# Ordered Risk Minimization: Learning More from Less Data

Peter Coppens and Panagiotis Patrinos<sup>†</sup>

Abstract—We consider the worst-case expectation of a permutation invariant ambiguity set of discrete distributions as a proxy-cost for data-driven expected risk minimization. For this framework, we coin the term ordered risk minimization to highlight how results from order statistics inspired the proxycost. Specifically, we show how such costs serve as pointwise high-confidence upper bounds of the expected risk. The confidence level can be determined tightly for any sample size. Conversely we also illustrate how to calibrate the size of the ambiguity set such that the high-confidence upper bound has some user specified confidence. This calibration procedure notably supports  $\phi$ -divergence based ambiguity sets. Numerical experiments then illustrate how the resulting scheme both generalizes better and is less sensitive to tuning parameters compared to the empirical risk minimization approach.

#### I. INTRODUCTION

The problem of *expected risk minimization* is ubiquitous in machine learning and statistics [1], [2]. It is based on the idea that the quality of a model can be assessed by measuring its expected error, quantified by some loss function. The expectation should be evaluated with respect to the datagenerating distribution. However, in practice, only samples are available. So the expectation needs to be replaced with a data-driven proxy, which aggregates the data. The common solution is *empirical risk minimization* or the *sample average approach (SAA)*, where one takes an average over the losses at the sampled data points.

Despite its advantages, the SAA often exhibits excessive sensitivity to the specific data realizations, particularly in high-dimensional settings [3, §8.H], leading to diminished generalization capabilities of the model. To address this, researchers have turned to *Distributionally Robust Optimization (DRO)*, aiming to robustify against disparities between the empirical and true data-generating distributions.

In DRO, a worst-case expectation with respect to distributions in an *ambiguity set* centered on the empirical distribution serves as a proxy for the true expected risk. This set can be based on the Wasserstein distance [4],  $\phi$ divergences [5], hypothesis tests [6], and others. See [7], [8] for recent surveys. However, DRO faces challenges when determining the ambiguity set's size. Current approaches rely on concentration inequalities [9] or asymptotic bounds [5], ensuring that true distribution is contained within the ambiguity set with high probability. Unfortunately, this often results in conservatism, caused by either loose constants in the concentration inequalities or the shape of the ambiguity set. As an alternative, bootstrapping or cross validation techniques are often employed (cf. [4]). These can be computationally expensive and lack statistical guarantees for finite samples, similarly to the asymptotic bounds. Such guarantees are a requirement in safety-critical applications like control (e.g. constraint tightening in tube-based MPC [10], [11]).

To address conservativeness issues, [12], [13] focus on bounding the expectation directly as an alternative to creating a confidence bound for the entire distribution. This mimics the focus on the expectation in the statistical learning framework of [2, §1]. However, their bounds are asymptotic and therefore also lack strong statistical guarantees. In this paper, we take the first steps towards a finite-sample version of their scheme. To achieve this, we draw inspiration from results in order statistics [14] and stochastic orders [15] to motivate the use of permutation invariant ambiguity sets. Notably, the  $\phi$ divergences used in [12], [13], [16] produce a specific case. We coin the term *ordered* risk minimization to emphasize our statistical motivations. We demonstrate how to calibrate the ambiguity set's size to upper bound the true expectation with high probability, even when the true distribution does not fall within the ambiguity set. This probability serves as an intuitive tuning parameter.

The remainder of the paper continues as follows. We first present some notation, before moving on the the problem statement in §II. There the proxy costs we use are presented as well as the calibration problem. We present the statistical interpretation of these proxy costs as high confidence upper bounds in §III and solve the calibration problem in §IV. Numerical experiments are then presented in §V.

*Notation:* Let  $\mathbb{R}$  denote the reals and  $\overline{\mathbb{R}}$  the extended reals. For some convex function  $\phi \colon \mathbb{R}^n \to \overline{\mathbb{R}}$ , let  $\phi^*$  denote the convex conjugate,  $\partial \phi$  the subgradient and dom  $\phi$  its domain. For a set  $\mathcal{X}$  let  $\iota_{\mathcal{X}}(x) = 0$  if  $x \in \mathcal{X}$  and  $+\infty$ otherwise be the indicator function of  $\mathcal{X}$ . For integers a, blet  $[a, b] = \{a, \dots, b\}$  and  $[b] = \{1, \dots, b\}$ . Let  $[x]_+ =$  $\max(0, x)$ . For real vectors  $x, y \in \mathbb{R}^n$  we use  $\langle x, y \rangle$  to denote the Euclidean inner product and  $\mathbb{1}_n \in \mathbb{R}^n$  is the vector of all ones. For a cone  $\mathcal{K}$  let  $\mathcal{K}^{\circ} := \{y : \langle x, y \rangle \leq$  $0, \forall x \in \mathcal{K}$  denote its polar cone. Let  $\Pi^n$  denote the permutations of [n] (i.e., all bijections  $[n] \rightarrow [n]$ ). We write  $\pi x = (x_{\pi(1)}, \dots, x_{\pi(n)})$  for  $\pi \in \Pi^n, x \in \mathbb{R}^n$ and similarly let  $\Pi^n y = \{\pi y : \pi \in \Pi^n\}$  denote the orbit of y under  $\Pi^n$ . Let  $\mathbb{R}^n_{\uparrow} := \{x \colon x_1 \leq x_2 \leq \cdots \leq x_n\}$ denote the monotone cone and  $\mathcal{M}^n$  its polar (cf. Lem. A.1 and (6)). Let  $\Delta^n := \{ \mu : \sum_{i=1}^n \mu_i = 1, \mu_i \ge 0, i \in [n] \}$ 

<sup>&</sup>lt;sup>†</sup>P. Coppens and P. Patrinos are with the Department of Electrical Engineering (ESAT-STADIUS), KU Leuven, Kasteelpark Arenberg 10, 3001 Leuven, Belgium. Email: peter.coppens@kuleuven.be, panos.patrinos@kuleuven.be

This work was supported by: the Research Foundation Flanders (FWO) PhD grant 11E5520N and research projects G081222N, G033822N, G0A0920N; European Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie grant agreement No. 953348.

denote the probability simplex. For a vector  $x \in \mathbb{R}^n$  let  $x_{(1)} \leq x_{(2)} \leq \cdots \leq x_{(n)}$  be the increasing permutation of the elements of x with  $x_{\uparrow} = (x_{(1)}, x_{(2)}, \dots, x_{(n)})$ . For sets A, B let  $A + B := \{a + b : a \in A, b \in B\}$  denote the Minkowski sum. For random variables X, Y we write  $X \cong Y$  to say X is identically distributed to Y. Let esssup[X] denote the essential supremum. Let  $X \cong U[\ell, u]$ imply X is uniformly distributed over  $[\ell, u]$ . For a set A, $X \cong U[A]$  is then uniformly distributed over A. Finally let  $X \cong \mathcal{N}(\mu, \Sigma)$  denote X that is normally distributed with mean  $\mu$  and covariance  $\Sigma$ .

#### II. PROBLEM STATEMENT

We consider expected risk minimization

$$\underset{\theta \in \Theta}{\text{minimize}} \quad \mathbb{E}[\ell(\theta, \xi)]. \tag{1}$$

Here  $\Theta \subseteq \mathbb{R}^{n_{\theta}}$  and  $\xi \colon \Omega \to \Xi \subseteq \mathbb{R}^{n_{\xi}}$  is a random vector on a probability space  $(\Omega, \mathscr{F}, \mathbb{P})$ . As is common [1], [2], we then assume access to *independent and identically distributed (iid) samples*  $\xi^{(1)}, \ldots, \xi^{(n-1)}$ . Let  $\ell_i(\theta) = \ell(\theta, \xi^{(i)})$  for i = $1, \ldots, n-1$  and take  $\ell_n(\theta)$  such that it upper bounds the cost almost surely. That is  $\mathbb{P}[\ell(\theta, \xi) \leq \ell_n(\theta)] = 1$ . It is assumed that  $\ell_n(\theta)$  is finite for any  $\theta$ . We do so for two reasons: *(i)* an assumption on the tail of the distribution of  $\ell(\theta, \xi)$  is required to find a confidence interval for the mean [17]; *(ii)* the scheme is simplified considerably. An example of a valid bound is  $\ell_n(\theta) = \sup_{\xi \in \Xi} \ell(\theta, \xi)$ .

We will use a data-driven proxy for the expectation by introducing *permutation invariant ambiguity sets*. These are subsets  $\mathcal{A}$  of the probability simplex  $\Delta^n$  such that, for each  $\mu \in \mathcal{A}$  any permutation of  $\mu$  is also in  $\mathcal{A}$ . The *ordered risk minimization* problem is then:

$$\underset{\theta \in \Theta}{\operatorname{minimize}} \quad \sup_{\mu \in \mathscr{A}} \sum_{i=1}^{n} \mu_{i} \ell_{i}(\theta).$$
(2)

This proxy cost interpolates between the robust case for  $\mathcal{A} = \Delta^n$  and the sample average (including a term associated with  $\ell_n(\theta)$ ) when  $\mathcal{A} = \{\mathbb{1}_n/n\}$ . The interpolation interpretation is also common in DRO [5]. To find a good balance, our goal is to select  $\mathcal{A}$  such that, for all  $\theta \in \Theta$ ,

$$\mathbb{P}\left[\sup_{\mu\in\mathscr{A}}\sum_{i=1}^{n}\mu_{i}\ell_{i}(\theta)\geq\mathbb{E}[\ell(\theta,\xi)]\right]\geq1-\delta,\qquad(3)$$

We refer to this problem as the *calibration problem*. It robustifies against disparities between the empirical and true distributions. The parameter  $\delta$  then serves as an intuitive, user-determined parameter that controls the conservativeness of the method. However, as illustrated by experiments, our method is relatively insensitive to the value of  $\delta$ .

