# Confidence regions in Wasserstein distributionally robust estimation

### By JOSE BLANCHET

Department of Management Science and Engineering, Stanford University, Huang Engineering Center, 475 Via Ortega, Stanford, California 94305, U.S.A. jose.blanchet@stanford.edu

# KARTHYEK MURTHY

Engineering Systems and Design pillar, Singapore University of Technology and Design, 8 Somapah Road, Singapore 487372 karthyek\_murthy@sutd.edu.sg

#### AND NIAN SI

Department of Management Science and Engineering, Stanford University, Huang Engineering Center, 475 Via Ortega, Stanford, California 94305, U.S.A. niansi@stanford.edu

#### SUMMARY

Estimators based on Wasserstein distributionally robust optimization are obtained as solutions of min-max problems in which the statistician selects a parameter minimizing the worst-case loss among all probability models within a certain distance from the underlying empirical measure in a Wasserstein sense. While motivated by the need to identify optimal model parameters or decision choices that are robust to model misspecification, these distributionally robust estimators recover a wide range of regularized estimators, including square-root lasso and support vector machines, among others. This paper studies the asymptotic normality of these distributionally robust estimators as well as the properties of an optimal confidence region induced by the Wasserstein distributionally robust optimization formulation. In addition, key properties of minmax distributionally robust optimization problems are also studied; for example, we show that distributionally robust estimators regularize the loss based on its derivative, and we also derive general sufficient conditions which show the equivalence between the min-max distributionally robust optimization problem and the corresponding max-min formulation.

Some key words: Asymptotic normality; Confidence region; Distributionally robust optimization; Wasserstein distance.

## 1. Introduction

In recent years, distributionally robust optimization formulations based on Wasserstein distances have sparked a substantial amount of interest. One reason for this interest, as demonstrated by a range of examples in statistical learning and operations research, is that these formulations

provide a flexible way to quantify and hedge against the impact of model misspecification. Motivated by those applications, this paper aims to understand their fundamental statistical properties, such as asymptotic normality of the distributionally robust estimators and the associated confidence regions deemed optimal in a suitable sense to be described shortly.

Before providing a review of Wasserstein distributionally robust optimization and its connections to several areas, such as artificial intelligence, machine learning and operations research, we set the stage by first introducing the elements of a typical data-driven distributionally robust estimation problem.

Suppose that  $\{X_k : 1 \le k \le n\} \subset \mathbb{R}^m$  are independent and identically distributed samples from an unknown distribution  $P_*$ . A typical nonrobust stochastic optimization formulation informed by  $P_n$  focuses on minimizing empirical expected loss of the form  $E_{P_n}$   $\{\ell(X;\beta)\} = n^{-1} \sum_{i=1}^n \ell(X_i;\beta)$  over the parameter choices  $\beta \in B \subseteq \mathbb{R}^d$ . In this paper we take B to be a closed, convex subset of  $\mathbb{R}^d$ . Let the empirical risk minimization estimators be

$$\beta_n^{\text{ERM}} \in \arg\min_{\beta \in B} E_{P_n} \{ \ell(X; \beta) \}.$$
 (1)

On the other hand, a distributionally robust formulation recognizes the distributional uncertainty inherent in  $P_n$  being a noisy representation of an unknown distribution. Therefore, it enriches the empirical risk minimization (1) by considering an estimator of the form

$$\beta_n^{\mathrm{DRO}}(\delta) \in \arg\min_{\beta \in B} \sup_{P \in \mathcal{U}_\delta(P_n)} E_P \{\ell(X; \beta)\},$$
 (2)

where the set  $\mathcal{U}_{\delta}(P_n)$  is called the distributional uncertainty region and  $\delta$  is the size of the distributional uncertainty. Here, given a measurable function  $f(\cdot)$ , the notation  $E_P\{f(X)\}$  denotes expectation with respect to a probability distribution P. Wasserstein distributionally robust formulations advocate choosing

$$\mathcal{U}_{\delta}(P_n) = \{ P \in \mathcal{P}(\Omega) : W(P_n, P) \leqslant \delta^{1/2} \},$$

where  $W(P_n, P)$  is the Wasserstein distance between distributions  $P_n$  and P defined below, and  $\mathcal{P}(\Omega)$  is the set of probability distributions supported on a closed set  $\Omega \subseteq \mathbb{R}^m$ .

Definition 1 (Wasserstein distances). Given a lower semicontinuous function  $c: \Omega \times \Omega \to [0, \infty]$ , the optimal transport cost  $D_c(P, Q)$  between any two distributions  $P, Q \in \mathcal{P}(\Omega)$  is defined as

$$D_c(P,Q) = \min_{\pi \in \Pi(P,Q)} E_{\pi} \{c(X,Y)\},$$

where  $\Pi(P,Q)$  denotes the set of all joint distributions of the random vector (X,Y) with marginal distributions P and Q, respectively. If we specifically take  $c(x,y) = d(x,y)^2$ , where  $d(\cdot)$  is a metric, we obtain the Wasserstein distance of order 2 by setting  $W(P,Q) = \{D_c(P,Q)\}^{1/2}$ .

The quantity  $W(P_n, P)$  may be interpreted as the cheapest way to transport mass from the distribution  $P_n$  to the mass of another probability distribution P while measuring the cost of transportation from location  $x \in \Omega$  to location  $y \in \Omega$  in terms of the squared distance between x and y. In this paper we shall work with Wasserstein distances of order 2, which explains why it is natural to use  $\delta^{1/2}$  to specify the distributional uncertainty region  $\mathcal{U}_{\delta}(P_n)$  as above. Since

 $W(P_n, P_n) = 0$ , the empirical risk minimizing estimator in (1) can be seen as a special case of the formulation (2) by setting  $\delta = 0$ .

The need for selecting model parameters or making decisions using a data-driven approach which is robust to model uncertainties has sparked a rapidly growing literature on Wasserstein distributionally robust optimization, via formulations such as (2); see, for example, Gao & Kleywegt (2016), Chen et al. (2018), Gao et al. (2018), Mohajerin Esfahani & Kuhn (2018), Zhao & Guan (2018) and Blanchet & Murthy (2019) for applications in operations research, and Yang (2017, 2020) for examples specifically in stochastic control.

In principle, the min-max formulation (2) is distributionally robust in the sense that its solution guarantees a uniform performance over all probability distributions in  $\mathcal{U}_{\delta_n}(P_n)$ . Roughly speaking, for every choice of parameter or decision  $\beta$ , the min-max game type formulation in (2) introduces an adversary that chooses the most adversarial distribution from a class of distributions  $\mathcal{U}_{\delta_n}(P_n)$ . The goal of the procedure is to then choose a decision that also hedges against these adversarial perturbations, thus introducing adversarial robustness into settings where the quality of optimal solutions are sensitive to incorrect model assumptions.

Interestingly, the min-max formulation (2), which is derived from the above robustness viewpoint, has been shown to recover many machine learning estimators when applied to suitable loss functions  $\ell(\cdot)$ ; some examples include the square-root lasso and support vector machines (Blanchet et al., 2019a), the group lasso (Blanchet & Kang, 2017), adaptive regularization (Volpi et al., 2018; Blanchet et al., 2019b), among others (Shafieezadeh-Abadeh et al., 2015; Chen & Paschalidis, 2018; Duchi et al., 2020) and Gao et al., 2020. The utility of the distributionally robust formulation (2) has also been explored in adversarial training of neural networks; see, for example, Staib & Jegelka (2017) and Sinha et al. (2018).

Generic formulations such as (2) are becoming increasingly tractable; see, for example, Luo & Mehrotra (2017) and Mohajerin Esfahani & Kuhn (2018) for convex programming based approaches, and Sinha et al. (2018) and Blanchet et al. (2020) for stochastic gradient descent based iterative schemes.

Motivated by this wide range of applications, we investigate the asymptotic behaviour of the optimal value and optimal solutions of (2). In order to specifically describe the contributions, we introduce the following notation. For any positive integer n and  $\delta_n > 0$ , let

$$\Psi_n(\beta) = \sup_{P \in \mathcal{U}_{\delta_n}(P_n)} E_P \{\ell(X; \beta)\}$$

denote the distributionally robust objective function in (2). Suppose that  $\beta_*$  uniquely minimizes the population risk. According to (1) and (2), we have  $\beta_n^{\text{DRO}}$  and  $\beta_n^{\text{ERM}}$  minimize, respectively, the distributionally robust loss  $\Psi_n(\beta)$  and the empirical loss in (1). Next, let

$$\Lambda_{\delta_n}(P_n) = \left\{ \beta \in B : \beta \in \arg\min_{\beta \in B} E_P \left\{ \ell(X; \beta) \right\} \text{ for some } P \in \mathcal{U}_{\delta_n}(P_n) \right\}$$
 (3)

denote the set of choices of  $\beta \in B$  that are compatible with the distributional uncertainty region, in the sense that for every  $\beta \in \Lambda_{\delta_n}(P_n)$  there exists a probability distribution  $P \in \mathcal{U}_{\delta_n}(P_n)$  for which  $\beta$  is optimal. In other words, if  $\mathcal{U}_{\delta_n}(P_n)$  represents the set of probabilistic models which are, based on the empirical evidence, plausible representations of the underlying phenomena, then each such representation induces an optimal decision and  $\Lambda_{\delta_n}(P_n)$  encodes the set of plausible decisions. Let  $\Lambda_{\delta_n}^+(P_n)$  be the closure of  $\cap_{\epsilon>0}\Lambda_{\delta_n+\epsilon}(P_n)$ . Typically,  $\Lambda_{\delta_n}^+(P_n) = \Lambda_{\delta_n}(P_n)$ , but this is not always true as illustrated in Example 1. Asymptotically, as  $\delta_n$  decreases to zero, the distinction is negligible. However, choosing a set such as  $\Lambda_{\delta_n}^+(P_n)$  as a natural set of plausible

decisions is sensible because we guarantee that a distributionally robust solution belongs to this region. Our main result also implies that all distributionally robust solutions are asymptotically equivalent within  $o_p(n^{-1/2})$  distance from each other.

With the above notation, the key contributions of this article can be described as follows. We first establish the convergence in distribution of the triplet

$$\left[n^{1/2}\{\beta_n^{\text{ERM}} - \beta_*\}, \ n^{\bar{\gamma}/2}\{\beta_n^{\text{DRO}}(\delta_n) - \beta_*\}, \ n^{1/2}\{\Lambda_{\delta_n}^+(P_n) - \beta_*\}\right] \tag{4}$$

for a suitable  $\bar{\gamma} \in (0, 1/2]$  that depends on the rate at which the size of the distributional uncertainty,  $\delta_n$ , is decreased to zero; see Theorem 1. We identify the joint limiting distributions of the triplet (4). The third component of the triplet in (4),  $n^{1/2}\{\Lambda_{\delta_n}^+(P_n) - \beta_*\}$ , considers a suitably scaled and centred version of the choices of  $\beta \in B$  which are compatible with the respective distributional uncertainty region  $\mathcal{U}_{\delta_n}(P_n)$  in the sense described above. Therefore,  $\Lambda_{\delta_n}^+(P_n)$  is a natural choice of the confidence region. We further develop an approximation for  $\Lambda_{\delta_n}^+(P_n)$ .

