

1        **STOCHASTIC DECOMPOSITION METHOD FOR TWO-STAGE  
2        DISTRIBUTIONALLY ROBUST LINEAR OPTIMIZATION \***

3        HARSHA GANGAMMANAVAR<sup>†</sup> AND MANISH BANSAL<sup>‡</sup>

4        **Abstract.** In this paper, we present a sequential sampling-based algorithm for the two-stage  
5        distributionally robust linear program (2-DRLP) with general ambiguity set. The algorithm is a  
6        distributionally robust version of the well-known stochastic decomposition algorithm of Higle and  
7        Sen (Math. of OR 16(3), 650-669, 1991) that was designed for risk-neutral two-stage stochastic linear  
8        programs. We refer to the algorithm as the distributionally robust stochastic decomposition (DRSD)  
9        method. The algorithm works with data-driven approximations of ambiguity sets that are constructed  
10      during the course of the algorithm using samples of increasing size. It constructs statistical  
11      approximations of the worst-case expectation function by solving subproblems corresponding to the  
12      latest observation(s) in every iteration. We show that the DRSD method asymptotically identifies an  
13      optimal solution, with probability one, for a family of ambiguity sets that includes the moment-based  
14      and Wasserstein distance-based ambiguity sets. We also computationally evaluate the performance  
15      of the DRSD method for solving distributionally robust variants of instances considered in the sto-  
16      chastic programming literature. The numerical results corroborate our analysis of the DRSD method  
17      and illustrate the computational advantage over an external sampling-based decomposition approach  
18      and reformulation techniques known in the literature.

19        **Key words.** Distributionally robust optimization, stochastic programming, stochastic decom-  
20        position, sequential sampling, cutting-plane method.

21        **AMS subject classifications.** 90C15, 90C06, 90C47

22        **1. Introduction.** Stochastic programming (SP) is a well-known framework for  
23        decision-making under uncertainty that arises in applications such as finance, capac-  
24        ity expansion, manufacturing, wildfire planning, power systems, healthcare, and many  
25        more. The SP models with recourse, particularly in a two-stage setting, have gained  
26        wide acceptance across these application domains. In the two-stage SP models, the  
27        first-stage decision (referred to as the here-and-now decision) is taken before the re-  
28        alization of uncertainty. Following this, the second-stage decision (referred to as the  
29        wait-and-see decision) is taken in response to the first-stage decision and a realization  
30        of the uncertain data. In the classical setting of two-stage stochastic linear programs  
31        (2-SLPs), the decisions are solutions to linear programs in both stages [12].

32        The SP models are stated with an expectation-valued objective function. There-  
33        fore, stating an SP model either requires complete knowledge of the underlying prob-  
34        ability distribution or the ability to simulate observations from this distribution. The  
35        latter leads to the construction of a sample average approximation (SAA) of the  
36        problem. In many practical applications, the distribution associated with random  
37        parameters in the optimization model is not precisely known. It either has to be esti-  
38        mated from data or constructed by expert judgments, which tend to be subjective. In  
39        any case, identifying a distribution using available information may be cumbersome  
40        at best. Stochastic min-max programming that has gained significant attention in re-  
41        cent years under the name of *distributionally robust optimization* (DRO) is intended  
42        to address the ambiguity in distributional information.

---

\*Submitted to the editors on August 9, 2022.

**Funding:** M. Bansal is funded by National Science Foundation Grant CMMI-1824897, which is  
gratefully acknowledged

<sup>†</sup>Department of Engineering Management, Information, and Systems, Southern Methodist Uni-  
versity, Dallas, TX 75275 ([harsha@smu.edu](mailto:harsha@smu.edu)).

<sup>‡</sup>Department of Industrial and Systems Engineering, Virginia Tech, Blacksburg, VA 24061  
([bansal@vt.edu](mailto:bansal@vt.edu)).

43 In this paper, we study a particular manifestation of the DRO problem in the two-  
 44 stage setting, viz., the two-stage distributionally robust linear program (2-DRLP). We  
 45 state this problem as:

46 (1.1) 
$$\min \{f(x) = c^\top x + \mathbb{Q}(x) \mid x \in \mathcal{X}\}.$$

48 Here,  $c$  is the coefficient vector of a linear cost function and  $\mathcal{X}$  is the feasible set  
 49 of the first-stage decision vector. The feasible region takes the form of a compact  
 50 polyhedron, i.e.,  $\mathcal{X} = \{x \in \mathbb{R}^{d_x} \mid Ax \geq b, x \geq 0\}$ , where  $A \in \mathbb{R}^{m_1 \times d_x}$  and  $b \in \mathbb{R}^{m_1}$ .  
 51 The function  $\mathbb{Q}(x)$  is the worst-case expected recourse cost, that we define as:

52 (1.2) 
$$\mathbb{Q}(x) = \max_{P \in \mathfrak{P}} \left\{ \mathbb{Q}(x; P) := \mathbb{E}_P[\mathbb{Q}(x, \tilde{\omega})] \right\}.$$

54 We define the random vector  $\tilde{\omega} \in \mathbb{R}^d$  on a measurable space  $(\Omega, \mathcal{F})$ , where  $\Omega$  is  
 55 a continuous or discrete sample space equipped with the sigma-algebra  $\mathcal{F}$ .  $\mathfrak{P}$  is  
 56 a set of probability distributions defined on the measurable space  $(\Omega, \mathcal{F})$ . The set  
 57 of probability distributions  $\mathfrak{P}$  is referred to as the *ambiguity set*. The expectation  
 58 operation  $\mathbb{E}_P[\cdot]$  is taken with respect to a probability distribution  $P \in \mathfrak{P}$ . For a given  
 59  $x \in \mathcal{X}$ , we refer to the optimization problem in (1.2) as the *distribution separation*  
 60 *problem*. For a given realization  $\omega$  of the random vector  $\tilde{\omega}$  and a first-stage solution  
 61  $x$ , the recourse value in (1.2) is the optimal value of the following second-stage linear  
 62 program:

63 (1.3) 
$$\begin{aligned} Q(x, \omega) := \min & \quad g(\omega)^\top y \\ 64 \quad \text{s.t.} & \quad y \in \mathcal{Y}(x, \omega) := \{W(\omega)y = r(\omega) - T(\omega)x, y \in \mathbb{R}_+^{d_y}\}. \end{aligned}$$

66 Here, for each  $\omega \in \Omega$ , we have uncertain second-stage parameters:  $g(\omega)$ , the recourse  
 67 matrix  $W(\omega)$ , the right-hand side vector  $r(\omega)$ , and the technology matrix  $T(\omega)$  of  
 68 appropriate dimensions. A special case of 2-DRLP is the 2-SLP where  $\mathfrak{P}$  is a singleton,  
 69 i.e.,  $\mathfrak{P} = \{P^*\}$ , resulting in the following optimization problem:

70 (1.4) 
$$\min \{c^\top x + \mathbb{E}_{P^*}[\mathbb{Q}(x, \tilde{\omega})] \mid x \in \mathcal{X}\}.$$

72 Most data-driven and SAA-based approaches to solve 2-SLPs tackle the problem  
 73 in two steps. In the first simulation/sampling step, an uncertainty representation is  
 74 generated using a finite set of observations that serves as an approximation of  $\Omega$  and  
 75 the corresponding empirical distribution serves as an approximation of  $P^*$ . For a given  
 76 uncertainty representation, one obtains a deterministic approximation of (1.4). In  
 77 the second optimization step, the approximate problem is solved using deterministic  
 78 optimization methods. Such a two-step approach may lead to poor out-of-sample  
 79 performance, forcing the entire process to be repeated from scratch with an improved  
 80 uncertainty representation. Since sampling is performed prior to the optimization  
 81 step, this two-step approach is also referred to as the *external sampling procedure*.  
 82 This procedure has also been utilized for solving 2-DRLPs where in the first step, an  
 83 approximation of the ambiguity set  $\mathfrak{P}$  is obtained using a finite set of observations.  
 84 Then, in the second step, a deterministic min-max problem, i.e., Problem (1.1) where  
 85 expectation operator is replaced by summation over the finite sample, is solved. Once  
 86 again, using a finite sample to approximate the original sample space may result in  
 87 similar out-of-sample performance as in the case for 2-SLP.

88 **1.1. Contributions.** In light of the above observations regarding the two-step  
 89 external sampling procedure, the main contributions of this manuscript are as follows.

90 1. *A Sequential Sampling Algorithm:* We present a sequential sampling approach  
 91 for solving a 2-DRLP. We refer to this algorithm as the *distributionally robust*  
 92 *stochastic decomposition* (DRSD) algorithm following its risk-neutral prede-  
 93 cessor, the two-stage stochastic decomposition (SD) method [24] that was  
 94 designed for 2-SLPs. The DRSD algorithm concurrently performs the sim-  
 95 ulation and optimization steps in every iteration. In the simulation step, new  
 96 observation(s) are included to improve the representation of the ambiguity  
 97 set. The sequential inclusion of observations results in approximate ambi-  
 98 guity sets that evolve over the course of the algorithm. In the optimization  
 99 step, the solution is updated in an online manner by solving second-stage  
 100 programs for only the new observation(s) in each iteration. In this sense, the  
 101 DRSD method is an *internal sampling procedure*. Moreover, the algorithmic  
 102 design of the DRSD does not depend on any specific ambiguity set descrip-  
 103 tion. Hence, this method is suitable for any (general) ambiguity set for which  
 104 the distribution separation problem (1.2) can be solved efficiently.

105 2. *Convergence Analysis:* The DRSD method is an inexact bundle method that  
 106 creates outer linearization for the dynamically evolving approximation of the  
 107 first-stage problem. We provide the asymptotic analysis of DRSD and identify  
 108 conditions on ambiguity sets under which the sequential sampling approach  
 109 identifies an optimal solution to the 2-DRLP in (1.1) with probability one.

110 3. *Computational Evidence of Performance:* We provide the first set of exper-  
 111 iments that illustrates the advantages of a sequential sampling approach to  
 112 solving 2-DRLPs. We demonstrate these advantages through computational  
 113 experiments conducted on well-known problems in the SP literature. These  
 114 problems are modified to create distributionally robust variants with moment-  
 115 based,  $\ell_1$ -type Wasserstein, and  $\ell_\infty$ -type Wasserstein ambiguity sets.

116 **1.2. Related work.** For 2-SLPs with finite support, including the SAA prob-  
 117 lems, the L-shaped method due to Van Slyke and Wets [46] has proven to be very  
 118 effective. Other algorithms for 2-SLPs such as the Dantzig-Wolfe decomposition [13]  
 119 and the progressive hedging (PH) algorithm [36] also operate on problems with finite  
 120 support. The well-established theory of SAA (see Chapter 5 in [43]) supports the ex-  
 121 ternal sampling procedure for 2-SLP. The quality of the solution obtained by solving  
 122 an SAA problem is assessed using the procedures developed, e.g., in [5]. When the  
 123 quality of the SAA solution is not acceptable, a new SAA is constructed with a larger  
 124 number of observations. Prior works, such as [6] and [38], provide rules on how to  
 125 choose the sample sizes in a sequential SAA procedure.

126 In contrast to the above, SD incorporates one new observation in every iteration  
 127 to create approximations of the dynamically updating SAAs of (1.4). First proposed  
 128 in [24], this method has seen significant development in the past three decades with  
 129 the introduction of the quadratic proximal term [25], statistical optimality rules [27],  
 130 and extensions to multistage stochastic linear programs [21]. The DRSD method  
 131 extends the notion of sequential sampling of SD to DRO problems.

132 The concept of DRO dates back to the work of Scarf [40], and has gained sig-  
 133 nificant attention in recent years. Readers can refer to [34] for a comprehensive  
 134 treatment on various aspects of the DRO. The algorithmic works on DRO are either  
 135 decomposition-based or reformulation-based approaches. The decomposition-based  
 136 methods for 2-DRLP mimic the two-stage SP approach of using a deterministic repre-

137 sentation of the sample space using a finite number of observations. As a consequence,  
 138 the SP solution methods with suitable adaptation can be applied to solve the 2-DRLP  
 139 problems. For instance, Breton and El Hachem [11] applied the PH algorithm for a  
 140 2-DRLP model with a moment-based ambiguity set. Riis and Anderson [35] extended  
 141 the L-shaped method for 2-DRLP with continuous recourse and moment-based ambi-  
 142 guity set. Bansal et al. [1] extended the algorithm in [35], which they refer to as the  
 143 distributionally robust (DR) L-shaped method, to solve 2-DRLPs, with an ambiguity  
 144 set defined by a polytope. Further extensions of this decomposition approach are  
 145 presented in [1] and [2] for two-stage DRO problems with mixed-binary recourse and  
 146 disjunctive programs, respectively. Lately in [3], the authors considered two-stage  
 147 DRO problems with  $p$ -order conic mixed-integer programs in the second stage and  
 148 utilized scenario-based cuts to obtain linear programming equivalent (for  $p = 1$ ) and  
 149 convex approximations (for  $p \geq 2$ ) of the second-stage problems. We discuss key  
 150 differences of DRSD with SD and DR L-shaped method in Remark 3.3 at the end of  
 151 §3.

152 Another predominant approach to solve 2-DRLP problems is to reformulate the  
 153 distribution separation problem in (1.2) as a minimization problem, pose the problem  
 154 in (1.1) as a single deterministic optimization problem, and use off-the-shelf deter-  
 155 ministic optimization tools to solve the reformulation. For example, Shapiro and  
 156 Kleywegt [44] and Shapiro and Ahmed [42] used such an approach for a 2-DRLP with  
 157 moment matching set. They derived an equivalent stochastic program defined with a  
 158 reference distribution. Bertsimas et al. [8] provided tight semidefinite programming  
 159 reformulations for 2-DRLP where the ambiguity set is defined using multivariate dis-  
 160 tributions with known first and second moments. Likewise, Hanasusanto and Kuhn  
 161 [22] provided a conic programming reformulation for 2-DRLP where the ambiguity  
 162 set comprises of a  $\ell_2$ -type Wasserstein ball centered at a discrete distribution. Xie [47]  
 163 provided similar reformulations to tractable convex programs for 2-DRLP problems  
 164 with ambiguity set defined using  $\ell_\infty$  Wasserstein metric. By taking the dual of the in-  
 165 ner maximization problem, Love and Bayraksan [4] demonstrated that a 2-DRLP with  
 166 the ambiguity set defined using  $\phi$ -divergence and finite sample space is equivalent to  
 167 2-SLP with a coherent risk measure. A similar reformulation approach is employed in  
 168 [16] for ambiguity sets defined using Wasserstein and quadratic transport function on  
 169 unbounded and hyper-rectangle support. Jiang and Guan [29] reduced the worst-case  
 170 expectation in 2-DRLP, where the ambiguity set is defined using the  $\ell_1$ -norm on the  
 171 space of all (continuous and discrete) probability distributions, to a convex combina-  
 172 tion of CVaR and an essential supremum. Under the assumption of finite support,  
 173 [28] showed that a 2-DRLP with CVaR objective can be reformulated into a linear  
 174 program. On the other hand, the two-stage DRO problem with a linear recourse was  
 175 reformulated as a conic optimization problem under an assumption that second-stage  
 176 decisions are affine functions of the random vector in [30]. When reformulations result  
 177 in equivalent stochastic programs (as in [4, 28, 29, 42], for instance), an SAA of the  
 178 reformulation is used to obtain an approximate problem. This approximate problem  
 179 is amenable to standard cutting plane or bundle type methods prevalent in SP.

180 Data-driven approaches for DRO have been presented for specific ambiguity sets.  
 181 In [14], problems with ellipsoidal moment-based ambiguity set whose parameters are  
 182 estimated using sampled data are addressed. Esfahani et al. [32] tackled data-driven  
 183 problems with Wasserstein metric-based ambiguity sets with convex reformulations.  
 184 In both these works, the authors provide finite-sample performance guarantees that  
 185 probabilistically bound the gap between approximate and true DRO problems. Sun  
 186 and Xu presented asymptotic convergence analysis of DRO problems with ambigu-

187 ity sets that are based on moments and mixture distributions constructed using a  
 188 finite set of observations in [45]. A practical approach to incorporate the results of  
 189 these works to identify a high-quality DRO solution is similar to the sequential SAA  
 190 procedure for SP in [6]. Such an approach involves the following steps performed in  
 191 series – obtaining a deterministic representation of ambiguity set using sampled obser-  
 192 vations, applying appropriate reformulation, and solving the resulting deterministic  
 193 optimization problem. If the quality of the solution is deemed insufficient, then the  
 194 entire series of steps is repeated with an improved representation of the ambiguity set  
 195 (possibly with a larger number of observations).

196 **Organization.** We organize the remainder of the paper as follows. In §2, we  
 197 present the two key ideas of the DRSD- the sequential approximation of the ambiguity  
 198 set and the recourse function. We provide a detailed description of the DRSD method  
 199 in §3. We show the convergence of the value functions and solutions generated by the  
 200 DRSD method in §4. We present results of our computational experiments in §5, and  
 201 finally we conclude and discuss potential future directions in §6.