To find an ambiguity set  $\mathcal{A}$  satisfying (3) we need to somehow parametrize it. A well known class of permutation invariant ambiguity sets uses  $\phi$ -divergences [5]. Let  $\phi: \mathbb{R}_+ \to \mathbb{R}$  be lower semicontinuous, convex and  $\phi(1) =$ 0. Also, let<sup>1</sup>  $I_{\phi}(\mu, \nu) := \sum_{i=1}^{n} \nu_i \phi(\mu_i/\nu_i)$  for all  $\mu, \nu \in \Delta^n$ .

<sup>1</sup>We take the lower semicontinuous envelope of the terms inside the sum [18, Def. 6] to handle cases where  $\nu_i$  equals zero.

A (centered)  $\phi$ -divergence ambiguity set is then

$$\mathcal{A}_{\alpha} := \left\{ \mu \in \Delta^{n} \colon I_{\phi}\left(\mu, \frac{\mathbb{1}_{n}}{n}\right) = \sum_{i=1}^{n} \frac{\phi(n\mu_{i})}{n} \leq \alpha \right\}.$$
(4)

In this work we consider two examples: *total variation* (*TV*) for which  $\phi(t) = |t - 1|$  and *Kullback Leibler (KL)* divergence for which  $\phi(t) = t \log t - t + 1$ . However, our method works for any divergence. See [5] for more examples.

To calibrate  $\mathcal{A}_{\alpha}$  the radius  $\alpha \in \mathbb{R}$  should then be the smallest value such that (3) still holds. We also provide an alternative parametrization, related to a well known bound by Anderson [19] and the conditional value-at-risk.

It is important to note that the constraint in (3) is less stringent compared to DRO, which guarantees that the supremum in (3) acts as a high-confidence upper bound, uniformly over  $\theta^2$ . After all, we never require that the true distribution is contained within  $\mathcal{A}$  (as is the case in [20]). The gap between the point-wise (3) and the uniform equivalent is examined for  $\phi$ -divergences in [12], [13] in the asymptotic regime.

We numerically approximate the calibration problem without samples from  $\xi$ . So the parameters of the set  $\mathcal{A}$  only need to be computed once and can be tabulated afterwards. This contrasts the complex derivations and the resulting conservative constants associated with analytical approaches used to compute the radius of an ambiguity set in DRO [4], [5], [9]. We show experimentally how our calibration of  $\mathcal{A}$ according to (3) greatly improves generalization.

#### **III. STATISTICAL FRAMEWORK**

The analysis of this section investigates upper bounds for the mean of a scalar random variable  $(rv) Z: \Omega \to \mathbb{R}$ , defined on some probability space  $(\Omega, \mathscr{F}, \mathbb{P})$ . These findings remain applicable to the previous section, when considering  $Z = \ell(\theta, \xi)$  for fixed  $\theta$ . Let  $F(z) = \mathbb{P}[Z \leq z]$  denote the cumulative distribution function (cdf) of Z. We assume access to iid samples  $Z_1, \ldots, Z_{n-1}$  and an upper bound denoted as  $Z_n = \text{esssup}[Z]$  for notational convenience, which satisfies  $F(Z_n) = 1$ .

We then introduce the coverages

$$W_i = F(Z_{(i)}) - F(Z_{(i-1)}), \quad \forall i \in [n],$$
(5)

where  $-\infty = Z_{(0)} \leq Z_{(1)} \leq Z_{(2)} \leq \cdots \leq Z_{(n)}$  denotes an increasing permutation of  $Z_1, \ldots, Z_n$  called the *order statistics* of Z and with  $Z_{(0)}$  added for convenience.

The coverages will be used to bound the expectation. To do so, we require the following cone in  $\mathbb{R}^n$ :

$$\mathcal{M}^{n} = \left\{ x \colon \sum_{i=1}^{k} x_{i} \ge 0, \forall k \in [n-1], \sum_{i=1}^{n} x_{i} = 0 \right\}.$$
 (6)

In Lem. A.1 we prove that it corresponds to the polar of the monotone cone  $\mathbb{R}^n_{\uparrow}$ . It has a history in isotonic regression

<sup>2</sup>The mean bound will be uniform when

$$\mathbb{P}\left[\sup_{\mu \in \mathcal{A}} \sum_{i=1}^{n} \mu_{i} \ell_{i}(\theta) \geq \mathbb{E}[\ell(\theta, \xi)], \, \forall \theta \in \Theta\right] \geq 1 - \delta$$

[21], majorization [22] and also describes a stochastic order between discrete distributions [15, p. 4] as we illustrate later.

The expectation bound is then as follows:

**Proposition III.1.** Suppose that  $\mathcal{A} \subseteq \Delta^n$  is permutation invariant<sup>3</sup>. Take W as in (5). Then, for any  $rv Z: \Omega \to \mathbb{R}$  with iid samples  $\{Z_i\}_{i=1}^{n-1}$  and  $Z_n = \text{esssup}[Z] < +\infty$ ,

$$\mathbb{P}\left[\sup_{\mu\in\mathscr{A}}\sum_{i=1}^{n}\mu_{i}Z_{i}\geq\mathbb{E}[Z]\right]\geq\mathbb{P}[W\in\mathscr{A}+\mathcal{M}^{n}].$$

*Proof.* We split up the expectation<sup>4</sup> as follows:

$$\mathbb{E}[Z] = \int_{-\infty}^{Z_{(n)}} z \mathrm{d}F(z) = \sum_{i=1}^{n} \int_{Z_{(i-1)}}^{Z_{(i)}} z \mathrm{d}F(z), \quad (7)$$

with  $Z_{(1)} \leq \cdots \leq Z_{(n-1)}$  the order statistics and  $Z_{(n)} = \text{esssup}[Z]$ . The first equality follows from  $\mathbb{P}[Z \leq Z_{(n)}] = 1$  and the second from [23, Thm. 16.9]. For each term  $\int_{Z_{(i-1)}}^{Z_{(i)}} z dF(z) \leq Z_{(i)} \int_{Z_{(i-1)}}^{Z_{(i)}} dF(z) = Z_{(i)}(F(Z_{(i)}) - F(Z_{(i-1)}))$ . Hence  $\mathbb{E}[Z] \leq \sum_{i=1}^{n} Z_{(i)}W_i$ .

Condition on  $W \in \mathcal{A} + \mathcal{M}^n$ . Then,

$$\mathbb{E}[Z] \le \sum_{i=1}^{n} W_i Z_{(i)} \le \sup_{\mu \in \mathcal{A} + \mathcal{M}^n} \sum_{i=1}^{n} \mu_i Z_{(i)}.$$

The expression on the right can be simplified by noting that

$$\sup_{\mu \in \mathcal{A} + \mathcal{M}^n} \sum_{i=1}^n \mu_i Z_{(i)} = \sup_{\mu \in \mathcal{A}} \sum_{i=1}^n \mu_i Z_{(i)},$$

where we use  $(Z_{(1)}, \ldots, Z_{(n)}) \in \mathbb{R}^n_{\uparrow}$ , the definition of the polar cone and Lem. A.1, which implies  $\sum_{i=1}^n s_i Z_{(i)} \leq 0$  for any  $s \in \mathcal{M}^n$  with equality for s = 0.

Finally, let  $\pi \in \Pi^n$  be the permutation such that  $Z_{\pi(i)} = Z_{(i)}$  for  $i \in [n]$  and let  $\pi^{-1}$  denote its inverse (which exists, since permutations are bijections). Then

$$\sup_{\mu \in \mathscr{A}} \sum_{i=1}^{n} \mu_i Z_i = \sup_{\mu \in \mathscr{A}} \sum_{i=1}^{n} \mu_{\pi^{-1}(i)} Z_{(i)}$$
$$= \sup_{\pi \mu \in \mathscr{A}} \sum_{i=1}^{n} \mu_i Z_{(i)} \stackrel{(i)}{=} \sup_{\mu \in \mathscr{A}} \sum_{i=1}^{n} \mu_i Z_{(i)}$$

where (i) uses permutation invariance. Hence, we showed  $\mathbb{P}[\sup_{\mu \in \mathcal{A}} \sum_{i=1}^{n} \mu_i Z_{(i)} \ge \mathbb{E}[Z]] \ge \mathbb{P}[W \in \mathcal{A} + \mathcal{M}^n].$ 

From Prop. III.1, it is clear that the distribution of W is important. Interestingly, when F is continuous, then W is always uniformly distributed over  $\Delta^n$  [24, Thm. 8.7.4]. For general distributions however, we can still establish a type of *stochastic order* between the two distributions using  $\mathcal{M}^n$ :

**Lemma III.2.** Take  $W = (W_1, \ldots, W_n) \in \Delta^n$  as in (5). Then, for any (Lebesgue measurable)  $\mathcal{A} \subseteq \Delta^n$ ,

$$\mathbb{P}[W \in \mathcal{A} + \mathcal{M}^n] \ge \mathbb{P}[\nu \in \mathcal{A} + \mathcal{M}^n],$$

with  $\nu \cong U[\Delta^n]$  uniformly distributed over  $\Delta^n$ . For continuous cdf we have  $\mathbb{P}[W \in \mathcal{A}] = \mathbb{P}[\nu \in \mathcal{A}]$ .

*Proof.* For continuous cdf we refer to [24, Thm. 8.7.4]. For discontinuous cdf we first introduce a construction of the



Fig. 1: Lower bounds of the cumulative distribution for nonnegative Z with cdf F. The tightest lower bound supported on the samples is depicted in blue, while a feasible lower bound is depicted in red. The green area is the expectation.

joint distribution of random vectors W' and  $\nu'$ , such that  $W' \cong W$  and  $\nu' \cong \nu$  (i.e., the marginals are as specified in the lemma). For this construction, we show that  $\nu' \in \mathcal{A} + \mathcal{M}^n$  implies  $W' \in \mathcal{A} + \mathcal{M}^n$  almost surely. So  $\mathbb{P}[W' \in \mathcal{A} + \mathcal{M}^n] \geq \mathbb{P}[\nu' \in \mathcal{A} + \mathcal{M}^n]$ . From  $W' \cong W$  and  $\nu' \cong \nu$ we then get the required result.