Second, we utilize the limiting result of (4) to examine how the choice of the size of distributional ambiguity,  $\delta_n$ , affects the qualitative properties of the distributionally robust estimators and the induced confidence regions. Specifically, choosing  $\delta_n = \eta n^{-\gamma}$ , we characterize the behaviour of the solutions for different choices of  $\eta, \gamma \in (0, \infty)$  as  $n \to \infty$ . It emerges that the canonical,  $O(n^{-1/2})$ , rate of convergence is achieved only if  $\gamma \leq 1$ , and the limiting distribution corresponding to the distributionally robust estimator and that of the empirical risk minimizer are different only if  $\gamma \geq 1$ . Hence, to both obtain the canonical rate and tangible benefits from the distributionally robust optimization formulation we must choose  $\gamma = 1$ , which corresponds to the resulting  $\bar{\gamma}$  in (4) being equal to 1. Moreover, given any  $\alpha \in (0, 1)$ , utilizing the limiting distribution of the triplet in (4) we are able to identify a positive constant  $\eta_{\alpha} \in (0, +\infty)$  such that whenever  $\eta \geq \eta_{\alpha}$  in the choice  $\delta_n = \eta/n$ , the set  $\Lambda_{\delta_n}^+(P_n)$  is an asymptotic  $(1 - \alpha)$ -confidence region for  $\beta_*$ .

Finally, we establish the existence of an equilibrium game value. The distributionally robust optimization formulation assumes that the adversary selects a probability model after the statistician chooses a parameter. The equilibrium value of the game is attained if inf sup equals sup inf in (2), namely, if we allow the statistician to choose a parameter optimally after the adversary selects a probability model. We show in great generality that the equilibrium value of the game exists.

We end the introduction with a discussion of related statistical results. The asymptotic normality of M-estimators which minimize an empirical risk of the form  $E_{P_n}\{\ell(X;\beta)\}$  was first established in the pioneering work of Huber (1967). Subsequent asymptotic characterizations in the presence of constraints on the choices of parameter vector  $\beta$  have been developed in Dupacova & Wets (1988) and Shapiro (1989, 1991, 1993, 2000), again in the standard M-estimation setting. Our work here is different because of the presence of the adversarial perturbation to the loss represented by the inner maximization in (2).

Asymptotic normality in the related context of regularized estimators for least squares regression was established in Knight & Fu (2000). As mentioned earlier, distributionally robust estimators of the form (2) recover lasso-type estimators as particular examples (Blanchet et al., 2019a). In these cases, the inner max problem involving the adversary can be solved in closed form, resulting in the presence of regularization. However, our results can be applied even in the general context in which no closed-form solution to the inner maximization can be obtained. Therefore, our results in this paper can be seen as extensions of the results by Knight & Fu (2000), from a distributionally robust optimization perspective.

We comment that some of our results involving convergence of sets may be of interest to applications in the area of empirical likelihood (Owen, 1988, 1990, 2001). This is because  $\Lambda_{\delta_n}(P_n)$  can be characterized in terms of a function, namely, the robust Wasserstein profile function, which resembles the definition of the empirical likelihood profile function. We refer the reader to Blanchet et al. (2019a) for more discussion on the robust Wasserstein profile function and its connections to empirical likelihood. We also refer to Cisneros-Velarde et al. (2020) for additional applications, including graphical lasso, which could benefit from our results.

#### 2. Preliminaries and assumptions

# 2.1. Convergence of closed sets

For a sequence  $\{A_k : k \ge 1\}$  of closed subsets of  $\mathbb{R}^d$ , the inner and outer limits are defined, respectively, by

$$\operatorname{Li}_{n \to \infty} A_n = \{ z \in \mathbb{R}^d : \text{ there exists a sequence } (a_n)_{n \geqslant 1} \text{ with } a_n \in A_n \text{ convergent to } z \}, \text{ and } \operatorname{Ls}_{n \to \infty} A_n = \{ z \in \mathbb{R}^d : \text{ there exist positive integers } n_1 < n_2 < n_3 < \cdots \text{ and } a_k \in A_{n_k} \text{ such that the sequence } (a_k)_{k \geqslant 1} \text{ is convergent to } z \}.$$

We clearly have  $\text{Li}_{n\to\infty} A_n \subseteq \text{Ls}_{n\to\infty} A_n$ . The sequence  $\{A_n : n \geqslant 1\}$  is said to converge to a set A in the Painlevé–Kuratowski sense if

$$A = \operatorname{Li}_{n \to \infty} A_n = \operatorname{Ls}_{n \to \infty} A_n$$
.

Since  $\mathbb{R}^d$  is a locally compact Hausdorff space, the topology induced by Painlevé–Kuratowski convergence on the space of closed subsets of  $\mathbb{R}^d$  is completely metrizable, separable, and coincides with the well-known topology of closed convergence, also known as Fell topology; see Molchanov (2005, Ch. 1). The notion of convergence of sets we utilize here will be the above-defined Painlevé–Kuratowski convergence. After equipping the space of closed subsets with the Borel  $\sigma$ -algebra, we are able to define probability measures and further define the usual weak convergence of measures; see, for example, Billingsley (2013, Ch. 1).

# 2.2. Notation and assumptions

Throughout the paper, we use A > 0 to denote that a given symmetric matrix A is positive definite, and the notation  $C^{\circ}$  and cl(C) to denote the interior and closure of a subset C of Euclidean space, respectively. In the case of taking expectations with respect to the data-generating distribution  $P_*$ , we drop the subindex in the expectation operator as in  $E_{P_*}\{f(X)\} = E\{f(X)\}$ . We use  $\Rightarrow$  to denote weak convergence and  $\rightarrow$  to denote convergence in probability. We let  $\mathbb{I}(\cdot)$  be the indicator function. Let  $\|\cdot\|_p$  be the dual norm of  $\|\cdot\|_q$ , where 1/p + 1/q = 1 for  $q \in (1, \infty)$ , and  $p = \infty$  or 1 for q = 1 or  $\infty$ , respectively.

As mentioned in § 1, suppose that  $\Omega$  is a closed subset of  $\mathbb{R}^m$  and B is a closed, convex subset of  $\mathbb{R}^d$ . Assumptions 1 and 2 below are taken to be satisfied throughout the development, unless indicated otherwise.

Assumption 1. The transportation cost  $c: \Omega \times \Omega \to [0, \infty]$  is of the form  $c(u, w) = \|u - w\|_q^2$ 

Assumption 2. The function  $\ell: \Omega \times B \to \mathbb{R}$  satisfies the following properties:

- (a) The loss function  $\ell(\cdot)$  is twice continuously differentiable, and for each x,  $\ell(x, \cdot)$  is convex
- (b) Let  $h(x,\beta) = D_{\beta}\ell(x,\beta)$ , and assume there exists  $\beta_* \in B^{\circ}$  satisfying the optimality condition  $E\{h(X,\beta_*)\} = 0$ . In addition,  $E\{\|h(X,\beta_*)\|_2^2\} < \infty$ , the matrix  $C = E\{D_{\beta}h(X,\beta_*)\} > 0$ ,  $E\{D_xh(X,\beta_*)D_xh(X,\beta_*)^T\} > 0$ , and  $P\{\|D_x\ell(X,\beta_*)\|_p > 0\} > 0$ .
- (c) For every  $\beta \in \mathbb{R}^d$ ,  $\|D_{xx}\ell(\cdot;\beta)\|_p$  is uniformly continuous and bounded by a continuous function  $M(\beta)$ . Further, there exists a positive constant  $M' < \infty$  such that  $\|D_x h(x,\beta)\|_q \leq M'(1+\|x\|_q)$  for  $\beta$  in a neighbourhood of  $\beta_*$ . In addition,  $D_x h(\cdot)$  and  $D_\beta h(\cdot)$  satisfy the following locally Lipschitz continuity:

$$||D_x h(x + \Delta, \beta_* + u) - D_x h(x, \beta_*)||_q \le \kappa'(x) (||\Delta||_q + ||u||_q),$$
  
$$||D_\beta h(x + \Delta, \beta_* + u) - D_\beta h(x, \beta_*)||_q \le \bar{\kappa}(x) (||\Delta||_q + ||u||_q),$$

for  $\|\Delta\|_q + \|u\|_q \le 1$ , where  $\kappa', \bar{\kappa} : \mathbb{R}^m \to [0, \infty)$  are such that  $E[\{\kappa'(X_i)\}^2] < \infty$  and  $E[\bar{\kappa}^2(X_i)] < \infty$ .

Assumption 1 covers most of the cases in the literature described in  $\S$  1. One exception that does not immediately satisfy Assumption 1, but which can be easily adapted after a simple change of variables, is the weighted  $l_2$  norm, also known as the Mahalanobis distance, namely  $c(x,y) = (x-y)^{\mathrm{T}}A(x-y)$ , where A > 0; see Blanchet et al. (2020). The requirement that  $\ell(\cdot)$ is twice differentiable in Assumption 2(a) is useful in the analysis to identify a second-order expansion for the objective in (2), which helps quantify the impact of adversarial perturbations. The convexity of  $\ell(x,\cdot)$ , together with C being positive definite in Assumption 2(b), implies the uniqueness of  $\beta_*$ . The uniqueness of  $\beta_*$  is a standard assumption in the derivation of rates of convergence for estimators; see, for example, Huber (1967) and van der Vaart & Wellner (1996, § 3.2.2). Assumption 2(b) also allows us to rule out redundancies in the underlying source of randomness, e.g., collinearity in the setting of linear regression. The first part of Assumption 2(c) ensures that the inner maximization in (2) is finite by controlling the magnitude of the adversarial perturbations. The local Lipschitz continuity requirement in x arises with the optimal transportation analysis technique in Blanchet et al. (2019a, cf. their Assumption A6). Analogous regularity in  $\beta$  is useful in proving the confidence region limit theorem; see the discussion following Theorem 3. Limiting results which study the impact of relaxing some of these assumptions are given immediately after describing the main result in § 3.1.

### 3. Main results

# 3.1. *The main limit theorem*

In order to state our main results we introduce more definitions. Define

$$\varphi(\xi) = 4^{-1} E[\|\{D_x h(X, \beta_*)\}^T \xi\|_p^2],$$

and its convex conjugate,  $\varphi^*(\zeta) = \sup_{\xi \in \mathbb{R}^d} \left\{ \xi^{\mathsf{T}} \zeta - \varphi(\xi) \right\}$ . In addition, define

$$S(\beta) = \left[ E\left\{ \|D_x \ell(X; \beta)\|_p^2 \right\} \right]^{1/2}, \tag{5}$$

$$f_{\eta,\gamma}(x) = x \mathbb{I}(\gamma \geqslant 1) - \eta^{1/2} D_{\beta} S(\beta_*) \mathbb{I}(\gamma \leqslant 1)$$
(6)

for  $\eta \ge 0$ ,  $\gamma \ge 0$ . By Assumption 2(b), we have that  $S(\beta)$  is differentiable at  $\beta_*$ . Recalling the matrix  $C = E\{D_{\beta}h(X, \beta_*)\}$  introduced in Assumption 2(b), we define

$$\Lambda_{\delta_n}^+(P_n) = \operatorname{cl} \left\{ \cap_{\epsilon > 0} \Lambda_{\delta_n + \epsilon}(P_n) \right\},\,$$

which is the right limit of  $\Lambda_{\delta_n}(P_n)$  defined in (3). Finally, define the sets

$$\Lambda_{\eta} = \left\{ u : \varphi^*(Cu) \leqslant \eta \right\}, \qquad \Lambda_{\eta, \gamma} = \begin{cases} \Lambda_{\eta} & \text{if } \gamma = 1, \\ \mathbb{R}^d & \text{if } \gamma < 1, \\ \{0\} & \text{if } \gamma > 1. \end{cases}$$
 (7)

We now state our main result.