202 **Notations and Definitions.** We define the ambiguity sets over  $\mathcal{M}$ , the set of  
 203 all finite signed measures on the measurable space  $(\Omega, \mathcal{F})$ . A nonnegative measure  
 204 that satisfies  $P(\Omega) = 1$  is a probability distribution. For probability distributions  
 205  $P, P' \in \mathfrak{P}$ , we define

$$206 \quad (1.5) \quad \text{dist}(P, P') := \sup_{F \in \mathcal{F}} |\mathbb{E}_P[F(\tilde{\omega})] - \mathbb{E}_{P'}[F(\tilde{\omega})]|$$

208 as the uniform distance of expectation, where  $\mathcal{F}$  is a class of measurable functions.  
 209 The above is the distance with  $\zeta$ -structure that is used for stability analysis in SP  
 210 [37]. The distance between a single probability distribution  $P$  to a set of distributions  
 211  $\mathfrak{P}$  is given as  $\text{dist}(P, \mathfrak{P}) = \inf_{P' \in \mathfrak{P}} \text{dist}(P, P')$ . The distance between two sets of  
 212 probability distributions  $\mathfrak{P}$  and  $\widehat{\mathfrak{P}}$  is given as

$$213 \quad (1.6) \quad \mathbb{D}(\mathfrak{P}, \widehat{\mathfrak{P}}) := \sup_{P \in \widehat{\mathfrak{P}}} \text{dist}(P, \mathfrak{P}).$$

215 Finally, the Hausdorff distance between  $\mathfrak{P}$  and  $\widehat{\mathfrak{P}}$  is defined as

$$216 \quad (1.7) \quad \mathbb{H}(\mathfrak{P}, \widehat{\mathfrak{P}}) := \max\{\mathbb{D}(\mathfrak{P}, \widehat{\mathfrak{P}}), \mathbb{D}(\widehat{\mathfrak{P}}, \mathfrak{P})\}.$$

218 With suitable definitions for the set  $\mathcal{F}$ , the distance in (1.5) accepts the bounded  
 219 Lipschitz, the Kantorovich and the  $p$ -th order Fortet-Mourier metrics (see [37]).

220 **2. Approximating Ambiguity Set and Recourse Function.** In this section,  
 221 we present the building blocks that we embed within a sequential sampling setting  
 222 of the DRSD method. Specifically, we present procedures to approximate ambiguity  
 223 set  $\mathfrak{P}$  and recourse function  $Q(x, \omega)$  in an iteration of the DRSD. Going forward we  
 224 make the following assumptions on the 2-DRLP models:

- 225 (A1) The first-stage feasible region  $\mathcal{X}$  is a non-empty and compact set.
- 226 (A2)  $Q(\cdot)$  satisfies relatively complete recourse. The dual feasible region of the  
 227 recourse problem is a nonempty compact polyhedral set. The transfer (or  
 228 technology) matrix satisfies  $\sup_{P \in \mathfrak{P}} \mathbb{E}_P[T(\tilde{\omega})] < \infty$ .
- 229 (A3) The randomness only affects the right-hand sides of constraints in (1.3).
- 230 (A4) The sample space  $\Omega$  is a compact metric space and the ambiguity set  $\mathfrak{P}$  is  
 231 nonempty.

232 As a consequence of (A2), the recourse function satisfies  $Q(x, \tilde{\omega}) < \infty$  with probability  
 233 one for all  $x \in \mathcal{X}$ . It also implies that the second-stage feasible region, i.e.,  $\{y : Wy =$   
 234  $r(\omega) - T(\omega)x, y \geq 0\}$ , is non-empty for all  $x \in \mathcal{X}$  and every  $\omega \in \Omega$ . The non-  
 235 empty dual feasible region  $\mathcal{D}$  implies that there exists a constant  $L > -\infty$  such that  
 236  $Q(x, \tilde{\omega}) > L$ , almost surely. Without loss of generality, we assume that  $L = 0$ . As a  
 237 consequence of (A3), the cost coefficient vector  $g$  and the recourse matrix  $W$  are not  
 238 affected by uncertainty. Problems that satisfy (A3) are said to have a fixed recourse.  
 239 Finally, the compactness of the support  $\Omega$  guarantees that every probability measure  
 240  $P \in \mathfrak{P}$  is tight.

241 **2.1. Approximating the Ambiguity Set.** The DRO approach assumes only  
 242 partial knowledge about the underlying uncertainty that is captured by a suitable  
 243 description of the ambiguity set. An ambiguity set must capture the true distribution  
 244 with an absolute or high degree of certainty and must be computationally manageable.  
 245 The description of the ambiguity set involves parameters that are determined based  
 246 on a practitioner's risk preferences. The ambiguity set descriptions that are prevalent  
 247 in the literature include moment-based ambiguity sets with linear constraints (e.g.,  
 248 [15]) or conic constraints (e.g., [14]); Kantorovich distance or Wasserstein metric-  
 249 based ambiguity sets [31];  $\zeta$ -structure metrics [48],  $\phi$ -divergences such as  $\chi^2$  distance  
 250 and Kullback-Leibler divergence [7]; Prokhorov metrics [17], among others. In this  
 251 section, we present steps to construct approximate ambiguity sets in a data-driven  
 252 manner. We use moment-based and Wasserstein distance-based ambiguity sets to  
 253 illustrate these steps.

254 In a data-driven setting, the parameters used in the description of ambiguity sets  
 255 are estimated using a finite set of independent observations which can either be past  
 256 realizations of the random variable  $\tilde{\omega}$  or generated using computer simulations. We  
 257 will denote such a sample by  $\Omega^k \subseteq \Omega$ . When one observation is added to the sample  
 258 in every iteration, we obtain  $\Omega^k = \{\omega^j\}_{j=1}^k$ . Naturally, we can view  $\Omega^k$  as a random  
 259 sample and define the empirical frequency

$$260 \quad (2.1) \quad \hat{p}^k(\omega) = \frac{\kappa(\omega)}{k} \quad \text{for all } \omega \in \Omega^k,$$

262 where  $\kappa(\omega)$  denotes the number of times observation  $\omega$  is observed in the sample. Since  
 263 in the sequential sampling setting, the sample set is updated within the optimization  
 264 algorithm, it is worthwhile to note that the empirical frequency can be updated using  
 265 the following recursive equations:

$$266 \quad (2.2) \quad \hat{p}^k(\omega) = \begin{cases} \theta^k \hat{p}^{k-1}(\omega) & \text{if } \omega \in \Omega^{k-1}, \omega \neq \omega^k \\ \theta^k \hat{p}^{k-1}(\omega) + (1 - \theta^k) & \text{if } \omega \in \Omega^{k-1}, \omega = \omega^k \\ (1 - \theta^k) & \text{if } \omega \notin \Omega^{k-1}, \omega = \omega^k. \end{cases}$$

268 where  $\theta^k = \frac{k-1}{k}$ . In general, when more than one observation is added to the sample  
 269 in every iteration, we have  $\theta^k \in (0, 1)$ . We will succinctly denote the above using the  
 270 operator  $\Theta^k : \mathbb{R}^{|\Omega^{k-1}|} \rightarrow \mathbb{R}^{|\Omega^k|}$ .

271 In this paper, we focus on a setting where the ambiguity set  $\mathfrak{P}$  is replaced by a  
 272 sequence of *approximate ambiguity sets*  $\{\hat{\mathfrak{P}}^k\}_{k>0}$  such that the following properties  
 273 are satisfied: (B1) for any  $P \in \hat{\mathfrak{P}}^{k-1}$ , there exists  $\theta^k \in (0, 1)$  such that  $\Theta^k(P) \in \hat{\mathfrak{P}}^k$   
 274 and (B2)  $\mathbb{H}(\hat{\mathfrak{P}}^k, \mathfrak{P}) \rightarrow 0$  as  $k \rightarrow \infty$ , almost surely. We show that approximate  
 275 ambiguity sets for the moment-based ambiguity set  $\mathfrak{P}_{\text{mom}}$  and Wasserstein distance-  
 276 based ambiguity set  $\mathfrak{P}_{\text{w}}$  can be constructed such that these properties are satisfied  
 277 (Propositions 2.1 and 2.3, respectively).

278 Let  $\mathcal{F}^k = \sigma(\omega^j \mid j \leq k)$  be the  $\sigma$ -algebra generated by the observations in the  
 279 sample  $\Omega^k$ . Notice that  $\mathcal{F}^{k-1} \subseteq \mathcal{F}^k$ , and hence,  $\{\mathcal{F}^k\}_{k \geq 1}$  is a filtration. We will  
 280 define the approximate ambiguity sets over the measurable space  $(\Omega^k, \mathcal{F}^k)$ . These  
 281 sets should be interpreted to include all distributions that could have been generated  
 282 using the sample  $\Omega^k$ , which share a certain relationship with sample statistics. We  
 283 will use  $\mathcal{M}^k$  to denote the finite signed measures on  $(\Omega^k, \mathcal{F}^k)$ .

284 **2.1.1. Moment-based Ambiguity Sets.** Given the first  $q$  moments associated  
 285 with the random variable  $\tilde{\omega}$ , the moment-based ambiguity set can be defined as

$$286 \quad (2.3) \quad \mathfrak{P}_{\text{mom}} = \left\{ P \in \mathcal{M} \mid \begin{array}{l} \int_{\Omega} dP(\tilde{\omega}) = 1, \\ \int_{\Omega} \psi_i(\tilde{\omega}) dP(\tilde{\omega}) = b_i \quad i = 1, \dots, q \end{array} \right\}.$$

288 While the first constraint ensures the definition of a probability measure, the moment  
 289 requirements are guaranteed by the second constraints. Here,  $\psi_i(\tilde{\omega})$  denotes a real  
 290 valued measurable function on  $(\Omega, \mathcal{F})$  and  $b_i \in \mathbb{R}$  is a scalar for  $i = 1, \dots, q$ . Existence  
 291 of moments ensures that  $b_i < \infty$  for all  $i = 1, \dots, q$ . Notice that the description of  
 292 the ambiguity set requires explicit knowledge of the following statistics: the support  
 293  $\Omega$  and the moments  $b_i$  for  $i = 1, \dots, q$ . In the data-driven setting, the support is  
 294 approximated by  $\Omega^k$  and the sample moments  $\hat{b}_i^k = (1/k) \sum_{j=1}^k \psi_i(\omega^j)$  are used to  
 295 define the following approximate ambiguity set

$$296 \quad (2.4) \quad \hat{\mathfrak{P}}_{\text{mom}}^k = \left\{ P \in \mathcal{M}^k \mid \begin{array}{l} \sum_{\omega \in \Omega^k} p(\omega) = 1, \\ \sum_{\omega \in \Omega^k} p(\omega) \psi_i(\omega) = \hat{b}_i^k \quad i = 1, \dots, q \end{array} \right\}.$$

298 The following result characterizes the relationship between distributions drawn  
 299 from the above approximate ambiguity set, as well as asymptotic behavior of the  
 300 sequence  $\{\hat{\mathfrak{P}}_{\text{mom}}^k\}_{k \geq 1}$ .

301 **PROPOSITION 2.1.** *For any  $P \in \hat{\mathfrak{P}}_{\text{mom}}^{k-1}$ , we have  $\Theta^k(P) \in \hat{\mathfrak{P}}_{\text{mom}}^k$ . Further, suppose  $\hat{\mathfrak{P}}_{\text{mom}}^k \neq \emptyset$  for all  $k \geq 1$ ,  $\mathbb{H}(\hat{\mathfrak{P}}_{\text{mom}}^k, \mathfrak{P}_{\text{mom}}) \rightarrow 0$  as  $k \rightarrow \infty$ , almost surely.*

303 *Proof.* See Appendix §A. □

304 In the context of DRO, similar ambiguity sets have been studied in [9, 15] where  
 305 only the first moment (i.e.,  $q = 1$ ) is considered. The above form of ambiguity set  
 306 also relates to those used in [14, 35, 40, 45] where constraints were imposed only on  
 307 the mean and covariance. In the data-driven setting of [14] and [45], the statistical  
 308 estimates are used in constructing the approximate ambiguity set as in the case of  
 309 (2.4). However, the ambiguity sets in these previous works are defined over the original  
 310 sample space  $\Omega$ , as opposed to  $\Omega^k$  that is used in (2.4). This marks a critical deviation  
 311 in the way the approximate ambiguity sets are constructed.

312 **Remark 2.2.** When moment information is available about the underlying distri-  
 313 bution  $P^*$ , an approximate moment-based ambiguity set with constant parameters in  
 314 (2.4) (i.e., with  $\hat{b}_i^k = b_i$  for all  $k$ ) can be constructed. Such an approximate ambiguity  
 315 set defined over  $\Omega^k$  is studied in [35]. Notice that these approximate ambiguity sets  
 316 satisfy  $\cup_{k \geq 1} \hat{\mathfrak{P}}^k \subseteq \mathfrak{P}$  and  $\hat{\mathfrak{P}}^k \subseteq \hat{\mathfrak{P}}^{k+1}$ , for all  $k \geq 1$ . Therefore, they satisfy the  
 317 properties (i) and (ii) necessary for approximate ambiguity sets.

318 **2.1.2. Wasserstein distance-based Ambiguity Sets.** We next present ap-  
 319 proximations of another class of ambiguity sets that has gained significant attention  
 320 in the DRO literature, viz., the Wasserstein distance-based ambiguity sets. Consider  
 321 probability distributions  $\mu_1, \mu_2 \in \mathcal{M}$ , and a function  $\nu : \Omega \times \Omega \rightarrow \mathbb{R}_+ \cup \{\infty\}$  such that

322  $\nu$  is symmetric,  $\nu^{\frac{1}{r}}(\cdot)$  satisfies triangle inequality for  $1 \leq r < \infty$ , and  $\nu(\omega_1, \omega_2) = 0$   
 323 whenever  $\omega_1 = \omega_2$ . If  $\mathcal{J}(\mu_1, \mu_2)$  denotes the joint distribution of random vectors  $\omega_1$   
 324 and  $\omega_2$  with marginals  $\mu_1$  and  $\mu_2$ , respectively, then the Wasserstein metric of order  
 325  $r$  is given by

$$326 \quad (2.5) \quad d_w(\mu_1, \mu_2) = \left[ \inf_{\eta \in \mathcal{J}(\mu_1, \mu_2)} \left\{ \int_{\Omega \times \Omega} \nu^r(\omega_1, \omega_2) \eta(d\omega_1, d\omega_2) \right\} \right]^{1/r}.$$

328 In the above definition, the decision variable  $\eta \in \mathcal{J}$  can be viewed as a plan to trans-  
 329 port goods/mass from an entity whose spatial distribution is given by the measure  $\mu_1$   
 330 to another entity with spatial distribution  $\mu_2$ . Therefore, the  $d_w(\mu_1, \mu_2)$  measures the  
 331 optimal transport cost between the measures. Notice that an arbitrary norm  $\|\bullet\|^r$  on  
 332  $\mathbb{R}^d$  satisfies the requirement of the function  $\nu(\cdot)$ . In our presentation, we will use the  
 333  $\ell_1$  Wasserstein metric. However, the definition of the approximate ambiguity sets and  
 334 their use within the solution method are applicable to ambiguity sets defined using  
 335 Wasserstein metric of higher orders. Using the  $\ell_1$  Wasserstein metric, we define an  
 336 ambiguity set as follows:

$$337 \quad (2.6) \quad \mathfrak{P}_w = \{P \in \mathcal{M} \mid d_w(P, P^*) \leq \epsilon\}$$

339 for a given  $\epsilon > 0$  and a reference distribution  $P^*$ . In practice, the value of  $\epsilon$  is chosen  
 340 based on user's risk preferences; a smaller value indicates lower risk aversion. As done  
 341 in §2.1.1, we present approximate Wasserstein distance-based ambiguity sets defined  
 342 over the measurable space  $(\Omega^k, \mathcal{F}^k)$  as follows:

$$343 \quad (2.7) \quad \widehat{\mathfrak{P}}_w^k = \{P \in \mathcal{M}^k \mid d_w(P, \widehat{P}^k) \leq \epsilon\},$$

345 where  $\widehat{P}^k = (\widehat{p}^k(\omega))_{\omega \in \Omega^k}$ . For this approximate ambiguity set, the distribution sepa-  
 346 ration problem in (1.2) is a finite dimensional linear program:

$$347 \quad (2.8a) \quad \max \sum_{\omega \in \Omega^k} p(\omega) Q(x, \omega)$$

(2.8b)

$$348 \quad \text{subject to } P \in \widehat{\mathfrak{P}}_w^k = \left\{ P \in \mathcal{M}^k \left| \begin{array}{l} \sum_{\omega \in \Omega^k} p(\omega) = 1 \\ \sum_{\omega' \in \Omega^k} \eta(\omega, \omega') = p(\omega) \quad \forall \omega \in \Omega^k, \\ \sum_{\omega \in \Omega^k} \eta(\omega, \omega') = \widehat{p}^k(\omega') \quad \forall \omega' \in \Omega^k, \\ \sum_{(\omega, \omega') \in \Omega^k \times \Omega^k} \|\omega - \omega'\| \eta(\omega, \omega') \leq \epsilon \\ \eta(\omega, \omega') \geq 0 \quad \forall \omega, \omega' \in \Omega^k \end{array} \right. \right\}.$$

349 Note that when Wasserstein metric of order  $r > 1$  is used in the definition of the  
 350 ambiguity sets, the foregoing optimization problem remains a finite dimensional linear  
 351 program. In this case, the coefficients  $\|\omega - \omega'\|$  and right-hand side  $\epsilon$  in the fourth  
 352 set of constraints in (2.8b) must be replaced by  $\|\omega - \omega'\|^r$  and  $\epsilon^r$ , respectively. The  
 353 following result characterizes the distributions drawn from the approximate ambiguity  
 354 sets of the form in (2.7), or equivalently (2.8b).

355 **PROPOSITION 2.3.** *Under compactness of the support set  $\Omega \subset \mathbb{R}^d$ , i.e., (A4),  
 356 with  $d > 2$ , the sequence of Wasserstein distance-based approximate ambiguity sets  
 357 satisfies the following properties (i) for any  $P \in \widehat{\mathfrak{P}}_w^{k-1}$ , we have  $\Theta^k(P) \in \widehat{\mathfrak{P}}_w^k$ , and  
 358 (ii)  $\mathbb{H}(\widehat{\mathfrak{P}}_w^k, \mathfrak{P}_w) \rightarrow 0$  as  $k \rightarrow \infty$ , almost surely.*

360 *Proof.* See appendix §A. □

361 Note that, as in the case of moment-based ambiguity set, we also define Wasserstein  
 362 distance-based approximate ambiguity set over an approximation of the true sample  
 363 space, i.e.,  $\Omega^k$ . This approach precludes the need for exact knowledge of the sample  
 364 space and allows us to depend only on what is known until iteration  $k$ .