The construction starts by taking  $U_i \cong U[0, 1]$  as uniform random variables, for  $i \in [n-1]$ , with  $U_{(1)} \leq U_{(2)} \leq \cdots \leq U_{(n-1)} \leq U_{(n)} = 1$  the uniform order statistics. Let  $\sum_{i=1}^{k} \nu'_i = U_{(k)}$  for  $k \in [n]$ . Then, by [14, §6.4],  $\nu' \cong U[\Delta^n]$  or  $\nu' \cong \nu$ . Meanwhile, the quantile transform [25, Lem. 1.2.4(i)] states that  $Z'_i = F^{-1}(U_i)$  has cdf F, where  $F^{-1}$  is the quantile function of Z. Since F (and therefore  $F^{-1}$ ) is nondecreasing, we can apply [25, Lem. 1.2.1] to claim that  $F^{-1}(U_{(i)})$  is distributed as the *i*-th order statistic  $Z_{(i)}$ . So with  $Z'_{(i)} = F^{-1}(U_{(i)})$  we take W' analogously to (5). From  $Z'_{(i)} \cong Z_{(i)}$  we then have  $W' \cong W$ .

Both marginals are related as, for  $k \in [n]$ ,

$$\sum_{i=1}^{k} W'_{i} = F(Z'_{(k)})$$
  
=  $F(F^{-1}(U_{(k)})) = F(F^{-1}(\sum_{i=1}^{k} \nu'_{i})),$  (8)

where the first equality follows by summing (5) for  $i \in [k]$ with  $W'_i, Z'_{(i)}$  in place of  $W_i, Z_i$ , the second by construction of  $Z'_{(k)}$  and the third by construction of  $\nu'$ . Note that, for general distributions, we have  $F(F^{-1}(p)) \ge p$  for all  $p \in$ [0, 1] [26, Ex. 3.2], with strict inequality iff  $p \in (0, 1)$  is not in the range of F. Applying this to (8) gives

$$\sum_{i=1}^{k} \nu_i' \le \sum_{i=1}^{k} W_i', \, \forall k \in [n-1] \text{ and } \sum_{i=1}^{n} \nu_i' = \sum_{i=1}^{n} W_i'.$$
(9)

Observe how (9) corresponds to a conic inequality under  $\mathcal{M}^n$ . The inequalities are strict iff  $\sum_{i=1}^k \nu'_i$  is not in the range of F (i.e., it lies in a discontinuous jump of F). In that sense, (9) models the gap between  $U[\Delta^n]$  and the coverages W.

To complete the proof, assume that  $\nu' \in \mathcal{A} + \mathcal{M}^n$ . By definition of the Minkowski sum, this is equivalent to there being some  $\mu \in \mathcal{A}$  such that  $\nu - \mu \in \mathcal{M}^n$  or,  $\sum_{i=1}^k \mu_i \leq \sum_{i=1}^k \nu'_i$  for  $k \in [n-1]$  and  $\sum_{i=1}^n \mu_i = \sum_{i=1}^n \nu'_i$  (cf. (6)). Thus, from (9), we have  $\sum_{i=1}^k \mu_i \leq \sum_{i=1}^k W'_i$  for  $k \in [n-1]$  and  $\sum_{i=1}^n \mu_i = \sum_{i=1}^n W'_i$ . By definition of  $\mathcal{A} + \mathcal{M}^n$  and (6) this shows that  $W' \in \mathcal{A} + \mathcal{M}^n$ . So, by our arguments at the start of the proof, we showed the required result.

<sup>&</sup>lt;sup>3</sup>Specifically, for any  $\pi \in \Pi^n$  and  $\mu \in \mathcal{A}$ ,  $\pi \mu = (\mu_{\pi(1)}, \dots, \mu_{\pi(n)}) \in \mathcal{A}$ . Moreover we assume sets  $\mathcal{A}$  are Lebesgue measurable.

<sup>&</sup>lt;sup>4</sup>For details on the integral notation see [23, Eq. 17.22].

The result in Lem. III.2 can be interpreted in terms of lower bounding the cdf of Z. To illustrate this, we introduce the weighted empirical cdf

$$F^{n}_{\mu}(x) = \sum_{i=1}^{n} \mu_{i} \mathbb{1}_{[Z_{(i)}, +\infty)}, \tag{10}$$

where  $\mu \in \Delta^n$  and  $\mathbb{1}_{[Z_{(i)},+\infty)}(z) = 1$  when  $z \ge Z_{(i)}$  and zero otherwise. Note that  $W \in \{\mu\} + \mathcal{M}^n$  holds iff

$$F^{n}_{\mu}(Z_{(i)}) = \sum_{i=1}^{k} \mu_{i} \le \sum_{i=1}^{k} W_{i} = F(Z_{(i)}), \, \forall k \in [n],$$

with  $\sum_{i=1}^{n} \mu_i = \sum_{i=1}^{n} W_i = 1$ . This relationship relates to (9) and (6) and implies that  $F_{\mu}^n$  should lower bound the cdf F everywhere, as depicted in Fig. 1. This inequality between cdfs is the usual stochastic order [15, §1.A.1]. The event  $W \in \mathcal{A} + \mathcal{M}^n$  is then equivalent to the existence of a cdf  $F_{\mu}^n$  with weights  $\mu \in \mathcal{A}$  that lower bounds the true cdf. In terms of this interpretation, Prop. III.1 follows from the fact that a lower bound on the cdf implies an upper bound on the expectation (cf. [19], [27, Eq. 1.A.5]).

By combining Prop. III.1 and Lem. III.2 we directly prove the main contribution of this paper:

**Theorem III.3.** Assuming the setting of Prop. III.1 and taking  $\nu \cong U[\Delta^n]$ , then

$$\mathbb{P}\left[\sup_{\mu\in\mathscr{A}}\sum_{i=1}^{n}\mu_{i}Z_{i}\geq\mathbb{E}[Z]\right]\geq\mathbb{P}[\nu\in\mathscr{A}+\mathcal{M}^{n}].$$

Note that, the important requirement of  $\mathcal{A}$  in Thm. III.3 is that it is permutation invariant and a subset of the probability simplex. The resulting support functions are related to *lawinvariant, coherent risk measures* in literature [28, §6.3.5], [6], where the *law* in our case is a permutation of the random vector. The  $\phi$ -based ambiguity sets considered below are the most frequently studied case of such risk measures.

#### IV. CALIBRATION PROBLEM

This section considers calibrating ambiguity sets  $\mathcal{A}$  such that (3) holds. To do so, we first consider the  $\phi$ -divergence parametrization in (4) and try to upper bound the smallest  $\alpha$  such that (3) still holds for  $\mathcal{A} = \mathcal{A}_{\alpha}$ , which we denote as  $\alpha_{\star}$ . Later we also provide an alternative parametrization similar to the conditional value-at-risk (cf. Cor. IV.3 and the discussion below).

We can use the previous result to simplify (3). Note that, for a fixed  $\theta$ ,  $\ell(\theta, \xi)$  is simply a scalar random variable, which we will denote as Z. So we consider

$$\inf \left\{ \alpha \colon \mathbb{P}\left[ \sup_{\mu \in \mathscr{A}_{\alpha}} \sum_{i=1}^{n} \mu_{i} Z_{i} \ge \mathbb{E}[Z] \right] \ge 1 - \delta, \, \forall Z \right\}, \quad (11)$$

where we inherit the notation from the previous section. It's solution will upper bound  $\alpha_{\star}$ . Note that all previous results were distribution-free. They hold for all (bounded) random variables Z, invariant of their underlying distribution. As such, by using Thm. III.3, the constraint in (11) can be conservatively approximated by

$$\mathbb{P}\left[\nu \in \mathcal{A}_{\alpha} + \mathcal{M}^{n}\right] \ge 1 - \delta, \quad \text{with } \nu \cong \mathrm{U}[\Delta^{n}]. \tag{12}$$

We use this to approximate the calibration problem:

**Proposition IV.1.** Let  $I_{\phi}$  denote the  $\phi$ -divergence, with  $\mathcal{A}_{\alpha}$  the associated ambiguity set as in (4). Let  $\alpha_{\star}$  denote the smallest  $\alpha$  such that (3) holds for  $\mathcal{A}_{\alpha}$ . Then

$$\alpha_{\star} \le \inf \left\{ \alpha \in \mathbb{R} \colon \mathbb{P} \left[ I_{\phi}^{\diamond}(\nu) \le \alpha \right] \ge 1 - \delta \right\}, \qquad (13)$$

with  $\nu \cong U[\Delta^n]$  and, for  $\phi^*(s) = \sup_{t \ge 0} \{ts - \phi(t)\},\$ 

$$I_{\phi}^{\diamond}(\nu) := \sup_{\lambda \in \mathbb{R}^{n}_{\uparrow}} \left\{ \sum_{i=1}^{n} \nu_{i} \lambda_{i} - \frac{1}{n} \phi^{*}(\lambda_{i}) \right\}.$$
(14)

*Proof.* Note that  $\nu \in \mathcal{A}_{\alpha} + \mathcal{M}^{n}$  holds iff there is some  $\mu \in \mathcal{A}_{\alpha}$  such that  $\nu - \mu \in \mathcal{M}^{n}$ . Equivalently, there should exist a  $\mu \in \Delta^{n}$ , which satisfies  $I_{\phi}(\mu, \mathbb{1}_{n}/n) \leq \alpha$  and  $\nu - \mu \in \mathcal{M}^{n}$ . We already showed that the smallest  $\alpha$  satisfying (12) upper bounds  $\alpha_{\star}$ . The left-hand side of (12) equals

$$\mathbb{P}\left[\inf_{\mu\in\Delta^n}\left\{\sum_{i=1}^n\frac{1}{n}\phi(n\mu_i)\colon\nu-\mu\in\mathcal{M}^n\right\}\leq\alpha\right].$$

Since  $\nu \cong U[\Delta^n]$ , the constraint  $\nu - \mu \in \mathcal{M}^n$ , together with  $\sum_{i=1}^n \nu_i = 1$  implies (cf. (6))  $\sum_{i=1}^n \mu_i = 1$ . So the infimum can be taken over  $\mathbb{R}^n_+$ . We place the constraint inside of the cost by noting that the indicator function of a polar cone equals the support function of its dual [20, Ex. 2.26]. So we can rewrite (12) as follows:

$$\inf_{\mu \in \mathbb{R}^n_+} \sup_{\lambda \in \mathbb{R}^n_+} \left\{ \sum_{i=1}^n \frac{1}{n} \phi(n\mu_i) + \langle \lambda, \nu - \mu \rangle \right\}, \\
\leq \sup_{\lambda \in \mathbb{R}^n_+} \sum_{i=1}^n \inf_{t_i \ge 0} \left\{ \phi(t_i) - t_i \lambda_i \right\} / n + \lambda_i \nu_i,$$

The inequality follows by weak duality, the substitution  $t_i = n\mu_i$  and separability. The infima inside the sum are  $-\phi^*(\lambda_i)$ , completing the proof by the argument preceding (12).

