

35 $\mathcal{F} \times \mathcal{F}$. The set Ξ is partitioned as: $\Xi = \cup_{i=1}^{m+1} \Xi^i$, where $\Xi^i = \{\xi^i\}$ for $i \in [m]$ and
 36 $\Xi^{m+1} = \Xi \setminus (\cup_{i \in [m]} \Xi^i)$. We call (WRO) a Wasserstein-robust optimization problem
 37 from hereafter.

38 Wasserstein metric has been used to study the convergence of an empirical distri-
 39 bution from the i.i.d. samples to the true distribution. Specifically, Barrio et al. [1]
 40 showed that under the Wasserstein metric the empirical distribution converges to the
 41 true distribution almost surely as the number of samples go to infinity. Fournier and
 42 Guillin [22] have further shown that if the true distribution P_{true} has a light tail, then
 43 the probability $\Pr\{W(P_0, P_{\text{true}}) \geq r_0\}$ can be bounded from above by a function that
 44 decays exponentially with the sample size and r_0 .

45 The (WRO) modeling framework allows to protect against ambiguity in the dis-
 46 tribution when arriving at a decision. It also allows to perform sensitivity analysis
 47 with respect to the empirical distribution. This is useful particularly when the sample
 48 size m is small.

49 **Contributions and organization of this paper.** This paper makes the fol-
 50 lowing contributions:

- 51 1. It is common in robust optimization to dualize the inner problem to develop a
 52 reformulation of the original model [2, 4, 6, 7]. The definition of Wasserstein
 53 metric in (2) uses a semi-infinite number of equality constraints, therefore
 54 its direct use is not suitable for dualization of the inner problem. We prove
 55 in Section 3 that the inner problem in (WRO) is equivalent to a conic linear
 56 optimization problem. We show that this conic program can be dualized with
 57 no duality gap, thus obtaining a semi-infinite programming reformulation of
 58 (WRO). The constraint sets in this semi-infinite program decompose in the
 59 observed samples. The decomposition allows for an independent separation
 60 problem for each observed sample ξ^i for algorithms that use cuts to solve
 61 the reformulated semi-infinite program. The results in Section 3 only assume
 62 that $h(\theta, \xi)$ is a continuous and bounded function in θ , and Ξ is a compact set
 63 (See Assumption 3.1). These results, for example, are applicable for problems
 64 involving mixed-integer variables used in the definition of Ξ .
- 65 2. We adapt the exchange algorithm and the central cutting-surface algorithm
 66 from [29] and [43] in Sections 4.1 and 4.2 for the general and convex cases,
 67 respectively. We show finite convergence of the exchange method to a solution
 68 with a desirable accuracy. The linear rate of convergence for the cutting
 69 surface algorithm presented here exploits the structure of the semi-infinite
 70 program. Specifically, a global linear rate of convergence is proved.
- 71 3. In Section 5.2 we present results on the computational performance of the
 72 central cutting-surface algorithm for solving WRLR problems. We find that
 73 the number of oracle calls are typically 20 ~ 50, and the number of cuts added
 74 to the model are typically 3 ~ 10 times the number of training samples. The
 75 solution time is \lesssim 100 times that of solving the ordinary logistic regression
 76 model.
 77 In Section 5.3, we present performance results on the quality of model ob-
 78 tained by the Wasserstein-robust logistic regression (WRLR) approach and
 79 compare it with the performance of the ordinary logistic regression (LR). Our
 80 motivation is to study a setting in which the number of available samples is
 81 small, and robustness is used to understand and possibly improve the quality
 82 of the trained model. This is typically the case at an early stage of a study
 83 (e.g., in healthcare) when limited data is available due to data collection ex-

84 pens. We use m to be $\{50, 75, 100, 150\}$ in the numerical experiments to
 85 test the performance of the Wasserstein-robust logistic regression model (See
 86 Section 5). Eleven data sets from UCI Machine Learning Repository are used.
 87 We use area under the receiver operator characteristic curve (AUC) to evalu-
 88 ate the performance of the models [69]. We find that the Wasserstein-robust
 89 logistic regression WRLR model has a significantly better out of sample per-
 90 formance than logistic-regression model (with $\alpha = 0.05$) in 24 (55%) cases.
 91 The predictive performance of WRLR is worse in 7 (16%) cases ($\alpha = 0.05$),
 92 and for the remaining 13 (29%) cases the difference is not statistically signif-
 93 icant. The WRLR models also have smaller standard error when compared
 94 to logistic regression, suggesting that the model is more robust.