THEOREM 1. Suppose that Assumptions 1 and 2 are satisfied with  $q \in (1, \infty)$ ,  $\Omega = \mathbb{R}^m$  and  $E(\|X\|_2^2) < \infty$ . If  $H \sim \mathcal{N}[0, \text{cov}\{h(X, \beta_*)\}]$  and  $\delta_n = n^{-\gamma}\eta$  for some  $\gamma, \eta \in (0, \infty)$ , then we have the following joint convergence in distribution:

$$\left[ n^{1/2} \left( \beta_n^{\text{ERM}} - \beta_* \right), \ n^{\bar{\gamma}/2} \left\{ \beta_n^{\text{DRO}}(\delta_n) - \beta_* \right\}, \ n^{1/2} \left\{ \Lambda_{\delta_n}^+(P_n) - \beta_* \right\} \right] 
\Rightarrow \left[ C^{-1}H, \ C^{-1} f_{n,\gamma}(H), \ \Lambda_{n,\gamma} + C^{-1}H \right],$$

where  $\bar{\gamma} = \min{\{\gamma, 1\}}$  and  $\Lambda_{\eta, \gamma}$  is defined as in (7).

The proof of Theorem 1 is presented in § 5.3. For q=1 or  $\infty$ , which corresponds to  $p=\infty$  or 1,  $S(\beta)$  may not be differentiable at  $\beta_*$ , in which case the limited distribution presents a discontinuity which makes it difficult to use in practice. Hence, we prefer not to cover this here. Theorem 1 can be used as a powerful conceptual tool. For example, let us examine how a sensible choice for the parameter  $\delta_n$  can be obtained as an application of Theorem 1 by considering the following cases.

Case 1,  $\gamma > 1$ : If  $n\delta_n \to 0$ , corresponding to the case  $\gamma > 1$ , we have  $f_{0,\gamma}(H) = H$  from the definition of the parametric family in (6). Therefore, from Theorem 1,

$$\left[n^{1/2}(\beta_n^{\text{ERM}} - \beta_*), \ n^{\bar{\gamma}/2} \{\beta_n^{\text{DRO}}(\delta_n) - \beta_*\}, \ n^{1/2} \{\Lambda_{\delta_n}^+(P_n) - \beta_*\}\right] 
\Rightarrow \left[C^{-1}H, C^{-1}H, \{C^{-1}H\}\right],$$

which implies that the influence of the robustification vanishes in the limit when  $\delta_n = o(n^{-1})$ . Case 2,  $\gamma < 1$ : If  $n\delta_n \to \infty$ , corresponding to the case  $\gamma < 1$ , the rate of convergence for the distributionally robust estimator is slower than the canonical  $O(n^{-1/2})$  rate:

$$\beta_n^{\text{DRO}}(\delta_n) = \beta_* - \eta^{1/2} n^{-\gamma/2} C^{-1} D_{\beta} S(\beta_*) + o_p(n^{-\gamma/2}), \tag{8}$$

where  $n^{\gamma/2}o_p(n^{-\gamma/2}) \to 0$  in probability as  $n \to \infty$ . The relationship (8) reveals an uninteresting limit,  $n^{1/2}\{\Lambda_{\delta_n}^+(P_n) - \beta_*\} \Rightarrow \mathbb{R}^d$ , exposing a slower than  $O(n^{-1/2})$  rate of convergence  $\Lambda_{\delta_n}^+(P_n)$ . In fact, (8) indicates that an  $O(n^{-\gamma/2})$  scaling will result in a nondegenerate limit.

Case 3,  $\gamma = 1$ : when  $\delta_n = \eta/n$ , we have that all components in the triplet in Theorem 1 have nontrivial limits.

Theorem 2 provides a geometric insight relating  $\beta_n^{\mathrm{DRO}}(\delta_n)$ ,  $\beta_n^{\mathrm{ERM}}$  and  $\Lambda_{\delta_n}^+(P_n)$ , which justifies a picture describing  $\Lambda_{\delta_n}^+(P_n)$  as a set containing both  $\beta_n^{\mathrm{DRO}}(\delta_n)$  and  $\beta_n^{\mathrm{ERM}}$ . The observation that

 $\beta_n^{\mathrm{ERM}} \in \Lambda_{\delta_n}(P_n)$  is immediate because  $\Lambda_{\delta}(P_n)$  is increasing in  $\delta$ , so  $\beta_n^{\mathrm{ERM}} \in \Lambda_0(P_n) \subset \Lambda_{\delta_n}^+(P_n)$ . On the other hand, the observation that  $\beta_n^{\mathrm{DRO}}(\delta_n) \in \Lambda_{\delta_n}^+(P_n)$  is nontrivial and it relies on the exchangeability of inf and sup in Theorem 2. An appropriate choice of  $\eta$  which results in the set  $\Lambda_{\delta_n}^+(P_n)$  also possessing desirable coverage for  $\beta_*$  is prescribed in § 3.2.

THEOREM 2. Suppose that Assumption 1 is enforced. We further assume that the loss function  $\ell(\cdot)$  is continuous and nonnegative,  $\ell(x, \cdot)$  is convex for each x, and  $E_{P_*}\{\ell(X, \beta)\}$  has a unique optimizer  $\beta_* \in \mathcal{B}^{\circ}$ . Then, for any  $\delta > 0$ ,

$$\inf_{\beta \in B} \sup_{P \in \mathcal{U}_{\delta}(P_n)} E_P \{\ell(X; \beta)\} = \sup_{P \in \mathcal{U}_{\delta}(P_n)} \inf_{\beta \in B} E_P \{\ell(X; \beta)\}, \tag{9}$$

and there exists a distributionally robust estimator choice  $\beta_n^{DRO}(\delta) \in \Lambda_{\delta}^+(P_n)$ .

The proof of Theorem 2 is presented in the Supplementary Material. Example 1 demonstrates that the set of minimizers of the distributionally robust formulation (2) is not necessarily unique, and that the set  $\Lambda_{\delta}(P_n)$  may not contain distributionally robust solutions. Theorem 2 indicates that the right limit  $\Lambda_{\delta}^+(P_n)$  contains a distributionally robust solution. Theorem 1 implies that the minimizers of (2) differ by at most  $o_p(n^{-1/2})$  in magnitude, which indicates that they are asymptotically equivalent and the inclusion of one solution of (2) in  $\Lambda_{\delta}^+(P_n)$  is sufficient for the scaling considered.

Example 1. Let the loss function be

$$\ell(x,\beta) = f(\beta) + \{x^2 - \log(x^2 + 1)\}f(\beta - 4),$$

where  $f(\beta)=3\beta^2/4-1/8\beta^4+3/8$  for  $\beta\in[-1,1]$ , and  $f(\beta)=|\beta|$  otherwise. We have that  $\ell(x,\beta)$  is twice differentiable and convex, satisfying Assumptions 1 and 2. Then, if the empirical measure  $P_n$  is a Dirac measure centred at zero with n=1, and  $\delta=1$ , we have the distributionally robust estimators  $\beta_n^{\mathrm{DRO}}(\delta)\in[1,3]$ . Further,  $[1,3]\subset\Lambda_\delta^+(P_n)$ , but  $[1,3]\cap\Lambda_\delta(P_n)=\varnothing$ .

Next, we turn to the relationship between  $\beta_n^{\rm ERM}$  and  $\beta_n^{\rm DRO}(\delta_n)$  when  $\delta_n = \eta/n$ . From the first two terms in the triplet, we have

$$\beta_n^{\text{DRO}}(\delta_n) = \beta_n^{\text{ERM}} - \eta^{1/2} C^{-1} D_{\beta} S(\beta_*) n^{-1/2} + o_p (n^{-1/2})$$

$$= \beta_n^{\text{ERM}} - \delta_n^{1/2} C^{-1} D_{\beta} S(\beta_n^{\text{ERM}}) + o_p (\delta_n). \tag{10}$$

The right-hand side of (10) points to the canonical  $O(n^{-1/2})$  rate of convergence of the Wasserstein distributionally robust estimator, and it can readily be used to construct confidence regions, as we shall explain in § 3.2.

The relation (10) also exposes the presence of an asymptotic bias term, namely,  $S(\beta) = [E\{\|D_x\ell(X;\beta)\|_p^2\}]^{1/2}$ , which points towards selection of optimizers possessing reduced sensitivity with respect to perturbations in data. A precise mathematical statement of this sensitivity-reduction property is given in Corollary 1, and its proof is in the Supplementary Material.

COROLLARY 1. Suppose that Assumptions 1 and 2 are in force, and consider

$$\bar{\beta}_n^{\text{DRO}} \in \arg\min_{\beta \in B} \left( E_{P_n} \left\{ \ell(X; \beta) \right\} + n^{-1/2} \left[ \eta E_{P_n} \left\{ \| D_X \ell(X; \beta) \|_p^2 \right\} \right]^{1/2} \right). \tag{11}$$

Then, if  $\delta_n = \eta/n$ ,  $\beta_n^{DRO}(\delta_n) = \bar{\beta}_n^{DRO} + o_p(n^{-1/2})$ .

While the formulation on the right-hand side of (11) is conceptually appealing, it may not be desirable from an optimization point of view due to the potentially nonconvex nature of the objective involved. On the other hand, under Assumption 2, the distributionally robust objective  $\Psi_n(\beta)$  is convex; see, for example, the reasoning in Blanchet et al. (2020, Theorem 2a), while also enjoying the sensitivity-reduction property of the formulation in (11).

A result of a similar type to Corollary 1 is given in Gao et al. (2020), but the focus there is on the objective function of (2) being approximated by a suitable regularization. The difference between this type of result and Corollary 1 is that our focus is on the asymptotic equivalence of the actual optimizers. Behind a result such as Corollary 1, it is key to have a more nuanced approximation which precisely characterizes the second-order term of size  $O(\delta_n)$ ; see the Supplementary Material.

We conclude this section with results which examine the effects of relaxing some assumptions made in the statement of Theorem 1. Proposition 1 asserts that convergence of the natural confidence region  $\Lambda_{\delta_n}^+(P_n)$ , as identified in Theorem 1 holds even if the support of the probability distributions in the uncertainty region  $\mathcal{U}_{\delta_n}(P_n)$  is constrained to be a strict subset  $\Omega$  of  $\mathbb{R}^d$ . For this purpose, we introduce the following notation. For any set  $C \in \mathbb{R}^m$ , let  $C^{\epsilon} = \{x \in C : B_{\epsilon}(x) \subset C\}$ , where  $B_{\epsilon}(x)$  is the neighbourhood around x defined as  $B_{\epsilon}(x) = \{y : \|y - x\|_2 \le \epsilon\}$ . Thus, for any probability measure P, we have  $\lim_{\epsilon \to 0} P(C^{\epsilon}) = P(C^{\circ})$ .

