 2004.
- Paul Michel, Tatsunori Hashimoto, and Graham Neubig. Distributionally robust models with parametric likelihood ratios. *arXiv preprint arXiv:2204.06340*, 2022.
- Peyman Mohajerin Esfahani and Daniel Kuhn. Data-driven distributionally robust optimization using the wasserstein metric: Performance guarantees and tractable reformulations. *Mathematical Programming*, 171(1):115–166, 2018.
- April Morton, Eman Marzban, Georgios Giannoulis, Ayush Patel, Rajender Aparasu, and Ioannis A Kakadiaris. A comparison of supervised machine learning techniques for predicting short-term in-hospital length of stay among diabetic patients. In *2014 13th International Conference on Machine Learning and Applications*, pages 428–431. IEEE, 2014.
- Hongseok Namkoong and John C Duchi. Stochastic gradient methods for distributionally robust optimization with f-divergences. *Advances in neural information processing systems*, 29, 2016.
- Whitney K Newey and Daniel McFadden. Large sample estimation and hypothesis testing. *Handbook of econometrics*, 4:2111–2245, 1994.
- Xinkun Nie and Stefan Wager. Quasi-oracle estimation of heterogeneous treatment effects. *Biometrika*, 108(2):299–319, 2021.
- Xinkun Nie, Guido Imbens, and Stefan Wager. Covariate balancing sensitivity analysis for extrapolating randomized trials across locations. *arXiv preprint arXiv:2112.04723*, 2021.
- Michael Oberst, Nikolaj Thams, Jonas Peters, and David Sontag. Regularizing towards causal invariance: Linear models with proxies. In *International Conference on Machine Learning*, pages 8260–8270. PMLR, 2021.
- Yonatan Oren, Shiori Sagawa, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust language modeling. *arXiv preprint arXiv:1909.02060*, 2019.
- R Tyrrell Rockafellar and Stanislav Uryasev. Optimization of conditional value-at-risk. *Journal of risk*, 2: 21–42, 2000.
- Shiori Sagawa, Pang Wei Koh, Tatsunori B Hashimoto, and Percy Liang. Distributionally robust neural networks for group shifts: On the importance of regularization for worst-case generalization. *arXiv preprint arXiv:1911.08731*, 2019.
- Mani Sotoodeh and Joyce C Ho. Improving length of stay prediction using a hidden markov model. *AMIA Summits on Translational Science Proceedings*, 2019:425, 2019.
- Jörg Stoye. Minimax regret treatment choice with finite samples. *Journal of Econometrics*, 151(1):70–81, 2009.
- Adith Swaminathan and Thorsten Joachims. Batch learning from logged bandit feedback through counterfactual risk minimization. *The Journal of Machine Learning Research*, 16(1):1731–1755, 2015.
- Zhiqiang Tan. A distributional approach for causal inference using propensity scores. *Journal of the American Statistical Association*, 101(476):1619–1637, 2006.
- Nikolaj Thams, Michael Oberst, and David Sontag. Evaluating robustness to dataset shift via parametric robustness sets. *arXiv preprint arXiv:2205.15947*, 2022.

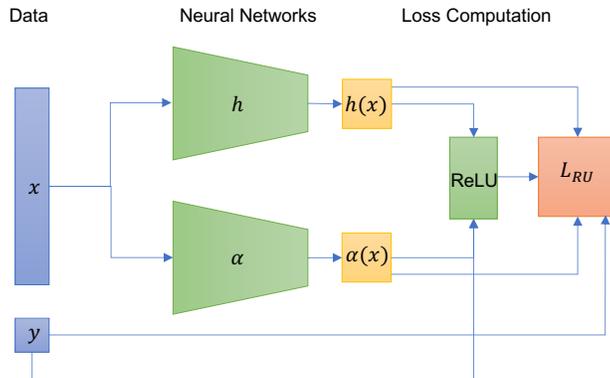


Figure 3: Model architecture for training with RU loss.

Aleksandr Filippovich Timan. *Theory of approximation of functions of a real variable*. Elsevier, 2014.

Sara A Van de Geer and Sara van de Geer. *Empirical Processes in M-estimation*, volume 6. Cambridge university press, 2000.

Sheng-Min Wang, Changsu Han, Soo-Jung Lee, Tae-Youn Jun, Ashwin A Patkar, Prakash S Masand, and Chi-Un Pae. Efficacy of antidepressants: bias in randomized clinical trials and related issues. *Expert Review of Clinical Pharmacology*, 11(1):15–25, 2018.

Shirly Wang, Matthew BA McDermott, Geeticka Chauhan, Marzyeh Ghassemi, Michael C Hughes, and Tristan Naumann. Mimic-extract: A data extraction, preprocessing, and representation pipeline for mimic-iii. In *Proceedings of the ACM conference on health, inference, and learning*, pages 222–235, 2020.

Steve Yadlowsky, Hongseok Namkoong, Sanjay Basu, John Duchi, and Lu Tian. Bounds on the conditional and average treatment effect with unobserved confounding factors. *arXiv preprint arXiv:1808.09521*, 2018.

Yingqi Zhao, Donglin Zeng, A John Rush, and Michael R Kosorok. Estimating individualized treatment rules using outcome weighted learning. *Journal of the American Statistical Association*, 107(499):1106–1118, 2012.

A Experiment Details

A.1 One-Dimensional Toy Example

A.1.1 Models

For the Standard ERM and Oracle ERM models, we train a neural network with 2 hidden layers and 128 units per layer and ReLU activation to learn the regression function h . For the RU Regression model, we jointly train two neural networks to learn the regression function h and the quantile function α , respectively. A visualization of the model architecture for RU Regression is provided in Figure 3. Each of the neural networks has 2 hidden layers and 64 units per layer and ReLU activation. We note that overall the Standard ERM and Oracle ERM models have 18.8K trainable parameters, and the RU Regression model has 10.6K trainable parameters.

A.1.2 Dataset Splits

For all methods, the train, validation, and test sets consists of 7000, 1400, and 10000 samples, respectively. For Standard ERM and RU Regression, the train and validation sets are generated via the data model specified in Equation 37 with $p = 0.2$. For Oracle ERM, the train and validation set is generated with the same data model with the parameter p matching that of the test distribution. All methods are evaluated on

the same test sets, which are generated via the data model in Equation 37 with parameter p taking value in $[0.1, 0.2, 0.5, 0.7, 0.9]$. For each of 6 random seeds $[0, 1, 2, 3, 4, 5]$, a new dataset (Standard ERM/RU Regression train and validation sets, Oracle ERM train and validation sets, and test sets) is generated.

A.1.3 Training Procedure

The models are trained for a maximum of 100 epochs with batch size equal to 1750 and we use the Adam optimizer with learning rate $1e-2$. Each epoch we check the loss obtained on the validation set and select the model that minimizes the loss on the validation set.

A.2 High-Dimensional Experiment

A.2.1 Models

We use the same models as in the one-dimensional experiment. See Section A.1.1 for details.

A.2.2 Dataset Splits

For all methods, the train, validation, and test sets consists of 100000, 20000, and 20000 samples, respectively. In the data model in Equation 38, we set

$$\mathbf{a} = [0.098, 0.430, 0.206, 0.090, -0.153, 0.292, -0.125, 0.784, \\ 0.927, -0.233, 0.583, 0.0578, 0.136, 0.851, -0.858, -0.826]$$

in all experiments. For Standard ERM and RU Regression, the train and validation sets are generated via Equation 38 with $p = 0.2$. For Oracle ERM, the train and validation set is generated with the same data model with the parameter p matching that of the test distribution. All methods are evaluated on the same test sets, which are generated via the data model in Equation 38 with parameter p taking value in $[0.1, 0.2, 0.5, 0.7, 0.9]$. For each of 6 random seeds $[0, 1, 2, 3, 4, 5]$, a new dataset (Standard ERM/RU Regression train and validation sets, Oracle ERM train and validation sets, and test sets) is generated.

A.2.3 Training Procedure

The models are trained for a maximum of 50 epochs with batch size equal to 25000 and we use the Adam optimizer with learning rate $1e-2$. Each step we check the loss obtained on the validation set and select the model that minimizes the loss on the validation set.

A.3 MIMIC-III Experiment

A.3.1 Dataset

Medical Information Mart for Intensive Care III (MIMIC-III) is a freely accessible medical database of critically ill patients admitted to the intensive care unit (ICU) at Beth Israel Deaconess Medical Center (BIDMC) from 2001 to 2012 [Johnson et al., 2016b, Goldberger et al., 2000]. During that time, BIDMC switched clinical information systems from Carevue (2001-2008) to Metavision (2008-2012). To ensure data consistency, only data archived via the Metavision system was used in the dataset.

A.3.2 Feature Selection and Data Preprocessing

We select the same patient features and imputed values as in Harutyunyan et al. [2019]. A total of 17 variables were extracted from the chartevents table to include in the dataset - capillary refill rate, blood pressure (systolic, diastolic, and mean), fraction of inspired oxygen, Glasgow Coma Score (eye opening response, motor response, verbal response, and total score), serum glucose, heart rate, respiratory rate, oxygen saturation, respiratory rate, temperature, weight, and arterial pH. For each unique ICU stay, values were extracted for the first 24 hours upon admission to the ICU and averaged. Normal values were imputed for missing variables as shown in Table 4.

Variable	MIMIC-III item ids from chartevents table	Imputed value
Capillary refill rate	(223951, 224308)	0
Diastolic blood pressure	(220051, 227242, 224643, 220180, 225310)	59.0
Systolic blood pressure	(220050, 224167, 227243, 220179, 225309)	118.0
Mean blood pressure	(220052, 220181, 225312)	77.0
Fraction inspired oxygen	(223835)	0.21
GCS eye opening	(220739)	4
GCS motor response	(223901)	6
GCS verbal response	(223900)	5
GCS total	(220739 + 223901 + 223900)	15
Glucose	(228388, 225664, 220621, 226537)	128.0
Heart Rate	(220045)	86
Height	(226707, 226730)	170.0
Oxygen saturation	(220227, 220277, 228232)	98.0
Respiratory rate	(220210, 224688, 224689, 224690)	19
Temperature	(223761, 223762)	97.88
Weight	(224639, 226512, 226531)	178.6
pH	(223830)	7.4

Table 4: Variables included in dataset

Following the cohort selection procedure in Wang et al. [2020], we further restrict to patients with covariates within physiologically valid range of measurements and length-of-stay less than or equal to 10 days.

A.3.3 Training Details

Models. For the Standard ERM model, we train a neural network with 2 hidden layers and 128 units per layer and ReLU activation to learn the regression function h . For the RU Regression model, we jointly train two neural networks to learn the regression function h and the quantile function α , respectively. Each of the neural networks has 2 hidden layers and 64 units per layer and ReLU activation. A visualization of the model architecture for RU Regression is provided in Figure 3. We note that overall the Standard ERM model has 18.8K trainable parameters, and the RU Regression model has 10.6K trainable parameters.

Dataset Splits. For all methods, the train, validation, and test sets consists of 7045, 4697, and 7829 samples, respectively.

B Standard Results

Definition 3. A function h on \mathcal{X} is said to satisfy a Holder condition with exponent β if there is a positive number γ such that $|h(x) - h(x_0)| \leq \gamma|x - x_0|^\beta$ for $x_0, x \in \mathcal{X}$.

Lemma 13. *If there is a function $Q_0(\theta)$ such that (i) $Q_0(\theta)$ is uniquely maximized at θ_0 ; (ii) θ_0 is an element of the interior of a convex set Θ and $\hat{Q}_n(\theta)$ is concave; and (iii) $\hat{Q}_n(\theta) \rightarrow Q_0(\theta)$ for all $\theta \in \Theta$, then $\hat{\theta}_n$ exists with probability approaching one and $\hat{\theta}_n \xrightarrow{P} \theta_0$ (Theorem 2.7, Newey and McFadden [1994]).*

Lemma 14. *If a functional $J : V \rightarrow \mathbb{R}$ is Gâteaux differentiable J' at $u_0 \in V$ and has a relative extremum at u_0 , then $J'(u_0; v) = 0$ for all $v \in V$.*

Lemma 15. *If $\{e_i\}$ is an orthonormal basis (a maximal orthonormal sequence) in a Hilbert space H then for any element $u \in H$ the ‘Fourier-Bessel series’ converges to u :*

$$u = \sum_{i=1}^{\infty} \langle u, e_i \rangle e_i.$$

Lemma 16. Let X be a Hilbert space, and suppose $f : X \rightarrow [-\infty, \infty]$ is lower semicontinuous and convex. If C is a closed, bounded, and convex subset of X , then f achieves its minimum on C ; i.e., there is some $x_0 \in C$ with $f(x_0) = \inf_{x \in C} f(x)$.

Lemma 17. Let A be a 2×2 symmetric matrix with $\text{tr}(A) > 0$ and $\det(A) \geq 0$. Then

$$\lambda_{\min}(A) \geq \frac{\det A}{\text{tr} A}, \quad \lambda_{\max}(A) \leq \text{tr} A.$$

Proof in Appendix D.1.

Lemma 18. Let $H(h, \alpha) = G(h) + F(h, \alpha)$, where G is strongly convex and Gâteaux differentiable in h and F is jointly convex in (h, α) , strictly convex in α , and Gâteaux differentiable in (h, α) . Then H is strictly convex in (h, α) . *Proof in Appendix D.2.*

C Proofs of Main Results

C.1 Notation

We introduce notation that is used in the proofs and technical lemmas.

$$L_{\text{RU},1}^\Gamma(z, y) := \Gamma^{-1}L(z, y) \tag{39}$$

$$L_{\text{RU},2}^\Gamma(a) := (1 - \Gamma^{-1})a \tag{40}$$

$$L_{\text{RU},3}^\Gamma(z, y, a) := (\Gamma - \Gamma^{-1}) \cdot (L(z, y) - a)_+. \tag{41}$$

Define

$$R_{f,c} := \{x \in \mathcal{X} \mid f(x) < c\} \tag{42}$$

$$S_{f,c} := \{x \in \mathcal{X} \mid f(x) > c\}. \tag{43}$$

When we consider loss functions L that satisfy Assumption 2, we define

$$\ell_1(y) := \begin{cases} \ell(y) & y > 0 \\ 0 & y \leq 0 \end{cases}, \quad \ell_2(y) := \begin{cases} 0 & y > 0 \\ \ell(y) & y \leq 0 \end{cases}, \tag{44}$$

$$T_{1,x}(c) := \mathbb{E}_{P_{Y|X=x}} [\ell(Y - c) \mid X = x], \tag{45}$$

$$T_{3,x}(c, d) := \begin{cases} \mathbb{E}_{P_{Y|X=x}} [(\ell(Y - c) - d)\mathbb{I}(\ell(Y - c) > d) \mid X = x] & d > 0 \\ \mathbb{E}_{P_{Y|X=x}} [\ell(Y - c) - d \mid X = x] & d \leq 0 \end{cases}. \tag{46}$$

C.2 Technical Lemmas

We prove lemmas about the transforms $T_{1,x}(c), T_{3,x}(c, d)$. These enable us to establish more general properties of the RU loss.

Lemma 19. Under Assumption 2, $T_{1,x}(c)$ is twice-differentiable in c and

$$\mathbb{E}_P [L_{\text{RU},1}^\Gamma(h(X), Y)] = \Gamma^{-1} \mathbb{E}_{P_X} [T_{1,X}(h(X))].$$

Proof in Appendix D.3.

Lemma 20. Under Assumption 2, 3, $T_{3,x}(c, d)$ is differentiable in c, d . In particular,

$$T_{3,x}^d(c, d) = \begin{cases} -\Pr(\ell(Y - c) > d \mid X = x) & d > 0 \\ -1 & d \leq 0 \end{cases}.$$

Equivalently,

$$T_{3,x}^d(c, d) = \begin{cases} -1 + P_{Y|X=x}(c + \ell_1^{-1}(d)) - P_{Y|X=x}(c + \ell_2^{-1}(d)) & d > 0 \\ -1 & d \leq 0 \end{cases}.$$

In addition, $T_{3,x}(c, d)$ is twice-differentiable in c, d when $d > 0$. The second derivatives are

$$\begin{aligned} T_{3,x}^{cc}(c, d) &= \sum_{i \in \{1,2\}} |\ell'(\ell_i^{-1}(d))| \cdot p_{Y|X=x}(c + \ell_i^{-1}(d)) + \mathbb{E}_{P_{Y|X}} [\ell''(Y - c) \mathbb{I}(\ell(Y - c) > d)], \\ T_{3,x}^{dd}(c, d) &= \sum_{i \in \{1,2\}} \frac{p_{Y|X=x}(c + \ell_i^{-1}(d))}{|\ell'(\ell_i^{-1}(d))|}, \\ T_{3,x}^{cd}(c, d) &= p_{Y|X=x}(c + \ell_1^{-1}(d)) - p_{Y|X=x}(c + \ell_2^{-1}(d)), \end{aligned}$$

where ℓ_1^{-1} is the inverse of $\ell(z)$ when $z > 0$ and ℓ_2^{-1} is the inverse of $\ell(z)$ when $z < 0$.

Also,

$$\mathbb{E}_P [L_{RU,3}^\Gamma(h(X), \alpha(X), Y)] = (\Gamma - \Gamma^{-1}) \mathbb{E}_{P_X} [T_{3,X}(h(X), \alpha(X))].$$

Proof in Appendix D.4.