95 **2. Literature Review.** We now provide a literature review of prior work on
 96 distributionally robust optimization (DRO) and semi-infinite programming. These
 97 topics are relevant because the (WRO) problem is reformulated as a semi-infinite
 98 program, and a separation oracle is needed in the algorithms.

99 **2.1. Distributionally robust optimization.** Distributional robust optimiza-
 100 tion is a generalization of robust optimization (RO) [10, 24, 64], where an ambiguity
 101 set is used to model problem data distribution. The use of an ambiguity set in distri-
 102 butionally robust optimization overcomes the weakness of traditional robust optimiza-
 103 tion framework, as the RO model is often considered to be very conservative. Even
 104 when a deterministic model is a convex optimization problem, its DR-counterpart
 105 is NP-hard in most cases [2]. Therefore, literature on RO and DRO either makes
 106 assumptions on the function form of the objective and constraints with respect to
 107 uncertain parameters to ensure the convexity of the model [5, 11], or purposes convex
 108 reformulations to approximate original problems [24, 62].

109 Studies of DRO focus on ways of defining the ambiguity set, reformulations of
 110 these models into computationally tractable problems, probability guarantee of the
 111 constraint satisfaction by the true distribution, and applications. We briefly review
 112 these aspects and the literature below.

113 **2.1.1. The ambiguity set.** A common approach to describe ambiguity sets is to
 114 use moments of the distribution followed by the random parameters [9, 12, 18, 19, 43,
 115 52, 56, 57]. Bertsimas and Popescu [11] discuss properties of probability distributions
 116 satisfying such constraints. It is also possible to use a *statistical distance* as a way
 117 of measuring the difference of two probability distributions. The statistical distances
 118 used in the DRO models are Wasserstein metric [21, 44, 49, 50, 54, 63], ϕ -divergence,
 119 χ^2 -distance, Kullback-Leibler divergence [3, 13, 32, 41, 61, 67], and the Prokhorov
 120 metric [20].

121 **2.1.2. Reformulation of DRO models.** Shapiro and Kleywegt [57] show that
 122 under mild regularity conditions, the DRO problem with a deterministic set of con-
 123 straints is equivalent to a stochastic programming problem based on a probability
 124 distribution that is a linear combination of distributions in the ambiguity set of the
 125 original DRO problem. It means that solving a DRO problem in this case is equiv-
 126 alent to solving a stochastic programming problem. Goh and Sim [24] investigate
 127 a two-stage DRO model whose objective has a linear structure, but can be used to
 128 express piece-wise linear utility functions and CVaR constraints. They show that by
 129 restricting the recourse variables to be affine mappings of uncertain parameters, this
 130 two-stage DRO model can be reformulated as a minimax linear programming prob-
 131 lem. Delage and Ye [18] show that DRO problems whose ambiguity sets are defined

132 by the first and second moment inequalities are polynomial-time solvable under the
 133 assumption that the objective is convex in the decision variables and concave in un-
 134 certain parameters. They also provide semidefinite formulations for the data-driven
 135 problems. As an often used objective, the least-square loss function is convex in both
 136 decision variables and model parameters, hence violating the assumption in [18]. To
 137 overcome this obstacle, Mehrotra and Zhang [44] give conic and semidefinite program-
 138 ming reformulations of DR-least-squares problems with ambiguity sets defined using
 139 the first two moments, Wasserstein metric, and bounds on the probability density
 140 functions, respectively. Mehrotra and Papp [43] use the central cutting-surface algo-
 141 rithm they developed to solve DRO models where the ambiguity set is specified using
 142 arbitrarily many generalized moments. Wiesemann et al. [62] propose a framework
 143 for modeling and solving distributionally robust convex optimization problems, in
 144 which the ambiguity set is conically representable and constraint functions are piece-
 145 wise affine in both decision variables and random parameters. They show that the
 146 reformulated problem is polynomial-time solvable under a strict nesting condition of
 147 the confidence sets. Esfahan and Kuhn [21] show that using the conic duality theo-
 148 ry, the data-driven DRO problem with a Wasserstein-metric ambiguity set can be
 149 reduced to finitely many tractable convex optimization problems, if the loss function
 150 can be expressed as the point-wise maximum of finitely many concave functions in
 151 the uncertain parameter. However, this assumption on the loss function is violated
 152 by many statistical learning models, such as the logistic regression model considered
 153 in the computational section of this paper.