PROPOSITION 1. Suppose that Assumptions 1 and 2 are satisfied with  $q \in [1, \infty]$  and  $E(\|X\|_2^2) < \infty$ . In addition, suppose that the data-generating measure  $P_*$  satisfies  $P_*(\Omega^\circ) = 1$ . If we take  $H \sim \mathcal{N}(0, \text{cov}\{h(X, \beta_*)\})$  and  $\delta_n = n^{-\gamma} \eta$  for some  $\gamma, \eta \in (0, \infty)$ , then the following convergence holds as  $n \to \infty$ :  $n^{1/2}\{\Lambda_{\delta_n}(P_n) - \beta_*\} \Rightarrow \Lambda_{\eta,\gamma} + C^{-1}H$ .

The steps involved in proving Proposition 1 are presented in § 5.1. A discussion on the validity of a central limit theorem for the estimator  $\beta_n^{\text{DRO}}$  in the presence of constraints restricting transportation within the support set  $\Omega$  is presented in § 6.

In the case where the unique minimizer  $\beta_*$  may not necessarily lie in the interior of the set B, as opposed to the requirement in Assumption 2(b), one may obtain the extension in Proposition 2 as the limiting result for the estimator  $\beta_n^{\text{DRO}}(\delta_n)$ . As in the previous results, we take  $h(x, \beta) = D_B \ell(x; \beta)$ . The proofs of all the subsequent propositions are given in the Supplementary Material.

PROPOSITION 2. Suppose that Assumptions 1, 2(a) and 2(c) are satisfied, and that  $\beta_*$  is the unique minimizer of  $\min_{\beta \in B} E\{\ell(X,\beta)\}$ . Suppose that the set B is compact, and there exist  $\varepsilon > 0$  and twice continuously differentiable functions  $g_i(\beta)$  such that

$$B \cap B_{\varepsilon}(\beta_*) = \{ \beta \in B_{\varepsilon}(\beta_*) : g_i(\beta) = 0, i \in I, g_j(\beta) \leq 0, j \in J \},\$$

where I,J are finite index sets and  $g_i(\beta_*) = 0$  for all  $i \in J$ . With this identification of the set B, suppose that the following so-called Mangasarian–Fromovitz constraint qualification is satisfied at  $\beta_*$ : the gradient vectors  $\{Dg_i(\beta_*): i \in I\}$  are linearly independent, and there exists a vector w such that  $w^TDg_i(\beta_*) = 0$  for all  $i \in I$  and  $w^TDg_j(\beta_*) < 0$  for all  $j \in J$ .

Suppose that  $\Lambda_0$  is the set of Lagrange multipliers satisfying the first-order optimality conditions and the following second-order sufficient conditions:  $\lambda \in \Lambda_0$  if and only if  $D_{\beta}L(\beta_*,\lambda) = 0$ ,  $\lambda_i \geq 0$  for  $i \in J$  and  $\max_{\lambda \in \Lambda_0} w^T D_{\beta\beta}L(\beta_*,\lambda)w > 0$  for all  $w \in C$ , where  $L(\beta,\lambda) = E\{\ell(X,\beta)\} + \sum_{i \in I \cup J} \lambda_i g_i(\beta)$  is the Lagrangian function associated with the minimization  $\min_{\beta \in B} E\{\ell(X,\beta)\}$ , and

$$C = [w : w^{\mathsf{T}} D g_i(\beta_*) = 0, i \in I, \ w^{\mathsf{T}} D g_j(\beta_*) \leqslant 0, j \in J, \ w^{\mathsf{T}} E \{h(X, \beta_*)\} \leqslant 0]$$

is the nonempty cone of critical directions. In addition, suppose that  $\omega(\xi)$  is the unique minimizer of  $\min_{u \in C} \{ \xi^T u + 2^{-1} q(u) \}$ , where  $q(u) = \max \{ u^T D_{\beta\beta} L(\beta_*, \lambda) u : \lambda \in \Lambda_0 \}$ . Then, if  $\delta_n = \eta n^{-1}$  for  $\eta \in (0, \infty)$ ,  $E\{\|h(X, \beta_*)\|_2^2\} < \infty$  and  $E\{D_{\beta}h(X, \beta_*)\} > 0$ , we have the following convergence as  $n \to \infty$ :

$$n^{1/2} \left\{ \beta_n^{\text{DRO}}(\delta_n) - \beta_* \right\} \Rightarrow \omega \left\{ -H + \eta^{1/2} D_{\beta} S(\beta_*) \right\},$$

where  $H \sim \mathcal{N}[0, \text{cov}\{h(X, \beta_*)\}]$ .

The Mangasarian–Fromovitz constraint qualification conditions and the necessary and sufficient conditions in the statement of Proposition 2 are standard in the literature if the optimal  $\beta_*$  lies on the boundary of the set B; see, for example, Shapiro (1989). Please refer to the discussion following Theorem 3.1 in Shapiro (1989) for sufficient conditions under which  $\omega(\xi)$  is unique.

Proposition 3 extends the sensitivity reduction property in Corollary 1 to settings where the minimizer for  $\min_{\beta \in B} E_{P_*} \{ \ell(X; \beta) \}$  is not unique.

PROPOSITION 3. Suppose that Assumptions 1, 2(a) and 2(c) are satisfied, the set B is compact and the choice of the radii  $(\delta_n : n \ge 1)$  is such that  $n\delta_n \to \eta \in (0, \infty)$ . Let the set  $B_*$  be arg  $\min_{\beta \in B} E_{P_*} \{\ell(X; \beta)\}$ . Then, the distributionally robust optimization objective  $\Psi_n(\beta)$  satisfies

$$n^{1/2} \left[ \Psi_n(\beta) - E\{\ell(X;\beta)\} \right] \Rightarrow Z(\beta) + \eta^{1/2} S(\beta),$$

where  $Z(\cdot)$  is a zero-mean Gaussian process with covariance function  $\operatorname{cov}\{Z(\beta_1), Z(\beta_2)\} = \operatorname{cov}\{\ell(X, \beta_1), \ell(X, \beta_2)\}$ . The above weak convergence holds, as  $n \to \infty$ , on the space of continuous functions equipped with the uniform topology on compact sets. Consequently, if  $\operatorname{arg\,min}_{\beta \in \mathcal{B}_*}\{Z(\beta) + \eta^{1/2}S(\beta)\}$  is singleton with probability one, we have, as  $n \to \infty$ ,

$$\beta_n^{\mathrm{DRO}}(\delta_n) \Rightarrow \arg\min_{\beta \in B_*} \{ Z(\beta) + \eta^{1/2} S(\beta) \}.$$

# 3.2. Construction of Wasserstein distributionally robust confidence regions

As mentioned in § 1, for suitably chosen  $\delta_n$ , the set  $\Lambda_{\delta_n}^+(P_n)$  represents a natural confidence region. In particular,  $\Lambda_{\delta_n}^+(P_n)$  possesses an asymptotically desired coverage, say at level at least  $1-\alpha$ , if and only if

$$1 - \alpha \leqslant \lim_{n \to \infty} \operatorname{pr} \left\{ \beta_* \in \Lambda_{\delta_n}^+(P_n) \right\} = \operatorname{pr}[-C^{-1}H \in \{u : \varphi^*(Cu) \leqslant \eta\}],$$

or, equivalently, if  $\eta \geqslant \eta_{\alpha}$ , where  $\eta_{\alpha}$  is the  $(1-\alpha)$ -quantile of the random variable  $\varphi^*(H)$ . Recall the earlier geometric insight describing  $\Lambda_{\delta_n}^+(P_n)$  as a set containing both  $\beta_n^{\mathrm{DRO}}(\delta_n)$  and  $\beta_n^{\mathrm{ERM}}$ , as a consequence of Theorem 2. Following this, if we let  $\eta \geqslant \eta_{\alpha}$ , we then have

$$\lim_{n\to\infty} \operatorname{pr}\left\{\beta_* \in \Lambda_{\delta_n}^+(P_n), \ \beta_n^{\operatorname{DRO}} \in \Lambda_{\delta_n}^+(P_n), \ \beta_n^{\operatorname{ERM}} \in \Lambda_{\delta_n}^+(P_n)\right\} = \lim_{n\to\infty} \operatorname{pr}\left\{\beta_* \in \Lambda_{\delta_n}^+(P_n)\right\}$$
$$\geqslant 1 - \alpha,$$

which presents the picture of  $\Lambda_{\delta_n}^+(P_n)$  as a confidence region simultaneously containing  $\beta_*$ ,  $\beta_n^{\text{ERM}}$  and  $\beta_n^{\text{DRO}}(\delta_n)$  with a desired level of confidence.

The function  $\varphi^*(H)$  can be computed in closed form in some settings. But, typically, computing  $\varphi^*(\cdot)$  may be challenging. We now describe how to obtain a consistent estimator for  $\eta_\alpha$ . Define the empirical version of  $\varphi(\xi)$ , namely

$$\varphi_n(\xi) = \frac{1}{4} E_{P_n} \left[ \left\| \{ D_x h(X, \beta_*) \}^{\mathsf{T}} \xi \right\|_p^2 \right] = \frac{1}{4n} \sum_{i=1}^n \left\| \{ D_x h(X, \beta_*) \}^{\mathsf{T}} \xi \right\|_p^2,$$

and the associated empirical convex conjugate,  $\varphi_n^*(\zeta) = \sup_{\xi \in \mathbb{R}^d} \{ \xi^T \zeta - \varphi_n(\xi) \}$ . Proposition 4 provides a basis for computing a consistent estimator for  $\eta_\alpha$ .

PROPOSITION 4. Let  $\Xi_n$  be any consistent estimator of  $\operatorname{cov}\{h(X,\beta)\}$ , and write  $\bar{\Xi}_n$  for any factorization of  $\Xi_n$  such that  $\bar{\Xi}_n \bar{\Xi}_n^{\mathsf{T}} = \Xi_n$ . Let Z be a d-dimensional standard Gaussian random vector independent of the sequence  $(X_n : n \ge 1)$ . Then, (i) the distribution of  $\varphi^*(Z)$  is continuous, (ii)  $\varphi^*_n(\cdot) \Rightarrow \varphi^*(\cdot)$  as  $n \to \infty$  uniformly on compact sets, and (iii)  $\varphi^*_n(\bar{\Xi}_n Z) \Rightarrow \varphi^*(H)$ .

Given the collection of samples  $\{X_i\}_{i=1}^n$ , we can generate independent and identically distributed copies of Z and use Monte Carlo to estimate the  $(1-\alpha)$ -quantile,  $\eta_\alpha$  (n), of  $\varphi_n^*(\bar{\Xi}_n Z)$ . The previous proposition implies that  $\eta_\alpha$   $(n) = \eta_\alpha + o_p(1)$  as  $n \to \infty$ . This is sufficient to obtain an implementable expression for  $\beta_n^{\mathrm{DRO}}\{\eta_\alpha(n)/n\}$  which is asymptotically equivalent to (10), as it differs only by an error of maginutude  $o_p(n^{-1/2})$ .