Observe that the right-hand side of (13) is essentially the  $1 - \delta$  quantile of the scalar random variable  $I_{\phi}^{\diamond}(\nu)$ . Unfortunately its distribution is unknown. However, we can sample from it. To do so note that  $\nu \cong U[\Delta^n]$  is Dirichlet distributed with parameters  $(1, \ldots, 1) \in \mathbb{R}^n$  [14, §6.4]. So we can sample it by either sampling from a Dirichlet distribution, or by sampling n - 1 uniform order statistics  $U_{(1)} \leq U_{(2)} \leq \cdots \leq U_{(n-1)}$  and using  $(U_{(1)}, U_{(2)} - U_{(1)}, \ldots, 1 - U_{(n-1)})$  as a sample from  $U[\Delta^n]$  (cf. [14, §6.4]). We can then evaluate  $I_{\phi}^{\diamond}(\nu)$  by solving (14), which is a special case of the optimization problem in [29]. That paper presents the *pool-adjacent violator (PAV)* algorithm. It solves (14) with complexity  $\mathcal{O}(n)$ .

The fact that we can sample  $I_{\phi}^{\diamond}(\nu)$  efficiently implies that (11) can be estimated through data-driven means.

**Theorem IV.2.** Let  $\alpha_1, \ldots, \alpha_m$  denote m iid samples from  $I_{\phi}^{\diamond}(\nu)$  with  $\nu \cong U[\Delta^n]$  and  $\alpha_{(1)} \leq \cdots \leq \alpha_{(m)}$  the associated order statistics. Then, for any  $\delta \in [0, 1]$  and  $k \in [m]$ ,

$$\mathbb{P}\left[\alpha_{\star} \le \alpha_{(k)}\right] \ge 1 - \beta \tag{15}$$

with  $\alpha_{\star}$  as in Prop. IV.1, and  $\beta = I_{1-\delta}(k, m-k+1)$  the regularized incomplete beta function (i.e., the cdf of a beta distribution) at level  $1 - \delta$ .

*Proof.* The result follows directly from Prop. IV.1 and a one-sided data-driven bound of a quantile stated in [30, §G.2.2]

applied to (13). Alternatively (13) can be interpreted as a *scenario program*, for which [31, Thm. 3.7] holds.  $\Box$ 

In practice, the user would select m (the number of samples from  $I_{\phi}^{\diamond}(\nu)$  computed using PAV). A larger value gives a tighter upper bound for  $\alpha_{\star}$  at the cost of additional computation time. The confidence level is then determined by fixing some  $\beta \in [0, 1]$  and then finding the smallest k such that  $I_{1-\delta}(k, m-k+1) \leq \beta$ . Such a k is determined with a scalar root finder<sup>5</sup>.

We can establish a connection with other results in literature, through the following corollary of Thm. III.3:

**Corollary IV.3.** Let  $\mu \in \mathbb{R}^n_{\uparrow} \cap \Delta^n$ . Then, for uniform order statistics  $U_{(1)} \leq U_{(2)} \leq \cdots \leq U_{(n-1)} \leq U_{(n)} = 1$  and  $Z_{(1)} \leq Z_{(2)} \leq \ldots Z_{(n-1)}$  and  $Z_{(n)} = \text{esssup}[Z]$  the order statistics of  $rv Z \colon \Omega \to \mathbb{R}$ ,

$$\mathbb{P}\left[\sum_{i=1}^{n} \mu_i Z_{(i)} \ge \mathbb{E}[Z]\right] \ge \mathbb{P}\left[\sum_{i=1}^{k} \mu_i \le U_{(k)}, \, \forall k \in [n]\right].$$
(16)

 $\begin{array}{l} \textit{Proof. Using [6] gives } \sum_{i=1}^{n} \mu_i Z_{(i)} = \sup_{\mu \in \mathscr{A}} \sum_{i=1}^{n} \mu_i Z_i, \\ \textit{with } \mathscr{A} = \operatorname{co} \Pi^n \mu \textit{ the convex hull of all permutations of } \mu. \\ \textit{Invoking Thm. III.3 then gives } \mathbb{P}\left[\sum_{i=1}^{n} \mu_i Z_{(i)} \geq \mathbb{E}[Z]\right] \geq \\ \mathbb{P}[\nu \in \mathscr{A} + \mathcal{M}^n], \textit{ with } \nu \textit{ distributed according to } U[\Delta^n]. \\ \textit{Assume that } \sum_{i=1}^{k} \mu_i \leq U_{(k)} \textit{ for all } k \in [n-1]. \textit{ Noting that} \\ \sum_{i=1}^{k} \nu_i \cong U_{(k)} \textit{ by } [14, \S 6.4] \textit{ and using (6), this implies} \\ \textit{that } \nu - \mu \in \mathcal{M}^n \textit{ and, since } \mu \in \mathscr{A}, \nu \in \mathscr{A} + \mathcal{M}^n. \textit{ So} \\ \mathbb{P}[\nu \in \mathscr{A} + \mathcal{M}^n] \geq \mathbb{P}[\sum_{i=1}^{k} \mu_i \leq U_{(k)}, \forall k \in [n]]. \end{array}$ 

Evaluating the right-hand side of (16) was studied in the context of lower bounding the cdf of a random variable (cf. [32] for an efficient and numerically stable algorithm).

The expression  $\sum_{i=1}^{n} \mu_i Z_{(i)}$  is called a *distortion risk* [6], which is a convex function of  $(Z_1, \ldots, Z_n)$  iff  $\mu \in \mathbb{R}^n_{\uparrow}$ . Since the uniform order statistics form a uniform distribution over a convex set (i.e.,  $\mathbb{R}^n_{\uparrow} \cap [0, 1]^n$ ), their density is quasiconcave [28, Ex. 4.10]. Hence we can leverage [28, Ex. 4.17, Cor. 4.42] to claim that the calibration problem

$$\inf_{\mu \in \mathbb{R}^n_{\uparrow} \cap \Delta^n} \left\{ \varphi(\mu) \colon \mathbb{P}\left[ \sum_{i=1}^k \mu_i \le U_{(k)}, \, \forall k \in [n] \right] \ge 1 - \delta \right\}$$

is a convex problem, whenever  $\varphi$  is convex. We will study the choice of  $\varphi$ , which should measure the size of  $\mathcal{A}$ , in future work as well as the tractability of the calibration problem. For now, we instead consider a specific parametrization:

$$\mu^{(\gamma)} := \left(0, \dots, 0, \frac{\lceil (n-1)\gamma \rceil}{n-1} - \gamma, \frac{1}{n-1}, \dots, \frac{1}{n-1}, \gamma\right), \quad (17)$$

with  $\gamma \geq 1/(n-1),$  such that  $\mu \in \Delta^n \cap {\rm I\!R}^n_{\uparrow}$  holds. Then

$$\sum_{i=1}^{n} \mu_i^{(\gamma)} Z_{(i)} = \left(\frac{d}{n-1} - \gamma\right) Z_{(d)} + \sum_{i=d+1}^{n-1} \frac{Z_{(i)}}{n-1} + \gamma Z_{(n)}, \quad (18)$$

with  $d = \lceil (n-1)\gamma \rceil$ . The final expression is a well-known bound for the expectation of Z due to Anderson [19]. The

re-interpretation in terms of a distortion risk is novel to the authors' knowledge. In §A-B we also show that the final expression is an affine transformation of the well-known conditional value-at-risk. Hence we denote it as  $\overline{\text{CV}@R}$  in the experiments. The value of  $\gamma$  in the context of [19] is usually determined using an asymptotic bound (cf. [24, Thm. 11.6.2]). We find an accurate value by solving:

$$\inf_{\gamma \ge 1/(n-1)} \left\{ \gamma \colon \mathbb{P}\left[ \sum_{i=1}^{k} \mu_i^{(\gamma)} \le U_{(k)}, \, \forall k \in [n] \right] \ge 1 - \delta \right\},\tag{19}$$

for a user-specified  $\delta \in [0, 1]$  with a scalar root finder<sup>5</sup>. The probability is evaluated numerically using [32]. The calibration problem (19) in [19] is interpreted as moving the empirical cdf down as little as possible, while still guaranteeing that it lower bounds the true cdf with high probability. It is also comparable to producing the tightest mean bound as in the calibration problem (11).

# V. CASE STUDIES

To illustrate the validity and potential of our method we provide several simple case studies. These are convex for maximum interpretability, as in the non-convex case a worse generalization performance might be caused by local optima. Nonetheless our method is also applicable in non-convex settings, where stochastic gradient descent methods can be used (cf. [33] for simple distortions and [18] for divergences). In the convex case we use duality to reformulate the proxy cost in (2). See [6], [34] for details.