154 Many DRO problems involve robust chance constraints, which are often non-
 155 convex. Jiang and Guan [33] study DR-chance constraints defined by the
 156 ϕ -divergences. They show that these constraints are equivalent to ordinary chance
 157 constraints based on some nominal probability measure. Hanasusanto et al. [26] show
 158 that if the ambiguity set in the robust chance constraint is defined by moments and
 159 satisfies a nested condition, the worst-case probability is an optimal solution of a
 160 conic optimization problem. For a recent review on tractable reformulations of robust
 161 chance constraints, see [51].

162 **2.1.3. The probability bound.** Studies on the probability that the true dis-
 163 tribution is contained in the ambiguity set are related to the ambiguous chance con-
 164 strained programming [14, 26, 70]. Erdoğan and Iyengar [20] show that the sampling
 165 of robust constraints is a good approximation for the DR-chance-constraint problem
 166 with a high probability. Delage and Ye [18] use the size of the ellipsoid confidence
 167 region using the second moment to satisfy a given level of probability guarantee. Es-
 168 fahan and Kuhn [21] give a result on the out-of-sample performance guarantee of
 169 the solution to the data-driven DRO problem with a Wasserstein-metric ambiguity
 170 set. Calafiore and Ghaoui [14] show that chance constraints of linear inequalities
 171 with respect to a radial distribution (i.e., Gaussian distribution and uniform distribu-
 172 tion on ellipsoidal support) can be converted explicitly into convex second-order cone
 173 constraints. Additionally, they show that distributionally robust chance constraints
 174 of linear inequalities under a few important distribution families (distributions with
 175 known mean and covariance, radially symmetric non-increasing distributions, etc.)
 176 can be guaranteed by some deterministic convex constraints. Bertsimas et al. [10]
 177 propose a novel scheme of constructing uncertainty sets for data-driven robust op-
 178 timization problems using hypothesis tests. The resulting model is computationally
 179 tractable and has application insights regarding using statistical estimate for chance
 180 constraint violation. Ben-Tal et al. [3] show that the robust counterpart of linear

181 optimization problems with uncertainty set defined by ϕ -divergence are tractable for
 182 most choices of ϕ . Constructing confidence sets using ϕ -divergence is also studied in
 183 [41, 67].

184 **2.1.4. Applications of DRO Models.** Modeling and solutions of the distri-
 185 butionally robust counterparts of deterministic optimization models are investigated
 186 in various areas in operations research, including but not limit to: inventory manage-
 187 ment [68], scheduling and logistics [31, 37], and risk management [39, 45]. In statistical
 188 learning, Lee and Mehrotra [36] study distributionally robust linear support vector
 189 machine (DR-SVM) models with a Wasserstein-metric ambiguity set. They find that
 190 the (DR-SVM) model can be reformulated as a semi-infinite program, in which the
 191 master problems are convex quadratic programs and separation problems are linear
 192 programs. They also find that (DR-SVM) models have improved generalization ca-
 193 pabilities than ordinary (SVM) models. Shaezadeh-Abadeh et al. [54] investigate a
 194 distributionally robust logistic regression model using a Wasserstein-metric ambiguity
 195 set. For the model in [54], the uncertainty set of the attribute vector is assumed to
 196 be the entire \mathbb{R}^n , and the authors show that the semi-infinite constraints in the dual
 197 problem are equivalent to a single constraint obtained using the conjugate function.
 198 The assumption in [54] is not practical in settings such as healthcare, where certain
 199 physiological variables (e.g., heart rate, blood pressure) must necessarily be bounded.
 200 In the framework of this paper, the uncertainty set is assumed to be compact. Addi-
 201 tionally, feature variables can be integer-valued.