Next, we provide rigorous support for the approximation  $\Lambda_{\delta_n}^+(P_n) \approx \beta_n^{\text{ERM}} + n^{-1/2}\Lambda_{\eta}$ , which can be used to approximate  $\Lambda_{\delta_n}^+(P_n)$ , providing we can estimate  $\Lambda_{\eta}$ .

COROLLARY 2. Under the assumptions of Theorem 1, and with  $\gamma = 1$ ,

$$n^{1/2}\left\{\Lambda_{\delta_n}^+(P_n) - \beta_n^{\mathrm{ERM}}\right\} \Rightarrow \Lambda_{\eta}.$$

Moreover, if  $\eta(n) = \eta + o(1)$  and  $C_n \to C$ , then  $\Lambda_{\eta(n)}^n = \{u : \varphi_n^*(C_n u) \leqslant \eta(n)\} \to \Lambda_{\eta}$ .

*Proof of Corollary* 2. Following directly from Theorem 1 and an application of the continuous mapping theorem,

$$n^{1/2} \left\{ \Lambda_{\delta_n}^+(P_n) - \beta_n^{\text{ERM}} \right\} = n^{1/2} \left\{ \Lambda_{\delta_n}^+(P_n) - \beta_* \right\} - n^{1/2} \left\{ \beta_n^{\text{ERM}} - \beta_* \right\}$$
$$\Rightarrow \Lambda_n + C^{-1}H - C^{-1}H.$$

The second part of the result follows from the regularity results in Proposition 4.

The next result, as we shall explain, allows us to obtain computationally efficient approximations of the set  $\Lambda_{\eta}$ . A completely analogous result can be used to estimate  $\Lambda_{\eta(n)}^n$ , simply replacing  $\varphi^*(\cdot)$ ,  $\varphi(\cdot)$  and C by  $\varphi_n^*(\cdot)$ ,  $\varphi_n(\cdot)$  and  $C_n$ .

PROPOSITION 5. The support function of the convex set  $\Lambda_{\eta} = \{u : \varphi^*(Cu) \leq \eta\}$  is  $h_{\Lambda_{\eta}}(v) = 2\{\eta\varphi(C^{-1}v)\}^{1/2}$ , where the support function of a convex set A is defined as  $h_A(x) = \sup\{x \cdot a : a \in A\}$ .

Remark 1. Proposition 5 can be used to obtain a tight envelope of the set  $\Lambda_{\eta}$  by evaluating an intersection of hyperplanes that enclose  $\Lambda_{\eta}$ . Recall from the definition of a support function

that  $\Lambda_{\eta} = \bigcap_{u} \{v : u \cdot v \leqslant h_{\Lambda_{\eta}}(u)\}$ . Therefore, for any  $u_1, \ldots, u_m$ , we have that  $\Lambda_{\eta}$  is contained in  $\bigcap_{u_1, \ldots, u_m} \{v : u_i \cdot v \leqslant h_{\Lambda_{\eta}}(u_i)\}$  and that  $\Lambda_{\eta(n)}^n$  is contained in  $\bigcap_{u_1, \ldots, u_m} \{v : u_i \cdot v \leqslant h_{\Lambda_{\eta(n)}^n}(u_i)\}$ .

# 4. NUMERICAL EXAMPLES: GEOMETRY AND COVERAGE PROBABILITIES

# 4.1. Distributionally robust linear regression

We first offer a brief introduction to the distributionally robust version of the linear regression problem considered in Blanchet et al. (2019a). Specifically, the data is generated by  $Y = \beta_*^T X + \epsilon$ , where  $X \in \mathbb{R}^d$  and  $\epsilon$  are independent,  $C = E(XX^T)$  and  $\epsilon \sim \mathcal{N}(0, \sigma^2)$ . We consider square loss  $\ell(x, y; \beta) = 1/2(y - \beta^T x)^2$  and take the cost function  $c : \mathbb{R}^{d+1} \times \mathbb{R}^{d+1} \to [0, \infty]$  to be

$$c\{(x,y),(u,v)\} = \begin{cases} \|x - u\|_q^2 & \text{if } y = v, \\ \infty & \text{otherwise.} \end{cases}$$
 (12)

Then, from Blanchet et al. (2019a, Theorem 1), we have

$$\min_{\beta \in \mathbb{R}^d} \sup_{P:D_c(P,P_n) \leq \delta_n} E_P \left[ \ell(X,Y;\beta) \right] = \frac{1}{2} \min_{\beta \in \mathbb{R}^d} \left[ E_{P_n} \left\{ (Y - \beta^T X)^2 \right\}^{1/2} + \delta_n^{1/2} \|\beta\|_p \right]^2, \quad (13)$$

where p satisfies 1/p + 1/q = 1. Following Corollary 2, an approximate confidence region is  $\Lambda_{\delta_n}^+(P_n) \approx n^{-1/2}\Lambda_{\eta_\alpha} + \beta_n^{\rm ERM}$ , where  $\Lambda_{\eta_\alpha} = \{\theta : \varphi^*(C\theta) \leqslant \eta_\alpha\}, \ \varphi(\xi) = 4^{-1}E\{\|e\xi - (\xi^T X)\beta_*\|_p^2\}$ , the constant  $\eta_\alpha$  is such that  $\Pr\{\varphi^*(H) \leqslant 1 - \alpha\} = \eta_\alpha$  for  $H \sim \mathcal{N}(0, C\sigma^2)$  and  $\delta_n = \eta_\alpha/n$ . By performing a change of variables via linear transformation in the analysis of the case  $c(x,y) = \|x - y\|_2^2$ , Theorem 1 can be directly adapted to the choice c(x,y) being a Mahalanobis metric as in

$$c(x,y) = (x-y)^{T} A(x-y),$$
 (14)

for some matrix A > 0. The respective  $\Lambda_{\eta_{\alpha}} = \{\theta : \varphi^*(C\theta) \leq \eta_{\alpha}\}$  is computed in terms of  $\varphi(\xi) = 4^{-1}E\{\|\xi^T D_x h(X, \beta_*) A^{-1/2}\|_2^2\}$ . For the choice  $c(x, y) = (x - y)^T A(x - y)$ , the relationship between distributionally robust and regularized estimators, as in (13), is

$$\min_{\beta \in \mathbb{R}^d} \sup_{P:D_c(P,P_n) \leqslant \delta_n} E_P \left\{ l(X,Y;\beta) \right\} = \frac{1}{2} \min_{\beta \in \mathbb{R}^d} \left[ E_{P_n} \left\{ (Y - \beta^\mathsf{T} X)^2 \right\}^{1/2} + \delta_n^{-1/2} \left\| A^{-1/2} \beta \right\|_2 \right]^2.$$

See Blanchet et al. (2019b) for an account of improved out-of-sample performance resulting from Mahalanobis cost choices.

# 4.2. Shape of confidence regions

The goal of this section is to provide some numerical implementations to gain intuition about the geometry of the set  $\Lambda_{\eta}$  for different transportation cost choices. We use the empirical set  $\Lambda_{\eta_{\alpha}}^{n} = \{\theta : \varphi_{n}^{*}(C_{n}\theta) \leq n^{-1/2}\tilde{\eta}_{\alpha}\}$ , to approximate the desired confidence region as in Corollary 2. In the above expression,  $\varphi_{n}(\xi) = 4^{-1}E_{P_{n}}\{\|e\xi - (\xi^{T}X)\beta_{n}^{\text{ERM}}\|_{p}^{2}\}$ ,  $\eta_{\alpha}(n)$  is such that  $\operatorname{pr}(\tilde{\varphi}_{n}^{*}(H) \leq 1 - \alpha) = \eta_{\alpha}(n)$  for  $H \sim \mathcal{N}(0, C_{n}\sigma_{n}^{2})$ ,  $C_{n} = E_{P_{n}}[XX^{T}]$ , and  $\sigma_{n}^{2} = E_{P_{n}}[\{(Y - (\beta_{n}^{\text{ERM}})^{T}X)\}^{2}]$ .

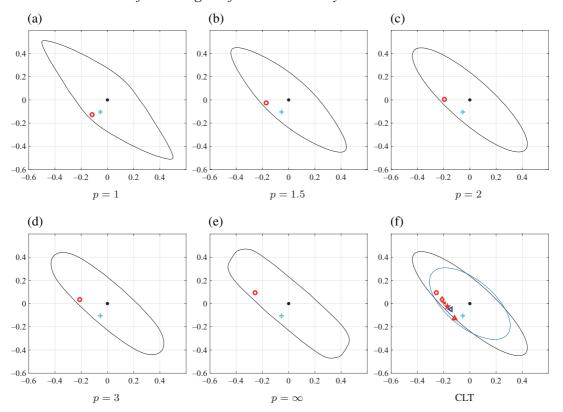


Fig. 1. Confidence regions for different norm choices and the central limit theorem based confidence region plotted together with the respective  $\beta_n^{\rm DRO}$  (red circle) estimators and  $\beta_n^{\rm ERM}$  (black dot). The blue cross represents the true value. In panel (f) the following estimators are included:  $\ell_1$  DRO (red triangle),  $\ell_3/2$  DRO (red star),  $\ell_2$  DRO (red cross) and  $\ell_3$  DRO (red diamond).

In the following numerical experiments, the data is sampled from a linear regression model with parameters  $\sigma^2 = 1$ ,  $\beta_* = [0.5, 0.1]^T$ , n = 100 and

$$X \sim \mathcal{N}\left(0, \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}\right),\tag{15}$$

with  $\rho=0.7$ . In Figs. 1(a)–1(e) we show the 95% confidence region corresponding to the choices  $p=1,3/2,2,3,\infty,q=\infty,3,2,3/2$  by means of support functions defined in Proposition 5. In addition, a confidence region for  $\beta_*$  resulting from the asymptotic normality of the least-squares estimator,  $n^{1/2}(\beta_n^{\rm ERM}-\beta^*) \Rightarrow \mathcal{N}(0,C^{-1}\sigma^2)$ , is

$$\Lambda_{\text{CLT}}(P_n) = n^{-1/2} \{ \theta : \theta^{\mathsf{T}} C \theta / \sigma^2 \leqslant \chi_{1-\alpha}^2(d) \} + \beta_n^{\text{ERM}},$$

where  $\chi^2_{1-\alpha}(d)$  is the  $(1-\alpha)$ -quantile of the chi-squared distribution with d degrees of freedom. One can select the matrix A in the Mahalanobis metric (14) such that the resulting confidence region coincides with  $\Lambda_{\text{CLT}}(P_n)$ . Namely, A is chosen by solving the equation

$$E\left[\left\{e\xi - \left(\xi^{\mathsf{T}}X\right)\beta_{*}\right\}A^{-1}\left\{e\xi - \left(\xi^{\mathsf{T}}X\right)\beta_{*}\right\}^{\mathsf{T}}\right] = C\sigma^{2}.$$

				· ·	
$oldsymbol{eta}_0$	ho	$\ell_2$ -confiden Coverage for $\beta_n^{\mathrm{DRO}}$	U	CLT confide Coverage for $\beta_n^{DRO}$	nce region Coverage for $\beta_*$
$\begin{bmatrix} 0.5 \\ 0.5 \end{bmatrix}$	0.95	100.0%	94.5%	99.4%	94.6%
	$0 \\ -0.95$	100.0% 100.0%	94.0% 94.8%	97.1% 75.8%	93.5% 94.4%
$\begin{bmatrix} 1.0 \\ 0.0 \end{bmatrix}$	0.95	100.0%	94.6%	93.7%	95.4%
	0	100.0%	94.6%	100%	94.1%
	-0.95	100.0%	95.3%	91.2%	94.9%

Table 1. Coverage probability

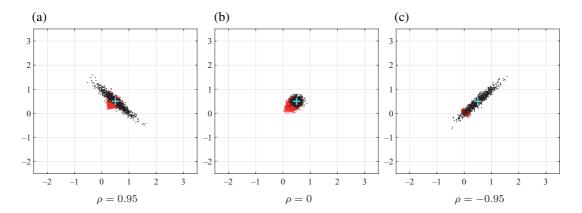


Fig. 2. Scatter plots of  $\beta_n^{\text{ERM}}$  (black circles) and  $\beta_n^{\text{DRO}}$  (red circles) for  $\beta_0 = [0.5, 0.5]^T$ . The blue cross represents the true value.