365 *Remark 2.4.* In [32], an approach that involves solving a sequence of DRO prob-  
 366 lems is used to tackle the risk-neutral 2-SLP problem (1.4). They use approximate  
 367 ambiguity set to be a ball constructed in the space of probability distributions that  
 368 are defined over the sample space  $\Omega$  and whose radius reduces with an increase in  
 369 the number of observations. Using Wasserstein balls of shrinking radii, the authors  
 370 of [32] show that the optimal value of the sequence of DRO problems converges to  
 371 the optimal value of the expectation-valued objective in (1.4) associated with the true  
 372 distribution  $P^*$ . A similar approach of involving a sequence of DRO problems is used  
 373 in [48] to solve (1.4), albeit using ambiguity sets with  $\zeta$ -structure. In contrast to  
 374 these works, our goal is to solve the DRO problem in (1.1). Therefore, we use a  
 375 constant radius for all  $k \geq 1$  to define the approximate ambiguity set in (2.7).

376 **2.2. Approximating the Recourse Problem.** Cutting plane methods for the  
 377 2-SLPs use an outer linearization-based approximation of the first-stage objective  
 378 function in (1.4). In such algorithms, the challenging aspect of computing the ex-  
 379 pectation is addressed by taking advantage of the structure of the recourse problem  
 380 (1.3). Specifically, for a given  $\omega$ , the recourse value  $Q(\cdot, \omega)$  is known to be convex in  
 381 the right-hand side parameters that includes the first-stage decision vector  $x$ . Addi-  
 382 tionally, if (A2) holds, then the function  $Q(\cdot, \omega)$  is polyhedral. Under assumptions  
 383 (A2) and (A4), this structural property of convexity extends to the expected recourse  
 384 value  $\mathcal{Q}(x)$ .

385 Due to the strong duality of linear programs, the recourse value is also equal to  
 386 the optimal value of the dual of (1.3), i.e.,

387 (2.9) 
$$Q(x, \omega) = \max \pi^\top [r(\omega) - T(\omega)x]$$
  
 388 subject to  $\pi \in \mathcal{D} := \{\pi \mid W^\top \pi \leq g\}$ .

390 Due to (A2) and (A4), the dual feasible region  $\mathcal{D}$  is a polytope that is not impacted  
 391 by the uncertainty. If  $\Pi \subseteq \mathcal{D}$  denotes the set of all extreme points of the polytope  
 392  $\mathcal{D}$ , then the recourse value can also be expressed as the pointwise maximum of affine  
 393 functions computed using elements of set  $\Pi$ :

394 (2.10) 
$$Q(x, \omega) = \max_{\pi \in \mathcal{D}} \pi^\top [r(\omega) - T(\omega)x].$$

396 The outer linearization approaches tend to approximate the above form of recourse  
 397 function by identifying the extreme points (optimal solutions to (2.9)) at a sequence of  
 398 candidate (or trial) solutions  $\{x^k\}$ , and generating the corresponding affine functions.  
 399 If  $\pi(x^k, \omega)$  is an optimal dual obtained by solving (2.9) with  $x^k$  as input, then the  
 400 affine function  $\alpha^k(\omega) + (\beta^k(\omega))^\top x$  is obtained by computing the coefficients  $\alpha^k(\omega) =$   
 401  $(\pi(x^k, \omega))^\top r(\omega)$  and  $\beta^k(\omega) = T(\omega)^\top \pi(x^k, \omega)$ . Following linear programming duality,  
 402 notice that this affine function is a supporting hyperplane to  $Q(x, \omega)$  at  $x^k$ , and lower  
 403 bounds the function at every other  $x \in \mathcal{X}$ .

404 If the support  $\Omega$  is finite, then one can solve a dual subproblem for all  $\omega \in \Omega$   
 405 with the candidate solution as input, generate the affine functions, and collate them  
 406 together to obtain an approximate first-stage objective function. This is the essence

407 of the L-shaped method applied to 2-SLP in (1.4). In each iteration of the L-shaped  
 408 method, the affine functions generated using a candidate solution  $x^k$  and information  
 409 gathered from individual observations are weighed by the probability density of the  
 410 observation to update the approximation of the first-stage objective function. The L-  
 411 shaped method can also be applied to the SAA of the 2-SLP with continuous sample  
 412 space  $\Omega$  that uses a sample  $\Omega_N \subset \Omega$  of finite size  $N$ . A similar approximation strategy  
 413 is used in the DR L-shaped method for 2-DRLP problems [1, 35].

414 Alternatively, we can consider the following approximation of the recourse func-  
 415 tion expressed in the form given in (2.10):

$$416 \quad (2.11) \quad Q^k(x, \omega) = \max_{\pi \in \Pi^k} \pi^\top [r(\omega) - C(\omega)x].$$

418 Notice that the above approximation is built using only a subset  $\Pi^k \subset \Pi$  of extreme  
 419 points, and therefore, satisfies  $Q^k(x, \omega) \leq Q(x, \omega)$ . Since  $Q(x, \omega) \geq 0$ , we begin  
 420 with  $\Pi^0 = \{0\}$ . Subsequently, we construct a sequence of sets  $\{\Pi^k\}$  such that  $\Pi^0 \subseteq$   
 421  $\dots \Pi^k \subseteq \Pi^{k+1} \subseteq \dots \subset \Pi$  that ensures  $Q^k(x, \omega) \geq 0$  for all  $k$ . The following result  
 422 from [24] captures the behavior of the sequence of approximation  $\{Q^k\}$ .

423 **PROPOSITION 2.5.** *The sequence  $\{Q^k(x, \omega)\}_{k \geq 1}$  converges uniformly to a contin-  
 424 uous function on  $\mathcal{X}$  for any  $\omega \in \Omega$ .*

425 *Proof.* See Appendix A. □

426 The approximation of the form in (2.11) is one of the principal features of the  
 427 SD algorithm (see [24, 25]). While the L-shaped and DR L-shaped methods require  
 428 finite support for  $\tilde{\omega}$ , SD is applicable even for problems with continuous support. The  
 429 algorithm uses an “incremental” SAA for the first-stage objective function by adding  
 430 one new observation in each iteration. Therefore, the first-stage objective function  
 431 approximation used in SD is built using the recourse problem approximation in (2.11)  
 432 and the incremental SAA. This approximation is given by:

$$433 \quad (2.12) \quad \mathcal{Q}^k(x) = c^\top x + \frac{1}{k} \sum_{j=1}^k Q^k(x, \omega^j).$$

435 The affine functions generated in SD provide an outer linearization for the approxi-  
 436 mation in (2.12). The sequence of sets that grow monotonically in size, viz.  $\{\Pi^k\}$ , is  
 437 generated by adding one new vertex to the previous set  $\Pi^{k-1}$  to obtain the updated  
 438 set  $\Pi^k$ . The newly added vertex is an optimal dual solution obtained by solving (2.9)  
 439 with the most recent observation  $\omega^k$  and candidate solution  $x^k$  as input.

440 We refer the reader to [10], [1, 35], and [24, 26] for the a detailed exposition of the  
 441 L-shaped, the DR L-Shaped, and the SD methods, respectively. Here, we only note the  
 442 key differences between these methods. Firstly, the sample used in the (DR) L-shaped  
 443 method is fixed before the optimization. In SD, this sample is updated dynamically  
 444 throughout the course of the algorithm. Secondly, in the (DR) L-shaped method,  
 445 subproblems corresponding to the current iterate and all observations in the sample  
 446 are solved exactly. The resulting optimal dual solutions are used to compute the affine  
 447 lower bounding functions (cuts). On the other hand, in SD, only two subproblems  
 448 corresponding to the latest observation are solved exactly, while the subproblems  
 449 corresponding to other observations in the sample use the approximation in (2.11).

450 **3. Distributionally Robust Stochastic Decomposition.** In this section, we  
 451 provide a detailed description of the DRSD algorithm. The pseudocode of the DRSD

**Algorithm 3.1** Distributionally Robust Stochastic Decomposition

---

1: **Input:** Incumbent solution  $\hat{x}^0 \in \mathcal{X}$ ; initial sample  $\Omega^0 \subseteq \Omega$ ; stopping tolerance  $\tau > 0$ ;  $\gamma \in (0, 1)$ ,  $\theta^1 = 0$ , and maximum and minimum iterations  $k^{\max} > k^{\min} > 1$ .  
2: **Initialization:** Set iteration counter  $k \leftarrow 1$ ;  $\Pi^0 = \emptyset$ ;  $\mathcal{L}^0 = \emptyset$ , and  $f^0(x) = 0$ .  
3: **while** ( $k \leq k^{\max}$ ) **do**  
4:     Solve the master problem (3.1) to obtain a candidate solution  $x^k$ .  
5:     **if**  $k > k^{\min}$  and  $f^{k-1}(\hat{x}^{k-1}) - f^{k-1}(x^k) < \tau f^{k-1}(\hat{x}^{k-1})$  **then**, Go to Line 28.  
6:     **end if**  
7:     Generate a scenario  $\omega^k \in \Omega$  to get sample  $\Omega^k \leftarrow \Omega^{k-1} \cup \{\omega^k\}$ .  
8:     Solve the second-stage linear program (1.3) with  $(x^k, \omega^k)$  as input;  
9:     Obtain the optimal value  $Q(x^k, \omega^k)$  and optimal dual solution  $\pi(x^k, \omega^k)$ ;  
10:    Update dual vertex set  $\Pi^k \leftarrow \Pi^{k-1} \cup \{\pi(x^k, \omega^k)\}$ .  
11:    **for**  $\omega \in \Omega^k \setminus \{\omega^k\}$  **do**  
12:       Use the argmax procedure (3.2) to identify dual vertex  $\pi(x^k, \omega^k)$ ;  
13:       Store  $Q^k(x^k, \omega) = (\pi(x^k, \omega))^{\top} [r(\omega) - T(\omega)x^k]$ .  
14:    **end for**  
15:    Solve the distribution separation problem using the ambiguity set  $\hat{\mathfrak{P}}^k$  and  $\{Q^k(x^k, \omega)\}_{\omega \in \Omega^k}$  to get an extremal distribution  $P^k := (p^k(\omega))_{\omega \in \Omega^k}$ .  
16:    Derive affine function  $\ell_k^k(x) = \alpha_k^k + (\beta_k^k)^{\top} x$  using  $\{\pi(x^k, \omega)\}_{\omega \in \Omega^k}$  and  $P^k$  to get lower bound approximation of  $\mathbb{Q}^k(x)$  as in (3.5);  
17:    Perform Steps 8-16 with  $\hat{x}^{k-1}$  (incumbent solution) to obtain  $\hat{\ell}_k^k(\cdot)$ .  
18:    **for**  $\ell_j^{k-1} \in \mathcal{L}^{k-1}$  **do**  
19:       Update previously generated affine functions  $\ell_j^{k-1}(x)$ :  

$$\alpha_j^k = \theta^k \alpha_j^{k-1} \text{ and } \beta_j^k = \theta^k \beta_j^{k-1};$$
  
20:       Set  $\ell_j^k(x) = \alpha_j^k + (\beta_j^k)^{\top} x$  that provides lower bound approx. of  $\mathbb{Q}^k(x)$ ;  
21:    **end for**  
22:    Build a collection of these affine functions, denoted by  $\mathcal{L}^k$ ;  
23:    Update approximation of the first-stage objective function:  

$$c^{\top} x + \mathbb{Q}^k(x) \geq f^k(x) = c^{\top} x + \max_{j \in \mathcal{L}^k} \{\alpha_j^k + (\beta_j^k)^{\top} x\};$$
  
24:    If incumbent update rule (3.9) is satisfied, then set  $\hat{x}^k \leftarrow x^k$  and  $\hat{x}^k \leftarrow \hat{x}^{k-1}$ , otherwise.  
25:    Update the master problem (3.1) by replacing  $f^{k-1}(x)$  with  $f^k(x)$ ;  
26:     $k \leftarrow k + 1$ ;  $\theta^k \leftarrow (k - 1)/k$   
27: **end while**  
28: **return** Incumbent solution  $\hat{x}^k$  and objective function estimate  $f^k(\hat{x}^k)$ .

---

452 method is given in Algorithm 3.1. In the following, we discuss the main steps of the  
453 algorithm in iteration  $k$  (Steps 4-26 of Algorithm 3.1). At the beginning of iteration  
454  $k$ , we have a certain approximation of the first-stage objective function that we denote  
455 as  $f^{k-1}(x)$ , a finite set of observations  $\Omega^{k-1}$  and an incumbent solution  $\hat{x}^{k-1}$ . We use  
456 the term *incumbent solution* to refer to the best solution discovered by the algorithm  
457 until iteration  $k$ . The solution identified in the current iteration is referred to as the  
458 *candidate solution* and denoted as  $x^k$  (without  $\hat{\bullet}$ ).  
459 Iteration  $k$  begins by first identifying the candidate solution by solving the fol-

460 lowing the master problem (Step 4):

461 (3.1) 
$$x^k \in \arg \min \{f^{k-1}(x) \mid x \in \mathcal{X}\}.$$

463 Following this, a new observation  $\omega^k \in \Omega$  is obtained and added to the current sample  
464 of observations  $\Omega^{k-1}$  to get  $\Omega^k = \Omega^{k-1} \cup \{\omega^k\}$  (Step 7).

465 In order to build the first-stage objective function approximation, we rely upon  
466 the recourse function approximation presented in Section 2.2. For the most recent  
467 observation  $\omega^k$  and the candidate solution  $x^k$ , we evaluate the recourse function value  
468  $Q(x^k, \omega^k)$  by solving (1.3), and obtain the dual optimum solution  $\pi(x^k, \omega^k)$  in Steps  
469 8–10. These dual vectors are added to a set  $\Pi^{k-1}$  of previously discovered optimal  
470 dual vectors. In other words, we recursively update  $\Pi^k \leftarrow \Pi^{k-1} \cup \{\pi(x^k, \omega^k)\}$ . For all  
471 other observations  $(\omega \in \Omega^k, \omega \neq \omega^k)$ , we identify a dual vector in  $\Pi^k$  that provides the  
472 best lower bounding approximation at  $Q(x^k, \omega)$  using the following operation (Steps  
473 12–13):

474 (3.2) 
$$\pi(x^k, \omega) \in \arg \max \{\pi^\top [r(\omega) - T(\omega)x^k] \mid \pi \in \Pi^k\}.$$

Note that the calculations in (3.2) are carried out only for previous observations as  
 $\pi(x^k, \omega^k)$  provides the best lower bound at  $Q(x^k, \omega^k)$ . Further, notice that

$$\pi(x^k, \omega)^\top [r(\omega) - T(\omega)x^k] = Q^k(x^k, \omega),$$

475 the approximate recourse function value at  $x^k$  defined in (2.11), for all  $\omega \in \Omega^k$ , and  
476  $Q^k(x^k, \omega^k) = Q(x^k, \omega^k)$ .

477 Using  $\{Q^k(x^k, \omega^j)\}_{j=1}^k$ , we solve a *distribution separation problem* (in Step 15):

478 (3.3) 
$$\mathbb{Q}^k(x^k) = \max \left\{ \sum_{\omega \in \Omega^k} p(\omega)Q^k(x^k, \omega) \mid p(\omega) \in \widehat{\mathfrak{P}}^k \right\}.$$

480 Let  $P^k = (p^k(\omega))_{\omega \in \Omega^k}$  denote an optimal solution of the above problem which we  
481 identify as a maximal/extremal probability distribution. Since the problem is solved  
482 over measures  $\mathcal{M}^k$  that are defined only over the observed set  $\Omega^k$ , the maximal proba-  
483 bility distribution has weights  $p^k(\omega)$  for  $\omega \in \Omega^k$ , and  $p^k(\omega) = 0$  for  $\omega \in \Omega \setminus \Omega^k$ . Notice  
484 that the problem in (1.2) differs from the distribution separation problem (3.3) as the  
485 latter uses the recourse function approximation  $Q^k(\cdot)$  and approximate ambiguity set  
486  $\widehat{\mathfrak{P}}^k$  as opposed to the true recourse function  $Q(\cdot)$  and ambiguity set  $\mathfrak{P}$ , respectively.  
487 For the moment-based and Wasserstein distance-based ambiguity sets (discussed in  
488 Section 2.1), the distribution separation problem is a deterministic linear program. In  
489 general, the distribution separation problems associated with well-known ambiguity  
490 sets remain deterministic convex optimization problems [34], and off-the-shelf solvers  
491 can be used to obtain the extremal distribution.