Lemma 21. Under Assumption 2, 3, there are symmetric matrices $A_x(c, d), B_x(c, d)$ such that

$$A_x(c, d) \preceq \nabla^2 T_{3,x}(c, d) \preceq B_x(c, d)$$

when $d > 0$. The entries of $A_x(c, d)$ are given by

$$\begin{aligned} A_{x,11}(c, d) &= \sum_{i \in \{1,2\}} |\ell'(\ell_i^{-1}(d))| \cdot p_{Y|X=x}(c + \ell_i^{-1}(d)) + C_{L,l} \cdot \Pr(\ell(Y - c) > d \mid X = x) \\ A_{x,22}(c, d) &= \sum_{i \in \{1,2\}} \frac{p_{Y|X=x}(c + \ell_i^{-1}(d))}{|\ell'(\ell_i^{-1}(d))|}, \\ A_{x,12}(c, d) &= p_{Y|X=x}(c + \ell_1^{-1}(d)) - p_{Y|X=x}(c + \ell_2^{-1}(d)). \end{aligned}$$

The entries of $B_x(c, d)$ are given by

$$\begin{aligned} B_{x,11}(c, d) &= \sum_{i \in \{1,2\}} |\ell'(\ell_i^{-1}(d))| \cdot p_{Y|X=x}(c + \ell_i^{-1}(d)) + \mathbb{E}_{P_{Y|X=x}} [\ell''(Y - c) \mid X = x], \\ B_{x,22}(c, d) &= \sum_{i \in \{1,2\}} \frac{p_{Y|X=x}(c + \ell_i^{-1}(d))}{|\ell'(\ell_i^{-1}(d))|}, \\ B_{x,12}(c, d) &= p_{Y|X=x}(c + \ell_1^{-1}(d)) - p_{Y|X=x}(c + \ell_2^{-1}(d)). \end{aligned}$$

Proof in Appendix D.5.

We give a few additional lemmas related to the RU loss. These lemmas are used in proofs of many of the results from Section 3.1.

Lemma 22. Under Assumption 1, 2, 3, $\mathbb{E}_P [L_{RU}^\Gamma(h(X), \alpha(X), Y)]$ is Gâteaux differentiable in (h, α) on $L^2(P_X, \mathcal{X}) \times L^2(P_X, \mathcal{X})$ and twice-Gâteaux differentiable in (h, α) on \mathcal{C} , where

$$\mathcal{C} = \{(h, \alpha) \in \Theta \mid \alpha(x) > 0 \quad \forall x \in \mathcal{X}\}.$$

Proof in Appendix D.6.

Lemma 23. Under Assumption 2, $\mathbb{E}_P [L_{RU,1}^\Gamma(h(X), Y)]$ is strongly convex in h . *Proof in Appendix D.7.*

Lemma 24. Under Assumptions 1, 2, 3, $\mathbb{E}_P [L_{RU,3}^\Gamma(h(X), \alpha(X), Y)]$ is strictly convex in α on \mathcal{A} .

Proof in Appendix D.8.

C.3 Proof of Theorem 1

First, we verify the first claim of the theorem—the convexity of L_{RU}^Γ . We verify that $L_{\text{RU},3}^\Gamma$ is convex in (z, a) . We note that $L(z, y) - a$ is convex in (z, a) because L is convex in z and does not depend on a , so it is convex in (z, a) , and $-a$ is clearly convex in a and doesn't depend on z , so it is also convex in (z, a) . Thus, their sum must also be convex in (z, a) . In addition, we note the function $z \mapsto (z)_+$ is nondecreasing in z . Since the composition of a nondecreasing function a convex function is convex, we have that $(L(z, y) - a)_+$ is convex in (z, a) . Thus, the third term of the RU loss $L_{\text{RU},3}^\Gamma$ is convex in (z, a) . By an analogous argument, we can also show that $L_{\text{RU},1}^\Gamma + L_{\text{RU},2}^\Gamma$ is convex in (z, a) . Since the sum of convex functions is convex, L_{RU} is convex in (z, a) .

We proceed to the proof of the second claim of the theorem. The proof of this part relies on the following lemmas.

Lemma 25. *Solving the worst-case population risk minimization problem*

$$\operatorname{argmin}_{h \in L^2(Q_X, \mathcal{X})} \sup \left\{ \mathbb{E}_Q [L(h(X), Y)] : Q \in S_\Gamma(P, Q_X) \right\} \quad (47)$$

amounts to solving the worst-case risk minimization problem conditionally for each x

$$\operatorname{argmin}_{h(x) \in \mathbb{R}} \sup \left\{ \mathbb{E}_{Q_{Y|X}} [L(h(x), Y) \mid X = x] : Q \in S_\Gamma(P, Q_X) \right\}. \quad (48)$$

Proof in Appendix D.9.

Theorem 26 (Rockafellar and Uryasev [2000], Theorem 1). *For each x , let the loss $L(h, x)$ be a random variable with distribution on \mathbb{R} induced by $x \in \mathbb{R}^d$, which has density $p(x)$. Let*

$$F_\beta(h, \alpha) = \alpha + (1 - \beta)^{-1} \int_{x \in \mathbb{R}^d} (L(h, x) - \alpha)_+ p(x) dx.$$

As a function of α , $F_\beta(h, \alpha)$ is convex and continuously differentiable. The β -CVaR of the loss associated with any $h \in H$ can be determined from the formula

$$\phi_\beta(h) = \min_{\alpha \in \mathbb{R}} F_\beta(h, \alpha). \quad (49)$$

The set consisting of the values of α for which the minimum is $A_\beta(h) = \operatorname{argmin}_{\alpha \in \mathbb{R}} F_\beta(h, \alpha)$ and is a nonempty closed bounded interval (perhaps reducing to a single point).

Theorem 27 (Rockafellar and Uryasev [2000], Theorem 2). *Minimizing the β -CVaR of the loss with h over all $h \in H$ is equivalent to minimizing $F_\beta(h, \alpha)$ over all $(h, \alpha) \in H \times \mathbb{R}$, in the sense that*

$$\min_{h \in H} \phi_\beta(h) = \min_{(h, \alpha) \in H \times \mathbb{R}} F_\beta(h, \alpha),$$

where, moreover a pair (h^, α^*) achieves the second minimum iff h^* achieves the first minimum and $\alpha^* \in A_\beta(h) = \operatorname{argmin}_{\alpha \in \mathbb{R}} F_\beta(h, \alpha)$. The minimization over $(h, \alpha) \in H \times \mathbb{R}$ produces a pair (h^*, α^*) , not necessarily unique such that h^* minimizes the β -CVaR and α^* gives the corresponding β -VaR.*

First, we apply Lemma 25 to see that (7) is equivalent to minimizing the worst case loss conditionally for every $x \in \mathcal{X}$,

$$\min_{h(x) \in \mathbb{R}} \sup \left\{ \mathbb{E}_{Q_{Y|X}} [L(h(x), Y) \mid X = x] : Q \in S_\Gamma(P, Q_X) \right\}. \quad (50)$$

By the Neyman-Pearson lemma, we can verify for any decision rule h ,

$$\begin{aligned} & \sup \{ \mathbb{E}_{Q_{Y|X}} [L(h(X), Y) \mid X = x] : Q \in S_\Gamma(P, Q_X) \} \\ &= \mathbb{E}_{P_{Y|X}} \left[L(h(X), Y) \left(\Gamma^{-1} + (\Gamma - \Gamma^{-1}) \mathbb{I}(L(h(X), Y) \geq q_{\eta(\Gamma)}^L(X; h(X))) \right) \mid X = x \right], \end{aligned} \quad (51)$$

where $q_\eta^L(x; h(x))$ is as defined in (10) and $\eta(\Gamma)$ is as defined in (11). So, the problem in (50) can be rewritten as

$$\min_{h(x) \in \mathbb{R}} \mathbb{E}_{P_{Y|X}} \left[L(h(x), Y) \left(\Gamma^{-1} + (\Gamma - \Gamma^{-1}) \mathbb{I}(L(h(x), Y) \geq q_\eta^L(X; h(x))) \right) \mid X = x \right]. \quad (52)$$

Thus, we can focus on the optimization problem in (52). We realize that this optimization problem depends $\mathbb{E}_{P_{Y|X}} [L(h(x), Y) \mathbb{I}(L(h(x), Y) \geq q_\eta^L(X; h(x))) \mid X = x]$, which can be rewritten as follows

$$\begin{aligned} & \mathbb{E}_{P_{Y|X}} [L(h(x), Y) \mathbb{I}(L(h(x), Y) \geq q_\eta^L(X; h(x))) \mid X = x] \\ &= \text{CVaR}_{\eta(\Gamma)}(L(h(x), Y)) \cdot \Pr(L(h(x), Y) \geq q_\eta^L(X; h(x)) \mid X = x). \end{aligned}$$

By the definition of $q_\eta^L(X; h(x))$, we have that $\Pr(L(h(x), Y) \geq q_\eta^L(X; h(x)) \mid X = x) = 1 - \eta(\Gamma)$. Thus, we have that

$$\mathbb{E}_{P_{Y|X}} [L(h(X), Y) \mathbb{I}(L(h(X), Y) > q_\eta^L(X; h(X))) \mid X = x] = (1 - \eta(\Gamma)) \cdot \text{CVaR}_{\eta(\Gamma)}(L(h(x), Y)). \quad (53)$$

Now, the problem in (52) can be written as

$$\min_{h(x) \in \mathbb{R}} \mathbb{E}_{P_{Y|X}} \left[L(h(X), Y) \left(\Gamma^{-1} + (\Gamma - \Gamma^{-1}) \mathbb{I}(L(h(x), Y) \geq q_\eta^L(X; h(x))) \right) \mid X = x \right] \quad (54)$$

$$= \min_{h(x) \in \mathbb{R}} \Gamma^{-1} \cdot \mathbb{E}_{P_{Y|X}} [L(h(X), Y) \mid X = x] + (\Gamma - \Gamma^{-1}) \cdot (1 - \eta(\Gamma)) \text{CVaR}_{\eta(\Gamma)}(L(h(x), Y)) \quad (55)$$

$$= \min_{h(x) \in \mathbb{R}} \Gamma^{-1} \cdot \mathbb{E}_{P_{Y|X}} [L(h(X), Y) \mid X = x] + (1 - \Gamma^{-1}) \cdot \text{CVaR}_{\eta(\Gamma)}(L(h(x), Y)). \quad (56)$$

By Theorem 26, we have that

$$\text{CVaR}_{\eta(\Gamma)}(L(h(x), Y)) = \min_{\alpha(x) \in \mathbb{R}} \alpha(x) + \frac{1}{1 - \eta(\Gamma)} \mathbb{E}_{P_{Y|X}} [(L(h(x), Y) - \alpha(x))_+ \mid X = x]. \quad (57)$$

By Theorem 27, we can use (57) to write (56) as a joint optimization problem over both $h(x)$ and $\alpha(x)$ as follows

$$\begin{aligned} & \min_{h(x) \in \mathbb{R}} \Gamma^{-1} \cdot \mathbb{E}_{P_{Y|X}} [L(h(X), Y) \mid X = x] + (1 - \Gamma^{-1}) \cdot \text{CVaR}_{\eta(\Gamma)}(L(h(x), Y)) \\ &= \min_{h(x), \alpha(x) \in \mathbb{R}} \Gamma^{-1} \cdot \mathbb{E}_{P_{Y|X}} [L(h(x), Y) \mid X = x] \\ & \quad + (1 - \Gamma^{-1}) \cdot \left(\alpha(x) + \frac{1}{1 - \eta(\Gamma)} \mathbb{E}_{P_{Y|X}} [(L(h(x), Y) - \alpha(x))_+ \mid X = x] \right) \\ &= \min_{h(x), \alpha(x) \in \mathbb{R}} \Gamma^{-1} \cdot \mathbb{E}_{P_{Y|X}} [L(h(x), Y) \mid X = x] + (1 - \Gamma^{-1}) \alpha(x) \\ & \quad + (\Gamma - \Gamma^{-1}) \mathbb{E}_{P_{Y|X}} [(L(h(x), Y) - \alpha(x))_+ \mid X = x] \\ &= \min_{h(x), \alpha(x) \in \mathbb{R}} \mathbb{E}_{P_{Y|X=x}} [L_{\text{RU}}^\Gamma(h(x), \alpha(x), Y) \mid X = x], \end{aligned}$$

where the last line follows from the definition of L_{RU}^Γ in (8). So, (56) is equivalent to the augmented conditional risk minimization

$$\min_{h(x), \alpha(x) \in \mathbb{R}} \mathbb{E}_{P_{Y|X}} [L_{\text{RU}}^\Gamma(h(x), \alpha(x), Y) \mid X = x]. \quad (58)$$

Lastly, functions $h_\Gamma^*, \alpha_\Gamma^*$ that for every $x \in \text{supp}(P_X)$ solve (18) also solve (9).

Now, we can show that h_Γ^* solves (7) for any Q_X that is absolutely continuous to P_X . Let \mathcal{T} be any set with nonzero measure with respect to Q_X . Then \mathcal{T} must also have nonzero measure with respect to P_X because $Q_X \ll P_X$. So, for any $h \in L^2(Q_X, \mathcal{X})$

$$\sup_{Q_{Y|X}: Q \in \mathcal{S}_\Gamma(P, Q_X)} \mathbb{E}_{Q_{Y|X}} [L(h_\Gamma^*(X), Y) \mid X \in \mathcal{T}] \leq \sup_{Q_{Y|X}: Q \in \mathcal{S}_\Gamma(P, Q_X)} \mathbb{E}_{Q_{Y|X}} [L(h(X), Y) \mid X \in \mathcal{T}]$$

for any set \mathcal{T} with nonzero measure with respect to Q_X . This is sufficient to show that h_Γ^* is a solution to (7) for any $Q_X \ll P_X$.

C.4 Proof of Lemma 2

Suppose for the sake of contradiction (h, α) is a minimizer of the population RU risk and $(h, \alpha) \notin \Theta$. There are three cases

1. $(h, \alpha) \in \mathcal{H}^c \times \mathcal{A}$,
2. $(h, \alpha) \in \mathcal{H} \times \mathcal{A}^c$,
3. $(h, \alpha) \in \mathcal{H}^c \times \mathcal{A}^c$.

First, we focus on the case where $(h, \alpha) \in \mathcal{H}^c \times \mathcal{A}$. We consider \bar{h} ,

$$\bar{h}(x) = \begin{cases} h(x) & h(x) \in [-2B, 2B] \\ 2B & h(x) > 2B \\ -2B & h(x) < -2B \end{cases}.$$

We note that $(\bar{h}, \alpha) \in \Theta$. We define $R_{h,-2B}$ and $S_{h,2B}$ following (42) and (43).

$$\begin{aligned} & \mathbb{E}_P [L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y)] - \mathbb{E}_P [L_{\text{RU}}^\Gamma(\bar{h}(X), \alpha(X), Y)] \\ &= \mathbb{E}_P [(L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y) - L_{\text{RU}}^\Gamma(\bar{h}(X), \alpha(X), Y))\mathbb{I}(R_{h,-2B})] \\ & \quad + \mathbb{E}_P [(L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y) - L_{\text{RU}}^\Gamma(\bar{h}(X), \alpha(X), Y))\mathbb{I}(S_{h,2B})] \end{aligned}$$

because they only differ on $R_{h,-2B}$ and $S_{h,2B}$. Analyzing the second term on the right side above, we see that

$$\begin{aligned} & \mathbb{E}_P [(L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y) - L_{\text{RU}}^\Gamma(\bar{h}(X), \alpha(X), Y)) \cdot \mathbb{I}(S_{h,2B})] \\ &= \mathbb{E}_{P_X} \left[\left(\Gamma^{-1}T_{1,X}(h(X), \alpha(X)) + (\Gamma - \Gamma^{-1}) \cdot T_{3,X}(h(X), \alpha(X)) \right) \mathbb{I}(S_{h,2B}) \right], \end{aligned}$$

where $T_{1,X}, T_{3,X}$ are defined in Lemma 19 and Lemma 20, respectively. For $x \in S_{h,2B}$,

$$\begin{aligned} & \Gamma^{-1}T_{1,x}(h(x), \alpha(x)) + (\Gamma - \Gamma^{-1}) \cdot T_{3,x}(h(x), \alpha(x)) - \Gamma^{-1}T_{1,x}(\bar{h}(x), \alpha(x)) - (\Gamma - \Gamma^{-1}) \cdot T_{3,x}(\bar{h}(x), \alpha(x)) \\ &= (h(x) - \bar{h}(x)) \cdot \left(\Gamma^{-1}T_{1,x}^c(\tilde{h}(x), \alpha(x)) + (\Gamma - \Gamma^{-1}) \cdot T_{3,X}^c(\tilde{h}(x), \alpha(x)) \right) \quad \tilde{h}(x) \in [\bar{h}(x), h(x)] \end{aligned} \quad (59a)$$

$$= (h(x) - \bar{h}(x)) \cdot \mathbb{E}_{P_{Y|X=x}} \left[\Gamma^{-1} \cdot (-\ell'(Y - \tilde{h}(x))) + (\Gamma - \Gamma^{-1}) \cdot (-\ell'(Y - \tilde{h}(x))) \cdot \mathbb{I}(\ell(Y - \tilde{h}(x)) > \alpha(x)) \right] \quad (59b)$$

$$\geq (h(x) - \bar{h}(x)) \cdot \mathbb{E}_{P_{Y|X=x}} \left[\Gamma^{-1} \cdot (-\ell'(Y - \tilde{h}(x))) \right] \quad (59c)$$

$$> 0. \quad (59d)$$

(59a) follows from the Mean Value Theorem, the differentiability of $T_{1,x}$ (Lemma 19), and the differentiability of $T_{3,x}$ (Lemma 20). (59b) follows from Lemma 19 and Lemma 20. The inequality in (59c) comes from the observation that for $x \in S_{h,2B}$, we have that $Y - \tilde{h}(x) \leq -B$ because $Y \in [-B, B]$ and $\tilde{h}(x) \in [2B, h(x)]$. So, $-\ell'(Y - \tilde{h}(x)) > 0$. Meanwhile, $h(x) - \bar{h}(x) > 0$. So, the product of $-\ell'(Y - \tilde{h}(x)) \cdot (h(x) - \bar{h}(x)) > 0$. Since $\Pr(\ell(Y - \tilde{h}(x)) > \alpha(x) | X = x) \geq 0$, (59c) holds. For the same reason, (59d) holds as well. Thus, if $S_{h,2B}$ has positive measure, then

$$\mathbb{E}_P [(L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y) - L_{\text{RU}}^\Gamma(\bar{h}(X), \alpha(X), Y))\mathbb{I}(X \in S_{h,2B})] > 0.$$

An analogous argument can be used to show that for $R_{h,-2B}$ with positive measure,

$$\mathbb{E}_P [(L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y) - L_{\text{RU}}^\Gamma(\bar{h}(X), \alpha(X), Y))\mathbb{I}(X \in R_{h,-2B})] > 0.$$

Thus, as long as $R_{h,-2B} \cup S_{h,2B}$ has positive measure, which must be the case under our assumption that the minimizer $(h, \alpha) \in \mathcal{H}^c \times \mathcal{A}$, then there is $(\bar{h}, \alpha) \in \Theta$ that achieves lower population RU risk. This is a contradiction, so the minimizer cannot be in $\mathcal{H}^c \times \mathcal{A}$.