Figure 1(f) gives the confidence region for the choice p=2 and  $\Lambda_{\text{CLT}}(P_n)$  superimposed with various distributionally robust minimizers along with the empirical risk minimizer. It is evident from the figures that p=1 gives a diamond shape, p=2 gives an elliptical shape and  $p=\infty$  gives a rectangular shape. Furthermore, we see that the distributionally robust optimization solutions all reside in their respective confidence regions, but may lie outside of the confidence regions of other norms.

We find that the induced confidence regions constructed by the Wasserstein distributionally robust optimization formulations are somewhat similar across the various  $l_p$  norms, but they are all different to the standard central limit theorem based confidence region. The Mahalanobis cost can be calibrated to exactly match the standard central limit theorem confidence region.

# 4.3. Coverage probabilities and distributionally robust optimization solutions

We now test the scenario in which the covariates are highly correlated. Specifically, the data is sampled from a linear regression model with parameters  $\sigma^2 = 1$ , n = 100, p = 2. The random vector X is taken to be distributed in (15), considering three different values for  $\rho$ :  $\rho = 0.95, 0, -0.95$ . We consider the following two cases for the underlying parameter  $\beta_*$ :  $\beta_* = [0.5, 0.5]^T$  and  $\beta_* = [1, 0]^T$ . In Table 1 we report the coverage probabilities of the underlying  $\beta_*$  and  $\beta_n^{\mathrm{DRO}}(\delta_n)$  in both the  $\ell_2$  confidence region and the central limit theorem based confidence regions. Specifically, we report the following four probabilities:  $\mathrm{pr}\{\beta_n^{\mathrm{DRO}} \in \Lambda_{\delta_n}^+(P_n)\}$ ,  $\mathrm{pr}\{\beta_* \in \Lambda_{\mathrm{CLT}}(P_n)\}$  and  $\mathrm{pr}\{\beta_* \in \Lambda_{\mathrm{CLT}}(P_n)\}$ .

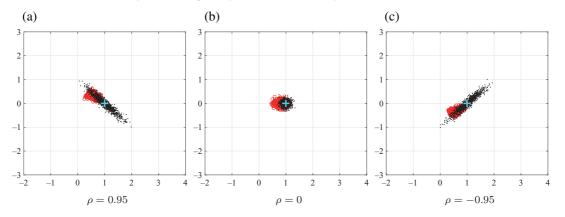


Fig. 3. Scatter plots of  $\beta_n^{\text{ERM}}$  (black circles) and  $\beta_n^{\text{DRO}}$  (red circles) for  $\beta_0 = [1.0, 0.0]^T$ . The blue cross represents the true value.

Figures 2 and 3 show scatterplots of the estimators  $\beta_n^{\text{ERM}}$  and  $\beta_n^{\text{DRO}}$  when the underlying  $\beta_*$  takes the values  $[0.5, 0.5]^{\text{T}}$  and  $[1, 0]^{\text{T}}$ , respectively. In the near-collinearity cases where  $\rho = 0.95$  or -0.95, the lower spreads for the distributionally robust estimators reveal their better performance over the empirical risk-minimizing solutions. The utility of the proposed  $\ell_2$  confidence region emerges in light of the better performance of the distributionally robust estimator  $\beta_n^{\text{DRO}}$  and its aforementioned lack of membership in  $\Lambda_{\text{CLT}}(P_n)$ .

We sample 1000 datasets and report the coverage probabilities in Table 1. We observe that for  $\beta_*$ , both the  $\ell_2$  confidence region and the central limit theorem based confidence region achieve the target 95% coverage. Furthermore, the coverage for the distributionally robust estimator of the  $\ell_2$  confidence region is 100%, which validates our theory. However, when  $\rho = -0.95$  and  $\beta_* = [0.5, 0.5]^T$ , the coverage for the distributionally robust estimator in the central limit theorem based confidence region is only 75.8%. In this example, the asymptotic results developed indicate that this coverage probability converges to zero, when n tends to infinity.

#### 5. Proofs of Main Results

#### 5.1. Preliminaries

Theorem 1 is obtained by considering appropriate level sets involving auxiliary functionals. Following Blanchet et al. (2019a), we define the robust Wasserstein profile function, associated with the estimation of  $\beta_*$  by solving  $E_{P_n}\{D_{\beta}h(X,\beta)\}=0$ , as follows:

$$R_n(\beta) = \inf_{P \in \mathcal{P}(\Omega)} \left[ D_c(P, P_n) : \beta \in \arg \min_{\beta \in B} E_P \left\{ \ell(X; \beta) \right\} \right].$$

This definition, as noted in Blanchet et al. (2019a), allows us to characterize the set  $\Lambda_{\delta}^+$  ( $P_n$ ) in terms of an associated level set; in particular, we have

$$\Lambda_{\delta}^{+}(P_n) = \operatorname{cl}\{\beta : R_n(\beta) \leqslant \delta\},\,$$

where  $cl(\cdot)$  denotes closure. Indeed, this is because

$$\Lambda_{\delta}^{+}(P_n) = \operatorname{cl} \Big[ \cap_{\epsilon > 0} \Big\{ \beta \in B : \beta \in \arg\min_{\beta \in B} E_P \{ \ell(X; \beta) \} \text{ for some } P \in \mathcal{U}_{\delta_n + \epsilon}(P_n) \Big\} \Big].$$

If 
$$\beta \in B^{\circ}$$
, we have  $R_n(\beta) = \inf_{P \in \mathcal{P}(\Omega)} [D_c(P, P_n) : E_P \{h(X, \beta)\} = 0]$ .

Next, for the sequence of radii  $\delta_n = n^{-\gamma} \eta$ , for some positive constants  $\eta, \gamma$ , define functions  $V_n^{\mathrm{DRO}} : \mathbb{R}^d \to \mathbb{R}$  and  $V_n^{\mathrm{ERM}} : \mathbb{R}^d \to \mathbb{R}$ , as below, by considering suitably scaled versions of the distributionally robust and empirical risk objective functions, namely

$$V_n^{\text{DRO}}(u) = n^{\bar{\gamma}} \{ \Psi_n (\beta_* + n^{-\bar{\gamma}/2} u) - \Psi_n(\beta_*) \} \text{ and }$$

$$V_n^{\text{ERM}}(u) = n [E_{P_n} \{ \ell(X; \beta_* + n^{-1/2} u) \} - E_{P_n} \{ \ell(X; \beta_*) \} ],$$

where  $\bar{\gamma} = \min \{ \gamma, 1 \}$  is defined in Theorem 1. Moreover, define  $V : \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$  via  $V(x, u) = x^T u + 2^{-1} u^T C u$ . The following result, as we shall see, can be used to establish Theorem 1 directly.

THEOREM 3. Suppose that the assumptions made in Theorem 1 hold. Then,

$$\left\{ V_n^{\text{ERM}}(\cdot), \ V_n^{\text{DRO}}(\cdot), \ nR_n \left( \beta_* + n^{-1/2} \times \cdot \right) \right\}$$
  
$$\Rightarrow \left\{ V(-H, \cdot), \ V\{-f_{\eta, \gamma}(H), \cdot\}, \ \varphi^*(H - C \times \cdot) \right\}$$

on the space  $C(\mathbb{R}^d;\mathbb{R})^3$  equipped with the topology of uniform convergence in compact sets.

Ensuring smoothness of  $D_{\beta}h(x+\Delta,\beta)$  and  $D_xh(x+\Delta,\beta)$  around  $\beta=\beta_*$  as in Assumption 2(c) is useful towards investigating the behaviour of  $nR_n(\cdot)$  in the neighbourhood of  $\beta^*$ , as required in the third component in the triplet in Theorem 3.

# 5.2. Proof of Theorem 3

Throughout this section we suppose that the assumptions imposed in Theorem 1 hold. Let  $H_n = n^{-1/2} \sum_{i=1}^n h(X_i, \beta_*)$ . The following sequence of results will be useful in proving Theorem 3 and Proposition 1. Propositions 6 and 7 hold true for  $\Omega = \mathbb{R}^d$ , while Propositions 8–12 hold true for general  $\Omega$  under the assumption  $P_*(\Omega^\circ) = 1$  in Proposition 1.

PROPOSITION 6. Fix  $\alpha \in [0, 1]$ . Given  $\varepsilon, \varepsilon', K > 0$ , there exists a positive integer  $n_0$  such that

$$\operatorname{pr}\left[\left|n^{\alpha-1}V_{n}^{\operatorname{ERM}}\left\{n^{(1-\alpha)/2}u\right\}-n^{\alpha/2}H_{n}^{\operatorname{T}}u-2^{-1}u^{\operatorname{T}}Cu\right|\leqslant\varepsilon'\right]\geqslant1-\varepsilon$$

for every  $n > n_0$  and  $||u||_2 \le K$ . Specifically, if  $\alpha = 1$ ,

$$\operatorname{pr}\left\{\left|V_{n}^{\operatorname{ERM}}(u) - H_{n}^{\operatorname{T}}u - 2^{-1}u^{\operatorname{T}}Cu\right| \leqslant \varepsilon'\right\} \geqslant 1 - \varepsilon. \tag{16}$$

PROPOSITION 7. Given  $\varepsilon, \varepsilon', K > 0$ , there exists a positive integer  $n_0$  such that

$$\operatorname{pr}\left\{\left|V_{n}^{\mathrm{DRO}}(u) + f_{n,\nu}(-H_{n})^{\mathsf{T}}u - 2^{-1}u^{\mathsf{T}}Cu\right| \leqslant \varepsilon'\right\} \geqslant 1 - \varepsilon \tag{17}$$

for every  $n > n_0$  and  $||u||_2 \leq K$ .