492 In Step 16 of Algorithm 3.1, we use the dual vectors  $\{\pi(x^k, \omega^j)\}_{j \leq k}$  and the  
493 maximal probability distribution  $P^k$  to generate a lower bounding affine function:

494 (3.4) 
$$\mathbb{Q}^k(x) = \max_{P \in \widehat{\mathfrak{P}}^k} \mathbb{E}_P[Q^k(x, \tilde{\omega})] \geq \sum_{\omega \in \Omega^k} p^k(\omega) \cdot (\pi(x^k, \omega))^\top [r(\omega) - T(\omega)x],$$

496 for the worst-case expected recourse function measured with respect to the maximal  
497 probability distribution  $P^k \in \widehat{\mathfrak{P}}^k$ . We denote the coefficients of the affine function on  
498 the right-hand side of (3.4) by

499 (3.5) 
$$\alpha_k^k = \sum_{\omega \in \Omega^k} p^k(\omega)\pi(x^k, \omega)^\top r(\omega) \text{ and } \beta_k^k = - \sum_{\omega \in \Omega^k} p^k(\omega)T(\omega)^\top \pi(x^k, \omega),$$

501 and succinctly write the affine function as  $\ell_k^k(x) = \alpha_k^k + (\beta_k^k)^\top x$ . Similar calculations  
 502 are carried out using the incumbent solution  $\hat{x}^{k-1}$  to identify a maximal probability  
 503 distribution and a lower bounding affine function resulting in the affine function  
 504  $\hat{\ell}_k^k(x) = \hat{\alpha}_k^k + (\hat{\beta}_k^k)^\top x$  (Step 17). Note that we use two indices for the cut coefficients  
 505  $(\alpha, \beta)$  and the affine function  $\ell$ . The superscript indicates the current iteration, while  
 506 the subscript indicates the iteration when the quantities were first computed. Since,  
 507 one observation is added to  $\Omega^k$  in every iteration, the subscript also indicates the  
 508 number of observations used in computing the quantities.

509 While the latest affine functions provide a lower bound for  $\mathbb{Q}^k$ , the affine functions  
 510 generated in previous iteration are not guaranteed to lower bound  $\mathbb{Q}^k$ . To see  
 511 this, let us consider the moment-based approximate ambiguity sets  $\{\hat{\mathfrak{P}}_{\text{mom}}^k\}_{k \geq 1}$ . Let  
 512  $P_{\text{mom}}^j \in \hat{\mathfrak{P}}_{\text{mom}}^j$  be the maximal distribution identified in an iteration  $j < k$  which was  
 513 used to compute the affine function  $\ell_j^j(x)$ . By assigning  $p^j(\omega) = 0$  for all new obser-  
 514 vations encountered after iteration  $j$ , i.e.,  $\omega \in \Omega^k \setminus \Omega^j$ , we can construct a probability  
 515 distribution  $\bar{P} = ((p^j(\omega))_{\omega \in \Omega^j}, (0)_{\omega \in \Omega^k \setminus \Omega^j}) \in \mathbb{R}_+^{|\Omega^k|}$ . This reconstructed distribution  
 516 satisfies  $\sum_{\omega \in \Omega^k} \bar{p}(\omega) = 1$ . However, it is easy to see that  $\sum_{\omega \in \Omega^k} \psi_i(\omega) \bar{p}(\omega) = \hat{b}_i^j \neq \hat{b}_i^k$   
 517 for all  $i = 1, \dots, q$ . Therefore,  $\bar{P} \notin \mathfrak{P}^k$ . In other words, while the coefficients  $(\alpha_j^j, \beta_j^j)$   
 518 are  $\mathcal{F}^j$ -measurable, the corresponding measure is not feasible to the approximate am-  
 519 biguity set  $\mathfrak{P}^k$ . Therefore,  $\ell_j^j(x)$  is not a valid lower bound to  $\mathbb{Q}^k$ . The arguments for  
 520 the Wasserstein-based approximate ambiguity set are more involved, but persistence  
 521 of a similar issue can be demonstrated.

522 To address this, we recursively update the previously generated affine functions  
 523  $\ell_j^{k-1}(x) = \alpha_j^{k-1} + (\beta_j^{k-1})^\top x$  for  $j < k$  as follows (Steps 18 - 21):

524 (3.6)  $\alpha_j^k = \theta^k \alpha_j^{k-1}, \beta_j^k = \theta^k \beta_j^{k-1}$ , and  $\ell_j^k(x) = \alpha_j^k + (\beta_j^k)^\top x$  for all  $j < k$ ,  
 525 such that  $\ell_j^k(x)$  provides lower bound approximation of  $\mathbb{Q}^k(x)$  for all  $j \in \{1, \dots, k-1\}$ .  
 526 Similarly, we update the affine functions  $\hat{\ell}_j^k(x)$ ,  $j < k$ , associated with incumbent  
 527 solution. The candidate and the incumbent affine functions ( $\ell_k^k(x)$  and  $\hat{\ell}_k^k(x)$ , respec-  
 528 tively), as well as the updated collection of previously generated affine functions are  
 529 used to build the set of affine functions which we denote by  $\mathcal{L}^k$  (Step 22). Using  
 530 this collection of affine functions  $\mathcal{L}^k$ , we update the approximation of the first-stage  
 531 objective function in Step 23, as follows:

533 (3.7) 
$$f^k(x) = c^\top x + \max_{\ell \in \mathcal{L}^k} \{\ell(x)\}.$$
  
 534

535 The lower bounding property of this first-stage objective function approximation is  
 536 captured in the following result.

537 **THEOREM 3.1.** *Under assumption (A2), the first-stage objective function approx-  
 538 imation in (3.7) satisfies*

539 
$$f^k(x) \leq c^\top x + \mathbb{Q}^k(x) \text{ for all } x \in \mathcal{X} \text{ and } k \geq 1.$$

541 *Proof.* For the non-empty approximate ambiguity set  $\hat{\mathfrak{P}}^1$  of ambiguity set  $\mathfrak{P}$ , the  
 542 construction of the affine function ensures that  $\ell_1^1(x) \leq \mathbb{Q}^1(x)$ . Now assume that  
 543  $\ell(x) \leq \mathbb{Q}^{k-1}(x)$  for all  $\ell \in \mathcal{L}^{k-1}$  and  $k > 1$ . The maximal nature of the probability  
 544 distribution  $P^k$  satisfies:

545 
$$\sum_{\omega \in \Omega^k} p^k(\omega) Q^k(x, \omega) \geq \sum_{\omega \in \Omega^k} p(\omega) Q^k(x, \omega) \quad \forall P \in \hat{\mathfrak{P}}^k.$$
  
 546

547 Using above and the monotone property of the approximate recourse function, we  
 548 have

$$549 \quad \sum_{\omega \in \Omega^k} p^k(\omega) Q^k(x, \omega) \geq \sum_{\omega \in \Omega^k} p(\omega) Q^{k-1}(x, \omega) \\ 550 \quad (3.8) \quad = \sum_{\omega \in \Omega^k \setminus \{\omega^k\}} p(\omega) Q^{k-1}(x, \omega) + p(\omega^k) Q^{k-1}(x, \omega^k), \\ 551$$

552 for all  $\{p(\omega)\}_{\omega \in \Omega^k} \in \widehat{\mathfrak{P}}^k$ . Based on the properties of  $\mathfrak{P}$  and  $\{\widehat{\mathfrak{P}}^k\}_{k \geq 1}$  (similar to  
 553 Propositions 2.1 and 2.3), we know that for every  $P \in \widehat{\mathfrak{P}}^{k-1}$  we can construct a  
 554 probability distribution in  $\mathfrak{P}^k$  using the mapping  $\Theta^k$  defined by (2.2). Considering a  
 555 probability distribution  $P' = \{p'(\omega)\}_{\omega \in \Omega^{k-1}} \in \widehat{\mathfrak{P}}^{k-1}$  we have  $\Theta^k(P') \in \widehat{\mathfrak{P}}^k$  and the  
 556 inequality (3.8) reduces to

$$557 \quad \sum_{\omega \in \Omega^k} p^k(\omega) Q^k(x, \omega) \geq \sum_{\omega \in \Omega^k \setminus \{\omega^k\}} [\theta^k p'(\omega) Q^{k-1}(x, \omega)] + [\theta^k p'(\omega^k) + (1 - \theta^k)] Q^{k-1}(x, \omega^k) \\ 558 \quad = \theta^k \left[ \sum_{\omega \in \Omega^{k-1}} p'(\omega) Q^{k-1}(x, \omega) \right] + (1 - \theta^k) Q^{k-1}(x, \omega^k) \\ 559 \quad \geq \theta^k \left[ \sum_{\omega \in \Omega^{k-1}} p'(\omega) Q^{k-1}(x, \omega) \right]. \\ 560$$

561 The last inequality is due to assumption (A2), i.e.,  $Q(x, \omega^k) \geq 0$  and the construction  
 562 of recourse function approximation  $Q^k$  described in §2.2. Since  $\ell(x)$  lower bounds the  
 563 term in bracket, we have

$$564 \quad \sum_{\omega \in \Omega^k} p^k(\omega) Q^k(x, \omega) \geq \theta^k \ell(x). \\ 565$$

566 Using the same arguments for all  $\ell \in \mathcal{L}^{k-1}$ , and the fact that the  $\ell_k^k(x)$  and  $\hat{\ell}_k^k(x)$  are  
 567 constructed as lower bounds to the  $\mathbb{Q}^k$ , we have  $f^k(x) \leq c^\top x + \mathbb{Q}^k(x)$ . This completes  
 568 the proof by induction.  $\square$

569 Once the approximation (3.7) is updated, the performance of the candidate so-  
 570 lution is compared relative to the incumbent solution (Step 24). This comparison is  
 571 performed by verifying if the following inequality is satisfied:

$$572 \quad (3.9) \quad f^k(x^k) - f^k(\hat{x}^{k-1}) < \gamma [f^{k-1}(x^k) - f^{k-1}(\hat{x}^{k-1})],$$

574 where parameter  $\gamma \in (0, 1)$ . If so, the candidate solution is designated to be the  
 575 next incumbent solution, i.e.,  $\hat{x}^k = x^k$ . If the inequality is not satisfied, the previous  
 576 incumbent solution is retained as  $\hat{x}^k = \hat{x}^{k-1}$ . This completes a DRSD iteration.

577 *Remark 3.2.* We can extend the algorithm design for 2-DRLPs where the rel-  
 578 atively complete recourse assumption of (A2) and/or assumption (A3) is not satisfied.  
 579 For problems where relatively complete recourse condition is not met, a candidate  
 580 solution may lead to one or more subproblems to be infeasible. In this case, the dual  
 581 extreme rays can be used to compute a feasibility cut that is included in the first-stage  
 582 approximation. The argmax procedure in (3.2) is only valid when assumption (A3)  
 583 is satisfied. In problems where the uncertainty also affects the cost coefficients, the  
 584 argmax procedure presented in [20] can be utilized. These algorithmic enhancements  
 585 can be incorporated without affecting the convergence properties of DRSD.

586 *Remark 3.3* (Relation between DRSD, SD, and DR L-shaped Method). We  
 587 close this section by identifying the key differences in the DRSD algorithm design  
 588 when compared to SD and DR L-shaped methods.

589 • There are two main differences between DRSD and the DR L-Shaped method.  
 590 Firstly, the DR L-shaped method operates with a deterministic representation of  
 591 the ambiguity set computed using a fixed sample of observations, an input to the  
 592 algorithm. In contrast, a new observation is added (Line 7) in every iteration  
 593 of DRSD to improve the approximation of the ambiguity set. Secondly, every  
 594 iteration of the DR L-shaped method involves solving a subproblem corresponding  
 595 to each observation used in the ambiguity set representation. On the other hand,  
 596 in DRSD, only two subproblems corresponding to latest observation  $\omega^k$  are solved  
 597 to optimality, and the argmax procedure is used for the other observations.  
 598 • While DRSD is designed to address the 2-DRLP problem (1.1), the SD and its vari-  
 599 ants [24, 25, 41] are for risk-neutral 2-SLP. This generalization introduces another  
 600 layer of approximation to SD, viz., the approximation of ambiguity sets. The algo-  
 601 rithmic enhancements necessary to address this new layer of approximations make  
 602 the DRSD significantly different from its risk-neutral predecessors. For instance,  
 603 we need to solve an approximate distribution separation problem in every iteration  
 604 (Line 15). The cut coefficients are computed and updated (in Lines 18-21) in a  
 605 manner that is consistent with the updates carried out to approximate the ambi-  
 606 guity sets (see propositions 2.1 and 2.3, and coefficient updates in (3.6)). The cut  
 607 updates that are undertaken in SD only need to be consistent with the updates in  
 608 empirical distribution. This critical difference in cut computations also introduces  
 609 significant differences in the convergence analysis of DRSD that we present next.

610 **4. Convergence Analysis.** In this section we provide the convergence result of  
 611 the sequential sampling-based approach to solve DRO problems. In order to facilitate  
 612 the exposition of our theoretical results, we will define certain quantities for notational  
 613 convenience that are not necessarily computed during the course of the algorithm. Our  
 614 convergence results are built upon stability analyses presented in [45] and convergence  
 615 analysis of the SD algorithm in [24].

616 We define a function over the approximate ambiguity set using the recourse func-  
 617 tion  $Q(\cdot, \cdot)$ , that is

$$618 \quad (4.1) \quad g^k(x) := c^\top x + \max_{P \in \hat{\mathfrak{P}}^k} \mathbb{E}_P[Q(x, \tilde{\omega})].$$

620 We begin by analyzing the behavior of the sequence  $\{g^k\}_{k \geq 1}$  as  $k \rightarrow \infty$ . In particular,  
 621 we will assess the sequence of function evaluations at a converging subsequence of  
 622 first-stage solutions. The result is captured in the following proposition.

623 **PROPOSITION 4.1.** *Suppose  $\{\hat{x}^{k_n}\}$  denotes a subsequence of  $\{\hat{x}^k\}$  such that  $\hat{x}^{k_n} \rightarrow$   
 624  $\bar{x}$ , then  $\lim_{n \rightarrow \infty} |g^{k_n}(\hat{x}^{k_n}) - f(\bar{x})| = 0$ , with probability one.*

625 *Proof.* Consider an approximate ambiguity set  $\hat{\mathfrak{P}}^k$ . For  $i = 1, 2$  and  $x_i \in \mathcal{X}$ , let  
 626  $P(x_i) \in \arg \max_{P \in \hat{\mathfrak{P}}^k} \{\mathbb{E}_P[Q(x_i, \tilde{\omega})]\}$ . Then,

$$\begin{aligned} 627 \quad g^k(x_1) &= c^\top x_1 + \mathbb{E}_{P(x_1)}[Q(x_1, \tilde{\omega})] \geq c^\top x_1 + \mathbb{E}_{P(x_2)}[Q(x_1, \tilde{\omega})] \\ 628 &= c^\top x_2 + \mathbb{E}_{P(x_2)}[Q(x_2, \tilde{\omega})] + c^\top(x_1 - x_2) + \\ 629 &\quad \mathbb{E}_{P(x_2)}[Q(x_1, \tilde{\omega})] - \mathbb{E}_{P(x_2)}[Q(x_2, \tilde{\omega})] \\ 630 &= g^k(x_2) + c^\top(x_1 - x_2) + \mathbb{E}_{P(x_2)}[Q(x_1, \tilde{\omega})] - \mathbb{E}_{P(x_2)}[Q(x_2, \tilde{\omega})]. \end{aligned}$$

632 The inequality in the above follows from optimality of  $P(x_1)$ . The above implies that

$$633 \quad g^k(x_2) - g^k(x_1) \leq c^\top(x_2 - x_1) + \mathbb{E}_{P(x_2)}[Q(x_2, \tilde{\omega})] - \mathbb{E}_{P(x_2)}[Q(x_1, \tilde{\omega})] \\ 634 \quad \leq |c^\top(x_2 - x_1)| + \left| \mathbb{E}_{P(x_2)}[Q(x_2, \tilde{\omega})] - \mathbb{E}_{P(x_2)}[Q(x_1, \tilde{\omega})] \right| \\ 635$$

636 The second relationship is due to the triangular inequality. Under assumption (A2),  
 637 the recourse function  $Q(x, \tilde{\omega})$  is a uniformly Lipschitz continuous function, with prob-  
 638 ability one (see Chapter 2 in [43] for details). This implies that there exists a constant  
 639  $C$  such that  $|\mathbb{E}_P[Q(x_1, \tilde{\omega})] - \mathbb{E}_P[Q(x_2, \tilde{\omega})]| \leq C\|x_1 - x_2\|$  for any probability distribu-  
 640 tion  $P$ . As a result,

$$641 \quad (4.2) \quad g^k(x_2) - g^k(x_1) \leq (\|c\| + C)\|x_2 - x_1\|. \\ 642$$

643 Starting with  $x_2$  and using the same arguments, we have

$$644 \quad (4.3) \quad g^k(x_1) - g^k(x_2) \leq (\|c\| + C)\|x_1 - x_2\|. \\ 645$$

646 Therefore, the function  $g^k(x)$  is equi-continuous on  $x \in \mathcal{X}$ . Now consider ambiguity  
 647 sets  $\mathfrak{P}$  and  $\widehat{\mathfrak{P}}^k$ . Note that for all  $x \in \mathcal{X}$ ,

$$648 \quad |f(x) - g^k(x)| = \left| \max_{P \in \mathfrak{P}} \mathbb{E}_P[Q(x, \tilde{\omega})] - \max_{P' \in \widehat{\mathfrak{P}}^k} \mathbb{E}_{P'}[Q(x, \tilde{\omega})] \right| \\ 649 \quad \leq \max_{P \in \mathfrak{P}} \min_{P' \in \widehat{\mathfrak{P}}^k} |\mathbb{E}_P[Q(x, \tilde{\omega})] - \mathbb{E}_{P'}[Q(x, \tilde{\omega})]| \\ 650 \quad \leq \max_{P \in \mathfrak{P}} \min_{P' \in \widehat{\mathfrak{P}}^k} \sup_{x \in \mathcal{X}} |\mathbb{E}_P[Q(x, \tilde{\omega})] - \mathbb{E}_{P'}[Q(x, \tilde{\omega})]|. \\ 651$$

652 Using the definition of deviation (1.6) and Hausdorff distance (1.7) between ambiguity  
 653 sets  $\mathfrak{P}$  and  $\widehat{\mathfrak{P}}^k$ , we have

$$654 \quad (4.4) \quad |f(x) - g^k(x)| \leq \mathbb{D}(\mathfrak{P}, \widehat{\mathfrak{P}}^k) \leq \mathbb{H}(\mathfrak{P}, \widehat{\mathfrak{P}}^k). \\ 655$$

656 For  $\hat{x}^{k_n}$  and  $\bar{x}$ , using the triangle inequality we have

$$657 \quad |f(\bar{x}) - g^{k_n}(\hat{x}^{k_n})| \leq |f(\bar{x}) - g^{k_n}(\bar{x})| + |g^{k_n}(\bar{x}) - g^{k_n}(\hat{x}^{k_n})| \\ 658 \quad \leq \mathbb{H}(\mathfrak{P}, \widehat{\mathfrak{P}}^{k_n}) + (\|c\| + C)\|\bar{x} - \hat{x}^{k_n}\|. \\ 659$$

660 The second inequality is justified by combining (4.2), (4.3), and (4.4). As  $n \rightarrow \infty$ ,  
 661  $\mathbb{H}(\mathfrak{P}, \widehat{\mathfrak{P}}^{k_n}) \rightarrow 0$  due to property (B2) of the considered family of ambiguity sets.  
 662 Furthermore, since  $\hat{x}^{k_n} \rightarrow \bar{x}$ , the right-hand side of the above inequality vanishes.  
 663 Therefore, we conclude that  $g^{k_n}(\hat{x}^{k_n}) \rightarrow f(\bar{x})$  as  $n \rightarrow \infty$ .  $\square$

664 Notice that the behavior of the approximate ambiguity sets defined in §2.1, in  
 665 particular, the condition  $\mathbb{H}(\mathfrak{P}, \widehat{\mathfrak{P}}^k) \rightarrow 0$  as  $k \rightarrow \infty$  plays a central role in the above  
 666 proof. Recall that for the moment and Wasserstein distance-based ambiguity sets, the  
 667 condition is established in propositions 2.1 and 2.3, respectively. It is also worthwhile  
 668 to note that under the foregoing conditions, (4.4) also implies uniform convergence of  
 669 the sequence  $\{g^k\}$  to  $f(x)$ , with probability one.