Now, we consider the next case that the minimizer $(h, \alpha) \in \mathcal{H} \times \mathcal{A}^c$. Consider $\bar{\alpha} \in \mathcal{A}$,

$$\bar{\alpha}(x) = \begin{cases} 0 & \alpha(x) < 0 \\ \alpha(x) & 0 \leq \alpha(x) \leq M_u \\ M_u & \alpha(x) > M_u \end{cases}$$

Note that $(h, \bar{\alpha}) \in \Theta$. We define $R_{\alpha,0}$ and S_{α, M_u} according to (42) and (43), respectively. We have that

$$\begin{aligned} & \mathbb{E}_P [L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y)] - \mathbb{E}_P [L_{\text{RU}}^\Gamma(h(X), \bar{\alpha}(X), Y)] \\ &= \mathbb{E}_P [(L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y) - L_{\text{RU}}^\Gamma(h(X), \bar{\alpha}(X), Y))\mathbb{I}(R_{\alpha,0})] \\ & \quad + \mathbb{E}_P [(L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y) - L_{\text{RU}}^\Gamma(h(X), \bar{\alpha}(X), Y))\mathbb{I}(S_{\alpha, M_u})]. \end{aligned}$$

because they only differ on $R_{\alpha,0}$ and S_{α, M_u} . We find that

$$\begin{aligned} & \mathbb{E}_P [(L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y) - L_{\text{RU}}^\Gamma(h(X), \bar{\alpha}(X), Y)) \cdot \mathbb{I}(R_{\alpha,0})] \\ &= (1 - \Gamma^{-1})\mathbb{E}_P [(\alpha(X) - \bar{\alpha}(X))\mathbb{I}(R_{\alpha,0})] + (\Gamma - \Gamma^{-1})\mathbb{E}_P [(L(h(X), Y) - \alpha(X))_+ \cdot \mathbb{I}(R_{\alpha,0})] \\ & \quad - (\Gamma - \Gamma^{-1})\mathbb{E}_P [(L(h(X), Y) - \bar{\alpha}(X))_+ \cdot \mathbb{I}(R_{\alpha,0})] \\ &= (1 - \Gamma^{-1})\mathbb{E}_X [\alpha(X)\mathbb{I}(R_{\alpha,0})] + (\Gamma - \Gamma^{-1})\mathbb{E}_P [(L(h(X), Y) - \alpha(X))\mathbb{I}(R_{\alpha,0})] \\ & \quad - (\Gamma - \Gamma^{-1})\mathbb{E}_P [L(h(X), Y)\mathbb{I}(R_{\alpha,0})] \\ &= (1 - \Gamma)\mathbb{E}_P [\alpha(X) \cdot \mathbb{I}(R_{\alpha,0})]. \end{aligned}$$

If $R_{\alpha,0}$ has positive measure, then

$$\mathbb{E}_P [(L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y) - L_{\text{RU}}^\Gamma(h(X), \bar{\alpha}(X), Y)) \cdot \mathbb{I}(R_{\alpha,0})] > 0$$

because on $R_{\alpha,0}$, we have that $\alpha(X) < 0$ and also $(1 - \Gamma) < 0$. In addition,

$$\begin{aligned} & \mathbb{E}_P [(L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y) - L_{\text{RU}}^\Gamma(h(X), \bar{\alpha}(X), Y)) \cdot \mathbb{I}(S_{\alpha, M_u})] \\ &= \mathbb{E}_{P_X} [\mathbb{E}_{P_{Y|X}} [L_{\text{RU},2}^\Gamma(\alpha(X)) - L_{\text{RU},2}^\Gamma(\bar{\alpha}(X)) + L_{\text{RU},3}^\Gamma(h(X), \alpha(X), Y) - L_{\text{RU},3}^\Gamma(h(X), \bar{\alpha}(X), Y) \mid X] \mathbb{I}(S_{\alpha, M_u})]. \end{aligned}$$

For $x \in S_{\alpha, M_u}$, we compute

$$\mathbb{E}_{P_{Y|X}} [L_{\text{RU},2}^\Gamma(\alpha(X)) - L_{\text{RU},2}^\Gamma(\bar{\alpha}(X)) + L_{\text{RU},3}^\Gamma(h(X), \alpha(X), Y) - L_{\text{RU},3}^\Gamma(h(X), \bar{\alpha}(X), Y) \mid X = x] \quad (60a)$$

$$= \mathbb{E}_{P_{Y|X=x}} [(1 - \Gamma^{-1})(\alpha(X) - \bar{\alpha}(X)) \mid X = x] \quad (60b)$$

$$+ \mathbb{E}_{P_{Y|X=x}} \left[(\Gamma - \Gamma^{-1}) \left(T_{3,X}(h(X), \alpha(X)) - T_{3,X}(h(X), \bar{\alpha}(X)) \right) \mid X = x \right] \quad (60c)$$

$$= (1 - \Gamma^{-1})(\alpha(x) - \bar{\alpha}(x)) + (\Gamma - \Gamma^{-1}) \left(T_{3,x}(h(x), \alpha(x)) - T_{3,x}(h(x), \bar{\alpha}(x)) \right) \quad (60d)$$

$$= (1 - \Gamma^{-1})(\alpha(x) - \bar{\alpha}(x)) + (\Gamma - \Gamma^{-1}) \cdot (\alpha(x) - \bar{\alpha}(x)) \cdot T_{3,x}^d(h(x), \tilde{\alpha}(x)) \quad \tilde{\alpha}(x) \in [\bar{\alpha}(x), \alpha(x)] \quad (60e)$$

$$= (1 - \Gamma^{-1})(\alpha(x) - \bar{\alpha}(x)) + (\Gamma - \Gamma^{-1}) \cdot (\alpha(x) - \bar{\alpha}(x)) \cdot (-1 + F_{x;h(x)}(\tilde{\alpha}(x))) \quad (60f)$$

$$> (1 - \Gamma^{-1})(\alpha(x) - \bar{\alpha}(x)) + (\Gamma - \Gamma^{-1}) \cdot (\alpha(x) - \bar{\alpha}(x)) \cdot (-1 + \eta(\Gamma)) \quad (60g)$$

$$= 0. \quad (60h)$$

In the above derivation, we have that (60d) follows from Lemma 20 and Assumption 2. Next, we apply the Mean Value Theorem to $T_{3,x}(c, d)$ to arrive at (60e). After that, we use the definition of $T_{3,x}^d(c, d)$ for $d > 0$ from Lemma 20, where $\tilde{\alpha}(x) > 0$. Finally, we recall that $F_{x;h(x)}$ is the distribution over $L(h(x), Y) = \ell(Y - h(x))$ when Y is distributed according to $P_{Y|X=x}$. We can show (60g) as follows. Since $\tilde{\alpha}(x) \in [\bar{\alpha}(x), \alpha(x)]$ and $x \in S_{\alpha, M_u}$, we have that

$$F_{x;h(x)}(\tilde{\alpha}(x)) \geq F_{x;h(x)}(\bar{\alpha}(x)) = F_{x;h(x)}(M_u),$$

and we have that

$$q_{\eta(\Gamma)}^L(x; h(x)) = F_{x;h(x)}^{-1}(\eta(\Gamma)) < M_u$$

by the definition of M_u (24). So, we see that $\eta(\Gamma) < F_{x;h(x)}(M_u)$. In addition, we note that $\alpha(x) - \bar{\alpha}(x) > 0$ for $x \in S_{\alpha, M_u}$ and $\Gamma - \Gamma^{-1} > 0$. We conclude that if S_{α, M_u} has positive measure, then

$$\mathbb{E}_P \left[(L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y) - L_{\text{RU}}^\Gamma(h(X), \bar{\alpha}(X), Y)) \cdot \mathbb{I}(S_{\alpha, M_u}) \right] > 0.$$

Thus, as long as $R_{\alpha, 0} \cup S_{\alpha, M_u}$ has positive measure, which must be the case because we assumed that $(h, \alpha) \in \mathcal{H} \times \mathcal{A}^c$, there is $(h, \bar{\alpha}) \in \Theta$ that achieves lower population RU risk than the minimizer (h, α) . This is a contradiction, so any minimizer cannot be in $\mathcal{H} \times \mathcal{A}^c$.

Combining the two previous arguments, we can show that any minimizer also cannot be in $\mathcal{H}^c \times \mathcal{A}^c$. Thus, any minimizer of the population RU risk must lie in Θ .

C.5 Proof of Lemma 3

The main goal of this proof is to apply Lemma 16 to the function $\mathbb{E}_P [L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y)]$ and set Θ . Clearly, the population RU risk is continuous. We have the RU loss is convex from the first part of Theorem 1, so the population RU risk is also convex in (h, α) . In addition, $\Theta \subset L^2(P_X, \mathcal{X}) \times L^2(P_X, \mathcal{X})$, which is a Hilbert space. In addition, since L^∞ balls are closed in $L^2(P_X, \mathcal{X})$, and Θ consists of a product of L^∞ balls (one of which is not centered at 0), so Θ is closed in $L^2(P_X, \mathcal{X})$. Also, Θ is convex and bounded. Thus Lemma 16 holds, so $\mathbb{E}_P [L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y)]$ must achieve a minimum on Θ .

C.6 Proof of Lemma 4

Let

$$\begin{aligned} F(h, \alpha) &= \mathbb{E}_P [L_{\text{RU},3}^\Gamma(h(X), \alpha(X), Y)] \\ G(h) &= \mathbb{E}_P [L_{\text{RU},1}^\Gamma(h(X), Y)] \\ H(h, \alpha) &= \mathbb{E}_P [L_{\text{RU},1}^\Gamma(h(X), Y)] + \mathbb{E}_P [L_{\text{RU},3}^\Gamma(h(X), \alpha(X), Y)]. \end{aligned}$$

Note that

$$\mathbb{E}_P [L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y)] = H(h, \alpha) + \mathbb{E}_P [L_{\text{RU},2}^\Gamma(\alpha(X))]. \quad (61)$$

Since the population RU risk is the sum of H and a function that is convex in (h, α) , then it suffices to show that H is strictly convex. The main goal of this proof is to show that the conditions of Lemma 18 hold so that we can conclude that H , as defined above, is strictly convex in (h, α) .

First, we note that F, G, H are all Gâteaux differentiable by Lemma 22.

Second, we show that G satisfies the conditions of Lemma 18. By Lemma 23, G is strongly convex with constant $\Gamma^{-1}C_{L,l}$.

Third, we show that F satisfies the conditions of Lemma 18. It follows from the first part of Theorem 1 that F is jointly convex in (h, α) . Also, F is strictly convex in α on \mathcal{A} by Lemma 24.

As a result, F, G satisfy the conditions of Lemma 18. So, we have that $H(h, \alpha)$ is strictly convex in (h, α) . Furthermore, because $\mathbb{E}_P [L_{\text{RU},2}^\Gamma(\alpha)]$ is convex in α and does not depend on h , it is also jointly convex in (h, α) . Due to the decomposition in (61), $\mathbb{E}_P [L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y)]$ is the sum of a strictly convex function and a convex function in (h, α) , and is thus strictly convex.

C.7 Proof of Theorem 5

First, by Lemma 2 we have that

$$\underset{(h, \alpha) \in L^2(P_X, \mathcal{X}) \times L^2(P_X, \mathcal{X})}{\operatorname{argmin}} \mathbb{E}_P [L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y)] = \underset{(h, \alpha) \in \Theta}{\operatorname{argmin}} \mathbb{E}_P [L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y)].$$

Second, we can show that $\mathbb{E}_P [L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y)]$ has a unique minimizer on Θ . By Lemma 4, we have that $\mathbb{E}_P [L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y)]$ is strictly convex on Θ , so it has at most one minimizer on the convex set Θ . From Lemma 3, we have that $\mathbb{E}_P [L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y)]$ has at least one minimizer on Θ . Thus, there is a unique minimizer $(h_\Gamma^*, \alpha_\Gamma^*)$ on Θ . Finally, $(h_\Gamma^*, \alpha_\Gamma^*)$ is also the unique minimizer over $L^2(P_X, \mathcal{X}) \times L^2(P_X, \mathcal{X})$.

C.8 Proof of Lemma 6

Let $L(h, \alpha) = \mathbb{E}_P [L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y)]$ as the population RU risk. Since $L(h, \alpha)$ is Gâteaux differentiable (Lemma 22) and has a unique minimizer at $(h_\Gamma^*, \alpha_\Gamma^*)$ (Theorem 5), we can use Lemma 14 to realize that the Gâteaux derivative in the direction ϕ is equal to 0 for all $\phi \in L^2(P_X, \mathcal{X})$, i.e.

$$L'_\alpha(h_\Gamma^*, \alpha_\Gamma^*, \phi) = 0, \quad \forall \phi \in L^2(P_X, \mathcal{X}).$$

Recall that from Lemma 22, we have that

$$L'_\alpha(h, \alpha; \phi) = (1 - \Gamma^{-1}) \mathbb{E}_{P_X} [\phi(X)] + (\Gamma - \Gamma^{-1}) \cdot \mathbb{E}_{P_X} [T_{3,X}^d(h_\Gamma^*(X), \alpha_\Gamma^*(X)) \phi(X)].$$

So, at $(h_\Gamma^*, \alpha_\Gamma^*)$, we have that

$$\mathbb{E}_{P_X} \left[\phi(X) \cdot \left(\frac{1 - \Gamma^{-1}}{\Gamma - \Gamma^{-1}} + T_{3,X}^d(h_\Gamma^*(X), \alpha_\Gamma^*(X)) \right) \right] = 0, \quad \forall \phi \in L^2(P_X, \mathcal{X}).$$

We note that by Lemma 20,

$$T_{3,x}^d(h(x), \alpha(x)) = -1 + F_{x;h(x)}(\alpha(x)),$$

where $F_{x;h(x)}$ is the distribution over $L(h(x), Y)$ where Y is distributed according to $P_{Y|X=x}$. So, we have that

$$\mathbb{E}_{P_X} [\phi(X) \cdot (-\eta(\Gamma) + F_{X,h_\Gamma^*(X)}(\alpha_\Gamma^*(X)))] = 0, \quad \forall \phi \in L^2(P_X, \mathcal{X}).$$

So, $-\eta(\Gamma) + F_{x,h_\Gamma^*(x)}(\alpha_\Gamma^*(x))$ must be equal to 0 almost everywhere for the above equation to hold for all ϕ . Therefore, we conclude that

$$\alpha_\Gamma^*(x) = F_{x;h_\Gamma^*(x)}^{-1}(\eta(\Gamma)) = q_{\eta(\Gamma)}^L(x; h_\Gamma^*(x)).$$

Now, with this definition of α_Γ^* , we can show that there exists $M_l > 0$ such that $\alpha_\Gamma^*(x) > M_l$ for all $x \in \mathcal{X}$. We aim to show that $\inf_{x \in \mathcal{X}, h \in \mathcal{H}} q_{\eta(\Gamma)}^L(x; h(x)) > 0$. For convenience, we define

$$m(x) := q_{\frac{\eta}{2}}^L(x; h(x)).$$

We note that $\eta(\Gamma) > \frac{1}{2}$. So,

$$q_{\eta(\Gamma)}^L(x; h(x)) \geq m(x).$$

We have that for any $x \in \mathcal{X}, h \in \mathcal{H}$,

$$\Pr(L(h(X), Y) \leq m(X) \mid X = x) = \frac{1}{2}.$$

Recall that under Assumption 2, $L(h(x), y) = \ell(y - h(x))$. We can apply Assumption 2 to see that

$$\Pr(Y \in [h(x) + \ell_2^{-1}(m(x)), h(x) + \ell_1^{-1}(m(x))] \mid X = x) = \frac{1}{2}.$$

Now, we can use the upper bound on the density of $p_{Y|X=x}$ from Assumption 3 to see that

$$C_{p,u} \cdot (\ell_1^{-1}(m(x)) - \ell_2^{-1}(m(x))) \geq \frac{1}{2}.$$

Rearranging, we have that

$$\ell_1^{-1}(m(x)) - \ell_2^{-1}(m(x)) \geq \frac{1}{2C_{p,u}}.$$

So,

$$\max\{\ell_1^{-1}(m(x)), -\ell_2^{-1}(m(x))\} \geq \frac{1}{4C_{p,u}}.$$

Applying ℓ to both sides, we conclude that

$$m(x) \geq \ell \left(\frac{1}{4C_{p,u}} \right).$$

Since $\frac{1}{4C_{p,u}} > 0$, we have that $m(x)$ is lower bounded by a positive constant for any choice of $h \in \mathcal{H}, x \in \mathcal{X}$. Thus, we have that

$$\alpha_{\Gamma}^*(x) = q_{\eta(\Gamma)}^L(x; h(x)) \geq \inf_{x \in \mathcal{X}, h \in \mathcal{H}} q_{\frac{1}{2}}^L(x; h(x)) \geq \ell \left(\frac{1}{4C_{p,u}} \right).$$

So, let $M_l = \ell(\frac{1}{4C_{p,u}})/2$. Then $\alpha^*(x) > M_l$ for all $x \in \mathcal{X}$.