PROPOSITION 8. Define the set  $\Theta \subset \mathbb{R}^d$  as  $\Theta = \{\beta \in B^\circ : 0 \in \text{conv}[\{h(x,\beta) \mid x \in \Omega\}]^\circ\}$ , where conv(S) denotes the convex hull of the set S. For  $\beta_* + n^{-1/2}u \in \Theta$ ,  $nR_n(\beta_* + n^{-1/2}u) = 0$ 

 $\max_{\xi} \left\{ -\xi^{\mathrm{T}} H_n - M_n(\xi, u) \right\}$ , where

$$M_n(\xi, u) = \frac{1}{n} \sum_{i=1}^n \max_{\Delta: X_i + n^{-1/2} \Delta \in \Omega} \left\{ \xi^{\mathsf{T}} \int_0^1 D_x h(X_i + n^{-1/2} t \Delta, \beta_* + n^{-1/2} t u) \Delta \, \mathrm{d}t + \xi^{\mathsf{T}} \int_0^1 D_\beta h(X_i + n^{-1/2} t \Delta, \beta_* + n^{-1/2} t u) u \, \mathrm{d}t - \|\Delta\|_q^2 \right\}.$$

Furthermore, there exists a neighbourhood of  $\beta *$ ,  $B_{\epsilon}(\beta_*)$ , such that  $B_{\epsilon}(\beta_*) \subset \Theta$ .

PROPOSITION 9. Consider any  $\varepsilon, \varepsilon', K > 0$ . Then there exists  $b_0 \in (0, \infty)$  such that, for any  $b \ge b_0$ ,  $c_0 > 0$ ,  $\epsilon_0 > 0$ , we have a positive integer  $n_0$  such that

$$\operatorname{pr}\left[\sup_{\|u\|_{2}\leqslant K}\left\{nR_{n}\left(\beta_{*}+n^{-1/2}u\right)-f_{\operatorname{up}}(H_{n},u,b,c)\right\}\leqslant\varepsilon'\right]\geqslant1-\varepsilon$$

for all  $n \ge n_0$ , and  $f_{up}(H_n, u, b, c)$  equals

$$\max_{\|\xi\|_p \le b} \left( -\xi^{\mathsf{T}} H_n - E \left[ 4^{-1} \| \{ D_x h(X, \beta_*) \}^{\mathsf{T}} \xi \|_p^2 + \xi^{\mathsf{T}} D_\beta h(X, \beta_*) u \right] \mathbb{I}(X \in C_0^{\epsilon_0}) \right),$$

with  $C_0 = \{x \in \Omega : ||x||_p \leqslant c_0\}.$ 

PROPOSITION 10. For any  $\varepsilon, \varepsilon', K, b > 0$ , there exists a positive integer  $n_0$  such that

$$\operatorname{pr}\left[\sup_{\|u\|_{2} \leqslant K} \left\{ nR_{n} \left(\beta_{*} + n^{-1/2} u\right) - f_{\operatorname{low}}(H_{n}, u, b) \right\} \geqslant -\varepsilon' \right] \geqslant 1 - \varepsilon$$

for all  $n > n_0$ , where

$$f_{\text{low}}(H_n, u, b) = \max_{\|\xi\|_p \le b} \left( -\xi^{\mathsf{T}} H_n - E\left[ 4^{-1} \| \{D_X h(X, \beta_*)\}^{\mathsf{T}} \, \xi \|_p^2 + \xi^{\mathsf{T}} D_\beta h(X, \beta_*) u \right] \right).$$

PROPOSITION 11. For any  $\varepsilon > 0$  there exist constants  $a, n_0 > 0$  such that, for every  $n \ge n_0$ , pr  $\{nR_n(\beta_*) \le a\} \ge 1 - \varepsilon$ .

PROPOSITION 12. For any  $\varepsilon, \varepsilon', K > 0$ , there exist positive constants  $n_0, \delta$  such that

$$\sup_{\substack{\|u_1-u_2\|_2 \leqslant \delta \\ \|u_i\|_2 \leqslant K}} \left| nR_n \left( \beta_* + n^{-1/2} u_1 \right) - nR_n \left( \beta_* + n^{-1/2} u_2 \right) \right| \leqslant \varepsilon'$$

with probability exceeding  $1 - \varepsilon$  for every  $n > n_0$ .

With the statements of these results, we proceed with the proof of Theorem 3.

*Proof of Theorem* 3. Since  $E\{h(X, \beta_*)\}=0$ , it follows from the central limit theorem that  $H_n \Rightarrow -H$ , where  $H \sim \mathcal{N}[0, E\{h(X, \beta_*)h(X, \beta_*)^T\}]$ . Since the inequalities (17) and (16) are associated with the same  $H_n$ , it follows from Propositions 6 and 7 that

$$V_n^{\mathrm{ERM}}(\cdot) \Rightarrow V^{\mathrm{ERM}}(\cdot) = V(-H, \cdot), \qquad V_n^{\mathrm{DRO}}(\cdot) \Rightarrow V^{\mathrm{DRO}}(\cdot) = V\{-f_{\eta,\gamma}(H), \cdot\}$$
 (18)

jointly on the space topologized by uniform convergence on compact sets.

To prove convergence of the third component of the triplet considered in Theorem 3, observe from the definitions of  $\varphi^*(\cdot)$  and C that

$$\varphi^*(H - Cu) = \max_{\xi} \left( \xi^{\mathsf{T}} [H - E\{D_{\beta}h(X, \beta_*)\}u] - 4^{-1}E \| \{D_X h(X, \beta_*)\}^{\mathsf{T}} \xi \|_p^2 \right). \tag{19}$$

Consider any fixed  $K \in (0, +\infty)$ . Due to the weak convergence  $H_n \Rightarrow -H$ , applications of the continuous mapping theorem to the bounds in Propositions 9 and 10 result in

$$f_{\text{up}}(H_n, u, b, c) \Rightarrow \max_{\|\xi\|_p \leqslant b} \left( \xi^{\mathsf{T}} H - E \left[ 4^{-1} \| \left\{ D_x h(X, \beta_*) \right\}^{\mathsf{T}} \xi \|_p^2 + \xi^{\mathsf{T}} D_\beta h(X, \beta_*) u \right] \mathbb{I}(X \in C_0^{\epsilon_0}) \right), \tag{20}$$

$$f_{\text{low}}(H_n, u, b) \Rightarrow \max_{\|\xi\|_p \leqslant b} \left( \xi^{\mathsf{T}} H - E \left[ 4^{-1} \| \{ D_x h(X, \beta_*) \}^{\mathsf{T}} \xi \|_p^2 + \xi^{\mathsf{T}} D_\beta h(X, \beta_*) u \right] \right)$$
(21)

for any u satisfying  $||u||_2 \le K$ . Since the bounds in Propositions 9 and 10 hold for arbitrarily large choices of the constants b, c, and an arbitrarily small choice for constant  $\epsilon_0$ , combining with the assumption that  $P_*(\Omega^\circ) = 1$ , we conclude from (19), (20) and (21) that

$$nR_n(\beta_* + n^{-1/2}u) \Rightarrow \varphi^*(H - Cu) \tag{22}$$

for any u satisfying  $||u||_2 \le K$ . Finally, from Propositions 11 and 12, the collection  $\{nR_n(\beta_* + n^{-1/2} \times \cdot)\}$  is tight; see, for example, Billingsley (2013, Theorem 7.4). As a consequence of this tightness and the finite-dimensional convergence in (22), we have  $nR_n(\beta_* + n^{-1/2} \times \cdot) \Rightarrow \varphi^*(H - C \times \cdot)$ . Combining this observation with those in (18), we obtain the desired convergence result in Theorem 3. Furthermore, since  $f_{\text{low}}(H_n, u, b)$  and  $f_{\text{up}}(H_n, u, b)$  are associated with the same  $H_n$  as inequalities (17) and (16), we have the three terms converging jointly.

# 5.3. Proof of Theorem 1

Theorem 1 is proved by considering suitable level sets of the component functions in the triplet  $\{V_n^{\text{ERM}}(\cdot), V_n^{\text{DRO}}(\cdot), nR_n(\beta_* + n^{-1/2} \times \cdot)\}$  considered in Theorem 3. To reduce the clutter in expressions, from here onwards we refer to the distributionally robust estimator (2) simply as  $\beta_n^{\text{DRO}}$ , with the dependence on the radius  $\delta_n$  to be understood from the context. To begin, consider the following tightness result.

PROPOSITION 13. The sequences  $\{\arg\min_u V_n^{\text{ERM}}(u) : n \ge 1\}$  and  $\{\arg\min_u V_n^{\text{DRO}}(u) : n \ge 1\}$  are tight.

Observe that  $V_n^{\rm ERM}(\cdot)$  and  $V_n^{\rm DRO}(\cdot)$  are minimized, respectively, at  $n^{1/2}(\beta_n^{\rm ERM}-\beta_*)$  and  $n^{\bar{\gamma}/2}(\beta_n^{\rm DRO}-\beta_*)$ . Furthermore, due to the positive definiteness of C in Assumption 2(b), we have that  $V^{\rm ERM}(\cdot)$  and  $V^{\rm DRO}(\cdot)$  are strongly convex with respect to u and have unique minimizers, with probability 1. Therefore, due to the tightness of the sequences  $\{n^{1/2}(\beta_n^{\rm ERM}-\beta_*)\}_{n\geqslant 1}$  and  $\{n^{\bar{\gamma}/2}(\beta_n^{\rm DRO}-\beta_*)\}_{n\geqslant 1}$ , see Proposition 13, and the weak convergence of  $V_n^{\rm ERM}(\cdot)$  and  $V_n^{\rm DRO}(\cdot)$  in Theorem 3, we have the following convergences:

$$n^{1/2}(\beta_n^{\text{ERM}} - \beta_*) \Rightarrow \arg\min_{u} V(-H, u) = C^{-1}H,$$
  

$$n^{\bar{\gamma}/2}(\beta_n^{\text{DRO}} - \beta_*) \Rightarrow \arg\min_{u} V^{\text{DRO}}(u) = C^{-1}f_{\eta,\gamma}(H).$$
(23)

Finally, to prove the convergence of the sets  $\Lambda_{\delta_n}^+(P_n)$ , we proceed as follows. Define  $G_n(u) = nR_n(\beta_* + n^{-1/2}u)$ ,  $G(u) = \varphi^*(H - Cu)$  and  $\alpha_n = n\delta_n$ . For any function  $f: B \to \mathbb{R}$  and  $\alpha \in [0, +\infty]$ , let lev $(f, \alpha)$  denote the level set  $\{x \in \mathbb{R}^d : f(x) \leq \alpha\}$ .

PROPOSITION 14. If  $\delta_n = n^{-1}\eta$ , then  $\operatorname{cl}\{\operatorname{lev}(G_n, \alpha_n)\} \Rightarrow \operatorname{lev}(G, \eta)$ .

PROPOSITION 15. If  $\delta_n = n^{-\gamma} \eta$  for some  $\gamma > 1$ , then  $\operatorname{cl}\{\operatorname{lev}(G_n, \alpha_n)\} \Rightarrow \{C^{-1}H\}$ .