We present problems of the form (1) and employ ordered risk minimization (2) or using (18), which we refer to as  $\overline{\text{CV}@R}$ . For divergences we use either the *total variation* (*TV*) or *Kullback-Leibler* (*KL*) and the radius is calibrated using Thm. IV.2 (with  $\beta = 0.005$  and  $m = 10\,000$ ). The value of  $\gamma$  in (18) is calibrated using (19).

## A. Newsvendor

We begin with a toy problem, illustrating the behavior of our method in low-sample settings. Let  $\xi: \Omega \to \mathbb{R}$  be Beta distributed with  $\alpha = 0.1$ ,  $\beta = 0.2$ , scaled by a factor  $\overline{D} := 100$ . Consider a newsvendor problem [28, §1.2.1]:

$$\underset{\theta \in \mathbb{R}}{\text{minimize}} \quad \mathbb{E}\left[\underline{c\theta + b[\xi - \theta]_+ + h[\theta - \xi]_+}\right],$$

with b = 14, h = 2 and c = 1. For samples  $\{\xi_i\}_{i=1}^{n-1}$  with n = 20 let  $\ell_i(\theta) = \ell(\theta, \xi_i)$  for  $i \in [n-1]$ ,  $\ell_n(\theta) = \max\{(c-b)\theta + b\overline{D}, (c+h)\theta\}$  a robust upper bound. We replace the expectation by a data-driven proxy as described at the start of the section. For the sample average approach (SAA) we take  $\sum_{i=1}^{n-1} \ell_i(\theta)/(n-1)$ .

The calibration problems are solved for  $\delta = 0.2$ . Their performance is compared over 200 sampled data sets in Fig. 2. The left plot shows the actual expected cost for the minimizers. The blue dashed line is the true optimum of §V-A. See [28, §1.2.1] for details on how to compute these values. Note how the SAA performs decently in the median, but has significantly more variance. The outliers above 240 were omitted, the largest of which was 428.2. Moreover,

<sup>&</sup>lt;sup>5</sup>We use brentq as implemented in scipy 1.10.0



Fig. 2: Box plots showing newsvendor expected cost (left); and difference between the predicted cost and expected cost (right). The colored area is the *inter-quartile range* (IQR), while the whiskers show the range of samples truncated to 1.5 times the IQR. Outliers outside of this range are depicted as diamonds. The red dashed lines depict the robust performance. The blue dashed line is the optimal cost.



Fig. 3: Regression using n = 50 samples with d = 20 and  $\lambda = 0.2$  for different risk measures.

the right plot depicts the difference between the optimum value of the proxy cost, and the true cost. The SAA often underestimates its true cost, while our methods overestimate it. The dashed red line depicts the behavior when taking  $\mathcal{A} = \Delta^n$  in (2) (cf. [28, Eq. 1.9]). As we almost never perform worse than this robust method, this shows that our methods learn from data without over-fitting on the sample.

In large sample cases, we can use the largest sample as an approximation of  $\ell_n(\theta)$ . This heuristic is similar to the one used in the scenario approach, the consequences of which have been studied in detail (cf. [35]). In combination with some regularization, this significantly boosts the performance of our method, as shown in the next examples.

## B. Regression

Let  $T_k \colon \mathbb{R} \to \mathbb{R}$  denote the Chebychev polynomials of the first kind for  $k \ge 0$  and  $f_d(x) = (T_k(x))_{k=0}^d \in \mathbb{R}^{d+1}$  a feature vector. Consider a lasso regression problem:

$$\underset{\theta \in \mathbb{R}^{d+1}}{\text{minimize}} \qquad \mathbb{E}\left[ \left( \langle f_d(X), \theta \rangle - Y \right)^2 \right] + \lambda \|\theta\|_1.$$
 (20)

Assuming access to samples  $\{(X_i, Y_i)\}_{i=1}^n$ , we replace the expectation with the proxy costs described above, where  $\ell_i(\theta) = (\langle f_d(X_i), \theta \rangle - Y_i)^2$  for  $i \in [n]$ . So we approximate the robust term with the largest sample.

For the parameters  $\theta_{\star} = (0, 0, 0.2, 0.5, 1.0)$  the data is generated as  $Y_i = \langle f_4(X_i), \theta_{\star} \rangle + E_i$  with  $X_i \cong \mathcal{U}(-1, 1)$ and  $E_i \cong \mathcal{U}(-0.2, 0.2)$  for  $i \in [n]$ . We over-parametrize the problem, taking d = 20, to illustrate the regularizing effect of our method. A fit is plotted for  $\lambda = 0.2$  and n = 50in Fig. 3. Note how the risk measures all perform similarly, while SAA has a worse fit.

			SAA	TV	$\overline{\mathrm{CV@R}}$	KL
ε	d	$\lambda$				
0.05	10	0.001	0.019(03)	0.018(02)	0.018(02)	0.018(02)
		0.01	0.018(05)	0.017(04)	0.017(04)	0.017(04)
		0.05	0.023(05)	0.017(05)	0.017(05)	0.018(05)
	20	0.001	0.023(03)	0.023(04)	0.023(04)	0.023(03)
		0.01	0.019(05)	0.019(04)	0.019(04)	0.019(04)
		0.05	0.024(05)	0.018(05)	0.018(05)	0.018(05)
0.2	10	0.001	0.019(03)	0.019(02)	0.019(02)	0.018(02)
		0.01	0.018(05)	0.017(04)	0.017(04)	0.017(04)
		0.05	0.023(05)	0.018(05)	0.018(05)	0.018(05)
	20	0.001	0.023(03)	0.023(03)	0.023(03)	0.023(03)
		0.01	0.019(05)	0.019(04)	0.019(04)	0.019(04)
		0.05	0.024(05)	0.018(04)	0.018(04)	0.018(04)

TABLE I: Regression generalization performance for various tuning parameters. Values are reported as *mean (standard deviation*  $\cdot$  10<sup>3</sup>) computed over 10 training sets. The same 10 sets were used for every selection of parameters and method. Note that  $\varepsilon$  does not affect SAA.

The methods are evaluated quantitatively by sampling an additional 100 000 data points and computing a sample approximation of the cost of (20). The resulting performance is compared for several tunings in Tbl. I, where any parameters not mentioned are kept as specified above. It is of note that our methods are significantly less sensitive to tuning parameters compared to the SAA. In fact, our methods outperform SAA for all tunings investigated.

#### C. Support Vector Machines

Consider a classification problem with  $X \cong \mathcal{N}(0, I_2)$ normally distributed and Y = 1 if  $X_1 X_2 \ge 0$  and Y = -1otherwise. A Support Vector Machine (SVM) solves:

$$\underset{(f,b)\in\mathcal{H}\times\mathbb{R}}{\text{minimize}} \qquad \frac{1}{2} \|f\|_{\mathcal{H}}^2 + \lambda \mathbb{E} \left[1 - Y(f(X) - b)\right]_+$$

with  $\lambda > 0$  and  $\mathcal{H}$  some *reproducing kernel Hilbert Space* (*RKHS*) [36, Def. 2.9]. The resulting classifier is then given by sign(f(X) - b). Henceforth  $\mathcal{H}$  is the RKHS associated with the *radial basis function* kernel [36, §2.3] with some standard deviation  $\sigma$ . Solving the primal problem is difficult for two reasons: (*i*) the true expectation is often unknown; (*ii*) optimizing over the infinite dimensional  $\mathcal{H}$  is intractable in general. We resolve (*i*) by replacing the expectation with a proxy-cost as described above and (*ii*) through the usual duality trick [36, §7.4]. Details are deferred to §B.

The proxy cost of three of the risks above – SAA, TV and  $\overline{\text{CV@R}}$  – is a maximum of linear functions and the dual problem is a QP. The sample average – C-SVC in [36, §7.5] – is the usual choice. We illustrate the superior performance our calibrated risks.

In Fig. 4, the three classifiers produced by the three proxy costs above are depicted. Note how both TV and  $\overline{\text{CV@R}}$  perform similarly and both visibly better than the usual SAA. Quantitative performance is compared through the fraction of incorrectly labeled samples in a test set of  $10^5$  samples, which we refer to as the misclassification rate. The performance is compared for several tunings in Tbl. II, where



Fig. 4: SVM classifiers trained using n = 250 samples with  $\sigma = 0.25$  and  $\lambda = 10^4$  for different risks. The red and blue markers are samples for Y = 1 and -1 respectively. The line is the decision boundary and the color axis depicts f(X) - b.



Fig. 5: SVM misclassification rates for varying sample counts n. The center line depicts the mean, while the intervals depicts the empirical 0.2-confidence interval.

any parameters not mentioned are kept as specified above. It is of note that our methods are significantly less sensitive to tuning parameters compared to the SAA. In fact, even for the tunings where SAA performs best, our methods perform better for the same tuning, for reasonable choices of  $\delta$ .

We can also examine the effect of varying the sample count *n*. For each such value we train the classifiers, again using the parameters used to produce Fig. 4, for 30 training sets. The resulting misclassification rates are depicted in Fig. 5. Again note that  $\overline{\text{CV}@R}$  and TV both outperform SAA.