C.9 Proof of Theorem 7

The constant for strong convexity depends on lower bounds on the conditional density $p_{Y|X=x}(\cdot)$ and on $|\ell'(\ell^{-1}(\cdot))|$ when these functions are evaluated over a particular region. To ensure that they can be lower bounded, we pick the radius of the $\|\cdot\|_{\infty}$ ball about the minimizer.

Define

$$g_i(x; h, \alpha) = h(x) + \ell_i^{-1}(\alpha(x)). \quad (62)$$

For $0 < \epsilon < \frac{1-\eta(\Gamma)}{2C_{p,u}}$, we pick $0 < \delta(\epsilon) < M_l$ to ensure that for $(h, \alpha) \in \mathcal{C}_{\delta(\epsilon)}$, we have that

$$\sup_{x \in \mathcal{X}, i \in \{1,2\}} |g_i(x; h, \alpha) - g_i(x; h_{\Gamma}^*, \alpha_{\Gamma}^*)| < \epsilon.$$

By Lemma 6, we have that $\alpha_{\Gamma}^*(x) > M_l$ for all $x \in \mathcal{X}$. We consider $(h, \alpha) \in \mathcal{C}_{\delta(\epsilon)}$. For such α , we have that $\|\alpha - \alpha_{\Gamma}^*\|_{\infty} \leq \delta(\epsilon)$, and so $\alpha(x) \geq M_l - \delta(\epsilon)$ for all $x \in \mathcal{X}$. Since $\delta(\epsilon) < M_l$, for $(h, \alpha) \in \mathcal{C}_{\delta(\epsilon)}$, we have that $\alpha(x) > 0$. Since the RU loss is twice-differentiable when $\alpha(x) > 0$ (Lemma 22), we have that it is twice-differentiable on $\mathcal{C}_{\delta(\epsilon)}$.

Let $L(h, \alpha), L_1(h, \alpha), L_3(h, \alpha)$ be shorthand for the population RU risk, the first term of the population RU risk, and the third term of the population RU risk, respectively.

$$\begin{aligned} L(h, \alpha) &= \mathbb{E}_P [L_{\text{RU}}^{\Gamma}(h(X), \alpha(X), Y)], \\ L_1(h, \alpha) &= \mathbb{E}_P [L_{\text{RU},1}^{\Gamma}(h(X), Y)], \\ L_3(h, \alpha) &= \mathbb{E}_P [L_{\text{RU},3}^{\Gamma}(h(X), \alpha(X), Y)]. \end{aligned}$$

We compute the second Gâteaux derivative of the population RU risk.

$$\langle L''(h, \alpha; \psi, \phi), (\psi, \phi) \rangle \quad (63a)$$

$$= \langle L_1''(h, \alpha; \psi, \phi) + L_3''(h, \alpha; \psi, \phi), (\psi, \phi) \rangle \quad (63b)$$

$$\geq \langle L_3''(h, \alpha; \psi, \phi), (\psi, \phi) \rangle \quad (63c)$$

$$= (\Gamma - \Gamma^{-1}) \cdot \mathbb{E}_{P_X} \left[\begin{bmatrix} \psi(X) & \phi(X) \end{bmatrix} \nabla^2 T_{3,X}(h(X), \alpha(X)) \begin{bmatrix} \psi(X) \\ \phi(X) \end{bmatrix} \right] \quad (63d)$$

$$\geq (\Gamma - \Gamma^{-1}) \mathbb{E}_{P_X} \left[\begin{bmatrix} \psi(X) & \phi(X) \end{bmatrix} A_X(h(X), \alpha(X)) \begin{bmatrix} \psi(X) \\ \phi(X) \end{bmatrix} \right] \quad (63e)$$

(63b) follows from Lemma 22. (63c) holds because $\langle L_1''(h, \alpha; \psi, \phi), (\psi, \phi) \rangle \geq 0$ because $L_1(h, \alpha)$ is strongly convex in h (Lemma 23) and does not depend on α . Next, (63d) follows from Lemma 20. Finally, $A_x(c, d)$ is the lower bound on the Hessian matrix of $T_{3,x}(c, d)$ defined in Lemma 21.

To develop a lower bound for $\langle L''(h, \alpha; \psi, \phi), (\psi, \phi) \rangle$, we aim to apply Lemma 17 to $A_x(h(x), \alpha(x))$. Before we verify the conditions of Lemma 17, we introduce the following notation

$$\begin{aligned} a_{i,x} &:= |\ell'(\ell_i^{-1}(\alpha(x)))| \quad i = 1, 2, \\ f_{i,x} &:= p_{Y|X=x}(h(x) + \ell_i^{-1}(\alpha(x))) \quad i = 1, 2, \end{aligned}$$

and we develop upper and lower bounds on $a_{i,x}$ for $i \in \{1, 2\}$, $\sum_{i \in \{1,2\}} f_{i,x}$, and $1 - F_{x;h(x)}(\alpha(x))$.

First, we focus on $a_{i,x}$. By the definition of Θ , we have that $\alpha(x) \leq M_u$. Since $|\ell'(\ell_i^{-1}(y))|$ is strictly increasing in y and on $\mathcal{C}_{\delta(\epsilon)}$, $\alpha(x) \geq M_l - \delta(\epsilon)$ for all $x \in \mathcal{X}$, we can recall the definition of $C_{a,l,\delta}, C_{a,u}$ from (29), (25) to see that

$$0 < C_{a,l,\delta(\epsilon)} \leq a_{i,x} \leq C_{a,u} < \infty \quad i = 1, 2, x \in \mathcal{X}.$$

Second, we aim to show that $\sum_{i \in \{1,2\}} f_{i,x}$ is similarly upper and lower bounded. The upper bound is straightforward from Assumption 3. To obtain the lower bound, we first analyze $\sum_{i \in \{1,2\}} p_{Y|X=x}(h_\Gamma^*(x) + \ell_i^{-1}(\alpha_\Gamma^*(x)))$, which can be written as $\sum_{i \in \{1,2\}} p_{Y|X=x}(g_i(x; h_\Gamma^*, \alpha_\Gamma^*))$ using the definition of g in (62).

Let $\ell_i^{-1}(q_{\eta(\Gamma)}^L(x; h_\Gamma^*))$ corresponds to the $c_{i,x}$ -th quantile of Y , where Y is distributed following $P_{Y|X=x}$. We realize that

$$\begin{aligned} \sum_{i \in \{1,2\}} p_{Y|X=x}(g_i(x; h_\Gamma^*, \alpha_\Gamma^*)) &= \sum_{i \in \{1,2\}} p_{Y|X=x}(h_\Gamma^*(x) + \ell_i^{-1}(\alpha_\Gamma^*(x))) \\ &= \sum_{i \in \{1,2\}} p_{Y|X=x}(h_\Gamma^*(x) + \ell_i^{-1}(q_{\eta(\Gamma)}^L(x; h_\Gamma^*(x)))) \\ &= \sum_{i \in \{1,2\}} p_{Y|X=x}(h_\Gamma^*(x) + q_{c_{i,x}}^Y(x) - h_\Gamma^*(x)) \\ &= \sum_{i \in \{1,2\}} p_{Y|X=x}(q_{c_{i,x}}^Y(x)). \end{aligned}$$

Furthermore, we realize that either $c_{1,x}$ or $c_{2,x}$ lies in $[1 - \frac{\eta(\Gamma)}{2}, 1 + \frac{\eta(\Gamma)}{2}]$. First, because $q_{\eta(\Gamma)}^L(x; h_\Gamma^*(x))$ corresponds to the $\eta(\Gamma)$ -th quantile of the conditional losses, we must have that

$$c_{1,x} - c_{2,x} = \eta(\Gamma). \quad (64)$$

In addition, $c_{1,x} \leq 1$, so $c_{2,x} \leq 1 - \eta(\Gamma)$. So, $c_{2,x} \in [0, 1 - \eta(\Gamma)]$. Suppose that $c_{2,x} \in [1 - \frac{\eta(\Gamma)}{2}, 1 - \eta(\Gamma)]$, then clearly the desired claim holds. If $c_{2,x} \notin [1 - \frac{\eta(\Gamma)}{2}, 1 - \eta(\Gamma)]$, this means that $c_{2,x} \in [0, 1 - \frac{\eta(\Gamma)}{2})$. So, we must have that $c_{1,x} \in [\eta(\Gamma), 1 + \frac{\eta(\Gamma)}{2})$. Thus, we have that at least one of $c_{1,x}, c_{2,x}$ lies in the interval $[1 - \frac{\eta(\Gamma)}{2}, 1 + \frac{\eta(\Gamma)}{2}]$.

Now, we have that

$$\sum_{i \in \{1,2\}} f_{i,x} = \sum_{i \in \{1,2\}} p_{Y|X=x}(g_i(x; h, \alpha)),$$

and $\delta(\epsilon)$ was chosen so that for $(h, \alpha) \in \mathcal{C}_{\delta(\epsilon)}$

$$\sup_{x \in \mathcal{X}, i \in \{1,2\}} |g_i(x; h, \alpha) - g_i(x; h_\Gamma^*, \alpha_\Gamma^*)| = \sup_{x \in \mathcal{X}, i \in \{1,2\}} |g_i(x; h, \alpha) - p_{Y|X=x}(q_{c_{i,x}}^Y(x))| < \epsilon.$$

Thus, we realize that for $(h, \alpha) \in \mathcal{C}_{\delta(\epsilon)}$,

$$g_i(x; h, \alpha) = q_{c_{i,x}}^Y(x) + b_i(x), \quad b_i(x) \in (-\epsilon, \epsilon), i \in \{1, 2\}, x \in \mathcal{X}. \quad (65)$$

So, for $(h, \alpha) \in \mathcal{C}_{\delta(\epsilon)}$, we realize that a lower bound on $\sum_{i \in \{1,2\}} f_{i,x} = \sum_{i \in \{1,2\}} p_{Y|X=x}(g_i(x; h, \alpha))$ is given by $C_{p,l,\epsilon}$ from (30). Thus, we have that

$$0 < C_{p,l,\epsilon} \leq \sum_{i \in \{1,2\}} f_{i,x} \leq 2C_{p,u} < \infty \quad i = 1, 2, x \in \mathcal{X},$$

and clearly each $f_{i,x}$ must be nonnegative.

Third, we aim to show that $1 - F_{x;h(x)}(\alpha(x))$ is similarly upper and lower bounded on $\mathcal{C}_{\delta(\epsilon)}$. Clearly, an upper bound on this quantity is 1. To compute the lower bound, we see that for $(h, \alpha) \in \mathcal{C}_{\delta(\epsilon)}$,

$$\begin{aligned} 1 - F_{x;h(x)}(\alpha(x)) &= 1 - P_{Y|X=x}(g_1(x; h, \alpha)) + P_{Y|X=x}(g_2(x; h, \alpha)) \\ &= 1 - P_{Y|X=x}(q_{c_{1,x}}^Y(x) + b_1(x)) + P_{Y|X=x}(q_{c_{2,x}}^Y(x) + b_2(x)) \quad b_1(x), b_2(x) \in (-\epsilon, \epsilon) \\ &\geq 1 - c_{1,x} - C_{p,u} \cdot \epsilon + c_{2,x} - C_{p,u} \cdot \epsilon \\ &= 1 - \eta(\Gamma) - 2C_{p,u}\epsilon \\ &> 0. \end{aligned}$$

The first line follows from the definition of F and g_i from (62). In the second line, we apply (65). In the third line, we note that the c.d.f. of $P_{Y|X=x}$ at $q_{c_{i,x}}^Y(x) + b_i(x)$ can be closely approximated by the value of the c.d.f. at $q_{c_{i,x}}^Y(x)$. Next, we apply (64). The last line follows because $\epsilon < \frac{1-\eta(\Gamma)}{2C_{p,u}}$. Thus, we have that

$$1 - F_{x;h(x)}(\alpha(x)) \geq 1 - \eta(\Gamma) - 2C_{p,u}\epsilon > 0. \quad (66)$$

Now, we finally verify the conditions of Lemma 17. We note that $A_x(h(x), \alpha(x))$ is a symmetric matrix by definition. We realize that $\text{tr } A_x(h(x), \alpha(x)) \geq 0$ because

$$\text{tr } A_x(h(x), \alpha(x)) = A_{x,11}(h(x), \alpha(x)) + A_{x,22}(h(x), \alpha(x)) \quad (67)$$

$$= \sum_{i \in \{1,2\}} a_{i,x} \cdot f_{i,x} + \sum_{i \in \{1,2\}} \frac{f_{i,x}}{a_{i,x}} + C_{L,\ell}(1 - F_{x;h(x)}(\alpha(x))) \quad (68)$$

$$\geq C_{L,\ell}(1 - F_{x;h(x)}(\alpha(x))) \quad (69)$$

$$> 0. \quad (70)$$

(69) follows from the observation that $f_{i,x}, a_{i,x} \geq 0$. (70) follows from (66). In addition, we see that $\det A_x(h(x), \alpha(x)) \geq 0$ because

$$\det A_x(h(x), \alpha(x)) \quad (71a)$$

$$= A_{x,11}(h(x), \alpha(x)) \cdot A_{x,22}(h(x), \alpha(x)) - (A_{x,12}(h(x), \alpha(x)))^2 \quad (71b)$$

$$= \left(\sum_{i \in \{1,2\}} a_{i,x} \cdot f_{i,x} + C_{L,\ell}(1 - F_{x;h(x)}(\alpha(x))) \right) \cdot \left(\sum_{i \in \{1,2\}} \frac{f_{i,x}}{a_{i,x}} \right) - (f_{1,x} - f_{2,x})^2 \quad (71c)$$

$$= \left(\frac{a_{1,x}}{a_{2,x}} + \frac{a_{2,x}}{a_{1,x}} + 2 \right) \cdot f_{1,x} \cdot f_{2,x} + C_{L,\ell}(1 - F_{x;h(x)}(\alpha(x))) \cdot \left(\sum_{i \in \{1,2\}} \frac{f_{i,x}}{a_{i,x}} \right) \quad (71d)$$

$$\geq C_{L,\ell} \cdot \left(\sum_{i \in \{1,2\}} \frac{f_{i,x}}{a_{i,x}} \right) \cdot (1 - F_{x;h(x)}(\alpha(x))) \quad (71e)$$

$$\geq C_{L,\ell} \cdot \frac{1}{C_{a,u}} \cdot \left(\sum_{i \in \{1,2\}} f_{i,x} \right) \cdot (1 - \eta(\Gamma) - 2C_{p,u}\epsilon) \quad (71f)$$

$$\geq C_{L,\ell} \cdot \frac{1}{C_{a,u}} \cdot C_{p,l,\epsilon} \cdot (1 - \eta(\Gamma) - 2C_{p,u}\epsilon) \quad (71g)$$

$$> 0. \quad (71h)$$

Thus, we can apply Lemma 17 to $A_x(h(x), \alpha(x))$ to see that

$$\lambda_{\min}(A_x(h(x), \alpha(x))) \geq \frac{\det A_x(h(x), \alpha(x))}{\text{tr } A_x(h(x), \alpha(x))}.$$

We can combine the lower bound on $\det A$ from (71g) with the following upper bound on $\text{tr } A$ to find a lower bound on $\lambda_{\min}(A_x(h(x), \alpha(x)))$ that does not depend on the choice of $x \in \mathcal{X}$ and $(h, \alpha) \in \mathcal{C}_{\delta(\epsilon)}$.

$$\text{tr } A_x(h(x), \alpha(x)) = \sum_{i \in \{1,2\}} a_{i,x} \cdot f_{i,x} + \sum_{i \in \{1,2\}} \frac{f_{i,x}}{a_{i,x}} + C_{L,\ell}(1 - F_{x;h(x)}(\alpha(x))) \quad (72a)$$

$$\leq 2C_{p,u} \left(C_{a,u} + \frac{1}{C_{a,l,\delta(\epsilon)}} \right) + C_{L,\ell} \quad (72b)$$

$$= \frac{2C_{p,u}(C_{a,u} \cdot C_{a,l,\delta(\epsilon)} + 1) + C_{L,\ell} \cdot C_{a,l,\delta(\epsilon)}}{C_{a,l,\delta(\epsilon)}} \quad (72c)$$

Therefore, applying (72c) and (71g), we find that

$$\begin{aligned}\lambda_{\min}(A_x(h(x), \alpha(x))) &\geq \frac{\det A_x(h(x), \alpha(x))}{\text{tr } A_x(h(x), \alpha(x))} \\ &\geq \left(C_{L,l} \cdot (1 - \eta(\Gamma) - 2C_{p,u}\epsilon) \cdot \frac{1}{C_{a,u}} \cdot C_{p,l,\epsilon} \right) \cdot \left(\frac{C_{a,l,\delta(\epsilon)}}{2C_{p,u}(C_{a,u} \cdot C_{a,l,\delta(\epsilon)} + 1) + C_{L,l} \cdot C_{a,l,\delta(\epsilon)}} \right) \\ &\geq \frac{C_{L,l} \cdot (1 - \eta(\Gamma) - 2C_{p,u} \cdot \epsilon) \cdot C_{p,l,\epsilon}}{2C_{p,u} \cdot (C_{a,u} \cdot C_{a,l,\delta(\epsilon)} + 1) + C_{L,l} \cdot C_{a,l,\delta(\epsilon)}} \cdot \frac{C_{a,l,\delta(\epsilon)}}{C_{a,u}}.\end{aligned}$$

Recall the definition of $\kappa_{1,\epsilon}$ from (31). We realize that for all $x \in \mathcal{X}$, $(h, \alpha) \in \mathcal{C}_{\delta(\epsilon)}$, we have that

$$(\Gamma - \Gamma^{-1}) \cdot A_x(h(x), \alpha(x)) \succeq \kappa_{1,\epsilon} \cdot I_2.$$