Proposition 16. If  $\delta_n = n^{-\gamma} \eta$  for some  $\gamma < 1$ , then  $\operatorname{cl}\{\operatorname{lev}(G_n, \alpha_n)\} \Rightarrow \mathbb{R}^d$ .

Propositions 14–16 allow us to complete the proof of Theorem 1 as follows. From the definition of  $R_n(\beta)$ ,  $\Lambda_{\delta_n}^+(P_n) = \{\beta : R_n(\beta) \le \delta_n\} = \beta_* + n^{-1/2} \{u : G_n(u) \le \alpha_n\}$ . From Propositions 14–16,

$$n^{1/2}\left\{\Lambda_{\delta_n}^+(P_n) - \beta_*\right\} = \{u : G_n(u) \leqslant \alpha_n\} \Rightarrow \begin{cases} \operatorname{lev}(G, \eta) & \text{if } \gamma = 1, \\ \mathbb{R}^d & \text{if } \gamma < 1, \\ \{C^{-1}H\} & \text{if } \gamma > 1. \end{cases}$$

Observe that  $\varphi^*(u) = \varphi^*(-u)$ . Therefore,  $\operatorname{lev}(G, \eta) = \{u : \varphi^*(H - Cu) \leq \eta\} = C^{-1}H + \{u : \varphi^*(Cu) \leq \eta\}$ . Since the three terms in Theorem 3 converge jointly, we have the three terms in Theorem 1 also converging jointly. This completes the proof of Theorem 1.

Proposition 1 follows by adopting exactly the same steps used to establish the convergence of  $n^{1/2} \left\{ \Lambda_{\delta_n}^+(P_n) - \beta_* \right\}$  in the proof of Theorem 1.

## 6. DISCUSSION

We discuss the subtleties in deriving a limit theorem for the distributionally robust estimator  $\beta_n^{\mathrm{DRO}}$  when the support of the random vector X is a strict subset of  $\mathbb{R}^m$ . Suppose that the support of X is constrained to be contained in the set  $\Omega = \{x \in \mathbb{R}^m : Ax \leq b\}$  specified in terms of linear constraints involving an  $l \times m$  matrix A and  $b \in \mathbb{R}^l$ . For the sake of clarity, we discuss here only the nondegenerate case where  $\delta_n = \eta/n$ .

Considering the transportation cost  $c(x,y) = \|x-y\|_2^2$  in Definition 1, we demonstrate in the Supplementary Material that the central limit theorem,  $n^{1/2}\{\beta_n^{\mathrm{DRO}}(\delta_n) - \beta_*\} \Rightarrow C^{-1}H - \eta^{1/2}C^{-1}D_{\beta}S(\beta_*)$ , continues to hold, for example, in the elementary case where the matrix A has linearly independent rows, X has a probability density which is absolutely continuous with respect to the Lebesgue measure on  $\mathbb{R}^m$  and the support  $\Omega$  is compact. A key element which emerges in the verification is that the fraction of samples which get transported to the boundary of the set  $\Omega$  stays at  $O_p(n^{-1/2})$  as  $n \to \infty$ .

On the other hand, when the set  $\Omega = \{x \in \mathbb{R}^m : Ax \leq b\}$  has equality constraints, as in, for example,  $\Omega = \{(x_1, x_2, \dots, x_m) \in \mathbb{R}^2 : x_1 - x_2 = 0\}$ , the bias term in the limit theorem gets affected due to the constraint binding all the samples  $\{X_1, \dots, X_n\}$ , and the fraction of samples which get transported to the boundary of the set  $\Omega$  is 1. This can be easily seen in the linear regression example in § 4, where  $\ell(x, y; \beta) = (y - \beta^T x)^2$  and the support is taken as  $\Omega = \{(x_1, x_2) \in \mathbb{R}^2 : x_1 = x_2\}$ . For this elementary example, we instead have

$$n^{1/2} \{ \beta_n^{\text{DRO}}(\delta_n) - \beta_* \} \Rightarrow C^{-1}H - \eta^{1/2}C^{-1}D_{\beta}\tilde{S}(\beta_*),$$
 (24)

where  $\tilde{S}(\beta)$  is different from the term  $S(\beta)$ , as in  $\tilde{S}(\beta_*) = 2^{1/2-1/q} |\beta^T 1| \|\beta\|_p^{-1} S(\beta)$ . Here, recall the earlier definition  $S(\beta) = [E\{\|D_x\ell(X;\beta)\|_p^2\}]^{1/2}$  in (5) for the unconstrained support case. The computations required to arrive at the above conclusion are presented in the Supplementary Material. In the presence of general support constraints of the form  $\Omega = \{x \in \mathbb{R}^m : Ax = b\}$ , we show that (24) holds with  $\tilde{S}(\beta) = \|P_{\mathcal{N}(A)}\beta\|_2$  for quadratic losses of the form  $\ell(x;\beta) = a + \beta^T x + \beta^T C\beta$ ; here,  $\ell(x;\beta) = a + \beta^T x +$ 

#### ACKNOWLEDGEMENT

This work was supported by the Air Force Office of Scientific Research (FA9550-20-1-0397), the National Science Foundation (1915967, 1820942 and 1838576) and the Ministry of Education Singapore (2018 134).

#### SUPPLEMENTARY MATERIAL

Supplementary Material available at *Biometrika* online includes the technical proofs.

#### REFERENCES

BILLINGSLEY, P. (2013). Convergence of Probability Measures. New York: John Wiley and Sons.

BLANCHET, J. & KANG, Y. (2017). Distributionally robust groupwise regularization estimator. *Proc. Mach. Learn. Res.* 77, 97–112.

BLANCHET, J., KANG, Y. & MURTHY, K. (2019a). Robust Wasserstein profile inference and applications to machine learning. *J. Appl. Prob.* **56**, 830–57.

BLANCHET, J., KANG, Y., MURTHY, K. & ZHANG, F. (2019b). Data-driven optimal transport cost selection for distributionally robust optimization. In *Proc. Winter Simulation Conf. (WSC)*, pp. 3740–51.

BLANCHET, J. & MURTHY, K. (2019). Quantifying distributional model risk via optimal transport. *Math. Oper. Res.* 44, 565–600.

BLANCHET, J., MURTHY, K. & ZHANG, F. (2020). Optimal transport based distributionally robust optimization: Structural properties and iterative schemes. *arXiv*:1810.02403v2.

CHEN, Ř. & PASCHALIDIS, I. C. (2018). A robust learning approach for regression models based on distributionally robust optimization. *J. Mach. Learn. Res.* **19**, 517–64.

CHEN, Z., KUHN, D. & WIESEMANN, W. (2018). Data-driven chance constrained programs over Wasserstein balls. arXiv:1809.00210.

CISNEROS-VELARDE, P., PETERSEN, A. & OH, S.-Y. (2020). Distributionally robust formulation and model selection for the graphical lasso. *Proc. Mach. Learn. Res.* **108**, 756–65.

Duchi, J. C., Hashimoto, T. & Namkoong, H. (2020). Distributionally robust losses against mixture covariate shifts. arXiv:2007.13982.

DUPACOVA, J. & WETS, R. (1988). Asymptotic behavior of statistical estimators and of optimal solutions of stochastic optimization problems. *Ann. Statist.* **16**, 1517–49.

GAO, R., CHEN, X. & KLEYWEGT, A. J. (2020). Wasserstein distributional robustness and regularization in statistical learning. *arXiv*:1712.06050v3.

GAO, R. & KLEYWEGT, A. J. (2016). Distributionally robust stochastic optimization with Wasserstein distance. arXiv:1604.02199v2.

GAO, R., XIE, L., XIE, Y. & XU, H. (2018). Robust hypothesis testing using Wasserstein uncertainty sets. In *Proc. 32nd Int. Conf. Neural Information Processing Systems*, pp. 7913–23.

HUBER, P. J. (1967). The behavior of maximum likelihood estimates under nonstandard conditions. In *Proc. 5th Berkeley Symp. Mathematical Statistics and Probability*, vol. 1. Berkeley, CA: University of California Press.

KNIGHT, K. & Fu, W. (2000). Asymptotics for lasso-type estimators. Ann. Statist. 28, 1356–78.

Luo, F. & Mehrotra, S. (2017). Decomposition algorithm for distributionally robust optimization using Wasserstein metric. *arXiv*:1704.03920.

Mohajerin Esfahani, P. & Kuhn, D. (2018). Data-driven distributionally robust optimization using the Wasserstein metric: Performance guarantees and tractable reformulations. *Math. Program.* **171**, 115–66.

- Molchanov, I. S. (2005). Theory of Random Sets, vol. 19. Berlin: Springer.
- OWEN, A. (1988). Empirical likelihood ratio confidence intervals for a single functional. Biometrika 75, 237-49.
- OWEN, A. (1990). Empirical likelihood ratio confidence regions. Ann. Statist. 18, 90–120.
- OWEN, A. B. (2001). Empirical Likelihood. London: Chapman and Hall/CRC.
- Shafiezadeh-Abadeh, S., Esfahani, P. & Kuhn, D. (2015). Distributionally robust logistic regression. In *Proc. 28th Conf. Advances in Neural Information Processing Systems*, pp. 1576–84.
- SHAPIRO, A. (1989). Asymptotic properties of statistical estimators in stochastic programming. *Ann. Statist.* 17, 841–58.
- SHAPIRO, A. (1991). Asymptotic analysis of stochastic programs. *Ann. Oper. Res.* **30**, 169–86.
- Shapiro, A. (1993). Asymptotic behavior of optimal solutions in stochastic programming. *Math. Oper. Res.* **18**, 829–45.
- SHAPIRO, A. (2000). On the asymptotics of constrained local *m*-estimators. *Ann. Statist.* **28**, 948–60.
- SINHA, A., NAMKOONG, H. & DUCHI, J. (2018). Certifiable distributional robustness with principled adversarial training. In *Proc. Int. Conf. Learning Representations*.
- STAIB, M. & JEGELKA, S. (2017). Distributionally robust deep learning as a generalization of adversarial training. In *Proc. NIPS Workshop on Machine Learning and Computer Security*.
- VAN DER VAART, A. & WELLNER, J. (1996). Weak Convergence and Empirical Processes: With Applications to Statistics. New York: Springer.
- VOLPI, R., NAMKOONG, H., SENER, O., DUCHI, J., MURINO, V. & SAVARESE, S. (2018). Generalizing to unseen domains via adversarial data augmentation. In *Proc. 32nd Int. Conf. Neural Information Processing Systems*, pp. 5339–49.
- YANG, I. (2017). A convex optimization approach to distributionally robust Markov decision processes with Wasserstein distance. *IEEE Control Syst. Lett.* **1**, 164–9.
- YANG, I. (2020). Wasserstein distributionally robust stochastic control: A data-driven approach. *IEEE Trans. Automatic Control*, doi:10.1109/TAC.2020.3030884.
- ZHAO, C. & GUAN, Y. (2018). Data-driven risk-averse stochastic optimization with Wasserstein metric. *Oper. Res. Lett.* **46**, 262–7.

[Received on 5 June 2019. Editorial decision on 19 March 2021]