#### REFERENCES

- S. Shalev-Shwartz and S. Ben-David, Understanding Machine Learning. Cambridge University Press, 2014.
- [2] V. Vapnik, Statistical Learning Theory. John Wiley & Sons, 1998.
- [3] J. Royset and R. Wets, An optimization primer. Springer, 2022.
- [4] P. Mohajerin Esfahani and D. Kuhn, "Data-driven distributionally robust optimization using the Wasserstein metric: performance guarantees and tractable reformulations," *Mathematical Programming*, vol. 171, no. 1-2, pp. 115–166, 2018.
- [5] A. Ben-Tal, D. Den Hertog, A. De Waegenaere, B. Melenberg, and G. Rennen, "Robust solutions of optimization problems affected by uncertain probabilities," *Management Science*, vol. 59, no. 2, pp. 341– 357, 2013.
- [6] D. Bertsimas and D. B. Brown, "Constructing Uncertainty Sets for Robust Linear Optimization," *Operations Research*, vol. 57, no. 6, pp. 1483–1495, 2009.
- [7] H. Rahimian and S. Mehrotra, "Distributionally Robust Optimization: A Review," 2019, arXiv: 1908.05659.
- [8] F. Lin, X. Fang, and Z. Gao, "Distributionally Robust Optimization: A review on theory and applications," *Numerical Algebra, Control & Optimization*, vol. 12, no. 1, p. 159, 2022.
- [9] E. Delage and Y. Ye, "Distributionally Robust Optimization Under Moment Uncertainty with Application to Data-Driven Problems," *Operations Research*, vol. 58, no. 3, pp. 595–612, 2010.
- [10] M. Lorenzen, F. Dabbene, R. Tempo, and F. Allgower, "Constraint-Tightening and Stability in Stochastic Model Predictive Control," *IEEE TAC*, vol. 62, no. 7, pp. 3165–3177, 2017.
- [11] L. Aolaritei, M. Fochesato, J. Lygeros, and F. Dörfler, "Wasserstein tube MPC with exact uncertainty propagation," 2023, arXiv: 2304.12093.

δ	$\sigma$	$\lambda$	SAA	TV	$\overline{\mathrm{CV@R}}$
0.05	0.05	$10^{4}$	0.474 (0.052)	0.126 (0.100)	0.111 (0.054)
		$10^{6}$	0.170 (0.057)	0.052 (0.035)	0.052 (0.041)
		$10^{8}$	0.053 (0.020)	0.020 (0.008)	0.022 (0.011)
		$10^{10}$	0.039 (0.029)	0.027 (0.013)	0.022 (0.017)
	0.25	$10^{4}$	0.097 (0.047)	0.045 (0.038)	0.038 (0.035)
		$10^{6}$	0.042 (0.012)	0.019 (0.004)	0.016(0.003)
		$10^{8}$	0.022 (0.010)	0.019 (0.004)	0.018 (0.003)
		$10^{10}$	0.067 (0.082)	0.061 (0.114)	0.043 (0.045)
	0.5	$10^{4}$	0.055 (0.021)	0.023 (0.007)	0.023 (0.006)
		$10^{6}$	0.030 (0.008)	0.020 (0.003)	0.020(0.004)
		$10^{8}$	0.023 (0.008)	0.024 (0.007)	0.024 (0.007)
		$10^{10}$	0.202 (0.140)	0.145 (0.110)	0.114 (0.110)
0.1	0.05	$10^{4}$	0.474 (0.052)	0.154 (0.093)	0.145 (0.075)
		$10^{6}$	0.170 (0.057)	0.059 (0.048)	0.060 (0.046)
		$10^{8}$	0.053 (0.020)	0.022 (0.011)	0.021 (0.012)
		$10^{10}$	0.039 (0.029)	0.030 (0.021)	0.030 (0.020)
	0.25	$10^{4}$	0.097 (0.047)	0.046 (0.042)	0.041 (0.036)
		$10^{6}$	0.042 (0.012)	0.015 (0.003)	0.015 (0.003)
		$10^{8}$	0.022 (0.010)	0.020 (0.005)	0.018 (0.004)
		$10^{10}$	0.067 (0.082)	0.076 (0.119)	0.037 (0.039)
	0.5	$10^{4}$	0.055 (0.021)	0.025 (0.009)	0.024 (0.009)
		$10^{6}$	0.030 (0.008)	0.021 (0.004)	0.021 (0.005)
		$10^{8}$	0.023 (0.008)	0.024 (0.007)	0.024 (0.007)
		$10^{10}$	0.202 (0.140)	0.145 (0.097)	0.112 (0.110)
0.2	0.05	$10^{4}$	0.474 (0.052)	0.225 (0.142)	0.197 (0.100)
		$10^{6}$	0.170 (0.057)	0.074 (0.050)	0.075 (0.051)
		$10^{8}$	0.053 (0.020)	0.025 (0.015)	0.026(0.015)
		$10^{10}$	0.039 (0.029)	0.033 (0.023)	0.036 (0.026)
	0.25	$10^{4}$	0.097 (0.047)	0.049 (0.044)	0.045 (0.040)
		$10^{6}$	0.042 (0.012)	0.019 (0.005)	0.016(0.003)
		$10^{8}$	0.022 (0.010)	0.017 (0.004)	0.017 (0.004)
		$10^{10}$	0.067 (0.082)	0.070 (0.113)	0.039 (0.033)
	0.5	$10^{4}$	0.055 (0.021)	0.027 (0.014)	0.027 (0.015)
		$10^{6}$	0.030 (0.008)	0.021 (0.005)	0.020(0.005)
		$10^{8}$	0.023 (0.008)	0.024 (0.007)	0.024 (0.007)
		$10^{10}$	0.202 (0.140)	0.186 (0.144)	0.112 (0.101)

TABLE II: SVM misclassification rates for various tuning parameters. Reported values are the *mean (standard devia-tion)* over 10 training sets. The same 10 sets were used for every parameter selection and method. Note that  $\delta$  does not affect SAA. The lowest values in a column are bold. Observe that in the rows where SAA achieves its best performance, our methods still perform better.

- [12] J. C. Duchi, P. W. Glynn, and H. Namkoong, "Statistics of robust optimization: A generalized empirical likelihood approach," *Mathematics* of Operations Research, vol. 46, no. 3, pp. 946–969, 2021.
- [13] H. Lam, "Recovering best statistical guarantees via the empirical divergence-based distributionally robust optimization," *Operations Research*, vol. 67, no. 4, pp. 1090–1105, 2019.
- [14] H. A. David and H. N. Nagaraja, Order Statistics. John Wiley & Sons, Inc., 2003.
- [15] M. Shaked and J. G. Shanthikumar, *Stochastic Orders*. Springer Series in Statistics, Springer New York, 2007.
- [16] B. P. Van Parys, P. M. Esfahani, and D. Kuhn, "From data to decisions: Distributionally robust optimization is optimal," *Management Science*, vol. 67, no. 6, pp. 3387–3402, 2021.
- [17] R. R. Bahadur and L. J. Savage, "The Nonexistence of Certain Statistical Procedures in Nonparametric Problems," *The Annals of Mathematical Statistics*, vol. 27, no. 4, pp. 1115 – 1122, 1956.
- [18] É. Chouzenoux, H. Gérard, and J.-C. Pesquet, "General risk measures for robust machine learning," *Foundations of Data Science*, vol. 1, no. 3, pp. 249–269, 2019.
- [19] T. W. Anderson, "Confidence limits for the expected value of an arbitrary bounded random variable with a continuous distribution function," tech. rep., Stanford University, 1969.

- [20] A. Beck, First-Order Methods in Optimization. SIAM, 2017.
- [21] R. E. Barlow and H. D. Brunk, "The Isotonic Regression Problem and its Dual," *Journal of the American Statistical Association*, vol. 67, no. 337, pp. 140–147, 1972.
- [22] A. Steerneman, "G-Majorization, group-induced cone orderings, and reflection groups," *Linear Algebra and its Applications*, vol. 127, pp. 107–119, 1990.
- [23] P. Billingsley, Probability and Measure. John Wiley & Sons, Inc., 1995.
- [24] S. S. Wilks, Mathematical Statistics. John Wiley & Sons, Inc., 1964.
- [25] R.-D. Reiss, Approximate Distributions of Order Statistics. Springer Series in Statistics, Springer New York, 1989.
- [26] G. R. Shorack, *Probability for statisticians*. Springer, 2nd ed., 2017.
  [27] G. R. Shorack and J. A. Wellner, *Empirical Processes with Applications to Statistics*. SIAM, 2009.
- [28] A. Shapiro, D. Dentcheva, and A. Ruszczynski, Lectures on Stochastic Programming: Modeling and Theory, Third Edition. SIAM, 2021.
- [29] M. J. Best, N. Chakravarti, and V. A. Ubhaya, "Minimizing Separable Convex Functions Subject to Simple Chain Constraints," *SIAM Journal* on Optimization, vol. 10, no. 3, pp. 658–672, 2000.
- [30] L. A. E. William Q. Meeker, Gerald J. Hahn, *Statistical Intervals: A Guide for Practitioners and Researchers*. Wiley, 2ed. ed., 2017.
- [31] M. C. Campi and S. Garatti, *Introduction to the Scenario Approach*. SIAM, 2018.
- [32] A. Moscovich, "Fast calculation of p-values for one-sided Kolmogorov-Smirnov type statistics," 2020, arXiv: 2009.04954.
- [33] R. Mehta, V. Roulet, K. Pillutla, L. Liu, and Z. Harchaoui, "Stochastic optimization for spectral risk measures," in *Proceedings of The* 26th International Conference on Artificial Intelligence and Statistics, vol. 206, pp. 10112–10159, PMLR, 2023.
- [34] M. Schuurmans and P. Patrinos, "A general framework for learningbased distributionally robust MPC of Markov jump systems," *IEEE Transactions on Automatic Control*, pp. 1–16, 2023.
- [35] F. A. Ramponi, "Consistency of the scenario approach," SIAM Journal on Optimization, vol. 28, no. 1, pp. 135–162, 2018.
- [36] B. Schölkopf, A. J. Smola, F. Bach, et al., Learning with kernels: support vector machines, regularization, optimization, and beyond. MIT press, 2002.
- [37] S. Boyd and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.
- [38] R. T. Rockafellar and S. Uryasev, "Optimization of conditional valueat-risk," *The Journal of Risk*, vol. 2, no. 3, pp. 21–41, 2000.
- [39] H. H. Bauschke and P. L. Combettes, Convex Analysis and Monotone Operator Theory in Hilbert Spaces. Springer New York, 2011.