Revisiting (63e), we have that

$$\langle L''(h, \alpha; (\psi, \phi)), (\psi, \phi) \rangle \geq (\Gamma - \Gamma^{-1}) \cdot \mathbb{E}_{P_X} \left[[\psi(X) \quad \phi(X)] A_X(h(X), \alpha(X)) \begin{bmatrix} \psi(X) \\ \phi(X) \end{bmatrix} \right] \quad (73)$$

$$\geq \mathbb{E}_{P_X} \left[[\psi(X) \quad \phi(X)] \kappa_{1,\epsilon} I_2 \begin{bmatrix} \psi(X) \\ \phi(X) \end{bmatrix} \right] \quad (74)$$

$$= \kappa_{1,\epsilon} \mathbb{E}_{P_X} [\psi(X)^2 + \phi(X)^2] \quad (75)$$

$$= \kappa_{1,\epsilon} \|(\psi, \phi)\|^2. \quad (76)$$

Thus, $\mathbb{E}_P [L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y)]$ is $\kappa_{1,\epsilon}$ strongly convex in (h, α) on $\mathcal{C}_{\delta(\epsilon)}$. We note that as $\epsilon \rightarrow 0$, then $\delta(\epsilon) \rightarrow 0$, as well. So, $C_{a,l,\delta(\epsilon)} \rightarrow C_{a,l}$, where $C_{a,l}$ is defined in (26) and $C_{p,l,\epsilon} \rightarrow C_{p,l}$, where $C_{p,l}$ is defined in (27), and $\epsilon \cdot C_{p,u} \rightarrow 0$. So, we have that

$$\begin{aligned}\lim_{\epsilon \rightarrow 0} (\Gamma - \Gamma^{-1}) \frac{C_{L,l} \cdot (1 - \eta(\Gamma) - 2C_{p,u} \cdot \epsilon) \cdot C_{p,l,\epsilon}}{2C_{p,u} \cdot (C_{a,u} \cdot C_{a,l,\delta(\epsilon)} + 1) + C_{L,l} \cdot C_{a,l,\delta(\epsilon)}} \cdot \frac{C_{a,l,\delta(\epsilon)}}{C_{a,u}} \\ = (\Gamma - \Gamma^{-1}) \cdot \frac{C_{L,l} \cdot (1 - \eta(\Gamma)) \cdot C_{p,l}}{2C_{p,u} \cdot (C_{a,u} \cdot C_{a,l} + 1) + C_{L,l} \cdot C_{a,l}} \cdot \frac{C_{a,l}}{C_{a,u}} \\ = (1 - \Gamma^{-1}) \cdot \frac{C_{L,l} \cdot C_{p,l}}{2C_{p,u} \cdot (C_{a,u} \cdot C_{a,l} + 1) + C_{L,l} \cdot C_{a,l}} \cdot \frac{C_{a,l}}{C_{a,u}}\end{aligned}$$

Thus, as $\epsilon \rightarrow 0$, then $\kappa_{1,\epsilon} \rightarrow \kappa_1$, where κ_1 is defined in (28).

C.10 Proof of Theorem 8

Let $L(h, \alpha)$, $L_1(h, \alpha)$, $L_3(h, \alpha)$, $a_{i,x}$, $f_{i,x}$ be defined as in the proof of Theorem 7. To show that the population RU risk is κ_2 -smooth on $\mathcal{C}_{\delta(\epsilon)}$, we show that

$$\langle L''_{h\alpha}(h, \alpha; \psi, \phi), (\psi, \phi) \rangle \leq \kappa_2 \|(\psi, \phi)\|_{L^2(P_X, \mathcal{X})}^2.$$

We have that

$$\begin{aligned}\langle L''(h, \alpha; \psi, \phi), (\psi, \phi) \rangle &= \langle L''_1(h, \alpha; \psi, \phi) + L''_3(h, \alpha; \psi, \phi), (\psi, \phi) \rangle \\ &\leq \mathbb{E}_{P_X} \left[[\psi(X) \quad \phi(X)] \cdot \left(\Gamma^{-1} \nabla^2 T_{1,X}(h(X), \alpha(X)) + (\Gamma - \Gamma^{-1}) \nabla^2 T_{3,X}(h(X), \alpha(X)) \right) \begin{bmatrix} \psi(X) \\ \phi(X) \end{bmatrix} \right] \\ &\leq \mathbb{E}_{P_X} \left[[\psi(X) \quad \phi(X)] \cdot \left(\Gamma^{-1} \nabla^2 T_{1,X}(h(X), \alpha(X)) + (\Gamma - \Gamma^{-1}) \nabla^2 B_X(h(X), \alpha(X)) \right) \begin{bmatrix} \psi(X) \\ \phi(X) \end{bmatrix} \right],\end{aligned}$$

The second line follows from Lemma 22. The matrix $B_x(h(x), \alpha(x))$ is as defined in Lemma 21. It suffices to show that there is $\kappa_{2,\epsilon}$ such that

$$\Gamma^{-1}\nabla^2 T_{1,x}(h(x), \alpha(x)) + (\Gamma - \Gamma^{-1})B_x(h(x), \alpha(x)) \preceq \kappa_{2,\epsilon} I_2 \quad \forall x \in \mathcal{X}.$$

Applying Lemma 19 and Assumption 4, we have that

$$\nabla^2 T_{1,x}(h(x), \alpha(x)) = \begin{bmatrix} \mathbb{E}_{P_{Y|X=x}} [\ell''(Y - h(x))] & 0 \\ 0 & 0 \end{bmatrix} \preceq C_{L,u} I_2. \quad (77)$$

From the proof of Theorem 7, for $0 < \epsilon < \frac{1-\eta(\Gamma)}{2C_{p,u}}$, there exists $0 < \delta(\epsilon) < M_l$ so that for $(h, \alpha) \in \mathcal{C}_{\delta(\epsilon)}$, $(\Gamma - \Gamma^{-1})\nabla^2 T_{3,x}(h(x), \alpha(x))$ is positive definite. So, on this set $\mathcal{C}_{\delta(\epsilon)}$, $B_x(h(x), \alpha(x))$ is also certainly positive definite. So, $B_x(h(x), \alpha(x))$ satisfies the conditions of Lemma 17, so we can conclude that $\lambda_{\max}(B_x(h(x), \alpha(x))) \leq \text{tr} B_x(h(x), \alpha(x))$. We can compute an upper bound on $\text{tr} B_x(h(x), \alpha(x))$. We note that for $(h, \alpha) \in \mathcal{C}_{\delta(\epsilon)}$, $\alpha(x) \geq M_l - \delta(\epsilon)$ for all $x \in \mathcal{X}$ because $\alpha^*(x) > M_l$ by Lemma 6 and $\|\alpha - \alpha^*\|_\infty < \delta(\epsilon)$. So, we have that

$$\begin{aligned} \text{tr} B_x(h(x), \alpha(x)) &= \sum_{i \in \{1,2\}} a_{i,x} \cdot f_{i,x} + \sum_{i \in \{1,2\}} \frac{f_{i,x}}{a_{i,x}} + \mathbb{E}_{P_{Y|X=x}} [\ell''(Y - h(x))] \\ &\leq 2C_{p,u} \left(C_{a,u} + \frac{1}{C_{a,l,\delta(\epsilon)}} \right) + C_{L,u}. \end{aligned}$$

We arrive at the second inequality by recalling the definition of $C_{p,u}$ from Assumption 3, $C_{a,u}$ from (25), $C_{a,l,\delta}$ from (29), and $C_{L,u}$ from Assumption 4. So, we have that

$$B_x(h(x), \alpha(x)) \preceq \left(2C_{p,u} \left(C_{a,u} + \frac{1}{C_{a,l}} \right) + C_{L,u} \right) I_2. \quad (78)$$

Combining the constants from (77) and (78), we have that for $(h, \alpha) \in \mathcal{C}_{\delta(\epsilon)}$

$$\Gamma^{-1}\nabla^2 T_{1,x}(h(x), \alpha(x)) + (\Gamma - \Gamma^{-1})B_x(h(x), \alpha(x)) \preceq \kappa_{2,\epsilon} I_2 \quad \forall x \in \mathcal{X},$$

where $\kappa_{2,\epsilon}$ is defined as in (33). Thus, we conclude that $\mathbb{E}_P [L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y)]$ is $\kappa_{2,\epsilon}$ -smooth in (h, α) on $\mathcal{C}_{\delta(\epsilon)}$. As, $\epsilon \rightarrow 0$, $\delta(\epsilon) \rightarrow 0$. So, $C_{a,l,\delta(\epsilon)} \rightarrow C_{a,l}$. This implies that $\kappa_{2,\epsilon} \rightarrow \kappa_2$ as the radius of the $\|\cdot\|_\infty$ -ball shrinks.

C.11 Proof of Lemma 9

We note that Θ_m is a convex subset of Θ . By Lemma 4, the population RU risk is strictly convex on Θ . So, it is strictly convex on Θ_m , which means that it has at most one minimizer on Θ_m . In addition, by an analogous argument as the proof of Lemma 3, the population RU risk has at least one minimizer on Θ_m . Combining these two facts, it has a unique minimizer on Θ_m called θ_m^* .

C.12 Proof of Theorem 10

In this proof, we use the following lemma.

Lemma 28. *Define $\pi_m : \Theta \rightarrow \Theta_m$ to be the projection of θ^* onto Θ_m . Under Assumptions 1, 2, 3,*

$$\|\pi_m(\theta^*) - \theta^*\|_{L^2(P_X, \mathcal{X})} \rightarrow 0.$$

Proof in Appendix D.10.

To simplify notation, let $L(\theta) = \mathbb{E}_P [L_{\text{RU}}^\Gamma(\theta(X), Y)]$. For the sake of contradiction, assume that θ_m^* does not limit to θ^* . This means that there exists $\delta_1 > 0$ such that for every $m \in \mathbb{N}$, there is $A_m \geq m$ such that

$$\|\theta_{A_m}^* - \theta^*\|_{L^2(P_X, \mathcal{X})} > \delta_1.$$

We have that $\theta^* \in \Theta_m \subset \Theta$. In addition, $\theta^* \in \Theta$ by Lemma 2. By the strict convexity of the population RU risk on Θ (Lemma 4), for some $\epsilon > 0$, we have that

$$L(\theta_{A_m}) > L(\theta^*) + \epsilon$$

because by strict convexity, $\|\theta - \theta^*\|_{L^2(P_X, \mathcal{X})} > \delta_1$ implies that $L(\theta) > L(\theta^*) + \epsilon$ for some $\epsilon > 0$.

Note that $L(\theta)$ is continuous at θ^* , so there exists $\delta_2 > 0$ such that $\|\theta - \theta^*\|_{L^2(P_X, \mathcal{X})} < \delta_2$ implies that $|L(\theta) - L(\theta^*)| < \epsilon$. Since θ^* is the unique minimizer of the population RU risk, we have that $L(\theta) < L(\theta^*) + \epsilon$ in particular. Since $\pi_m(\theta^*) \rightarrow \theta^*$, there exists $M \in \mathbb{N}$ such that $\|\pi_m(\theta^*) - \theta^*\|_{L^2(P_X, \mathcal{X})} < \delta_2$ for $m \geq M$. By continuity of the population RU risk, we have that

$$L(\pi_m(\theta^*)) < L(\theta^*) + \epsilon \quad \text{for } m \geq M. \quad (79)$$

In addition, there exists $A_M \geq M$ so that $\|\theta_{A_M}^* - \theta^*\|_{L^2(P_X, \mathcal{X})} > \delta$, implying that

$$L(\theta_{A_M}^*) > L(\theta^*) + \epsilon.$$

However, this is a contradiction because $\theta_{A_M}^*$ is by definition the unique minimizer of the population RU risk over Θ_{A_M} , but we find that $\pi_{A_M}(\theta^*) \in \Theta_{A_M}$ satisfies

$$L(\pi_{A_M}(\theta^*)) < L(\theta_{A_M}^*).$$

Thus, we must have that $\|\theta_m^* - \theta^*\|_{L^2(P_X, \mathcal{X})} \rightarrow 0$ as $m \rightarrow \infty$.

C.13 Proof of Lemma 11

The goal of the proof is to verify that the conditions of Lemma 13 hold so that we can conclude that $\hat{\theta}_{m,n}$ exists with probability approaching 1 and $\hat{\theta}_{m,n} \xrightarrow{P} \theta_m^*$. First, we note that over the sieve space Θ_m , the population RU risk is uniquely minimized at θ_m^* by Theorem 9. To check the second condition, we observe that for m sufficiently large, $\theta_m^* \in \text{Int}(\Theta_m)$ because $\theta_m^* \rightarrow \theta^*$ by Theorem 10 and $\theta^* = (h^*, \alpha^*)$ where $0 < M_l \leq \alpha^*(x) < M_u$ for all $x \in \mathcal{X}$. Furthermore, it follows from the first part of Theorem 1 that $\theta \mapsto L_{\text{RU}}^\Gamma(\theta(x), y)$ is convex, which implies that the empirical risk $\hat{\mathbb{E}}_P [L_{\text{RU}}(\theta(X), Y)]$ is also convex. Third, by the Weak Law of Large Numbers, we have the following pointwise convergence

$$\hat{\mathbb{E}}_P [L_{\text{RU}}^\Gamma(\theta(X), Y)] \xrightarrow{P} \mathbb{E}_P [L_{\text{RU}}^\Gamma(\theta(X), Y)].$$

Thus, $\hat{\mathbb{E}}_P [L_{\text{RU}}(\theta(X), Y)]$ and $\mathbb{E}_P [L_{\text{RU}}^\Gamma(\theta(X), Y)]$ satisfy the conditions of Lemma 13. So, we have that $\hat{\theta}_{m,n}$ exists with probability approaching 1 and $\hat{\theta}_{m,n} \xrightarrow{P} \theta_m^*$.

C.14 Proof of Theorem 12

The main goal of this proof is to show that the following theorem applies to our setting.

Theorem 29 (Chen [2007], Theorem 3.2). *Let Z_i be distributed i.i.d. following a distribution P . Let $\theta^* \in \Theta$ be the population risk minimizer*

$$\theta^* = \underset{\theta \in \Theta}{\operatorname{argmin}} \mathbb{E}_P [l(\theta, Z_i)].$$

Let $\hat{\theta}_n$ be the empirical risk minimizer given by

$$\frac{1}{n} \sum_{i=1}^n l(\hat{\theta}_n, Z_i) \leq \inf_{\theta \in \Theta_n} \frac{1}{n} \sum_{i=1}^n l(\theta, Z_i) + O_P(\epsilon_n^2).$$

Let $\|\cdot\|$ be a norm on Θ such that $\|\hat{\theta}_n - \theta^*\| = o_P(1)$. Let $\mathcal{F}_n = \{l(\theta, Z_i) - l(\theta^*, Z_i) : \|\theta - \theta^*\| \leq \delta, \theta \in \Theta_n\}$. For some constant $b > 0$, let

$$\delta_n = \inf \left\{ \delta \in (0, 1) : \frac{1}{\sqrt{n}\delta^2} \int_{b\delta^2}^{\delta} \sqrt{H_{\square}(w^{1+\frac{d}{2p}}, \mathcal{F}_n, \|\cdot\|)} dw \leq 1 \right\},$$

where $H_{\square}(w, \mathcal{F}_n, \|\cdot\|_r)$ is the $L^r(P)$ metric entropy with bracketing of the class \mathcal{F}_n .

Assume that the following conditions hold.

1. In a neighborhood of θ^* , $\mathbb{E} [l(\theta, Z_i) - l(\theta^*, Z_i)] \asymp \|\theta - \theta^*\|^2$.

2. There is $C_1 > 0$ s.t. for all small $\epsilon > 0$

$$\sup_{\theta \in \Theta_n: \|\theta - \theta^*\| \leq \epsilon} \text{Var} [l(\theta, Z_i) - l(\theta^*, Z_i)] \leq C_1 \epsilon^2.$$

3. For any $\delta > 0$, there exists a constant $s \in (0, 2)$ such that

$$\sup_{\theta \in \Theta_n: \|\theta - \theta^*\| \leq \delta} |l(\theta, Z_i) - l(\theta^*, Z_i)| \leq \delta^s U(Z_i)$$

with $\mathbb{E} [U(Z_i)^\gamma] \leq C_2$ for some $\gamma \geq 2$.

Then $\|\hat{\theta}_n - \theta^*\| = O_P(\epsilon_n)$, where

$$\epsilon_n = \max\{\delta_n, \inf_{\theta \in \Theta_n} \|\theta^* - \theta\|\}.$$

We will use the following lemmas to show that the conditions of the above theorem are satisfied for our setting.

Lemma 30 (Chen and Shen [1998], Lemma 2). For $\theta \in \Lambda_c^p(\mathcal{X})$, we have that $\|\theta\|_\infty \leq 2c^{1 - \frac{2p}{2p+d}} \|\theta\|_{L^2(\lambda, \mathcal{X})}^{\frac{2p}{2p+d}}$, where λ is the Lebesgue measure.

Lemma 31. Under Assumptions 2, 4 5, 6, for any $h \in \Lambda_c^p(\mathcal{X})$, there exists $\bar{L}(X, Y)$ such that

$$|L(h(x), y) - L(h_\Gamma^*(x), y)| \leq \bar{L}(x, y) \cdot |h(x) - h_\Gamma^*(x)|,$$

where $\sup_{x \in \mathcal{X}} \mathbb{E}_{P_{Y|X}} [\bar{L}(x, Y)^2 | X = x] \leq M < \infty$. *Proof in Appendix D.11.*

For the metric, we will use $\|\cdot\|_{L^2(P_X, \mathcal{X})}$. Since any function $\theta \in \Theta$ only depends on X , $\|\cdot\|_{L^2(P_X, \mathcal{X})} = \|\cdot\|_{L^2(P, \mathcal{X} \times \mathcal{Y})}$. From Theorem 10, Lemma 11, we have that $\hat{\theta}_n \xrightarrow{P} \theta^*$ with respect to the $L^2(P_X, \mathcal{X})$ norm. So, $\|\theta^* - \hat{\theta}_n\|_{L^2(P_X, \mathcal{X})} = o_P(1)$.