670 The above result applies to any algorithm that generates a converging sequence of  
 671 iterates  $\{x^k\}$  and a corresponding sequence of extremal distributions. Such an algo-  
 672 rithm is guaranteed to exhibit convergence to the optimal distributionally robust ob-  
 673 jective function value. Therefore, this result is applicable to the sequence of instances

674 constructed using external sampling and solved, for example, using reformulation-based methods. Such an approach was adopted in [35] and [45]. The analysis in [35] 675 relies upon two rather restrictive assumptions. The first assumption is that for all 676  $P \in \mathfrak{P}$ , there exists a sequence of measures  $\{P^k\}$  such that  $P^k \in \widehat{\mathfrak{P}}^k$  and converges 677 weakly to  $P$ . The second assumption requires the approximate ambiguity sets to be 678 strict subsets of the true ambiguity set, i.e.,  $\widehat{\mathfrak{P}}^k \subset \mathfrak{P}$ . Both of these assumptions are 679 very difficult to satisfy in a data-driven setting (also see Remark 2.2). 680

681 The analysis in [45], on the other hand, does not make the above assumptions. 682 Therefore, their analysis is more broadly applicable in settings where external sampling 683 is used to generate  $\Omega^k$ . DRO instances are constructed based on statistics 684 estimated using  $\Omega^k$  and solved to optimality for each  $k \geq 1$ . They show the convergence 685 of optimal objective function values and optimal solution sets of approximate 686 problems to the optimal objective function value and solutions of the true DRO problem, 687 respectively. In this regard, the result in Proposition 4.1 can alternatively be derived 688 using Theorem 1(i) in [45]. While the above function is not computed during 689 the course of the sequential sampling algorithm, it provides the necessary benchmark 690 for our convergence analysis. 691

691 One of the main point of deviation in our analysis stems from the fact that we 692 use the objective function approximations that are built based on the approximate 693 recourse function in (2.11). In order to study the piecewise affine approximation of 694 the first-stage objective function, we introduce another benchmark function 695

$$(4.5) \quad \phi^k(x) := c^\top x + \max_{P \in \widehat{\mathfrak{P}}^k} \mathbb{E}_P[Q^k(x, \tilde{\omega})].$$

697 Notice that the above function uses the approximations for the ambiguity set (as in the 698 case of (4.1)) as well as the approximation of the recourse function. This construction 699 ensures that  $\phi^k(x) \leq g^k(x)$  for all  $x \in \mathcal{X}$  and  $k \geq 1$ , which follows from the fact that 700  $Q^k(x, \tilde{\omega}) \leq Q(x, \tilde{\omega})$ , almost surely. Further, the result in Theorem 3.1 ensures that 701  $f^k(x) \leq \phi^k(x)$ . Putting these together, we obtain the following relationship:

$$(4.6) \quad f^k(x) \leq \phi^k(x) \leq g^k(x) \quad \forall x \in \mathcal{X}, k \geq 1.$$

702 While the previous proposition was focused on the upper limit in the above relationship, we present the asymptotic behavior of the  $\{f^k\}$  sequence in the following 703 results.

704 LEMMA 4.2. *Suppose  $\{\hat{x}^{k_n}\}$  denotes a subsequence of  $\{\hat{x}^k\}$  such that  $\hat{x}^{k_n} \rightarrow \bar{x}$ . 705 Then,  $\lim_{n \rightarrow \infty} f^{k_n}(\hat{x}^{k_n}) - f(\bar{x}) = 0$ , with probability one.*

706 *Proof.* From Proposition 4.1, we have  $\lim_{n \rightarrow \infty} |f(\bar{x}) - g^{k_n}(\hat{x}^{k_n})| = 0$ . Therefore, 707 there exists  $N_1 < \infty$  and  $\epsilon_1 > 0$  such that

$$(4.7) \quad \left| \max_{P \in \mathfrak{P}} \mathbb{E}_P[Q(\bar{x}, \tilde{\omega})] - \max_{P \in \widehat{\mathfrak{P}}^{k_n}} \mathbb{E}_P[Q(\hat{x}^{k_n}, \tilde{\omega})] \right| < \epsilon_1/2 \quad \forall n > N_1.$$

713 Now consider,

$$\begin{aligned} 714 \quad & \left| \max_{P \in \widehat{\mathfrak{P}}^{k_n}} \mathbb{E}_P[Q(\hat{x}^{k_n}, \tilde{\omega})] - \max_{P \in \widehat{\mathfrak{P}}^{k_n}} \mathbb{E}_P[Q^{k_n}(\hat{x}^{k_n}, \tilde{\omega})] \right| \\ 715 \quad & \leq \max_{P \in \widehat{\mathfrak{P}}^{k_n}} |\mathbb{E}_P[Q(\hat{x}^{k_n}, \tilde{\omega})] - \mathbb{E}_P[Q^{k_n}(\hat{x}^{k_n}, \tilde{\omega})]| \\ 716 \quad & = \max_{P \in \widehat{\mathfrak{P}}^{k_n}} \mathbb{E}_P[|Q(\hat{x}^{k_n}, \tilde{\omega}) - Q^{k_n}(\hat{x}^{k_n}, \tilde{\omega})|]. \end{aligned}$$

718 The last equality follows from the fact that  $Q(x, \tilde{\omega}) \geq Q^k(x, \tilde{\omega})$  for all  $x \in \mathcal{X}$  and  $k \geq 1$ ,  
 719 almost surely. Moreover, because of the uniform convergence of  $\{Q^k\}$  (Proposition  
 720 2.5), the sequence of approximate functions  $\{\phi^k\}$  also converges uniformly. This  
 721 implies that, there exists  $N_2 < \infty$  such that for all  $n > N_2$ ,

$$722 \quad (4.8) \quad \left| \max_{P \in \widehat{\mathfrak{P}}^{k_n}} \mathbb{E}_P[Q(\hat{x}^{k_n}, \tilde{\omega})] - \max_{P \in \widehat{\mathfrak{P}}^{k_n}} \mathbb{E}_P[Q^{k_n}(\hat{x}^{k_n}, \tilde{\omega})] \right| < \epsilon_1/2.$$

724 Let  $N = \max\{N_1, N_2\}$ . Using (4.7) and (4.8), we have for all  $n > N$ ,

$$725 \quad \left| \max_{P \in \mathfrak{P}} \mathbb{E}_P[Q(\bar{x}, \tilde{\omega})] - \max_{P \in \widehat{\mathfrak{P}}^{k_n}} \mathbb{E}_P[Q^{k_n}(x^{k_n}, \tilde{\omega})] \right| < \epsilon_1.$$

727 This implies that  $|f(\bar{x}) - \phi^{k_n}(\hat{x}^{k_n})| \rightarrow 0$  as  $n \rightarrow \infty$ . Based on (3.2), we have  
 728  $Q^{k_n}(\hat{x}^{k_n}, \omega) = (\pi(\hat{x}^{k_n}, \omega))^\top [r(\omega) - T(\omega)\hat{x}^{k_n}] \geq (\pi(\hat{x}^{k_n}, \omega))^\top [r(\omega) - T(\omega)x]$  for all  
 729  $x \in \mathcal{X}$  and  $\omega \in \Omega^{k_n}$ . Let

$$730 \quad \alpha_{k_n}^{k_n} = \sum_{\omega \in \Omega^{k_n}} p^{k_n}(\omega) (\pi(\hat{x}^{k_n}, \omega))^\top r(\omega) \text{ and } \beta_{k_n}^{k_n} = - \sum_{\omega \in \Omega^{k_n}} p^{k_n}(\omega) T(\omega)^\top \pi(\hat{x}^{k_n}, \omega),$$

732 where  $\{p^{k_n}(\omega)\}_{\omega \in \Omega^{k_n}}$  is an optimal solution of the distributional separation problem  
 733 (3.3) where index  $k$  is replaced by  $k_n$ . Then, the affine function  $\alpha_{k_n}^{k_n} + (c + \beta_{k_n}^{k_n})^\top x$   
 734 provides a lower bound approximation for function  $\phi^{k_n}(x)$ , i.e.,

$$735 \quad \phi^{k_n}(x) \geq \alpha_{k_n}^{k_n} + (c + \beta_{k_n}^{k_n})^\top x \quad \text{for all } x \in \mathcal{X},$$

737 with strict equality holding only at  $\hat{x}^{k_n}$ . Therefore, using the definition of  $f^k(x)$  we  
 738 have  $\lim_{n \rightarrow \infty} \alpha_{k_n}^{k_n} + (c + \beta_{k_n}^{k_n})^\top \hat{x}^{k_n} = \lim_{n \rightarrow \infty} f^{k_n}(\hat{x}^{k_n}) = \lim_{n \rightarrow \infty} \phi^{k_n}(\hat{x}^{k_n}) = f(\bar{x})$ ,  
 739 almost surely. This completes the proof.  $\square$

740 The above result characterizes the behavior of the sequence of affine functions  
 741 generated during the course of the algorithm. In particular, the sequence  $\{f^k(\hat{x}^k)\}_{k \geq 1}$   
 742 accumulates at the objective value of the original DRO problem (1.1). Recall that the  
 743 candidate solution  $x^k$  is a minimizer of  $f^{k-1}(x)$  and an affine function is generated  
 744 at this point such that  $f^k(x^k) = \phi^k(x^k)$  in all iterations  $k \geq 1$ . However, due to  
 745 the update procedure in (3.6) the quality of the estimates at  $x^k$  gradually diminishes  
 746 leading to a large value for  $(\phi^k(x^k) - f^k(x^k))$  as  $k$  increases. This emphasizes the role  
 747 of the incumbent solution and computing the incumbent affine function  $\hat{\ell}(x)$  during  
 748 the course of the algorithm. By updating the incumbent solution and frequently  
 749 reevaluating the affine functions at the incumbent solution, we can ensure that the  
 750 approximation is “sufficiently good” in the neighborhood of the incumbent solution.  
 751 In order to assess the improvement of approximation quality, we define

$$752 \quad (4.9) \quad \delta^k := f^{k-1}(x^k) - f^{k-1}(\hat{x}^{k-1}) \leq 0 \quad \forall k \geq 1.$$

754 The inequality follows from the optimality of  $x^k$  with respect to the objective func-  
 755 tion  $f^{k-1}$ . The quantity  $\delta^k$  measures the error in objective function estimate at the  
 756 candidate solution with respect to the estimate at the current incumbent solution.  
 757 The following result captures the asymptotic behavior of this error term.

758 LEMMA 4.3. *Let  $\mathcal{K}$  denote a sequence of iterations where the incumbent solution  
 759 changes. There exists a subsequence of iterations, denoted as  $\mathcal{K}^* \subseteq \mathcal{K}$ , such that  
 760  $\lim_{k \in \mathcal{K}^*} \delta^k = 0$ .*

761 *Proof.* We will consider two cases depending on whether the set  $\mathcal{K}$  is finite or not.  
 762 First, suppose that  $|\mathcal{K}|$  is not finite. By the incumbent update rule and (4.9),

763 
$$f^{k_n}(x^{k_n}) - f^{k_n}(\hat{x}^{k_n-1}) < \gamma[f^{k_n-1}(x^{k_n}) - f^{k_n-1}(\hat{x}^{k_n-1})] = \gamma\delta^{k_n} \leq 0 \quad \forall k_n \in \mathcal{K}.$$

765 Subsequently, we have  $\limsup_{n \rightarrow \infty} \delta^{k_n} \leq 0$ . Since  $x^{k_n} = \hat{x}^{k_n}$  and  $\hat{x}^{k_n-1} = \hat{x}^{k_n-1}$ , we  
 766 have

767 
$$f^{k_n}(\hat{x}^{k_n}) - f^{k_n}(\hat{x}^{k_n-1}) \leq \gamma\delta^{k_n} \leq 0.$$

769 The left-hand side of the above inequality captures the improvement in the objective  
 770 function value at the current incumbent solution over the previous incumbent solution.  
 771 Using the above, we can write the average improvement attained over  $n$  incumbent  
 772 changes as

773 
$$\frac{1}{n} \sum_{j=1}^n \left[ f^{k_j}(\hat{x}^{k_j}) - f^{k_j}(\hat{x}^{k_{j-1}}) \right] \leq \frac{1}{n} \sum_{j=1}^n \gamma\delta^{k_j} \leq 0 \quad \text{for all } n.$$

775 This implies that

776 
$$\underbrace{\frac{1}{n} \left( f^{k_n}(\hat{x}^{k_n}) - f^{k_1}(\hat{x}^{k_0}) \right)}_{(a)} + \underbrace{\frac{1}{n} \left[ \sum_{j=1}^{n-1} \left( f^{k_j}(\hat{x}^{k_j}) - f^{k_{j+1}}(\hat{x}^{k_j}) \right) \right]}_{(b)} \leq \frac{1}{n} \sum_{j=1}^n \gamma\delta^{k_j} \leq 0,$$

777 for all  $n$ . Under the assumption that the dual feasible region is non-empty and  
 779 bounded (this is ensured by relatively complete recourse, (A2)),  $\{f^k\}$  is a sequence  
 780 of Lipschitz continuous functions. This, along with compactness of  $\mathcal{X}$  (A1), implies  
 781 that  $f^{k_n}(\hat{x}^{k_n}) - f^{k_1}(\hat{x}^{k_0})$  is bounded from above. Hence, the term (a) reduces to zero  
 782 as  $n \rightarrow \infty$ . The term (b) converges to zero, with probability one, due to uniform  
 783 convergence of  $\{f^k\}$ . Since  $\gamma \in (0, 1)$ , we have

784 
$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n \delta^{k_j} = 0,$$

785 with probability one. Further,

787 
$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{j=1}^n \delta^{k_j} \leq \limsup_{n \rightarrow \infty} \delta^{k_n} \leq 0.$$

789 Thus, there exists a subsequence indexed by the set  $\mathcal{K}^*$  such that  $\lim_{k \in \mathcal{K}^*} \delta^k = 0$ ,  
 790 with probability one.

791 Now if  $|\mathcal{K}|$  is finite, then there exists  $\hat{x}$  and  $K < \infty$  such that for all  $k \geq K$ , we  
 792 have  $\hat{x}^k = \hat{x}$ . Notice that, if  $\lim_{k \in \mathcal{K}^*} x^k = \bar{x}$ , uniform convergence of the sequence  
 793  $\{f^k\}$  and Lemma 4.2 ensure that

794 (4.10a) 
$$\lim_{k \in \mathcal{K}^*} f^k(x^k) = \lim_{k \in \mathcal{K}^*} f^{k-1}(x^k) = f(\bar{x})$$

795 (4.10b) 
$$\lim_{k \in \mathcal{K}^*} f^k(\hat{x}) = \lim_{k \in \mathcal{K}^*} f^{k-1}(\hat{x}) = f(\hat{x}).$$

797 Further, since the incumbent is not updated in iterations  $k \geq K$ , we must have from  
 798 the update rule in (3.9) that

799 
$$f^k(x^k) - f^k(\hat{x}) \geq \gamma[f^{k-1}(x^k) - f^{k-1}(\hat{x})] = \gamma\delta^k \quad \text{for all } k \geq K.$$

801 Using (4.10), we have

$$802 \quad \lim_{k \in \mathcal{K}^*} (f^k(x^k) - f^k(\hat{x})) \geq \gamma \lim_{k \in \mathcal{K}^*} (f^{k-1}(x^k) - f^{k-1}(\hat{x})),$$

$$803 \quad \Rightarrow \quad f(\bar{x}) - f(\hat{x}) \geq \gamma(f(\bar{x}) - f(\hat{x})).$$

805 Noting that  $\gamma \in (0, 1)$ , the above inequality reduces to  $f(\bar{x}) - f(\hat{x}) \geq 0$ . Further,  
 806 using (4.9) in the limit as  $k \rightarrow \infty$  and the fact that  $\hat{x}^k = \hat{x}$  for all  $k \geq K$ , we  
 807 have  $f(\bar{x}) - f(\hat{x}) \leq 0$ . Therefore, we have  $f(\bar{x}) - f(\hat{x}) = 0$ . Hence,  $\lim_{k \in \mathcal{K}^*} \delta^k =$   
 808  $f(\bar{x}) - f(\hat{x}) = 0$ , with probability one.  $\square$

809 Equipped with the results in lemmas 4.2 and 4.3, we state the main theorem  
 810 which establishes the existence of a subsequence of the incumbent sequence generated  
 811 by the algorithm for which every accumulation point is an optimal solution to (1.1).