202 **2.2. Theory and numerical methods for semi-infinite programming.**
 203 The study of distributional robust optimization models considered in this paper ben-
 204 efit from the known literature on semi-infinite programming. Semi-infinite program-
 205 ming (SIP) problems are optimization problems with constraints induced by a con-
 206 tinuous parameter. The study of SIP is initialized by the work of Haar [25] and
 207 Charnes [15, 16, 17] focusing on linear-SIP problems. Later, the first and second
 208 order optimality conditions of general SIP were given in [27, 28, 30, 46, 47, 58]. For
 209 reviews of the theory and methods for SIP, see [29, 40, 53].

210 Numerical methods for solving convex SIP problems include primal methods [60],
 211 dual methods [29], penalty methods [38, 66], smooth approximation and projection
 212 methods [65], and cutting-plane methods [34]. Primal methods are based on searching
 213 for feasible descent directions, while dual methods are based on finding a solution of
 214 the system of KKT optimality conditions. In a penalty method, constraints are pe-
 215 nalized in the objective and the penalty term is an integral of the constraint function
 216 over the continuous parameter. In a smooth approximation and projection method,
 217 infinitely many functions are replaced by a integral entropy function as an approxi-
 218 mation. The SIP is solved by the smoothing projected gradient method. In a cutting-
 219 plane algorithm, a typical iteration involves solving a master problem with finitely
 220 many constraints and adding violated constraints obtained from solving a separation
 221 problem. Mehrotra and Papp [43] developed a central-cutting-surface algorithm with
 222 a linear rate of convergence for solving convex SIP problems. They demonstrate that
 223 adding cutting surfaces, as compared to cutting planes, can be computationally ef-
 224 fective for problems in high dimension and ensure greater stability in the algorithm's
 225 performance. The central-cutting surface method is related to the exchange method
 226 [29], however it uses a centrality parameter in the algorithm.

227 **3. Reformulation of the Wasserstein-robust Optimization Problem.** In
 228 this section we show that (WRO) is equivalent to a semi-infinite program. This semi-

229 infinite program decomposes in scenarios. We note that the definition of $\mathcal{W}(P, P_0) \leq$
 230 r_0 involves infinitely many equality constraints of the form $K(A \times \Omega) = P, \forall A \in \mathcal{F}$.
 231 Since this form is not suitable for dualization of the inner problem, Theorem 3.6 below
 232 reformulates the inner problem in (WRO) as a conic linear program with finitely many
 233 constraints. Proposition 3.3 provides an intermediate result, and shows the continuity
 234 of the objective function of the inner problem in (WRO) with respect to the probability
 235 measure. Lemma 3.5 is a technical result that generalizes similar convergent sequence
 236 existence results for the finite dimensional sets to a set defined by the Wasserstein
 237 metric. We make the following general assumption throughout this paper. Additional
 238 assumptions are made at appropriate places as results are developed. Note that results
 239 in this section make no assumption on the metric $d(\cdot, \cdot)$ used in defining (2).

240 ASSUMPTION 3.1. *We assume that $r_0 > 0, \forall \theta \in \Theta$ in (2). The feasible region*
 241 $\Theta \subseteq \mathbb{R}^n$ *and the domain Ξ are compact. The function $h(\cdot, \cdot)$ is bounded on $\Theta \times \Xi$,*
 242 *and for every $\theta \in \Theta$, there exists $C(\theta) > 0$ such that the function $h(\theta, \cdot)$ satisfies*
 243 $|h(\theta, s_1) - h(\theta, s_2)| \leq C(\theta)d(s_1, s_2), \forall s_1, s_2 \in \Xi$.