First, we note that our observed data (X_i, Y_i) is i.i.d.

We aim to verify the second condition. We note that by Theorems 7 and 8, the population RU risk is strongly convex and smooth in a $\|\cdot\|_\infty$ -ball about the minimizer θ^* . We note that all θ in this $\|\cdot\|_\infty$ -ball about θ^* also must lie in a $\|\cdot\|_{L^2(P_X, \mathcal{X})}$ -ball about θ^* . So, in a $L^2(P_X, \mathcal{X})$ -neighborhood of θ^* , we have that

$$\mathbb{E}_P [L_{\text{RU}}^\Gamma(\theta(X), Y)] - \mathbb{E}_P [L_{\text{RU}}^\Gamma(\theta^*(X), Y)] \asymp \|\theta - \theta^*\|_{L^2(P_X, \mathcal{X})}^2.$$

We aim to verify the third condition. First, we show the following three intermediate results.

$$\mathbb{E}_P [(L(h(X), Y) - L(h_\Gamma^*(X), Y))^2] \lesssim \|h - h_\Gamma^*\|_{L^2(P_X, \mathcal{X})}^2. \quad (80)$$

$$\mathbb{E}_P [(\alpha(X) - \alpha_\Gamma^*(X))^2] \asymp \|\alpha - \alpha_\Gamma^*\|_{L^2(P_X, \mathcal{X})}^2. \quad (81)$$

$$\mathbb{E}_P [((L(h(X), Y) - \alpha(X))_+ - (L(h_\Gamma^*(X), Y) - \alpha_\Gamma^*(X))_+)^2] \lesssim \|\theta - \theta^*\|_{L^2(P_X, \mathcal{X})}^2. \quad (82)$$

(80) can be shown by apply Lemma 31.

$$\begin{aligned} \mathbb{E}_P [(L(h(X), Y) - L(h_\Gamma^*(X), Y))^2] &= \mathbb{E}_P [\bar{L}(X, Y)^2 \cdot (h(X) - h_\Gamma^*(X))^2] \\ &= \mathbb{E}_{P_X} [\mathbb{E}_{P_{Y|X}} [\bar{L}(X, Y)^2 \cdot (h(X) - h_\Gamma^*(X))^2 | X = x]] \\ &\leq \sup_{x \in \mathcal{X}} \mathbb{E}_{P_{Y|X}} [\bar{L}(x, Y)^2 | X = x] \cdot \|h - h_\Gamma^*\|_{L^2(P_X, \mathcal{X})}^2 \\ &\asymp \|h - h_\Gamma^*\|_{L^2(P_X, \mathcal{X})}^2. \end{aligned}$$

(81) is true by definition. So, we proceed to show (82). We use (80), (81).

$$\begin{aligned}
& \mathbb{E}_P \left[\left((L(h(X), Y) - \alpha(X))_+ - (L(h_\Gamma^*(X), Y) - \alpha_\Gamma^*(X))_+ \right)^2 \right] \\
& \leq \mathbb{E}_P \left[\left((L(h(X), Y) - \alpha(X)) - (L(h_\Gamma^*(X), Y) - \alpha_\Gamma^*(X)) \right)^2 \right] \\
& = \mathbb{E}_P \left[\left((L(h(X), Y) - L(h_\Gamma^*(X), Y)) - (\alpha(X) - \alpha_\Gamma^*(X)) \right)^2 \right] \\
& \leq 2\mathbb{E}_P \left[(L(h(X), Y) - L(h_\Gamma^*(X), Y))^2 \right] + 2\mathbb{E}_P \left[(\alpha(X) - \alpha_\Gamma^*(X))^2 \right] \\
& \lesssim \|h - h_\Gamma^*\|_{L^2(P_X, \mathcal{X})}^2 + \|\alpha - \alpha_\Gamma^*\|_{L^2(P_X, \mathcal{X})}^2 \\
& = \|\theta - \theta^*\|_{L^2(P_X, \mathcal{X})}^2.
\end{aligned}$$

Now, we consider $\theta \in \mathcal{B}_\epsilon$ where

$$\mathcal{B}_\epsilon = \{\theta \in \Theta_n \mid \|\theta - \theta^*\|_{L^2(P_X, \mathcal{X})} \leq \epsilon\}.$$

We aim to show that $\text{Var}_P [L_{\text{RU}}^\Gamma(\theta(X), Y) - L_{\text{RU}}^\Gamma(\theta^*(X), Y)] \lesssim \epsilon^2$ when $\|\theta - \theta^*\|_{L^2(P_X, \mathcal{X})} \leq \epsilon$.

$$\begin{aligned}
& \text{Var}_P [L_{\text{RU}}^\Gamma(\theta(X), Y) - L_{\text{RU}}^\Gamma(\theta^*(X), Y)] \\
& \leq \mathbb{E}_P [(L_{\text{RU}}^\Gamma(\theta(X), Y) - L_{\text{RU}}^\Gamma(\theta^*(X), Y))^2] \\
& \leq 3\mathbb{E}_P [(L(h(X), Y) - L(h_\Gamma^*(X), Y))^2] + 3\mathbb{E}_P [(\alpha(X) - \alpha_\Gamma^*(X))^2] \\
& \quad + 3\mathbb{E}_P \left[\left((L(h(X), Y) - \alpha(X))_+ - (L(h_\Gamma^*(X), Y) - \alpha_\Gamma^*(X))_+ \right)^2 \right] \\
& \lesssim \|h - h_\Gamma^*\|_{L^2(P_X, \mathcal{X})}^2 + \|\alpha - \alpha_\Gamma^*\|_{L^2(P_X, \mathcal{X})}^2 + \|\theta - \theta^*\|_{L^2(P_X, \mathcal{X})}^2 \\
& \lesssim \|\theta - \theta^*\|_{L^2(P_X, \mathcal{X})}^2.
\end{aligned}$$

The second line comes from the Cauchy-Schwarz inequality and the second last line comes from (80), (81), and (82). This prove the third condition.

Finally, we verify the fourth condition. We consider $\theta \in \mathcal{B}_\delta$, where

$$\mathcal{B}_\delta = \{\theta \in \Theta_n \mid \|\theta - \theta^*\|_{L^2(P_X, \mathcal{X})}^2 \leq \delta\}.$$

Using a similar argument as in the previous condition, we apply Lemma 31.

$$|L_{\text{RU}}^\Gamma(\theta(x), y) - L_{\text{RU}}^\Gamma(\theta^*(x), y)| \lesssim |\bar{L}(x, y) \cdot (h(x) - h^*(x))| + |\alpha(x) - \alpha^*(x)| \quad (83)$$

$$\lesssim |\bar{L}(x, y)| \cdot \|\theta - \theta^*\|_\infty \quad (84)$$

$$\lesssim |\bar{L}(x, y)| \cdot \|\theta - \theta^*\|_{L^2(\lambda)}^{\frac{2p}{2p+d}} \quad (85)$$

$$\lesssim |\bar{L}(x, y)| \cdot \|\theta - \theta^*\|_{L^2(P_X, \mathcal{X})}^{\frac{2p}{2p+d}}. \quad (86)$$

Since Assumption 5 holds, we can apply Lemma 30 to see that for $\theta \in \Theta$, $\|\theta\|_\infty \lesssim \|\theta\|_{L^2(\lambda)}^{\frac{2p}{2p+d}}$, where λ is the Lebesgue measure. This gives (85). Under Assumption 7, $\|\theta - \theta'\|_{L^2(P_X, \mathcal{X})} \asymp \|\theta - \theta'\|_{L^2(\lambda)}$, which gives (86).

Therefore, the fourth condition holds with $s = \frac{2p}{2p+d}$ and $U(X_i, Y_i) = |\bar{L}(X_i, Y_i)|$. So, by Theorem 29, we have that $\|\hat{\theta}_n - \theta^*\|_{L^2(P_X, \mathcal{X})} = O_P(\max\{\delta_n, \inf_{\theta \in \Theta_n} \|\theta - \theta^*\|_{L^2(P_X, \mathcal{X})}\})$.

Let $\mathcal{F}_n = \{L_{\text{RU}}^\Gamma(\theta(X_i), Y_i) - L_{\text{RU}}^\Gamma(\theta^*(X_i), Y_i) : \|\theta - \theta^*\|_{L^2(P_X, \mathcal{X})} \leq \delta, \theta \in \Theta_n\}$. Let $H_{\square}(w, \mathcal{F}_n, \|\cdot\|_{L^2(P_X, \mathcal{X})})$ be the $L^2(P_X, \mathcal{X})$ -metric entropy with bracketing of the class \mathcal{F}_n .

Since in our setting, we satisfy the fourth condition of Theorem 29 with $s = \frac{2p}{2p+d}$,

$$H_{\square}(w, \mathcal{F}_n, \|\cdot\|_2) \leq \log N(w^{1+\frac{d}{2p}}, \Theta_n, \|\cdot\|_{L^2(P_X, \mathcal{X})}).$$

Recall that $\tilde{\Theta}_n$ is the sieve space *without truncation*. We note that the covering number of Θ_n is upper bounded by the covering number of $\tilde{\Theta}_n$, so we have that

$$H_{\square}(w, \mathcal{F}_n, \|\cdot\|_2) \leq \log N(w^{1+\frac{d}{2p}}, \tilde{\Theta}_n, \|\cdot\|_{L^2(P_X, \mathcal{X})}).$$

For the finite-dimensional linear sieves, such as those in Example 1 and 2 without truncation, we have that

$$\log N(w^{1+\frac{d}{2p}}, \tilde{\Theta}_n, \|\cdot\|_{L^2(P_X, \mathcal{X})}) \lesssim \dim(\tilde{\Theta}_n) \log\left(\frac{1}{w}\right)$$

from Van de Geer and van de Geer [2000]. Then, we have that

$$\frac{1}{\sqrt{n}\delta^2} \int_{b\delta^2}^{\delta} \sqrt{\log N(w^{1+\frac{d}{2p}}, \tilde{\Theta}_n, \|\cdot\|_{L^2(P_X, \mathcal{X})})} dw \lesssim \frac{1}{\delta} \sqrt{\frac{\dim(\tilde{\Theta}_n)}{n} \log \frac{1}{\delta}}.$$

We realize that

$$\delta_n \asymp \sqrt{\frac{\dim(\tilde{\Theta}_n) \log n}{n}}.$$

We note that $\tilde{\Theta}_n = \tilde{\mathcal{H}}_n \times \tilde{\mathcal{A}}_n$. We have that $\dim(\tilde{\Theta}_n) = 2J_n^d = O(J_n^d)$. Plugging this in, we have that

$$\delta_n \asymp \sqrt{\frac{J_n^d \log n}{n}}.$$

Now, we can bound the approximation error $\inf_{\theta \in \Theta_n} \|\theta^* - \theta\|_{L^2(P_X, \mathcal{X})}$. Since the truncation of the sieve space is a contraction map to the true minimizer, we have that

$$\inf_{\theta \in \Theta_n} \|\theta^* - \theta\|_{L^2(P_X, \mathcal{X})} \leq \inf_{\theta \in \tilde{\Theta}_n} \|\theta^* - \theta\|_{\infty} \leq O(J_n^{-p}),$$

where the last inequality follows from Timan [2014]. So, we can set $J_n = \left(\frac{n}{\log n}\right)^{\frac{1}{2p+d}}$. Thus, we have that $\|\hat{\theta} - \theta^*\|_{L^2(P_X, \mathcal{X})} = O_P\left(\left(\frac{\log n}{n}\right)^{\frac{p}{2p+d}}\right)$.

D Proofs of Technical Lemmas

D.1 Proof of Lemma 17

We note that the eigenvalues of a 2x2 matrix must satisfy

$$\lambda^2 - (\text{tr } A)\lambda + \det A = 0.$$

Since $\text{tr } A \geq 0$ and $\det A \geq 0$, the minimum eigenvalue is given by So,

$$\lambda_{\min}(A) = \frac{\text{tr } A - \sqrt{(\text{tr } A)^2 - 4 \det A}}{2}.$$

Let $x = \text{tr } A$ and $y = \sqrt{(\text{tr } A)^2 - 4 \det A}$. Note that $y \leq x$ because $\det A \geq 0$. Then we have that

$$\lambda_{\min}(A) = \frac{x - y}{2} = \frac{x^2 - y^2}{2(x + y)} \geq \frac{x^2 - y^2}{2(x + x)} = \frac{x^2 - y^2}{4x} = \frac{\det A}{\text{tr } A}.$$

In addition, we have that

$$\lambda_{\max}(A) = \frac{x + y}{2} \leq \frac{x + x}{2} = x = \text{tr } A.$$

D.2 Proof of Lemma 18

Since H is Gâteaux differentiable with derivative equal to $H'_{h,\alpha}$, we aim to show that, for any $h, \tilde{h}, \alpha, \tilde{\alpha}$,

$$\langle H'_{h,\alpha}(h, \alpha) - H'_{\tilde{h},\tilde{\alpha}}(0, 0), (h - \tilde{h}, \alpha - \tilde{\alpha}) \rangle > 0$$

to establish strict convexity. Without loss of generality, we assume that $(\tilde{h}, \tilde{\alpha}) = (0, 0)$. Define the Gâteaux derivative of F and G in (h, α) to be $F'_{h,\alpha}, G'_{h,\alpha}$, respectively. Let the Gâteaux derivative of G with respect to h be G'_h . We have that

$$\begin{aligned} \langle H'_{h,\alpha}(h, \alpha) - H'_{h,\alpha}(0, 0), (h, \alpha) \rangle &= \langle F'_{h,\alpha}(h, \alpha) + G'_{h,\alpha}(h, \alpha) - F'_{h,\alpha}(0, 0) - G'_{h,\alpha}(0, 0), (h, \alpha) \rangle \\ &= \langle F'_{h,\alpha}(h, \alpha) - F'_{h,\alpha}(0, 0), (h, \alpha) \rangle + \langle G'_h(h, \alpha) - G'_h(0, 0), h \rangle. \end{aligned}$$

Note that G does not depend on α , so the Gâteaux derivative is $G'_\alpha(h, \alpha) = G'_\alpha(0, 0) = 0$. Since F is jointly convex in (h, α) and G is strongly convex in h and does not depend on α , both terms above are nonnegative.

If $h \neq 0$, then we have that

$$\begin{aligned} \langle H'_{h,\alpha}(h, \alpha) - H'_{h,\alpha}(0, 0), (h, \alpha) \rangle &\geq \langle G'_h(h, \alpha) - G'_h(0, 0), h \rangle \\ &\geq \mu_1 \|h\|^2 > 0, \end{aligned}$$

where the last line follows from G 's strong convexity in h . If $h = 0$ and $\alpha \neq 0$, then we have that

$$\begin{aligned} \langle H'_{h,\alpha}(h, \alpha) - H'_{h,\alpha}(0, 0), (h, \alpha) \rangle &= \langle H'_{h,\alpha}(0, \alpha) - H'_{h,\alpha}(0, 0), (0, \alpha) \rangle \\ &= \langle F'_{h,\alpha}(0, \alpha) - F'_{h,\alpha}(0, 0), (0, \alpha) \rangle \\ &> 0, \end{aligned}$$

where the last inequality follows due to the strict convexity of F in α . Thus, H is strictly convex in (h, α) .

D.3 Proof of Lemma 19

We have that $T_{1,x}^c(c) = -\mathbb{E}_{P_{Y|X=x}}[\ell'(Y - c)]$. In addition, $T_{1,x}^{cc}(c) = \mathbb{E}_{P_{Y|X=x}}[\ell''(Y - c)]$. So, $T_{1,x}$ is twice differentiable in c . In addition, we realize that

$$\begin{aligned} \mathbb{E}_P [L_{\text{RU},1}^\Gamma(h(X), Y)] &= \mathbb{E}_{P_X} [\mathbb{E}_{P_{Y|X=x}} [L_{\text{RU},1}^\Gamma(h(X), Y)]] \\ &= \Gamma^{-1} \mathbb{E}_{P_X} [T_{X,1}(h(X))]. \end{aligned}$$

D.4 Proof of Lemma 20

First, we compute the first derivatives of $T_{3,x}(c, d)$. Second, we compute the second derivatives of $T_{3,x}(c, d)$ when $d > 0$. Finally, we show that $T_{3,x}$ can be used to express $\mathbb{E}_P [L_{\text{RU},3}^\Gamma(h(X), \alpha(X), Y)]$.

Computing derivatives for the $d \leq 0$ case is straightforward. When $d \leq 0$, we have that

$$\begin{aligned} T_{3,x}^c(c, d) &= -\mathbb{E}_{P_{Y|X}}[\ell'(Y - c) | X = x] \\ T_{3,x}^d(c, d) &= -1. \end{aligned}$$

Now, we consider the $d > 0$ case. To compute the derivatives of

$$T_{3,x}(c, d) = \mathbb{E}_{P_{Y|X}} [(\ell(Y - c) - d)\mathbb{I}(\ell(Y - c) > d) | X = x],$$

we first identify when the condition $\ell(Y - c) > d$ is satisfied. The strong convexity of ℓ given by Assumption 2 implies that ℓ is strictly increasing on $y > 0$ and ℓ is strictly decreasing on $y < 0$. We define ℓ_1^{-1} to be the inverse of $\ell(y)$ on $y > 0$. We define ℓ_2^{-1} to be the inverse of $\ell(y)$ on $y < 0$. By the Inverse Function Theorem, we have that

$$(\ell_i^{-1})'(z) = \frac{1}{\ell'(\ell_i^{-1}(z))} \quad i = 1, 2. \quad (87)$$

We note that $\ell_1^{-1}(z) > 0$, and $\ell(y)$ strictly increasing on $y > 0$, so $\ell'(\ell_1^{-1}(z)) > 0$. By (87), we have that $(\ell_1^{-1})'(z) > 0$. This means that ℓ_1^{-1} is strictly increasing on its domain. By an analogous argument, we have that $(\ell_2^{-1})'(z) < 0$ and ℓ_2^{-1} is strictly decreasing on its domain.