812 **THEOREM 4.4.** *Let  $\{x^k\}_{k=1}^\infty$  and  $\{\hat{x}^k\}_{k=1}^\infty$  be the sequence candidate and incumbent  
 813 solutions generated by the DRSD algorithm. There exists a subsequence  $\{\hat{x}^k\}_{k \in \mathcal{K}}$   
 814 for which every accumulation point is an optimal solution of 2-DRLP (1.1), with  
 815 probability one.*

816 *Proof.* Let  $x^* \in \mathcal{X}$  be an optimal solution of (1.1). Consider a subsequence  
 817 indexed by  $\mathcal{K}$  such that  $\lim_{k \in \mathcal{K}} \hat{x}^k = \bar{x}$ . Compactness of  $\mathcal{X}$  ensures the existence of  
 818 accumulation point  $\bar{x} \in \mathcal{X}$  and therefore,

$$820 \quad (4.11) \quad f(x^*) \leq f(\bar{x}).$$

821 From Theorem 3.1, we have for all  $k, x \in \mathcal{X}$

$$822 \quad f^k(x) \leq c^\top x + \mathbb{Q}^k(x) \leq c^\top x + \max_{P \in \mathfrak{P}^k} \mathbb{E}_P[Q(x, \tilde{\omega})] = g^k(x).$$

824 Thus, using the uniform convergence of  $\{g^k\}$  (Proposition 4.1) we have

$$825 \quad (4.12) \quad \limsup_{k \in \mathcal{K}'} f^k(x^*) \leq \lim_{k \in \mathcal{K}'} g^k(x^*) = f(x^*)$$

827 for all subsequences indexed by  $\mathcal{K}' \subseteq \{1, 2, \dots\}$ , with probability one. Recall that,

$$828 \quad \delta^k = f^{k-1}(x^k) - f^{k-1}(\hat{x}^{k-1}) \leq f^{k-1}(x^*) - f^{k-1}(\hat{x}^{k-1}) \quad \text{for all } k \geq 1.$$

830 The inequality in the above follows from the optimality of  $x^k$  with respect to  $f^{k-1}(x)$ .  
 831 Taking limit over  $\mathcal{K}$ , we have

$$832 \quad \lim_{k \in \mathcal{K}} \delta^k \leq \lim_{k \in \mathcal{K}} (f^{k-1}(x^*) - f^{k-1}(\hat{x}^{k-1}))$$

$$833 \quad \leq \limsup_{k \in \mathcal{K}} f^{k-1}(x^*) - \liminf_{k \in \mathcal{K}} f^{k-1}(\hat{x}^{k-1}) \leq f(x^*) - f(\bar{x}).$$

835 The last inequality follows from (4.12) and  $\lim_{k \in \mathcal{K}} f^{k-1}(\hat{x}^{k-1}) = f(\bar{x})$  (Lemma  
 836 4.2). From Lemma 4.3, there exists a subsequence indexed by  $\mathcal{K}^* \subseteq \mathcal{K}$  such that  
 837  $\lim_{k \in \mathcal{K}^*} \delta^k = 0$ . Therefore, if  $\{\hat{x}^k\}_{k \in \mathcal{K}^*} \rightarrow \bar{x}$ , we have

$$839 \quad f(x^*) - f(\bar{x}) \geq 0.$$

840 Using (4.11) and the above inequality, we conclude that  $\bar{x}$  is an optimal solution with  
 841 probability one.  $\square$

842     **5. Computational Experiment.** In this section, we evaluate the effectiveness  
 843 and efficiency of the DRSD method in solving 2-DRLPs. For our preliminary experiments,  
 844 we consider 2-DRLPs with moment-based ambiguity set  $\mathfrak{P}_{\text{mom}}$  for the first  
 845 two moments ( $q = 2$ ). We also consider 2-DRLPs with Wasserstein ambiguity set  $\mathfrak{P}_w$   
 846 with  $\ell_1$  and  $\ell_\infty$  distance metrics.

847     We report results from the computational experiments conducted on four well-  
 848 known SP test problems: the capacity expansion planning (CEP) [26], the power generation  
 849 planning (PGP) [26], multilocation transshipment (RETAIL) [23], and cargo flight scheduling (STORM) [33]. In Table 1, we provide the number of variables (#Var)

TABLE 1  
*Details of CEP, PGP, RETAIL, and STORM Test Problems, and Computational Results for the SD Algorithm*

Problem	Stage I		Stage II				SD Results		
	#Var	#Cons	#Var	#Cons	#RV	\Omega	#Iter	ObjEst	Time
PGP	4	2	16	7	3	576	215( $\pm 8$ )	446( $\pm 2.4$ )	0.43( $\pm 0.04$ )
CEP	8	5	15	7	3	216	153( $\pm 7$ )	343886( $\pm 12783$ )	0.18( $\pm 0.02$ )
RETAIL	7	0	70	22	7	$10^{11}$	721( $\pm 44$ )	154( $\pm 1.92$ )	4.20( $\pm 0.76$ )
STORM	121	185	1259	528	117	$10^{81}$	238( $\pm 17$ )	15173494( $\pm 657272$ )	2.83( $\pm 0.21$ )

850 and constraints (#Cons) in the first- and second-stage of the test problems. Notice  
 851 that the PGP and CEP have relatively smaller supports (216 and 576, respectively),  
 852 while RETAIL and STORM have a support size of  $10^{11}$  and  $10^{81}$ , respectively. In the  
 853 table, we also provide computational results from solving the risk-neutral versions of  
 854 these problems using the SD algorithm [41]. For these results, we report the number  
 855 of iterations (#Iter), objective function estimate (ObjEst) at termination, and total  
 856 time (in seconds) taken by the SD algorithm. We refer the readers interested in a  
 857 computational comparison between SD and an external sampling-based approach for  
 858 risk-neutral 2-SLPs to [41] and [20].

859     Following the rule of thumb adopted in experiments involving sampling-based  
 860 SP, we conduct 30 independent replications for each problem instance. The choice  
 861 of 30 replications is the same as in previous experiments with SD (see [20] and [41],  
 862 for example). Each replication uses a different seed for the random number generator.  
 863 The algorithms are implemented in the C programming language, and the  
 864 experiments are conducted on a 64-bit Intel core i7 - 4770 CPU at 3.4GHz  $\times$  8  
 865 machine with 32 GB memory. All linear programs, i.e., master problem, subproblems,  
 866 and distribution separation problem, are solved using CPLEX 12.10 callable subroutines.  
 867 For DRSD, we use  $\tau = 0.001$  and  $\gamma = 0.2$  in our experiments. We add one  
 868 new observation to the sample in every iteration and therefore,  $\theta^k = \frac{k-1}{k}$  is used  
 869 for the updates in (3.6). The source code for the DR L-shaped, DRSD algorithms,  
 870 and the reformulation techniques are available under the GNU general public license  
 871 at <https://github.com/SMU-SODA/distributionallyRobust.git>. The repository also  
 872 includes the test problems in SMPS file format.

873     **5.1. Results for 2-DRLPs with Moment-based Ambiguity Set.** The first  
 874 set of experiments concerns the 2-DRLP problems with a moment-based ambiguity  
 875 set  $\mathfrak{P}_{\text{mom}}$  for which we use an external sampling-based approach as a benchmark for  
 876 comparison with DRSD. The external sampling-based instances involve constructing  
 877 approximate problems of the form (4.1) with a pre-determined number of observations  
 878  $N \in \{100, 250, 500, 1000\}$ . The resulting instances are solved using the DR L-Shaped  
 879 method. For a fair comparison, the DRSD method is run for a maximum of  $N$  itera-  
 880

TABLE 2  
*Computational Results for 2-DRLP Instances with Moment-based Ambiguity Set*

N	DRSD Algorithm			DR L-Shaped Algorithm		
	#Iter	ObjEst	Time	#Iter	ObjEst	Time
PGP						
100	100 ( $\pm 0$ )	460.89 ( $\pm 3.76$ )	0.04 ( $\pm 0.00$ )	18 ( $\pm 0.9$ )	457.61 ( $\pm 3.28$ )	0.052 ( $\pm 0.00$ )
250	250 ( $\pm 0$ )	466.91 ( $\pm 2.52$ )	0.13 ( $\pm 0.00$ )	20 ( $\pm 0.7$ )	462.92 ( $\pm 2.28$ )	0.077 ( $\pm 0.00$ )
500	500 ( $\pm 0$ )	471.40 ( $\pm 3.49$ )	0.32 ( $\pm 0.00$ )	20 ( $\pm 0.6$ )	464.70 ( $\pm 1.95$ )	0.096 ( $\pm 0.00$ )
1000	504 ( $\pm 687$ )	463.19 ( $\pm 16.28$ )	0.35 ( $\pm 0.70$ )	20 ( $\pm 0.8$ )	466.10 ( $\pm 1.78$ )	0.121 ( $\pm 0.00$ )
CEP						
100	100 ( $\pm 0$ )	658831 ( $\pm 14453$ )	0.04 ( $\pm 0.00$ )	3 ( $\pm 0.2$ )	658817 ( $\pm 14457$ )	0.015 ( $\pm 0.00$ )
250	250 ( $\pm 0$ )	680795 ( $\pm 10524$ )	0.12 ( $\pm 0.00$ )	2 ( $\pm 0.2$ )	680736 ( $\pm 10511$ )	0.024 ( $\pm 0.00$ )
500	256 ( $\pm 0$ )	683300 ( $\pm 5955$ )	0.30 ( $\pm 0.00$ )	20 ( $\pm 0.6$ )	683252 ( $\pm 5949$ )	0.028 ( $\pm 0.00$ )
1000	256 ( $\pm 0$ )	683300 ( $\pm 5955$ )	0.30 ( $\pm 0.00$ )	2 ( $\pm 0$ )	679665 ( $\pm 4926$ )	0.028 ( $\pm 0.00$ )
RETAIL						
100	100 ( $\pm 0$ )	326.21 ( $\pm 15.35$ )	0.07 ( $\pm 0.00$ )	46 ( $\pm 1$ )	327.26 ( $\pm 14.79$ )	0.370 ( $\pm 0.01$ )
250	250 ( $\pm 0$ )	365.00 ( $\pm 19.03$ )	0.27 ( $\pm 0.01$ )	45 ( $\pm 2$ )	365.54 ( $\pm 19.45$ )	0.839 ( $\pm 0.03$ )
500	500 ( $\pm 0$ )	387.98 ( $\pm 17.10$ )	0.86 ( $\pm 0.02$ )	45 ( $\pm 1$ )	388.84 ( $\pm 17.48$ )	1.587 ( $\pm 0.05$ )
1000	625 ( $\pm 31$ )	396.67 ( $\pm 15.99$ )	1.13 ( $\pm 0.12$ )	45 ( $\pm 1$ )	401.71 ( $\pm 14.38$ )	3.176 ( $\pm 0.09$ )
STORM						
100	100 ( $\pm 0$ )	15755337 ( $\pm 12314$ )	0.74 ( $\pm 0.02$ )	12 ( $\pm 0.51$ )	15742456 ( $\pm 12192$ )	0.434 ( $\pm 0.02$ )
250	250 ( $\pm 0$ )	15795815 ( $\pm 8493$ )	4.66 ( $\pm 0.13$ )	11 ( $\pm 0.52$ )	15781725 ( $\pm 8754$ )	1.008 ( $\pm 0.05$ )
500	500 ( $\pm 0$ )	15811923 ( $\pm 5233$ )	20.54 ( $\pm 0.51$ )	12 ( $\pm 0.59$ )	15797020 ( $\pm 5346$ )	2.117 ( $\pm 0.10$ )
1000	516 ( $\pm 108$ )	15786865 ( $\pm 9155$ )	30.44 ( $\pm 15.28$ )	12 ( $\pm 0.52$ )	15806575 ( $\pm 3772$ )	4.318 ( $\pm 0.19$ )

881 tions to have an estimate based upon a sample of size not greater than  $N$ . Specifically,  
882 it terminates when conditions in Step 5 of Algorithm 3.1 are satisfied. We compare  
883 the solution quality provided by these methods along with the computational time.  
884 The results from the experiments are presented in Table 2.

885 Table 2 shows the number of iterations, objective function estimate  $f^k(\hat{x}^k)$  at  
886 termination, and solution time (in seconds) averaged over 30 replications. The val-  
887 ues in the parenthesis are the half-widths of the corresponding confidence intervals.  
888 Similar to SD, the number of iterations for DRSD is also equal to the number of  
889 observations used to approximate the ambiguity set. To begin, observe the increase  
890 in the objective function estimates of distributionally robust variants when compared  
891 to the risk-neutral results from SD in Table 1.

892 The objective function estimate obtained using the DRSD is comparable to the  
893 objective function estimate obtained using the DR L-shaped method. Notice that for  
894 instances with  $N = 1000$ , DRSD took less than 1000 iterations because the termi-  
895 nation conditions were satisfied. The same is true for CEP instances with  $N = 500$ .  
896 This shows the potential ability of DRSD to dynamically determine the number of  
897 observations by assessing the progress made during the algorithm. For instance, the  
898 DRSD objective function estimate for STORM that is based upon a sample of size 516  
899 (on average) is within 0.1% and 0.12% of the objective function value estimate pro-  
900 vided by the DR L-shaped method for  $N = 500$  and  $N = 1000$ , respectively. These  
901 results show that the optimal objective function estimate obtained from DRSD are  
902 comparable to those obtained using an external sampling-based approach.

903 The results for small scale instances (PGP and CEP) show that both DRSD and

904 the DR L-shaped method take a fraction of a second, but the computational time  
 905 for DRSD is higher than the DR L-shaped method for all  $N$ . We attribute this  
 906 behavior to two reasons. (i) The computational effort to solve all the subproblems  
 907 in each iteration does not increase significantly with  $N$  as they are easy to solve.  
 908 This observation is in-line with our computational experience with the SD method  
 909 for 2-SLPs (see [20]). (ii) The DRSD takes a larger number of iterations, resulting  
 910 in an increased number of master and distribution separation problems solved. It is  
 911 important to note that, while the computational time for the DR L-shaped method  
 912 on an individual instance may be lower, the iterative procedure necessary to identify  
 913 a sufficient sample size may require solving several instances with increasing sample  
 914 size. This may result in a significantly higher cumulative computational time.

915 On the other hand, for large-scale problems (**RETAIL** and **STORM**), we observe a  
 916 noticeable increase in the computational time for the DR L-shaped method with an  
 917 increase in  $N$ . A significant portion of this time is spent solving the subproblems.  
 918 Since the DRSD solves only two subproblems in each iteration, the time taken to solve  
 919 the subproblems is significantly less in comparison to the DR L-shaped method where  
 920 all subproblems corresponding to unique observations are solved in each iteration.  
 921 Notice that for **RETAIL**, the average number of iterations taken by DRSD is at least  
 922 8.2 times the average number of iterations taken by DR L-shaped for any  $N$ . This  
 923 increases the computational time spent for solving the master and distributional sepa-  
 924 ration problems. However, the reduction in the overall computational time is a direct  
 925 consequence of solving only two subproblems in each iteration. The results for **STORM**  
 926 also show similar behavior in terms of computational time associated with solving  
 927 master and subproblems. However, the overall increase in the computational time is  
 928 due to a significant computational expense ( $\sim 78\%$ ) in naively solving the distribution  
 929 separation problem. This computational time associated with solving the distribution  
 930 separation problem can be reduced by using column-generation procedures that take  
 931 advantage of the problem structure. Such an implementation is not undertaken for  
 932 our current experiments and is a fruitful future research avenue.

933 **5.2. Results for 2-DRLPs with  $\ell_1$ -type Wasserstein Ambiguity Set.** For  
 934 the Wasserstein distance-based ambiguity sets, we benchmark against the reformation  
 935 techniques proposed by [49]. Specifically, in [49], it has been shown that a 2-DRLP  
 936 (1.1) with Wasserstein ambiguity set can be reformulated as a two-stage robust optimi-  
 937 zation problem. This reformulation is given by

$$938 \quad (5.1) \quad \min_{x \in X, \eta \geq 0} \left\{ c^\top x + \eta \epsilon + \frac{1}{N_s} \sum_{n=1}^{N_s} \max_{\omega \in \Omega} \{Q(x, \omega) - \eta \|\omega - \bar{\omega}_n\|\} \right\},$$

939 where  $\{\bar{\omega}_1, \dots, \bar{\omega}_{N_s}\}$  is a finite set of observations obtained using true distribution.  
 940 Notice that the reformulation (5.1) can be written as the following semi-infinite pro-  
 941 gram:

$$943 \quad \min_{x \in X, \eta \geq 0} c^\top x + \eta \epsilon + \frac{1}{N_s} \sum_{n=1}^{N_s} \nu_n : \\ 944 \quad \text{s.t. } Q(x, \omega) - \eta \|\omega - \bar{\omega}_n\| \leq \nu_n, \quad n \in \{1, \dots, N_s\}, \omega \in \Omega.$$

945 For problem instances with  $\ell_1$ -type Wasserstein ambiguity set, we solve the foregoing  
 946 program using a Benders decomposition approach.