244 DEFINITION 3.2. *Let $(\mathcal{X}, d_{\mathcal{X}})$ be a metric space formed by defining a metric $d_{\mathcal{X}}$*
 245 *on the topological space \mathcal{X} . A function $f : \mathcal{X} \mapsto \mathbb{R}$ is continuous in $(\mathcal{X}, d_{\mathcal{X}})$ if*
 246 *for a given $x \in \mathcal{X}$, and $\varepsilon > 0, \exists \delta > 0$ such that $\forall x' \in \mathcal{X}, d_{\mathcal{X}}(x', x) < \delta$ we have*
 247 $|f(x') - f(x)| < \varepsilon$.

248 Lemma 3.5 below shows that the interior of the Wasserstein ball has a sequence
 249 of distributions that converge to a chosen point on the boundary of the ball. We need
 250 the following two results, proved in Appendix A, in the proof of
 251 Lemma 3.5.

252 PROPOSITION 3.3. *Let Assumption 3.1 hold, then the function $f(\theta, \cdot) :$*
 253 $\mathcal{M}(\Xi, \mathcal{F}) \mapsto \mathbb{R}$ *defined by $f(\theta, P) := \mathbb{E}_P[h(\theta, \xi)]$ is continuous in $(\mathcal{M}(\Xi, \mathcal{F}), \mathcal{W})$.*

254 PROPOSITION 3.4. *For any probability measure $P \in \mathcal{M}(\Omega, \mathcal{F})$, there exists a $K \in$*
 255 $\mathcal{M}(\Xi \times \Xi, \mathcal{F} \times \mathcal{F})$ *such that $K(A \times A) = P(A)$ and $K(\Xi \times A) = K(A \times \Xi) = P(A)$*
 256 *for any $A \in \mathcal{F}$. For such a joint probability measure K , we have*

$$257 \quad (3) \quad \int_{(s_1 \times s_2) \in \Xi \times \Xi} d(s_1, s_2) K(ds_1 \times ds_2) = 0.$$

258 LEMMA 3.5. *Let Assumption 3.1 hold, and \mathcal{P} be defined as in (1). Let $\mathcal{P}' :=$*
 259 $\{P \in \mathcal{M}(\Xi, \mathcal{F}) : \mathcal{W}(P, P_0) < r_0\}$ *be the interior of \mathcal{P} . Then for any $P \in \mathcal{P}$, there*
 260 *exists a sequence $\{P^n\}_{n=1}^{\infty} \subseteq \mathcal{P}'$ such that $\lim_{n \rightarrow \infty} \mathcal{W}(P^n, P) = 0$.*

261 *Proof.* For any given $P \in \mathcal{P}$ and for any $\varepsilon > 0$, by the definition of the Wasserstein
 262 metric, there exists $K^\varepsilon \in \mathcal{M}(\Xi \times \Xi, \mathcal{F} \times \mathcal{F})$ such that $K^\varepsilon(\Xi \times A) = P_0(A), K^\varepsilon(A \times \Xi) =$
 263 $P(A), \forall A \in \mathcal{F}$, and

$$264 \quad (4) \quad \int_{(s_1 \times s_2) \in \Xi \times \Xi} d(s_1, s_2) K^\varepsilon(ds_1 \times ds_2) \leq \mathcal{W}(P_0, P) + \varepsilon.$$

265 Define $K_0 \in \mathcal{M}(\Xi \times \Xi, \mathcal{F} \times \mathcal{F})$ such that $K_0(A \times B) = P_0(A \cap B), \forall A, B \in \mathcal{F}$. By
 266 Proposition 3.4, we have

$$267 \quad (5) \quad \int_{(s_1 \times s_2) \in \Xi \times \Xi} d(s_1, s_2) K_0(ds_1 \times ds_2) = 0.$$

268 Now define $\{P^n\}_{n=1}^{\infty} \subseteq \mathcal{M}(\Xi, \mathcal{F})$ as: $P^n := \lambda_n P + (1 - \lambda_n) P_0$ with $\lambda_n \in (0, 1)$ and
 269 $\lambda_n \rightarrow 1$. Define $K_n^\varepsilon := \lambda_n K^\varepsilon + (1 - \lambda_n) K_0$ as a probability measure in $\mathcal{M}(\Xi \times$