Based on the results above, we realize that for $d > 0$,

$$\{y \in \mathbb{R} \mid \ell(y - c) > d\} = \{y \in \mathbb{R} \mid y - c > \ell_1^{-1}(d)\} \cup \{y \in \mathbb{R} \mid y - c < \ell_2^{-1}(d)\}.$$

Thus, we can rewrite $T_{3,x}(c, d)$ for $d > 0$ as follows

$$T_{3,x}(c, d) = \mathbb{E}_{Y|X=x} [(\ell(Y - c) - d)\mathbb{I}(Y - c > \ell_1^{-1}(d))] + \mathbb{E}_{Y|X=x} [(\ell(Y - c) - d)\mathbb{I}(Y - c < \ell_2^{-1}(d))].$$

Now, we can compute the derivatives of $T_{3,x}(c, d)$ on $d > 0$ as follows.

$$\begin{aligned} T_{3,x}^d(c, d) &= \mathbb{E}_{P_{Y|X}} [-1 \cdot \mathbb{I}(\ell(Y - c) > d) \mid X = x] \\ &= \mathbb{E}_{P_{Y|X=x}} [-1 \cdot \mathbb{I}(Y - c > \ell_1^{-1}(d))] + \mathbb{E}_{P_{Y|X}} [-1 \cdot \mathbb{I}(Y - c < \ell_2^{-1}(d)) \mid X = x] \\ &= -\Pr(Y > c + \ell_1^{-1}(d) \mid X = x) - \Pr(Y < c + \ell_2^{-1}(d) \mid X = x) \\ &= -1 + P_{Y|X=x}(c + \ell_1^{-1}(d)) - P_{Y|X=x}(c + \ell_2^{-1}(d)). \end{aligned}$$

Another way to express $T_{3,x}^d = -\Pr(\ell(Y - c) > d \mid X = x)$.

We realize that

$$\lim_{d \rightarrow 0^+} T_{3,x}^d(c, d) = -1 + P_{Y|X=x}(-c) - P_{Y|X=x}(-c) = -1 = \lim_{d \rightarrow 0^-} T_{3,x}^d(c, d),$$

so $T_{3,x}(c, d)$ is differentiable at $d = 0$. Also,

$$\begin{aligned} T_{3,x}^c(c, d) &= -\mathbb{E}_{P_{Y|X}} [\ell'(Y - c)\mathbb{I}(\ell(Y - c) > d) \mid X = x] \\ &= -\mathbb{E}_{P_{Y|X}} [\ell'(Y - c) \cdot \mathbb{I}(Y - c > \ell_1^{-1}(d))] - \mathbb{E}_{Y|X=x} [\ell'(Y - c) \cdot \mathbb{I}(Y - c < \ell_2^{-1}(d)) \mid X = x] \end{aligned}$$

We realize that

$$\lim_{d \rightarrow 0^+} T_{3,x}^c(c, d) = -\mathbb{E}_{P_{Y|X}} [\ell'(Y - c) \mid X = x] = \lim_{d \rightarrow 0^-} T_{3,x}^c(c, d),$$

so $T_{3,x}(c, d)$ is differentiable with respect to c .

Second, we compute the second derivatives of $T_{3,x}(c, d)$ when $d > 0$. It is straightforward to see that

$$T_{3,x}^{dc}(c, d) = p_{Y|X=x}(c + \ell_1^{-1}(d)) - p_{Y|X=x}(c + \ell_2^{-1}(d)).$$

In addition, we have that

$$\begin{aligned} T_{3,x}^{dd}(c, d) &= p_{Y|X=x}(c + \ell_1^{-1}(d)) \cdot \frac{1}{\ell'(\ell_1^{-1}(d))} - p_{Y|X=x}(c + \ell_2^{-1}(d)) \cdot \frac{1}{\ell'(\ell_2^{-1}(d))} \\ &= \sum_{i \in \{1, 2\}} p_{Y|X=x}(c + \ell_i^{-1}(d)) \cdot \frac{1}{|\ell'(\ell_i^{-1}(d))|}. \end{aligned}$$

The second line follows because $\ell'(\ell_2^{-1}(y)) < 0$. Finally, we compute $T_{3,x}^{cc}(c, d)$. First, we recall $T_{3,x}^c(c, d)$ from Lemma 20 and simplify it as follows.

$$\begin{aligned} T_{3,x}^c(c, d) &= -\mathbb{E}_{P_{Y|X}} [\ell'(Y - c)\mathbb{I}(\ell(Y - c) > d) \mid X = x] \\ &= -\mathbb{E}_{P_{Y|X}} [\ell'(Y - c)\mathbb{I}(Y > \ell_1^{-1}(d) + c) \mid X = x] - \mathbb{E}_{P_{Y|X}} [\ell'(Y - c)\mathbb{I}(Y < \ell_2^{-1}(d) + c) \mid X = x] \\ &= -\int_{\ell_1^{-1}(d)+c}^{\infty} \ell'(y - c)p_{Y|X=x}(y)dy - \int_{-\infty}^{\ell_2^{-1}(d)+c} \ell'(y - c)p_{Y|X=x}(y)dy \\ &= -\int_{\ell_1^{-1}(d)}^{\infty} \ell'(y)p_{Y|X=x}(y + c)dy - \int_{-\infty}^{\ell_2^{-1}(d)} \ell'(y)p_{Y|X=x}(y + c)dy. \end{aligned}$$

Now, we compute $T_{3,x}^{cc}(c, d)$ by differentiating with respect to c and applying integration by parts.

$$T_{3,x}^{cc}(c, d) = -\int_{\ell_1^{-1}(d)}^{\infty} \ell'(y)p'_{Y|X=x}(y + c)dy - \int_{-\infty}^{\ell_2^{-1}(d)} \ell'(y)p'_{Y|X=x}(y + c)dy$$

$$\begin{aligned}
&= -\left(\ell'(y)p_{Y|X=x}(y+c)\right)\Big|_{\ell_1^{-1}(d)}^{\infty} - \int_{\ell_1^{-1}(d)}^{\infty} p_{Y|X=x}(y+c)\ell''(y)dy \\
&\quad - \left(\ell'(y)p_{Y|X=x}(y+c)\right)\Big|_{-\infty}^{\ell_2^{-1}(d)} - \int_{-\infty}^{\ell_2^{-1}(d)} p_{Y|X=x}(y+c)\ell''(y)dy \\
&= \ell'(\ell_1^{-1}(d))p_{Y|X=x}(\ell_1^{-1}(d)+c) + \int_{\ell_1^{-1}(d)}^{\infty} p_{Y|X=x}(y+c)\ell''(y)dy \\
&\quad - \ell'(\ell_2^{-1}(d))p_{Y|X=x}(\ell_2^{-1}(d)+c) + \int_{-\infty}^{\ell_2^{-1}(d)} p_{Y|X=x}(y+c)\ell''(y)dy \\
&= \ell'(\ell_1^{-1}(d))p_{Y|X=x}(\ell_1^{-1}(d)+c) + \int_{c+\ell_1^{-1}(d)}^{\infty} p_{Y|X=x}(y)\ell''(y-c)dy \\
&\quad - \ell'(\ell_2^{-1}(d))p_{Y|X=x}(\ell_2^{-1}(d)+c) + \int_{-\infty}^{c+\ell_2^{-1}(d)} p_{Y|X=x}(y)\ell''(y-c)dy \\
&= \sum_{i \in \{1,2\}} |\ell'(\ell_i^{-1}(d))| \cdot p_{Y|X=x}(\ell_i^{-1}(d)+c) + \mathbb{E}_{P_{Y|X}} [\ell''(Y-c)\mathbb{I}(\ell(Y-c) > d) | X=x].
\end{aligned}$$

Thus, when $d > 0$, $T_{3,x}(c, d)$ is twice differentiable in (c, d) .

Lastly, we find that

$$\begin{aligned}
&\mathbb{E}_P [L_{\text{RU},3}^{\Gamma}(h(X), \alpha(X), Y)] \\
&= (\Gamma - \Gamma^{-1}) \cdot \mathbb{E}_P [(\ell(Y - h(X)) - \alpha(X))_+] \\
&= (\Gamma - \Gamma^{-1}) \cdot \mathbb{E}_{P_X} [\mathbb{E}_{P_{Y|X}} [(\ell(Y - h(X)) - \alpha(X))_+ | X]] \\
&= (\Gamma - \Gamma^{-1}) \cdot \mathbb{E}_{P_X} \left[\begin{cases} \mathbb{E}_{P_{Y|X}} [(\ell(Y - h(X)) - \alpha(X))\mathbb{I}(\ell(Y - h(X)) > \alpha(X))] & \alpha(X) > 0 \\ \mathbb{E}_{P_{Y|X}} [(\ell(Y - h(X)) - \alpha(X))] & \alpha(X) \leq 0 \end{cases} \right] \\
&= (\Gamma - \Gamma^{-1})\mathbb{E}_{P_X} [T_{3,X}(h(X), \alpha(X))].
\end{aligned}$$

D.5 Proof of Lemma 21

Now, define a symmetric 2×2 matrix $A_x(c, d)$ where

$$\begin{aligned}
A_{x,11}(c, d) &= T_{3,x}^{cc}(c, d) - \mathbb{E}_{Y|X=x} [\ell''(Y-c)\mathbb{I}(\ell(Y-c) > d)] + C_{L,l} \cdot \Pr(\ell(Y-c) > d | X=x) \\
A_{x,22}(c, d) &= T_{3,x}^{dd}(c, d) \\
A_{x,12}(c, d) &= T_{3,x}^{dc}(c, d),
\end{aligned}$$

where F is the distribution over $\ell(Y-c)$ where Y follows $P_{Y|X=x}$. Under Assumption 2, we have that ℓ is $C_{l,l}$ -strongly convex, so

$$\mathbb{E}_{Y|X=x} [\ell''(Y-c)\mathbb{I}(\ell(Y-c) > d)] - C_{L,l} \cdot \Pr(\ell(Y-c) > d | X=x) \geq 0.$$

Thus, we have that

$$\begin{aligned}
\nabla^2 T_{3,x}(c, d) - A_x(c, d) &= \begin{bmatrix} \mathbb{E}_{Y|X=x} [(\ell''(Y-c) - C_{L,l})\mathbb{I}(\ell(Y-c) > d)] & 0 \\ 0 & 0 \end{bmatrix} \\
&= \begin{bmatrix} \mathbb{E}_{Y|X=x} [(\ell''(Y-c)\mathbb{I}(\ell(Y-c) > d)] - C_{L,l}\Pr(\ell(Y-c) > d | X=x) & 0 \\ 0 & 0 \end{bmatrix} \\
&\succeq 0.
\end{aligned}$$

So,

$$\nabla^2 T_{3,x}(c, d) \succeq A_x(c, d). \quad (88)$$

We can also define a symmetric 2×2 matrix $B_x(c, d)$ where

$$\begin{aligned} B_{x,11}(c, d) &= T_{3,x}^{cc}(c, d) + \mathbb{E}_{Y|X=x} [\ell''(Y - c)\mathbb{I}(\ell(Y - c) \leq d)] \\ B_{x,22}(c, d) &= T_{3,x}^{dd}(c, d) \\ B_{x,12}(c, d) &= T_{3,x}^{dc}(c, d). \end{aligned}$$

Under Assumption 2, we have that ℓ is strongly convex, so

$$\mathbb{E}_{Y|X=x} [\ell''(Y - c)\mathbb{I}(\ell(Y - c) \leq d)] \geq C_{\ell,l} \cdot \Pr(\ell(Y - c) \leq d|X = x) \geq 0.$$

Thus, we have that

$$B_x(c, d) - \nabla^2 T_{3,x}(c, d) = \begin{bmatrix} \mathbb{E}_{Y|X=x} [\ell''(Y - c)\mathbb{I}(\ell(Y - c) \leq d)] & 0 \\ 0 & 0 \end{bmatrix} \succeq 0.$$

So,

$$\nabla^2 T_{3,x}(c, d) \preceq B_x(c, d). \quad (89)$$

Combining (88) and (89) yields the desired result.

D.6 Proof of Lemma 22

Let $L(h, \alpha) = \mathbb{E}_P [L_{\text{RU}}^\Gamma(h(X), \alpha(X), Y)]$. First, we verify Gâteaux differentiability with respect to α . We show that the directional derivative of $L(h, \alpha)$ with respect to α in the direction ϕ exists for all $\phi \in L^2(P_X, \mathcal{X})$. We note that the directional derivative with respect to α in the direction ϕ is given by

$$L'_\alpha(h, \alpha; \phi) = \lim_{\theta \rightarrow 0^+} \frac{L(h, \alpha + \theta\phi) - L(h, \alpha)}{\theta}.$$

We simplify the numerator as follows

$$\begin{aligned} &L(h, \alpha + \theta\phi) - L(h, \alpha) \\ &= \mathbb{E}_P [L_{\text{RU},2}^\Gamma((\alpha + \theta\phi)(X))] + \mathbb{E}_P [L_{\text{RU},3}^\Gamma(h(X), (\alpha + \theta\phi)(X), Y)] \\ &\quad - \mathbb{E}_P [L_{\text{RU},2}^\Gamma(\alpha(X))] - \mathbb{E}_P [L_{\text{RU},3}^\Gamma(h(X), \alpha(X), Y)] \\ &= \theta(1 - \Gamma^{-1}) \cdot \mathbb{E}_{P_X} [\phi(X)] + (\Gamma - \Gamma^{-1}) \cdot \mathbb{E}_{P_X} [T_{3,X}(h(X), (\alpha + \theta\phi)(X)) - T_{3,X}(h(X), \alpha(X))]. \end{aligned}$$

The first line follows because only the second and third term of the RU loss depend on α . The second line follows by Lemma 20. We analyze the second term on the right side of the above equation. We note that by Lemma 20, the map $T_{3,x}(c, d)$ is differentiable with respect to c, d . So,

$$\lim_{\theta \rightarrow 0^+} \frac{T_{3,x}(h(x), \alpha(x) + \theta\phi(x)) - T_{3,x}(h(x), \alpha(x))}{\theta} = T_{3,x}^d(h(x), \alpha(x))\phi(x).$$

Therefore, we have that

$$\begin{aligned} L'_\alpha(h, \alpha; \phi) &= \lim_{\theta \rightarrow 0^+} (1 - \Gamma^{-1}) \cdot \frac{\theta \mathbb{E}_{P_X} [\phi(X)]}{\theta} \\ &\quad + \lim_{\theta \rightarrow 0^+} (\Gamma - \Gamma^{-1}) \cdot \frac{\mathbb{E}_{P_X} [T_{3,X}(h(X), \alpha(x) + \theta\phi(X)) - T_{3,X}(h(X), \alpha(X))]}{\theta} \\ &= (1 - \Gamma^{-1}) \cdot \mathbb{E}_{P_X} [\phi(X)] + (\Gamma - \Gamma^{-1}) \cdot \mathbb{E}_{P_X} [T_{3,X}^d(h(X), \alpha(X))\phi(X)] \\ &= \mathbb{E}_{P_X} [((1 - \Gamma^{-1}) + (\Gamma - \Gamma^{-1}) \cdot T_{3,X}^d(h(X), \alpha(X))) \cdot \phi(X)]. \end{aligned}$$

Since the directional derivative of $L(h, \alpha)$ with respect to α and in the direction ϕ exists for all $\phi \in \mathcal{A}$, then $L(h, \alpha)$ is Gâteaux differentiable in α .

We use a similar technique to verify Gâteaux differentiability with respect to h . We show that the directional derivative of $L(h, \alpha)$ with respect to h in the direction ψ exists for $\psi \in \mathcal{H}$. We recall that the directional derivative of $L(h, \alpha)$ with respect to h in the direction ψ is given by

$$L'_h(h, \alpha; \psi) = \lim_{\theta \rightarrow 0^+} \frac{L(h + \theta\psi, \alpha) - L(h, \alpha)}{\theta}. \quad (90)$$

We simplify the directional derivative in (90) as follows.

$$\begin{aligned} L'_h(h, \alpha; \psi) &= \lim_{\theta \rightarrow 0^+} \frac{L(h + \theta\psi, \alpha) - L(h, \alpha)}{\theta} \\ &= \frac{\mathbb{E}_P [L_{\text{RU},1}^\Gamma((h + \theta\psi)(X), Y) - L_{\text{RU},1}^\Gamma(h(X), Y)]}{\theta} \\ &\quad + \lim_{\theta \rightarrow 0^+} \frac{\mathbb{E}_P [L_{\text{RU},3}^\Gamma((h + \theta\psi)(X), \alpha(X), Y) - L_{\text{RU},3}^\Gamma(h(X), \alpha(X), Y)]}{\theta} \\ &= \lim_{\theta \rightarrow 0^+} \Gamma^{-1} \cdot \frac{\mathbb{E}_{P_X} [T_{1,X}(h(X) + \theta\psi(X)) - T_{1,X}(h(X))]}{\theta} \\ &\quad + \lim_{\theta \rightarrow 0^+} (\Gamma - \Gamma^{-1}) \cdot \frac{\mathbb{E}_{P_X} [T_{3,X}(h(X) + \theta\psi(X), \alpha(X)) - T_{3,X}(h(X), \alpha(X))]}{\theta} \\ &= \mathbb{E}_{P_X} [(\Gamma^{-1} \cdot T_{1,X}^c(h(X)) + (\Gamma - \Gamma^{-1}) \cdot T_{3,X}^c(h(X), \alpha(X))) \cdot \psi(X)]. \end{aligned}$$

The first line follows because only the first and third terms of the RU loss depend on h . The second line follows because of Lemma 19 and Lemma 20. The third line follows from the differentiability of $T_{1,x}, T_{3,x}$, which is given by Lemmas 19 and 20. Since the directional derivative of $L(h, \alpha)$ with respect to h and in the direction ψ exists for all $\psi \in L^2(P_X, \mathcal{X})$, and the directional derivative can be expressed as a continuous linear function (given the inner product on $L^2(P_X, \mathcal{X})$), then $L(h, \alpha)$ is Gâteaux differentiable in h .