947 For  $\ell_1$ -norm, the reformulation in (5.1) admits the application of Benders decom-  
 948 position algorithm. To address the semi-infinite nature of the linear program, [19]

TABLE 3  
*Computational Results for 2-DRLP Instances with Wasserstein-1 Ambiguity Set*

Problem	N	DRSD Algorithm		Reformulation Approach [49]	
		ObjEst	Time	ObjEst	Time
PGP	100	447.04 ( $\pm 3.34$ )	0.05 ( $\pm 0.00$ )	444.85 ( $\pm 3.26$ )	0.94 ( $\pm 0.07$ )
	250	454.06 ( $\pm 2.64$ )	0.25 ( $\pm 0.04$ )	449.85 ( $\pm 2.23$ )	7.02 ( $\pm 0.63$ )
	500	457.48 ( $\pm 2.76$ )	1.79 ( $\pm 0.26$ )	451.57 ( $\pm 1.92$ )	27.73 ( $\pm 1.81$ )
CEP	100	338295.71 ( $\pm 14430.81$ )	0.25 ( $\pm 0.01$ )	338295.71 ( $\pm 14430.81$ )	0.32 ( $\pm 0.02$ )
	250	355054.48 ( $\pm 11823.37$ )	3.12 ( $\pm 0.09$ )	355054.48 ( $\pm 11823.37$ )	2.29 ( $\pm 0.08$ )
	500	356757.34 ( $\pm 6917.77$ )	13.24 ( $\pm 0.26$ )	356757.34 ( $\pm 6917.77$ )	10.90 ( $\pm 0.15$ )
RETAIL	100	157.09 ( $\pm 4.00$ )	0.53 ( $\pm 0.04$ )	153.67 ( $\pm 3.89$ )	9.41 ( $\pm 0.67$ )
	250	155.32 ( $\pm 3.39$ )	7.75 ( $\pm 0.22$ )	154.06 ( $\pm 3.40$ )	331.15 ( $\pm 13.89$ )
	500	155.20 ( $\pm 2.39$ )	72.02 ( $\pm 2.05$ )	154.62 ( $\pm 2.38$ )	2189.66 ( $\pm 65.97$ )
STORM	100	15504501.91 ( $\pm 11397$ )	0.64 ( $\pm 0.04$ )	15498236.10 ( $\pm 11445$ )	21.51 ( $\pm 1.09$ )
	250	15508623.20 ( $\pm 7481$ )	8.86 ( $\pm 0.19$ )	15501074.50 ( $\pm 7571$ )	333.22 ( $\pm 11.80$ )
	500	15507815.12 ( $\pm 5059$ )	83.33 ( $\pm 2.55$ )	-	-

950 consider a special case where the sample space  $\Omega$  is defined by a bounded hyper-  
951 rectangle and derive a finite subset of the sample space (without loss of optimality)  
952 using extreme points of the hyper-rectangle. Since the test problems used in our  
953 experiments do not impose any restrictions on  $\Omega$ , we adopt a sampling-based dis-  
954cretization of the ambiguity set to tackle (5.1). Such a discretization satisfies the  
955 result in Proposition 2.3 and therefore, provides a suitable benchmark for DRSD.  
956 We use the reformulation corresponding to ambiguity set defined by the finite set of  
957 observations, i.e.  $\Omega := \{\omega_1, \dots, \omega_N\}$ ,  $N_s = N$ , and  $\bar{\omega}_i = \omega_i$  for  $i = 1, \dots, N$ .

958 In this second set of experiments, we consider  $N = 100, 150$ , and  $500$  observations.  
959 We use an external sampling approach to construct the instances of reformulation and  
960 solve these instances using the Benders decomposition method. We run the DRSD  
961 algorithm for the same number of iterations ( $N$ ) to have the same set of observations  
962 for approximating the ambiguity set (recall that we run replications of both algorithms  
963 with the same seed for random number generation). The results of this experiment  
964 are shown in Table 3 for  $\epsilon = 0.05$ . The table shows the average objective function  
965 estimates and computational time (in seconds) computed across 30 replications along  
966 with half widths of the corresponding confidence interval.

967 The results indicate that the estimates of the objective function obtained from  
968 the DRSD algorithm and the reformulation approach are comparable. For all the  
969 test problems, the computational time for both approaches increases with  $N$ . We  
970 attribute this to the increase in the size of the master problem. While the additional  
971 effort associated with solving distribution separation problems also contributes to  
972 the increased computation time in DRSD, the number of subproblems solved in each  
973 iteration of Benders decomposition increases with  $N$ . In any case, DRSD outperforms  
974 Benders decomposition applied to the reformulation across all test problems. Since we  
975 ran out of memory when solving the instances of STORM with  $N = 500$  using Benders  
976 decomposition, we do not report its results.

977 **5.3. Results for 2-DRLPs with  $\ell_\infty$ -type Wasserstein Ambiguity Set.** In  
978 contrast to the case of problems with  $\ell_1$ -type Wasserstein ambiguity sets, a problem  
979 with  $\ell_\infty$ -type Wasserstein ambiguity set (5.1) further reduces to a linear program  
980 (refer to Theorem 1 of [47]). We use this approach to benchmark the performance of

TABLE 4  
*Computational Results for 2-DRLP Instances with Wasserstein- $\infty$  Ambiguity Set*

Problem	$N$	DRSD Algorithm		Reformulation Approach [47]	
		ObjEst	Time	ObjEst	Time
PGP	100	448.94 ( $\pm 3.64$ )	0.05 ( $\pm 0.00$ )	447.10 ( $\pm 3.26$ )	0.00 ( $\pm 0.00$ )
	250	455.04 ( $\pm 2.49$ )	0.24 ( $\pm 0.04$ )	450.79 ( $\pm 2.03$ )	0.06 ( $\pm 0.00$ )
	500	458.40 ( $\pm 2.94$ )	1.58 ( $\pm 0.16$ )	451.31 ( $\pm 1.23$ )	0.21 ( $\pm 0.01$ )
CEP	100	338291.53 ( $\pm 14434.11$ )	0.25 ( $\pm 0.01$ )	338307.14 ( $\pm 14431.53$ )	0.00 ( $\pm 0.00$ )
	250	355061.07 ( $\pm 11823.37$ )	3.13 ( $\pm 0.08$ )	355066.83 ( $\pm 11823.82$ )	0.06 ( $\pm 0.01$ )
	500	356763.93 ( $\pm 6917.77$ )	12.65 ( $\pm 0.28$ )	356769.76 ( $\pm 6918.07$ )	0.22 ( $\pm 0.02$ )
RETAIL	100	156.40 ( $\pm 4.24$ )	0.44 ( $\pm 0.02$ )	153.49 ( $\pm 3.89$ )	0.18 ( $\pm 0.01$ )
	250	155.12 ( $\pm 3.43$ )	7.89 ( $\pm 0.16$ )	153.86 ( $\pm 3.40$ )	0.59 ( $\pm 0.02$ )
	500	155.12 ( $\pm 2.39$ )	60.39 ( $\pm 0.91$ )	154.42 ( $\pm 2.38$ )	1.36 ( $\pm 0.03$ )
STORM	100	15504413.45 ( $\pm 11396.84$ )	0.56 ( $\pm 0.01$ )	15502082.05 ( $\pm 11445.84$ )	20.25 ( $\pm 0.27$ )
	250	15508768.42 ( $\pm 7477.52$ )	8.65 ( $\pm 0.29$ )	15504919.06 ( $\pm 7571.13$ )	81.22 ( $\pm 1.30$ )
	500	15507936.99 ( $\pm 5058.71$ )	63.10 ( $\pm 0.87$ )	15503343.24 ( $\pm 5140.36$ )	220.65 ( $\pm 2.57$ )

981 DRSD for  $\ell_\infty$ -type Wasserstein ambiguity sets. As in the previous set of experiments,  
982 we use the empirical distribution with  $N$  observations as reference distribution for the  
983 ambiguity sets. We generate the observations using an external sampling approach  
984 and set up the linear programming reformulation. We solve this reformulation using  
985 an off-the-shelf solver (CPLEX 12.10). We summarize the results for  $N = 100, 250,$   
986 and 500 in Table 4.

987 For all the problems, the estimates of the objective function obtained from DRSD  
988 and the reformulation linear program are comparable. The results show that for  
989 instances test problems PGP, CEP, and RETAIL, the linear programming reformulation  
990 outperforms the DRSD algorithm. However, for larger problem STORM, the advantages  
991 of sequential sampling become prevalent resulting in a nearly 3.5 times decrease in  
992 the overall computational time for  $N = 500$ , for instance.

993 *Remark 5.1.* Overall, the computational experiments with all three ambiguity  
994 sets illustrate the advantages of the sequential sampling approach of DRSD to tackle  
995 large-scale 2-DRLP problems. Before we end this section, we note that the external  
996 sampling-based benchmark instances are set up and solved for a given  $N$ . Since we  
997 are dealing with sampling-based approximations, identifying a suitable  $N$  a priori is  
998 not a trivial task. A procedure to tackle this task involves solving several instances  
999 with progressively increasing sample sizes (see for e.g., [6] for risk-neutral SP). The  
1000 overall computational cost of identifying a high-quality solution is the cumulative cost  
1001 associated with individual instances. The DRSD method, and the sequential sampling  
1002 idea in general, mitigates the need for such an iterative process.

1003 **6. Conclusions and Future Work.** We presented a new decomposition ap-  
1004 proach for solving two-stage distributionally robust linear programs (2-DRLPs) with  
1005 a general ambiguity set defined using probability distributions with continuous or  
1006 discrete sample space. Since this approach extended the stochastic decomposition ap-  
1007 proach of Higle and Sen [24] for 2-DRLPs with a singleton ambiguity set, we referred  
1008 to it as Distributionally Robust Stochastic Decomposition (DRSD) method. The  
1009 DRSD method is a sequential sampling-based approach that allows sampling within  
1010 the optimization step where we solved second-stage subproblem(s) associated with

only the current observation in each iteration. While the design of DRSD accommodates general ambiguity sets, we provided its asymptotic convergence analysis for a family of ambiguity sets that includes the well-known moment-based and Wasserstein metric-based ambiguity sets. Furthermore, we performed computational experiments to evaluate the efficiency and effectiveness of solving distributionally robust variants of four well-known stochastic programming test problems that have supports of size ranging from 216 to  $10^8$ <sup>81</sup>. Based on our results, we observed that the objective function estimates obtained using the DRSD and the external sampling-based approaches are statistically comparable. These DRSD estimates are obtained while providing computational improvements on most problem instances. Such a computational edge will enable the application of DRO to critical applications that result in large-scale problem instances.

The preliminary computational experiments are encouraging. However, there are two components of the algorithm that require careful deliberation. Since DRSD is a randomized algorithm that simultaneously deals with the approximation of ambiguity sets and recourse function values, the deterministic stopping criteria are not applicable. Therefore, the development of reliable stopping criteria is a potential future research direction. Statistical approaches, similar to those developed in a series of papers for SD [26, 27, 41], could provide initial direction to address this issue. Another future research direction is to incorporate more efficient algorithms to solve the distribution separation problems. For example, instead of resolving distribution separation problem in every iteration, we can utilize a column generation procedure. Finally, we will explore a proximal point algorithm design to that will allow us to maintain a fixed-sized master problem.

**Appendix A. Proofs.** In this appendix, we provide the proofs for the propositions related to the asymptotic behavior of the approximate ambiguity sets defined in §2.1 and the recourse function approximation presented in §2.2.

*Proof.* (Proposition 2.1) For  $P = (p(\omega))_{\omega \in \Omega^{k-1}} \in \widehat{\mathfrak{P}}_{\text{mom}}^{k-1}$ , it is easy to verify that  $P' = (p'(\omega))_{\omega \in \Omega^k} = \Theta^k(P)$  satisfies the support constraint, viz.,  $\sum_{\omega \in \Omega^k} p'(\omega) = 1$ . Now consider for  $i = 1, \dots, q$ , we have

$$\begin{aligned} 1041 \quad \sum_{\omega \in \Omega^k} p'(\omega) \psi_i(\omega) &= \sum_{\omega \in \Omega^{k-1}, \omega \neq \omega^k} p'(\omega) \psi_i(\omega) + p'(\omega^k) \psi_i(\omega^k) \\ 1042 \quad &= \theta^k \sum_{\omega \in \Omega^{k-1}, \omega \neq \omega^k} p(\omega) \psi_i(\omega) + \theta^k p(\omega^k) \psi_i(\omega^k) + (1 - \theta^k) \psi_i(\omega^k) \\ 1043 \quad &= \theta^k \sum_{\omega \in \Omega^{k-1}} p(\omega) \psi_i(\omega) + (1 - \theta^k) \psi_i(\omega^k) = \hat{b}_i^{k-1} + (1 - \theta^k) \psi_i(\omega^k) = \hat{b}_i^k. \\ 1044 \end{aligned}$$

1045 This implies that  $\Theta^k(P) \in \widehat{\mathfrak{P}}_{\text{mom}}^k$ .

1046 Using Proposition 4 in [45], there exists a positive constant  $\chi$  such that

$$1047 \quad 0 \leq \mathbb{H}(\widehat{\mathfrak{P}}_{\text{mom}}^k, \mathfrak{P}_{\text{mom}}) \leq \chi \|\hat{\mathbf{b}}^k - \mathbf{b}\|.$$

1049 Here,  $\mathbf{b} = (b_i)_{i=1}^q$  and  $\hat{\mathbf{b}}^k = (\hat{b}_i^k)_{i=1}^q$ , and  $\|\cdot\|$  denotes the Euclidean norm. Since 1050 the approximate ambiguity sets are constructed using independent and identically 1051 distributed samples of  $\tilde{\omega}$ , using law of large numbers, we have  $\hat{b}_i^k \rightarrow b_i$  for all  $i = 1, \dots, q$ . This completes the proof.  $\square$

1053 *Proof.* (Proposition 2.3) Consider approximate ambiguity sets  $\widehat{\mathfrak{P}}_{\text{w}}^{k-1}$  and  $\widehat{\mathfrak{P}}_{\text{w}}^k$  of 1054 the form given in (2.8b). Let  $P = (p(\omega))_{\omega \in \Omega^{k-1}} \in \widehat{\mathfrak{P}}_{\text{w}}^{k-1}$ , and let the reconstructed

1055 probability distribution be denoted by  $P'$ . We can easily check that  $P' = \Theta^k(P)$  is  
 1056 indeed a probability distribution. With  $P' = (p'(\omega))_{\omega \in \Omega^k}$  fixed, it suffices now to  
 1057 show that the polyhedron

$$1058 \quad (A.1) \quad \mathcal{E}(P', \hat{P}^k) = \left\{ \eta' \in \mathbb{R}^{\Omega^k \times \Omega^k} \left| \begin{array}{ll} \sum_{\omega' \in \Omega^k} \eta'(\omega, \omega') = p'(\omega) & \forall \omega \in \Omega^k, \\ \sum_{\omega \in \Omega^k} \eta'(\omega, \omega') = \hat{p}^k(\omega') & \forall \omega' \in \Omega^k, \\ \sum_{(\omega, \omega') \in \Omega^k \times \Omega^k} \|\omega - \omega'\| \eta'(\omega, \omega') \leq \epsilon \end{array} \right. \right\}.$$

1060 is non-empty. Since  $P \in \widehat{\mathfrak{P}}_w^{k-1}$ , there exist  $\eta(\omega, \omega')$  for all  $(\omega, \omega') \in \Omega^{k-1} \times \Omega^{k-1}$  such  
 1061 that the constraints in the description of the approximate ambiguity set in (2.8b) are  
 1062 satisfied. We show that  $\mathcal{E}$  is non-empty by analyzing two possibilities,

1063 1. We encounter a previously seen observation, i.e.,  $\omega^k \in \Omega^{k-1}$  and  $\Omega^k = \Omega^{k-1}$ .  
 1064 Let  $\eta'(\omega, \omega') = \theta^k \eta(\omega, \omega')$  for  $\omega, \omega' \in \Omega^{k-1}$  and  $\omega \neq \omega' \neq \omega^k$ ; and  $\eta'(\omega^k, \omega^k) =$   
 1065  $\theta^k \eta(\omega^k, \omega^k) + (1 - \theta^k)$ . We verify the feasibility of this choice by checking the three  
 1066 sets of constraints in (A.1). For all  $\omega \in \Omega^k$

$$\begin{aligned} 1067 \quad & \sum_{\omega' \in \Omega^k} \eta'(\omega, \omega') = \sum_{\omega' \in \Omega^k \setminus \{\omega^k\}} \eta'(\omega, \omega') + \eta'(\omega, \omega^k) \\ 1068 \quad & = \sum_{\omega' \in \Omega^{k-1} \setminus \{\omega^k\}} \theta^k \eta(\omega, \omega') + \theta^k \eta(\omega, \omega^k) + \mathbf{1}_{\omega=\omega^k}(1 - \theta^k) \\ 1069 \quad & = \theta^k \left( \sum_{\omega' \in \Omega^{k-1}} \eta(\omega, \omega') \right) + \mathbf{1}_{\omega=\omega^k}(1 - \theta^k) = \theta^k p(\omega) + \mathbf{1}_{\omega=\omega^k}(1 - \theta^k) = p'(\omega). \\ 1070 \end{aligned}$$

1071 For all  $\omega' \in \Omega^k$ , we have

$$\begin{aligned} 1072 \quad & \sum_{\omega \in \Omega^k} \eta'(\omega, \omega') = \sum_{\omega \in \Omega^k \setminus \{\omega^k\}} \eta'(\omega, \omega') + \eta'(\omega^k, \omega') \\ 1073 \quad & = \sum_{\omega \in \Omega^{k-1} \setminus \{\omega^k\}} \theta^k \eta(\omega, \omega') + \theta^k \eta(\omega^k, \omega') + \mathbf{1}_{\omega'=\omega^k}(1 - \theta^k) \\ 1074 \quad & = \theta^k \sum_{\omega \in \Omega^{k-1}} \eta'(\omega, \omega') + \mathbf{1}_{\omega'=\omega^k}(1 - \theta^k) = \theta^k \hat{p}^{k-1}(\omega') + \mathbf{1}_{\omega'=\omega^k}(1 - \theta^k) = \hat{p}^k(\omega'). \\ 1075 \end{aligned}$$

1076 And finally,

$$\begin{aligned} 1077 \quad & \sum_{(\omega, \omega') \in \Omega^k \times \Omega^k} \|\omega - \omega'\| \eta'(\omega, \omega') \\ 1078 \quad & = \sum_{\substack{(\omega, \omega') \in \Omega^{k-1} \times \Omega^{k-1} \\ \omega \neq \omega' \neq \omega^k}} \theta^k \|\omega - \omega'\| \eta(\omega, \omega') + \|\omega^k - \omega^k\| \eta'(\omega^k, \omega^k) \\ 1079 \quad & = \theta^k \left( \sum_{(\omega, \omega') \in \Omega^{k-1} \times \Omega^{k-1}} \|\omega - \omega'\| \eta(\omega, \omega') \right) \leq \theta^k \epsilon \leq \epsilon. \\ 1080 \end{aligned}$$

1081 Since all the three constraints are satisfied, the chosen values for  $\eta$  is an element  
 1082 of the polyhedron  $\mathcal{E}$ , and therefore,  $\mathcal{E} \neq \emptyset$ .