## APPENDIX A PRELIMINARIES

#### A. Monotone Cone

Let  $\mathbb{R}^n_{\uparrow} := \{x \in \mathbb{R}^n : x_1 \leq x_2 \leq \cdots \leq x_n\}$  denote the monotone cone. This cone and its polar have a history in isotonic regression [21] and majorization [22].

We show that  $\mathcal{M}^n$  and  $\mathbb{R}^n_{\uparrow}$  are related.

**Lemma A.1.** Let  $\mathbb{R}^n_{\uparrow}$  be the monotone cone and  $\mathcal{M}^n$  as in (6). Then  $\mathcal{M}^n$  is the polar of  $\mathbb{R}^n_{\uparrow}$ .

*Proof.* The monotone cone is polyhedral with  $\mathbb{R}^n_{\uparrow} = \{x: Mx \leq 0\}$  for  $M \in \mathbb{R}^{n-1 \times n}$  with  $Mx = (x_1 - x_2, x_2 - x_3, \dots, x_{n-1} - x_n)$ . The definition of the polar cone is thus

$$(\mathbb{R}^n_{\uparrow})^{\circ} = \{ y \in \mathbb{R}^n \colon \langle x, y \rangle \le 0, \forall x \text{ s.t. } Mx \le 0 \}.$$

By Farkas' lemma [37, p. 263] we have either  $Mx \leq 0$  and  $\langle x, y \rangle > 0$  or  $M^{\top}\lambda = y$  and  $\lambda \geq 0$ . So

$$(\mathbb{R}^n_{\uparrow})^{\circ} = \{ y \in \mathbb{R}^n \colon y = M^{\top}\lambda, \lambda \ge 0 \}.$$

Note that  $\langle \lambda, Mx \rangle = \sum_{i=1}^{n-1} \lambda_i (x_i - x_{i+1}) = \sum_{j=1}^n x_j (\lambda_j - \lambda_{j-1}) = \langle M^\top \lambda, x \rangle$ , where  $\lambda_0 = \lambda_n = 0$ . Thus  $y \in (\mathbb{R}^n_{\uparrow})^\circ$ 

iff  $y = M^{\top}\lambda$  for  $\lambda \ge 0$ . Here  $y = M^{\top}\lambda$  holds iff

$$y_j = \lambda_j - \lambda_{j-1}, \qquad \forall j \in [n].$$
  
$$\Leftrightarrow \quad \sum_{j=1}^k y_j = \sum_{j=1}^k \lambda_j - \lambda_{j-1} = \lambda_k, \quad \forall k \in [n].$$

Since  $\lambda \ge 0$  we have  $\sum_{j=1}^{k} y_j = \lambda_k \ge 0$  for  $k \in [n-1]$  and  $\sum_{j=1}^{n} y_j = \lambda_n = 0$ . These are the constraints in (6).  $\Box$ 

## B. Details on $\overline{\mathrm{CV@R}}$

We first characterize the conditional value-at-risk in terms of order statistics as a *distortion risk* below.

**Lemma A.2.** Consider  $CV@R_n^{\gamma} \colon \mathbb{R}^n \to \mathbb{R}$  as

$$CV@R_n^{\gamma}[X] = \inf_{\tau} \left\{ \tau + \frac{1}{(1-\gamma)n} \sum_{i=1}^n [X_i - \tau]_+ \right\}, \quad (21)$$

for  $\gamma \in [0, 1]$ . Then

$$(1-\gamma)\mathrm{CV}@\mathbf{R}_n^{\gamma}[X] = \left(\frac{d}{n} - \gamma\right)X_{(d)} + \sum_{i=d+1}^n \frac{X_{(i)}}{n},$$

with  $d := \lceil n\gamma \rceil$ . So CV@R<sup> $\gamma$ </sup><sub>n</sub> is a distortion risk.

*Proof.* Consider the minimizers in the definition of CV@R:

$$\underset{\tau}{\operatorname{argmin}} \left\{ \tau + \frac{1}{(1-\gamma)n} \sum_{i=1}^{n} [X_i - \tau]_+ \right\}.$$

By [38, Thm. 1], this set is a closed bounded interval with the left endpoint being

$$V@\mathbf{R}_{n}^{\gamma}[X] := \inf_{x} \{ x \colon F_{n}(x) \ge \gamma \}$$
$$= \inf_{x} \left\{ x \colon \sum_{i=1}^{n} \mathbf{1}_{(-\infty,x]}(X_{i}) \ge \gamma n \right\} = X_{(d)},$$

with  $F_n$  the empirical cdf, the definition of which we plugged in for the second equality. For the third equality note that the left-hand side counts the number of values  $X_i$  smaller than or equal to x. Assume

$$X_{(d-k-1)} < X_{(d-k)} = X_{(d-k+1)} = \dots = X_{(d)},$$
 (22)

for  $k \ge 0$ . Then clearly there are at least  $d = \lceil n\gamma \rceil > n\gamma$ values smaller than or equal to  $X_{(d-k)}$ . For any  $z < X_{(d-k)}$ there are at most d - k - 1 samples values than or equal. Hence  $\operatorname{V}@R_n^{\gamma}[X] = X_{(d-k)} = X_{(d)}$ .

Plugging into the cost of (21) gives

$$X_{(d)} + \frac{1}{(1-\gamma)n} \sum_{i=1}^{n} [X_i - X_{(d)}]_+$$
  
=  $X_{(d)} + \frac{1}{(1-\gamma)n} \sum_{i=d+1}^{n} (X_{(i)} - X_{(d)}),$  (23)

where we used  $X_{(i)} \leq X_{(d)}$  for  $i \leq d$ . The stated result follows from some basic algebraic manipulation. Finally note that  $d/n - \gamma = (\lceil n\gamma \rceil - n\gamma)/n \leq 1/n$ . So the monotonicity constraint on  $\mu$  in Cor. IV.3 is also satisfied.

We can then rewrite (18) as

$$\overline{\text{CV@R}}_{n}^{\gamma}[X] := (1 - \gamma) \text{CV@R}_{n-1}^{\gamma}[(X_{(i)})_{i=1}^{n-1}] + \gamma X_{(n)}.$$

The associated distortion risk has weights  $\mu^{(\gamma)}$  as in (17). So we need  $\gamma \geq 1/(n-1)$  for  $\mu \in \Delta^n \cap \mathbb{R}^n_{\uparrow}$  to hold. The advantage of  $\overline{\text{CV@R}}$  as a well-calibrated risk measure over CV@R is that additional weight is placed on the largest sample. This often makes the mean bound associated with  $\overline{\text{CV@R}}$  less conservative compared to CV@R for the same confidence level.

# APPENDIX B SUPPORT VECTOR MACHINES

Let  $\mathcal{H}$  be some reproducing kernel Hilbert Space (RKHS) [36, Def. 2.9] with reproducing kernel  $\kappa \colon \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ for some  $d \in \mathbb{N}$ . Here  $\langle f, g \rangle_{\mathcal{H}}$  denotes the inner product associated with  $\mathcal{H}$  and  $||f||_{\mathcal{H}}^2 = \langle f, f \rangle$  for  $f, g \in \mathcal{H}$ . The primal problem for learning a support vector machine is usually given in terms of the hinge loss:

$$\underset{(f,b)\in\mathcal{H}\times\mathbb{R}}{\text{minimize}} \quad \frac{1}{2} \|f\|_{\mathcal{H}}^2 + \lambda \mathbb{E} \left[1 - Y(f(X) - b)\right]_+.$$

with  $\lambda > 0, X: \Omega \to \mathbb{R}^d$  and  $Y: \Omega \to \{-1, 1\}$ . Using the reproducing property of  $\kappa$  (cf. [36, Def. 2.9.1]) we have  $f(X) = \langle \kappa(X, \cdot), f \rangle$ . Given a sample  $\{(X_i, Y_i)\}_{i=1}^n$ , let  $\hat{Y} = \operatorname{diag}(Y_1, \ldots, Y_n) \in \mathbb{R}^{n \times n}$  and  $\hat{K}: \mathcal{H} \to \mathbb{R}^n$  a linear operator such that  $(\hat{K}f)_i = \langle \kappa(X_i, \cdot), f \rangle$ . We do not include a data-point modeling the esssup of the random loss. Instead the largest sample will act as a replacement.

We then replace the expectation with a proxy cost, as in (2). The robustified, data-driven problem then becomes

$$\underset{(f,b)\in\mathcal{H}\times\mathbb{R}}{\text{minimize}} \quad \frac{1}{2} \|f\|_{\mathcal{H}}^2 + \lambda \rho \left[\mathbb{1}_n - \hat{Y}\left(\hat{K}f - b\mathbb{1}_n\right)\right]_+, \quad (24)$$

where  $\rho(x) = \sup_{\mu \in \mathscr{A}} \langle \mu, x \rangle$  is a support function.

Proposition B.1. The value of (24) equals

$$\begin{array}{ll} \underset{(\alpha,\beta)\in\mathbb{R}^{n}_{+}\times\mathbb{R}^{n}_{+}}{\text{maximize}} & \sum_{i=1}^{n}\alpha_{i}-\frac{1}{2}\sum_{i,j=1}^{n}\alpha_{i}\alpha_{j}Y_{i}Y_{j}\kappa(X_{i},X_{j})\\ \text{subj. to} & \sum_{i=1}^{n}\alpha_{i}Y_{i}=0, \ \frac{\alpha+\beta}{\lambda}\in\mathscr{A}. \end{array}$$

Moreover, let  $\alpha^*, \beta^*$  denote the optimizers and  $\mathcal{J} := \{j \in [n]: \alpha_j > 0, \beta_j > 0\}$ . Then

$$\begin{split} f^{\star} &= \sum_{i=1}^{n} \alpha_{i}^{\star} Y_{i} \kappa(X_{i}, \cdot) \\ b^{\star} &= \sum_{i=1}^{n} \alpha_{i}^{\star} Y_{i} \kappa(X_{i}, X_{j}) - Y_{j}, \quad \forall j \in \mathcal{J} \text{ when } \mathcal{J} \neq \emptyset \end{split}$$

are the optimizers<sup>6</sup> of (24) when  $\mathcal{J} \neq \emptyset$ . When  $\mathcal{J} = \emptyset$ , then  $b^*$  can be determined by solving (24), keeping  $f = f^*$  fixed.