We can compute second derivatives of $L(h, \alpha)$ on \mathcal{C} by applying Lemma 19 and Lemma 20. Note that $T_{3,x}$ is twice-differentiable when $d > 0$. For $(h, \alpha) \in \mathcal{C}$, we have that $\alpha(x) \geq 0$. We note that the restriction of \mathcal{C} to the coordinate that corresponds to h is $L^2(P_X, \mathcal{X})$. Let \mathcal{A}' be the restriction of \mathcal{C} to the coordinate that corresponds to α . In the following result, we consider $\psi_1, \psi_2 \in L^2(P_X, \mathcal{X})$ and $\phi_1, \phi_2 \in \mathcal{A}'$. We find that

$$\begin{aligned} L''_{hh}(h, \alpha; \psi_1, \psi_2) &= L''_{1,hh}(h, \alpha; \psi_1, \psi_2) + L''_{3,hh}(h, \alpha; \psi_1, \psi_2) \\ &= \Gamma^{-1} \mathbb{E}_{P_X} [T_{1,X}^{cc}(h(X)) \psi_1(X) \psi_2(X)] + (\Gamma - \Gamma^{-1}) \mathbb{E}_{P_X} [T_{3,X}^{cc}(h(X), \alpha(X)) \psi_1(X) \psi_2(X)]. \\ L''_{h\alpha}(h, \alpha; \psi_1, \phi_1) &= L''_{3,h\alpha}(h, \alpha; \psi_1, \phi_1) \\ &= (\Gamma - \Gamma^{-1}) \mathbb{E}_{P_X} [T_{3,X}^{cd}(h(X), \alpha(X)) \psi_1(X) \phi_1(X)]. \\ L''_{\alpha\alpha}(h, \alpha; \phi_1, \phi_2) &= L''_{3,\alpha\alpha}(h, \alpha; \phi_1, \phi_2) \\ &= (\Gamma - \Gamma^{-1}) \mathbb{E}_{P_X} [T_{3,X}^{dd}(h(X), \alpha(X)) \phi_1(X) \phi_2(X)]. \end{aligned}$$

D.7 Proof of Lemma 23

Define $L_1(h, \alpha) = \mathbb{E}_P [L_{\text{RU},1}^\Gamma(h(X), Y)]$. From Lemma 22, we have that $L_1(h, \alpha)$ is twice Gâteaux differentiable in h with

$$\begin{aligned} L''_{1,h}(h, \alpha; \psi, \psi) &= \Gamma^{-1} \mathbb{E}_{P_X} [T_{1,X}^{cc}(h(X)) \cdot (\psi(X))^2] \\ &\geq \Gamma^{-1} \cdot C_{L,\ell} \|\psi\|_{L^2(\mathcal{X}, P_X)}^2 \end{aligned}$$

for $\psi \in L^2(P_X, \mathcal{X})$. The last line follows from Assumption 2, where we assume that ℓ is strongly convex. Thus, we have that $\mathbb{E}_P [L_{\text{RU},1}^\Gamma(h(X), Y)]$ is $\Gamma^{-1} \cdot C_{L,\ell}$ -strongly convex in h .

D.8 Proof of Lemma 24

Let $L_3(h, \alpha) = \mathbb{E}_P [L_{\text{RU},3}^\Gamma(h(X), \alpha(X), Y)]$. By Lemma 22, we have that $\mathbb{E}_P [L_{\text{RU},3}^\Gamma(h(X), \alpha(X), Y)]$ is Gâteaux differentiable in α with

$$L'_{3,\alpha}(h, \alpha; \phi) = (\Gamma - \Gamma^{-1}) \cdot \mathbb{E}_X [T_{3,X}^d(h(X), \alpha(X)) \cdot \phi(X)].$$

We aim to verify the strict convexity of $\alpha \mapsto \mathbb{E}_P [L_{\text{RU},3}^\Gamma(h(X), \alpha(X), Y)]$ via Lemma 24. So, we must show that for $\alpha_1, \alpha_2 \in \mathcal{A}$ that differ on a set of positive measure, we have that

$$\mathbb{E}_{P_X} [(T_{3,X}^d(h(X), \alpha_1(X)) - T_{3,X}^d(h(X), \alpha_2(X)) \cdot (\alpha_1(X) - \alpha_2(X))] > 0.$$

From Lemma 20, we have that

$$T_{3,x}^d(h(x), \alpha(x)) = \begin{cases} t_{3,x}^d(h(x), \alpha(x)) & \alpha(x) > 0 \\ -1 & \alpha(x) \leq 0 \end{cases},$$

where

$$t_{3,x}^d(h(x), \alpha(x)) = -1 + P_{Y|X=x}(h(x) + \ell_1^{-1}(\alpha(x))) - P_{Y|X=x}(h(x) + \ell_2^{-1}(\alpha(x))).$$

By the definition of ℓ_1^{-1}, ℓ_2^{-1} from Lemma 20, we have that

$$\ell_1^{-1}(\alpha(x)) > \ell_2^{-1}(\alpha(x)).$$

Under Assumption 3, we have that $P_{Y|X=x}$ is strictly increasing, so

$$P_{Y|X=x}(h(x) + \ell_1^{-1}(\alpha(x))) - P_{Y|X=x}(h(x) + \ell_2^{-1}(\alpha(x))) > 0,$$

which implies that

$$t_{3,x}^d(h(x), \alpha(x)) > -1. \quad (91)$$

Under Assumption 2, ℓ_1^{-1} is strictly increasing and ℓ_2^{-1} is strictly decreasing. We realize that if $\alpha_1(x) > \alpha_2(x)$, then

$$\begin{aligned} P_{Y|X=x}(h(x) + \ell_1^{-1}(\alpha_1(x))) &> P_{Y|X=x}(h(x) + \ell_1^{-1}(\alpha_2(x))) \\ P_{Y|X=x}(h(x) + \ell_2^{-1}(\alpha_1(x))) &< P_{Y|X=x}(h(x) + \ell_2^{-1}(\alpha_2(x))), \end{aligned}$$

so

$$t_{3,x}^d(h(x), \alpha_1(x)) > t_{3,x}^d(h(x), \alpha_2(x)). \quad (92)$$

Let $D = \{x \in \mathcal{X} \mid \alpha_1(x) \neq \alpha_2(x)\}$. Now, we compute

$$\mathbb{E}_{P_X} [(T_{3,X}^d(h(X), \alpha_1(X)) - T_{3,X}^d(h(X), \alpha_2(X)) \cdot (\alpha_1(X) - \alpha_2(X))] \quad (93a)$$

$$= \mathbb{E}_{P_X} [(T_{3,X}^d(h(X), \alpha_1(X)) - T_{3,X}^d(h(X), \alpha_2(X)) \cdot (\alpha_1(X) - \alpha_2(X)) \mathbb{I}(D)] \quad (93b)$$

$$= \mathbb{E}_{P_X} [((t_{3,X}^d(h(X), \alpha_1(X)) - t_{3,X}^d(h(X), \alpha_2(X))) \cdot (\alpha_1(X) - \alpha_2(X)) \cdot \mathbb{I}(S_{\alpha_1,0} \cap S_{\alpha_2,0} \cap D))] \quad (93c)$$

$$+ \mathbb{E}_{P_X} [(t_{3,X}^d(h(X), \alpha_1(X)) + 1)(\alpha_1(X) - \alpha_2(X)) \cdot \mathbb{I}(S_{\alpha_1,0} \cap S_{\alpha_2,0}^c \cap D)] \quad (93d)$$

$$+ \mathbb{E}_{P_X} [(-1 - t_{3,X}^d(h(X), \alpha_2(X)))(\alpha_1(X) - \alpha_2(X)) \cdot \mathbb{I}(S_{\alpha_1,0}^c \cap S_{\alpha_2,0} \cap D)]. \quad (93e)$$

The first line holds because $(T_{3,x}^d(h(x), \alpha_1(x)) - T_{3,x}^d(h(x), \alpha_2(x)) \cdot (\alpha_1(x) - \alpha_2(x)) = 0$ on D^c . The decomposition into (93c), (93d), (93e) holds because $T_{3,x}^d(h(x), \alpha_1(x)) - T_{3,x}^d(h(x), \alpha_2(x)) = 0$ when $\alpha_1(x) \leq 0$ and $\alpha_2(x) \leq 0$.

Since we have that $\alpha_1, \alpha_2 \in \mathcal{A}$ and D has positive measure, we can show that $P(S_{\alpha_1,0}^c \cap S_{\alpha_2,0}^c \cap D) < P(D)$. We consider two cases 1) $S_{\alpha_1,0} \cap D$ has positive measure and 2) $S_{\alpha_1,0} \cap D = \emptyset$. Suppose $S_{\alpha_1,0} \cap D$ has positive measure, then clearly

$$P(S_{\alpha_1,0}^c \cap S_{\alpha_2,0}^c \cap D) \leq P(S_{\alpha_1}^c \cap D) < P(D).$$

If $S_{\alpha_1,0} \cap D$ empty, this means that $\alpha_1(x) \leq 0$ for all $x \in D$. At the same time, we have that for all $\alpha \in \mathcal{A}$, $\alpha(x) \geq 0$ for every $x \in \mathcal{X}$. So, we must have that $\alpha_1 = 0$ on D . We must have that $\alpha_2(x) > 0$ on D , because α_1, α_2 must differ on D and $\alpha_2(x) \geq 0$ for all $x \in \mathcal{X}$. So, this means that $S_{\alpha_2,0} \cap D$ has positive measure, so

$$P(S_{\alpha_1,0}^c \cap S_{\alpha_2,0}^c \cap D) \leq P(S_{\alpha_2,0}^c \cap D) < P(D).$$

So, at least at least one of the sets $S_{\alpha_1,0} \cap S_{\alpha_2,0} \cap D$, $S_{\alpha_1,0} \cap S_{\alpha_2,0}^c \cap D$, $S_{\alpha_1,0}^c \cap S_{\alpha_2,0} \cap D$ has positive measure.

Suppose $S_{\alpha_1,0} \cap S_{\alpha_2,0} \cap D$ has positive measure. WLOG, if $\alpha_1(x) > \alpha_2(x)$, then $T_{3,x}^d(h(x), \alpha_1(x) - T_{3,x}^d(h(x), \alpha_2(x))) > 0$. In addition, if $\alpha_1(x) < \alpha_2(x)$, then $T_{3,x}^d(h(x), \alpha_1(x)) - T_{3,x}^d(h(x), \alpha_2(x)) < 0$. Then (93c) must be positive. We can use a similar argument to verify that (93d) will be positive if $S_{\alpha_1,0} \cap S_{\alpha_2,0}^c \cap D$ has positive measure and (93e) will be positive if $S_{\alpha_1,0}^c \cap S_{\alpha_2,0}$ has positive measure. Thus, we conclude that

$$\mathbb{E}_{P_X} [(T_{3,X}^d(h(X), \alpha_1(X)) - T_{3,X}^d(h(X), \alpha_2(X)) \cdot (\alpha_1(X) - \alpha_2(X))] > 0$$

so $\alpha \mapsto \mathbb{E}_P [L_{\text{RU},3}^\Gamma(h(X), \alpha(X), Y)]$ is strictly convex on \mathcal{A} .

D.9 Proof of Lemma 25

Let $h^* \in L^2(Q_X, \mathcal{X})$ be the solution to (47). Let the function \tilde{h} be minimizer of (48) at every x . Since \tilde{h} solves (48) for every $x \in \text{supp}(Q_X)$,

$$\sup_{Q \in S_\Gamma(P, Q_X)} \mathbb{E}_{Q_{Y|X}} [L(\tilde{h}(X), Y) | X = x] \leq \sup_{Q \in S_\Gamma(P, Q_X)} \mathbb{E}_{Q_{Y|X}} [L(h^*(X), Y) | X = x].$$

Given any marginal distribution Q_X , we can marginalize over X to see that

$$\mathbb{E}_{Q_X} \left[\sup_{Q_{Y|X}: Q \in S_\Gamma(P)} \mathbb{E}_{Q_{Y|X}} [L(\tilde{h}(X), Y) | X] \right] \leq \mathbb{E}_{Q_X} \left[\sup_{Q_{Y|X}: Q \in S_\Gamma(P)} \mathbb{E}_{Q_{Y|X}} [L(h^*(X), Y) | X] \right].$$

Based on our definition of $S_\Gamma(P, Q_X)$, we note that for any $h \in L^2(Q_X, \mathcal{X})$

$$\sup_{Q \in S_\Gamma(P, Q_X)} \mathbb{E}_Q [L(h(X), Y)] = \mathbb{E}_{Q_X} \left[\sup_{Q_{Y|X}: Q \in S_\Gamma(P, Q_X)} \mathbb{E}_{Q_{Y|X}} [L(h(X), Y) | X] \right].$$

Thus, we have that

$$\sup_{Q \in S_\Gamma(P, Q_X)} \mathbb{E}_Q [L(\tilde{h}(X), Y)] \leq \sup_{Q \in S_\Gamma(P, Q_X)} \mathbb{E}_Q [L(h^*(X), Y)].$$

Finally, by definition of h^* we must also have that

$$\sup_{Q \in S_\Gamma(P, Q_X)} \mathbb{E}_Q [L(h^*(X), Y)] \leq \sup_{Q \in S_\Gamma(P, Q_X)} \mathbb{E}_Q [L(\tilde{h}(X), Y)].$$

These last two inequalities yield the desired equivalence.

D.10 Proof of Lemma 28

Recall that we defined the sieve *with truncation* Θ_m and sieve *without truncation* $\tilde{\Theta}_m$. In addition, define $\tilde{\pi}_m : \Theta \rightarrow \tilde{\Theta}_m$ to be the projection of a function $\theta \in \Theta$ onto $\tilde{\Theta}_m$.

By Lemma 2, the truncation is a contraction map to the true minimizer, so

$$\|\pi_m(\theta^*) - \theta^*\|_{L^2(P_X, \mathcal{X})} \leq \|\tilde{\pi}_m(\theta^*) - \theta^*\|_{L^2(P_X, \mathcal{X})}.$$

Now, we verify the conditions of Lemma 15 to show that the right side of the above inequality converges to zero as $m \rightarrow \infty$. First, we note that $\tilde{\Theta}$ is a Hilbert space (with the $L^2(P_X, \mathcal{X})$ norm). Second, we note that

$$\tilde{\pi}_m(\theta^*) = \sum_{i=1}^m \langle \theta^*, \phi_i \rangle \phi_i,$$

where $\{\phi_i\}$ is an infinite-dimensional basis for Θ . Since $\tilde{\pi}_m(\theta^*)$ is a partial sum of the Fourier-Bessel series, we have that $\tilde{\pi}_m(\theta^*) \rightarrow \theta^*$. This implies that $\|\pi_m(\theta^*) - \theta^*\|_{L^2(P_X, \mathcal{X})}$, as well.

D.11 Proof of Lemma 31

First, by the Mean Value Theorem, we have that for any $z \in \mathbb{R}$,

$$|\ell'(z)| = |\ell'(z) - \ell'(0)| \leq |\ell''(\tilde{z})| \cdot |z|,$$

where \tilde{z} is between z and 0. By Assumption 4, $|\ell''(\tilde{z})| \leq C_{L,u}$, so

$$|\ell'(z)| \leq C_{L,u} \cdot |z|. \quad (94)$$

Again, by the Mean Value Theorem, we have that for any $h \in \Lambda_c^p(\mathcal{X})$ and $x \in \mathcal{X}$,

$$\begin{aligned} |L(h(x), y) - L(h^*(x), y)| &= |\ell(y - h(x)) - \ell(y - h^*(x))| \\ &= |\ell'(y - (\lambda(x) \cdot h(x) + (1 - \lambda(x)) \cdot h^*(x)))| \cdot |h(x) - h^*(x)| \quad \lambda(x) \in [0, 1]. \end{aligned}$$

We can define $\bar{L}(x, y) = |\ell'(y - (\lambda(x) \cdot h(x) + (1 - \lambda(x)) \cdot h^*(x)))|$. Now, we aim to verify that there exists some $0 < M < \infty$ such that

$$\sup_{x \in \mathcal{X}} \mathbb{E}_{P_{Y|X}} [\bar{L}(X, Y)^2 | X = x] < M.$$

We apply (94).

$$\begin{aligned} \mathbb{E}_{P_{Y|X}} [\bar{L}(x, Y)^2 | X = x] &= \mathbb{E}_{P_{Y|X}} [(\ell'(Y - (\lambda(x) \cdot h(x) + (1 - \lambda(x)) \cdot h^*(x))))^2 | X = x] \\ &= \mathbb{E}_{P_{Y|X}} [(\ell'(Y - (\lambda(x) \cdot h(x) + (1 - \lambda(x)) \cdot h^*(x)))) \cdot h^*(x))^2 | X = x] \\ &= \mathbb{E}_{P_{Y|X}} [C_{L,u}^2 \cdot ((Y - (\lambda(x) \cdot h(x) + (1 - \lambda(x)) \cdot h^*(x)))) \cdot h^*(x))^2 | X = x] \\ &\lesssim \mathbb{E}_{P_{Y|X}} [Y^2 | X = x] + h(x)^2 + h^*(x)^2 \\ &\lesssim \sup_{x \in \mathcal{X}} \mathbb{E}_{P_{Y|X}} [Y^2 | X = x] + c^2 \\ &< \infty. \end{aligned}$$

The last two lines follow from Assumption 5 and 6. Assumption 5 gives that $h, h^* \in \Lambda_c^p(x)$, so $|h(x)| \leq c$ and $|h^*(x)| \leq c$. Assumption 6 gives that $\sup_{x \in \mathcal{X}} \mathbb{E}_{P_{Y|X}} [Y^2 | X = x]$ is finite. Thus, $\sup_{x \in \mathcal{X}} \mathbb{E}_{P_{Y|X}} [\bar{L}(X, Y)^2 | X = x] < \infty$.