1083 2. We encounter a new observation, i.e.,  $\omega^k \notin \Omega^{k-1}$ . Let  $\eta'(\omega, \omega') = \theta^k \eta(\omega, \omega')$  for  
 1084  $\omega, \omega' \in \Omega^{k-1}$ ,  $\eta'(\omega^k, \omega') = 0$  for  $\omega' \in \Omega^{k-1}$ ,  $\eta'(\omega, \omega^k) = 0$  for  $\omega \in \Omega^{k-1}$ , and

1085  $\eta'(\omega^k, \omega^k) = (1 - \theta^k)$ . Let us again verify the three sets of constraints defining  
 1086 (A.1) with this choice for  $\eta'$ .

$$\begin{aligned} 1087 \quad \sum_{\omega' \in \Omega^k} \eta'(\omega, \omega') &= \sum_{\omega' \in \Omega^k \setminus \{\omega^k\}} \eta'(\omega, \omega') + \eta'(\omega, \omega^k) \\ 1088 \quad &= \sum_{\omega' \in \Omega^{k-1}} \theta^k \eta(\omega, \omega') + \mathbf{1}_{\omega=\omega^k} (1 - \theta^k) = \theta^k p(\omega) + \mathbf{1}_{\omega=\omega^k} (1 - \theta^k) = p'(\omega); \\ 1089 \quad \sum_{\omega \in \Omega^k} \eta'(\omega, \omega') &= \sum_{\omega \in \Omega^k \setminus \{\omega^k\}} \eta'(\omega, \omega') + \eta'(\omega^k, \omega') \\ 1090 \quad &= \sum_{\omega' \in \Omega^{k-1}} \theta^k \eta(\omega, \omega') + \mathbf{1}_{\omega'=\omega^k} (1 - \theta^k) = \theta^k \hat{p}^{k-1} + \mathbf{1}_{\omega'=\omega^k} (1 - \theta^k) = \hat{p}^k(\omega'); \\ 1091 \end{aligned}$$

1092 and finally,

$$\begin{aligned} 1093 \quad &\sum_{(\omega, \omega') \in \Omega^k \times \Omega^k} \|\omega - \omega'\| \eta'(\omega, \omega') \\ 1094 \quad &= \sum_{(\omega, \omega') \in \Omega^{k-1} \times \Omega^{k-1}} \theta^k \|\omega - \omega'\| \eta(\omega, \omega') + \|\omega^k - \omega^k\| \eta'(\omega^k, \omega^k) \\ 1095 \quad &+ \sum_{\omega \in \Omega^k} \|\omega - \omega^k\| \eta'(\omega, \omega^k) + \sum_{\omega' \in \Omega^k} \|\omega^k - \omega'\| \eta'(\omega^k, \omega') \leq \theta^k \epsilon \leq \epsilon. \\ 1096 \end{aligned}$$

1097 Therefore, the value of  $\eta'$  variables satisfies the constraints and  $\mathcal{E} \neq \emptyset$ . This implies  
 1098 that  $\Theta^k(P) \in \widehat{\mathfrak{P}}_w^k$ .

1099 Next, let us consider a distribution  $Q \in \widehat{\mathfrak{P}}_w^k$ . Then,

$$1100 \quad d_w(Q, P^*) \leq d_w(Q, \widehat{P}^k) + d_w(\widehat{P}^k, P^*) \leq \epsilon + d_w(\widehat{P}^k, P^*).$$

1102 The above inequality is a consequence of the triangle inequality of Wasserstein dis-  
 1103 tance. Since  $Q \in \widehat{\mathfrak{P}}_w^k$ , we have  $d_w(Q, \widehat{P}^k) \leq \epsilon$ . Under compactness assumption for  $\Omega$ ,  
 1104 we have  $\mathbb{E}_{P^*} [\exp(\|\tilde{\omega}\|^\alpha)] < \infty$ . Therefore, for  $d > 2$ , Theorem 2 in [18] guarantees

$$1105 \quad \text{Prob}[d_w(\widehat{P}^k, P^*) \leq \delta] \leq \begin{cases} C \exp(-ck\delta^d) & \text{if } \delta > 1 \\ C \exp(-ck\delta^a) & \text{if } \delta \leq 1 \end{cases}$$

1107 for all  $k \geq 1$ . This implies that the  $\lim_{k \rightarrow \infty} d_w(\widehat{P}^k, P^*) = 0$ , almost surely. Con-  
 1108 sequently, we obtain that  $d_w(Q, P^*) \leq \epsilon$  (or equivalently  $Q \in \mathfrak{P}_w$ ) as  $k \rightarrow \infty$ , almost  
 1109 surely. This completes the proof.  $\square$

1110 *Proof.* (Proposition 2.5) Recall that  $\mathcal{X} \times \Omega$  is a compact set because of Assump-  
 1111 tions (A1) and (A4), and  $\{Q^k\}$  is a sequence of continuous (piecewise linear and  
 1112 convex) functions. Further, the construction of the set of dual vertices satisfies  
 1113  $\Pi^0 = \{\mathbf{0}\} \subseteq \dots \subseteq \Pi^k \subseteq \Pi^{k+1} \subseteq \dots \subseteq \mathcal{D}$  which ensures that  $0 \leq Q^k(x, \omega) \leq$   
 1114  $Q^{k+1}(x, \omega) \leq Q(x, \omega)$  for all  $(x, \omega) \in (\mathcal{X}, \Omega)$  and  $k \geq 1$ . Since  $\{Q^k\}$  increases  
 1115 monotonically and is bounded by a finite function  $Q$  (due to (A2)), this sequence  
 1116 pointwise converges to some function  $\xi(x, \omega) \leq Q(x, \omega)$ . Once again due to (A2),  
 1117 we know that the set of dual vertices  $\mathcal{D}$  is finite and since  $\Pi^k \subseteq \Pi^{k+1} \subseteq \mathcal{D}$ , the set  
 1118  $\lim_{k \rightarrow \infty} \Pi^k := \overline{\Pi} (\subseteq \mathcal{D})$  is also a finite set. Clearly,

$$1119 \quad \xi(x, \omega) = \lim_{k \rightarrow \infty} Q^k(x, \omega) = \max \{\pi^\top [r(\omega) - T(\omega)x] \mid \pi \in \overline{\Pi}\}$$

1121 is the optimal value of a LP. Note that the right-hand side is a pointwise maximum of  
 1122 affine function and hence, is a continuous function. The compactness of  $\mathcal{X} \times \Omega$ , and  
 1123 continuity, monotonicity and pointwise convergence of  $\{Q^k\}$  to  $\xi$  guarantees that the  
 1124 sequence uniformly converges to  $\xi$  (implied by a slight modification of Theorem 7.13  
 1125 in [39]).  $\square$

## 1126 REFERENCES

1127 [1] M. Bansal, K. Huang, and S. Mehrotra. Decomposition Algorithms for Two-Stage Distributionally  
 1128 Robust Mixed Binary Programs. *SIAM Journal on Optimization*, pages 2360–2383,  
 1129 2018.

1130 [2] M. Bansal and S. Mehrotra. On solving two-stage distributionally robust disjunctive programs  
 1131 with a general ambiguity set. *European Journal of Operational Research*, 279(2):296–307,  
 1132 2019.

1133 [3] M. Bansal and Y. Zhang. Scenario-based cuts for structured two-stage stochastic and distri-  
 1134 butionally robust p-order conic mixed integer programs. *Journal of Global Optimization*,  
 1135 81:391–433, 2021.

1136 [4] G. Bayraksan and D. K. Love. Data-Driven Stochastic Programming Using Phi-Divergences.  
 1137 In *The Operations Research Revolution*, INFORMS TutORials in Operations Research,  
 1138 pages 1–19. INFORMS, 2015.

1139 [5] G. Bayraksan and D.P. Morton. Assessing solution quality in stochastic programs. *Mathematical  
 1140 Programming*, 108(2):495–514, 2006.

1141 [6] G. Bayraksan and D.P. Morton. A sequential sampling procedure for stochastic programming.  
 1142 *Operations Research*, 59(4):898–913, 2011.

1143 [7] A. Ben-Tal, D. den Hertog, A. De Waegenaere, B. Melenberg, and G. Rennen. Robust Solu-  
 1144 tions of Optimization Problems Affected by Uncertain Probabilities. *Management Science*,  
 1145 59(2):341–357, 2012.

1146 [8] D. Bertsimas, X. V. Doan, K. Natarajan, and C. Teo. Models for minimax stochastic linear  
 1147 optimization problems with risk aversion. *Mathematics of Operations Research*, 35(3):580–  
 1148 602, 2010.

1149 [9] D. Bertsimas and I. Popescu. Optimal inequalities in probability theory: A convex optimization  
 1150 approach. *SIAM Journal on Optimization*, 15(3):780–804, 2005.

1151 [10] J. R. Birge and F. Louveaux. *Introduction to Stochastic Programming*. Springer Series in  
 1152 Operations Research and Financial Engineering. Springer, 2011.

1153 [11] M. Breton and S. El Hachem. Algorithms for the solution of stochastic dynamic minimax  
 1154 problems. *Computational Optimization and Applications*, 4(4):317–345, 1995.

1155 [12] G. B. Dantzig. Linear programming under uncertainty. *Management Science*, 1(3-4):197–206,  
 1156 1955.

1157 [13] G. B. Dantzig and P. Wolfe. Decomposition principle for linear programs. *Operations Research*,  
 1158 8(1):101–111, 1960.

1159 [14] E. Delage and Y. Ye. Distributionally robust optimization under moment uncertainty with  
 1160 application to data-driven problems. *Operations Research*, 58:595–612, 2010.

1161 [15] J. Dupacová. The minimax approach to stochastic programming and an illustrative application.  
 1162 *Stochastics*, 20:73–88, 1987.

1163 [16] D. Dueque, S. Mehrotra, and D. Morton. Distributionally robust two-stage stochastic pro-  
 1164 gramming. Available at [http://www.optimization-online.org/DB\\_FILE/2020/09/8042.pdf](http://www.optimization-online.org/DB_FILE/2020/09/8042.pdf), 2020.

1165 [17] E. Erdogan and G. Iyengar. Ambiguous chance constrained problems and robust optimization.  
 1166 *Mathematical Programming*, 107(1-2):37–61, 2006.

1167 [18] N. Fournier and A. Guillin. On the rate of convergence in Wasserstein distance of the empirical  
 1168 measure. *Probability Theory and Related Fields*, 162(3-4):707–738, 2015.

1169 [19] Carlos Andrés Gamboa, Davi Michel Valladão, Alexandre Street, and Tito Homem-de Mello.  
 1170 Decomposition methods for wasserstein-based data-driven distributionally robust prob-  
 1171 lems. *Operations Research Letters*, 49(5):696–702, 2021.

1172 [20] H. Gangammanavar, Y. Liu, and S. Sen. Stochastic decomposition for two-stage stochastic  
 1173 linear programs with random cost coefficients. *INFORMS Journal on Computing*, 33(1):51–  
 1174 71, 2021.

1175 [21] H. Gangammanavar and S. Sen. Stochastic dynamic linear programming: A sequential sampling  
 1176 algorithm for multistage stochastic linear programming. *SIAM Journal on Optimization*,  
 1177 31(3):2111–2140, 2021.

[22] G. A. Hanusanto and D. Kuhn. Conic programming reformulations of two-stage distributionally robust linear programs over Wasserstein balls. *Operations Research*, 66(3):849–869, 2018.

[23] Y. T. Herer, M. Tzur, and E. Ycesan. The multilocation transshipment problem. *IIE Transactions*, 38(3):185–200, 2006.

[24] J. L. Higle and S. Sen. Stochastic decomposition: An algorithm for two-stage linear programs with recourse. *Mathematics of Operations Research*, 16(3):650–669, 1991.

[25] J. L. Higle and S. Sen. Finite master programs in regularized stochastic decomposition. *Mathematical Programming*, 67(1-3):143–168, 1994.

[26] J. L. Higle and S. Sen. *Stochastic Decomposition: A Statistical Method for Large Scale Stochastic Linear Programming*. Kluwer Academic Publishers, Boston, MA., 1996.

[27] J. L. Higle and S. Sen. Statistical approximations for stochastic linear programming problems. *Annals of Operations Research*, 85(0):173–193, 1999.

[28] R. Huang, S. Qu, Z. Gong, M. Goh, and Y. Ji. Data-driven two-stage distributionally robust optimization with risk aversion. *Applied Soft Computing*, 87:105978, 2020.

[29] R. Jiang and Y. Guan. Risk-averse two-stage stochastic program with distributional ambiguity. *Operations Research*, 66(5):1390–1405, 2018.

[30] B. Li, X. Qian, J. Sun, K. L. Teo, and C. Yu. A model of distributionally robust two-stage stochastic convex programming with linear recourse. *Applied Mathematical Modelling*, 58:86–97, 2018.

[31] S. Mehrotra and H. Zhang. Models and algorithms for distributionally robust least squares problems. *Mathematical Programming*, 148(1–2):123–141, 2014.

[32] P. Mohajerin Esfahani and D. Kuhn. Data-driven distributionally robust optimization using the Wasserstein metric: performance guarantees and tractable reformulations. *Mathematical Programming*, 171(1):115–166, 2018.

[33] J. M. Mulvey and A. Ruszczyński. A New Scenario Decomposition Method for Large-Scale Stochastic Optimization. *Operations Research*, 43(3):477–490, 1995. Publisher: INFORMS.

[34] H. Rahimian and S. Mehrotra. Distributionally robust optimization: A review. *arXiv preprint arXiv:1908.05659*, 2019.

[35] M. Riis and K. A. Andersen. Applying the minimax criterion in stochastic recourse programs. *European Journal of Operational Research*, 165(3):569–584, 2005.

[36] R. T. Rockafellar and R. J. B. Wets. Scenarios and policy aggregation in optimization under uncertainty. *Math. Oper. Res.*, 16(1):119–147, 1991.

[37] W. Römisch. Stability of stochastic programming problems. *Handbooks in operations research and management science*, 10:483–554, 2003.

[38] J. O. Royset and R. Szczytman. Optimal budget allocation for sample average approximation. *Operations Research*, 61(3):762–776, 2013.

[39] W. Rudin. *Principles of mathematical analysis*. McGraw-Hill Book Co., New York, third edition, 1976. International Series in Pure and Applied Mathematics.

[40] H. Scarf. A min-max solution of an inventory problem. In *Studies in the Mathematical Theory of Inventory and Production*, chapter 12, pages 201–209. RAND Corporation, Santa Monica CA, 1958.

[41] S. Sen and Y. Liu. Mitigating uncertainty via compromise decisions in two-stage stochastic linear programming: Variance reduction. *Operations Research*, 2016.

[42] A. Shapiro and S. Ahmed. On a class of minimax stochastic programs. *SIAM Journal on Optimization*, 14(4):1237–1249, 2004.

[43] A. Shapiro, D. Dentcheva, and A. Ruszczyński. *Lectures on Stochastic Programming: Modeling and Theory, Second Edition*. Society for Industrial and Applied Mathematics, Philadelphia, PA, USA, 2014.

[44] A. Shapiro and A. Kleywegt. Minimax analysis of stochastic problems. *Optimization Methods and Software*, 17(3):523–542, 2002.

[45] H. Sun and H. Xu. Convergence analysis for distributionally robust optimization and equilibrium problems. *Mathematics of Operations Research*, 41(2):377–401, 2016.

[46] R. M. Van Slyke and R. J. B. Wets. L-shaped linear programs with applications to optimal control and stochastic programming. *SIAM Journal on Applied Mathematics*, 17(4):638–663, 1969.

[47] W. Xie. Tractable reformulations of distributionally robust two-stage stochastic programs with  $\infty$ -Wasserstein distance. *arXiv preprint arXiv:1908.08454*, 2019.

[48] C. Zhao and Y. Guan. Data-driven risk-averse two-stage stochastic program with  $\zeta$ -structure probability metrics. Available at [http://www.optimization-online.org/DB\\_FILE/2015/07/5014](http://www.optimization-online.org/DB_FILE/2015/07/5014), 2015.

[49] C. Zhao and Y. Guan. Data-driven risk-averse stochastic optimization with wasserstein metric.