Proof. We write (24) with slack variables first

$$\begin{array}{ll} \underset{(f,b,s)\in\mathcal{H}\times\mathbb{R}\times\mathbb{R}^n}{\text{minimize}} & \frac{1}{2}\|f\|_{\mathcal{H}}^2 + \lambda\rho(s) \\ \text{subj. to} & \hat{Y}(\hat{K}f - \mathbb{1}_n b) - \mathbb{1}_n + s \ge 0 \\ & s \ge 0. \end{array}$$

<sup>6</sup>In practice, we pick  $b^*$  as the average of the values over all  $j \in \mathcal{J}$ .

The next step is to apply Lagrangian duality over Hilbert spaces as presented in [39, Prop. 19.18]. We first bring the problem in the standard form:

$$\underset{x \in G}{\text{minimize}} \qquad h(x) + g(Lx),$$

with  $\mathcal{G} := \mathcal{H} \times \mathbb{R} \times \mathbb{R}^n$  a Hilbert space, elements of which we partition as x = (f, b, s). Let  $L : \mathcal{G} \to \mathbb{R}^{2n}$  denote the linear operator defined as

$$Lx = \left(\hat{Y}\hat{K}f - \hat{Y}\mathbb{1}_n b + s, s\right).$$
(26)

Its adjoint – defined implicitly as  $L^* \colon \mathbb{R}^{2n} \to \mathcal{G}$  such that  $\langle Lx, v \rangle = \langle x, L^*v \rangle$  – is given as

$$L^*v = (\hat{K}^*\hat{Y}\alpha, -\mathbb{1}_n^T\hat{Y}\alpha, \alpha + \beta),$$

for  $v = (\alpha, \beta)$  and  $\hat{K}^* \hat{Y} \alpha = \sum_{i=1}^n \alpha_i Y_i \kappa(X_i, \cdot)$ . The functions  $h \colon \mathcal{G} \to \overline{\mathbb{R}}$  and  $g \colon \mathbb{R}^n \to \overline{\mathbb{R}}$  are given as

$$h(x) := \|f\|_{\mathcal{H}}^2/2 + \lambda \rho(s)$$
$$g(Lx) := \iota_{\mathbb{R}^{2n}_+}(Lx - \delta),$$

with  $\delta = (\mathbb{1}_n, 0) \in \mathbb{R}^{2n}$ .

First we prove that strong duality holds, a sufficient condition for which is  $\operatorname{int}(\operatorname{dom} g) \cap L \operatorname{dom} h$  (cf. [39, Thm. 15.23, Prop. 6.19(vii)]). Note that  $\operatorname{dom} g = \mathbb{R}^{2n}_+ + \delta$  and  $\operatorname{dom} h = \mathcal{H} \times \mathbb{R} \times \mathbb{R}^n$ , keeping in mind that  $\operatorname{dom} \rho = \mathbb{R}^n$  since  $\rho$  is coherent. Note that  $0 \in \mathcal{H}$  and  $0 \in \mathbb{R}$  and that L(0,0,s) = (s,s). So  $\mathfrak{D} = \{(s,s): s \in \mathbb{R}^n\} \subset L \operatorname{dom} h$ . Hence  $\operatorname{int}(\operatorname{dom} g) \cap L \operatorname{dom} h \supset (\operatorname{int}(\mathbb{R}^{2n}_+) + \delta) \cap \mathfrak{D} \neq \emptyset$ .

Consider the convex conjugates  $h^*$  and  $g^*$ . Then

$$h^*(\overline{x}) = \|f\|_{\mathcal{H}}^2 / 2 + \iota_{\{0\}}(\overline{b}) + \lambda \rho^*(\overline{s}/\lambda),$$

where we used the definition of the convex conjugate, the partitioning  $\overline{x} = (\overline{f}, \overline{b}, \overline{s}) \in \mathcal{G}$ , [39, Prop. 13.16, Prop. 13.20(i)] and seperability of *h*. Since  $\rho$  is a support function we have  $\rho^* = \iota_{\mathcal{A}}$ . Also, again letting  $v = (\alpha, \beta)$ ,

$$g^*(v) = \iota_{\mathbb{R}^{2n}}(v) + \mathbb{1}_n^\top \alpha$$

by [39, Prop. 13.20(ii)]. The dual problem, the value of which equals minus the primal by strong duality, is then given as [39, Prop. 19.18]:

$$\underset{(\alpha,\beta)\in\mathbb{R}^n\times\mathbb{R}^n}{\text{minimize}} \qquad \frac{1}{2} \|\hat{K}^*\hat{Y}\alpha\|_{\mathcal{H}}^2 - \mathbb{1}_n^\top \alpha$$
subj. to 
$$\mathbb{1}_n^\top \hat{Y}\alpha = 0, \ (\alpha+\beta)/\lambda \in \mathcal{A}, \ (\alpha,\beta) \ge 0$$

where we already integrated the indicator functions in the constraints. After adding the minus sign, this is equivalent to the problem in the theorem.

We next compute the subgradients. Note that,

$$\partial h^*(\overline{x}) = \{\overline{f}\} \times (\mathbb{R} \cap \{\overline{b}\}^{\perp}) \times \mathcal{N}_{\mathcal{A}}(\overline{s}/\lambda)/\lambda, \qquad (27)$$

with  $\mathcal{N}_{\mathcal{C}}$  denotes the *normal cone* for set  $\mathcal{C}$  [39, Def. 6.37] and  $\mathcal{C}^{\perp} = \{u: \langle x, u \rangle = 0, \forall x \in \mathcal{C}\}$  denotes the *orthogonal complement* of  $\mathcal{C}$ . The first term follows from differentiability. The second term follows from [39, Ex. 16.12, Ex. 6.39] which states  $\partial \iota_{\{0\}}(\bar{b}) = \mathcal{N}_{\{0\}}(\bar{b}) = \mathbb{R}^n \cap \{\bar{b}\}^{\perp}$ . The third follows from applying [39, Ex. 16.12, Cor. 16.42] to  $\iota_{\mathcal{A}} \circ \lambda^{-1} I_n$ . Similarly

$$\partial g^*(v) = \mathbb{R}^n_+ \cap \{v\}^\perp + \delta$$
$$= \{u + \delta \colon u \ge 0, \ u^\top v = 0\},$$
(28)

where we apply [39, Cor. 16.38], differentiability of the second term and again [39, Ex. 16.12, Ex. 6.39] to deal with the indicator.

By [39, Prop. 19.17(v), Prop. 19.18(v)] the optimizers  $v^{\star} = (\alpha^{\star}, \beta^{\star})$  and  $x^{\star} = (f^{\star}, b^{\star}, s^{\star})$  must satisfy

$$L^*\alpha^* \in \partial h(x^*) \quad \text{and} \quad -\alpha^* \in \partial g(Lx^*)$$
  
$$\Leftrightarrow \quad x^* \in \partial h^*(L^*\alpha^*) \quad \text{and} \quad Lx^* \in \partial g^*(-\alpha^*),$$

where we used [39, Cor. 16.24] to express the optimality conditions in terms of the subgradients of the conjugates. Plugging in the subgradient determined in (27) gives

$$f^* = \hat{K}^* \hat{Y} \alpha^* = \sum_{i=1}^n \alpha_i^* Y_i \kappa(X_i, \cdot).$$

Comparing with [36, Eq. 7.25] shows that this is the usual SVM solution. For the threshold  $b^*$  we use (28), giving

$$Lx^{\star} \in \partial g^{\star}(-\alpha^{\star})$$

$$\Leftrightarrow \sum_{i=1}^{n} \alpha_{i}^{\star} \left(Y_{i}(f^{\star}(X_{i})-b^{\star})-1+s_{i}^{\star}\right)+\sum_{i=1}^{n} \beta_{i}^{\star} s_{i}^{\star}=0$$

$$Y_{i}(f^{\star}(X_{i})-b^{\star})-1+s_{i} \geq 0, \quad \forall i \in [n],$$

$$s^{\star} \geq 0.$$

Let  $\mathcal{J} = \{j \in [n]: \alpha_j^* > 0, \beta_j^* > 0\}$ . From the above conditions we have  $s_j^* = 0$  for all  $j \in \mathcal{J}$ . Similarly we require

$$Y_j(f^{\star}(X_j) - b^{\star}) - 1 = 0, \quad \forall j \in \mathcal{J}.$$

Using  $Y_j \in \{-1, 1\}$  to mulyiply both sides with  $Y_j$  and the characterization of  $f^*$  above results in

$$b^{\star} = \sum_{i=1}^{n} \alpha_i^{\star} Y_i \kappa(X_i, X_j) - Y_j, \quad \forall j \in \mathcal{J}$$

We can again compare with the classical SVM setting (cf. [36, Eq. 47.32] and the discussion at [36, p. 206]), seeing that the condition is similar. When  $\mathcal{J} = \emptyset$ , we cannot generate trivial constraints on  $b^*$ . In that case, we note the  $f^*$  is still a valid minimizer, thus  $b^*$  must minimize (24) for  $f = f^*$ .  $\Box$ 

For  $\overline{\text{CV}@R}$  as in (18) we can characterize the ambiguity set  $\mathcal{A}$  efficiently in terms of a polyhedral set as in [6]. So (25) is a quadratic program. For divergence-based risk measures the ambiguity set  $\mathcal{A}_{\alpha}$  in (4) is convex so we can implement it directly. It is polyhedral for the total variation.