270 $\Xi, \mathcal{F} \times \mathcal{F}$). It is straightforward to verify that $K_n^\varepsilon(\Xi \times A) = P(A)$ and $K_n^\varepsilon(A \times$
 271 $\Xi) = P^n(A), \forall A \in \mathcal{F}$, using their definitions, which means the joint measure K_n^ε
 272 satisfies marginal conditions with respect to P and P^n . First, we need to verify that
 273 $\{P^n\}_{n=1}^\infty \subseteq \mathcal{P}'$. To see this, we have

$$\begin{aligned}
 274 \quad \mathcal{W}(P^n, P_0) &\leq \int_{(s_1 \times s_2) \in \Xi \times \Xi} d(s_1, s_2) K_n^\varepsilon(ds_1 \times ds_2) \\
 275 \quad &= \lambda_n \int_{\Xi \times \Xi} d(s_1, s_2) K^\varepsilon(ds_1 \times ds_2) + (1 - \lambda_n) \int_{\Xi \times \Xi} d(s_1, s_2) K_0(ds_1 \times ds_2) \\
 276 \quad (6) \quad &\leq \lambda_n [\mathcal{W}(P, P_0) + \varepsilon]. \quad \text{using (4-5)}
 \end{aligned}$$

278 Since ε can be chosen arbitrarily, we set

$$279 \quad (7) \quad \varepsilon = \min \left\{ 1, \frac{1}{2} \left(\frac{1}{\lambda_n} - 1 \right) \mathcal{W}(P, P_0) \right\}.$$

280 Substituting (7) into (6) yields $\mathcal{W}(P^n, P_0) \leq (1/2 + \lambda_n/2)\mathcal{W}(P, P_0) < r_0$, hence
 281 $\{P^n\}_{n=1}^\infty \subseteq \mathcal{P}'$.

282 It remains to verify that $\lim_{n \rightarrow \infty} \mathcal{W}(P^n, P) = 0$. To see this, define $K \in \mathcal{M}(\Xi \times$
 283 $\Xi, \mathcal{F} \times \mathcal{F})$ such that $K(A \times B) = P(A \cap B), \forall A, B \in \mathcal{F}$. By Proposition 3.4, we have

$$284 \quad (8) \quad \int_{(s_1 \times s_2) \in \Xi \times \Xi} d(s_1, s_2) K(ds_1 \times ds_2) = 0.$$

285 Let $\tilde{K}_n^\varepsilon := \lambda_n K + (1 - \lambda_n) K^\varepsilon$ be a joint probability measure in $\mathcal{M}(\Xi \times \Xi, \mathcal{F} \times \mathcal{F})$.
 286 Then we have $\tilde{K}^\varepsilon(\Xi \times A) = P(A)$ and $\tilde{K}^\varepsilon(A \times \Xi) = P^n(A), \forall A \in \mathcal{F}$, which means
 287 that \tilde{K}^ε satisfies marginal conditions with respect to P and P^n . It follows that

$$\begin{aligned}
 \mathcal{W}(P^n, P) &\leq \int_{(s_1 \times s_2) \in \Xi \times \Xi} d(s_1, s_2) \tilde{K}_n^\varepsilon(ds_1 \times ds_2) \\
 288 \quad (9) \quad &= \lambda_n \int_{\Xi \times \Xi} d(s_1, s_2) K(ds_1 \times ds_2) + (1 - \lambda_n) \int_{\Xi \times \Xi} d(s_1, s_2) K^\varepsilon(ds_1 \times ds_2) \\
 &\leq (1 - \lambda_n) [\mathcal{W}(P, P_0) + \varepsilon]. \quad \text{using (4),(8)}
 \end{aligned}$$

289 Since $\lambda_n \rightarrow 1$, we have $\mathcal{W}(P^n, P) \rightarrow 0$ as $n \rightarrow \infty$. \square

290 **THEOREM 3.6.** *Let Assumption 3.1 hold, and \mathcal{P} be defined as in (1). For a given*
 291 *θ , the inner problem $\sup_{P \in \mathcal{P}} \mathbb{E}_P[h(\theta, \xi)]$ has a finite optimal value, and it is equivalent*
 292 *to the following conic linear program (CLP):*

$$\begin{aligned}
 &\sup_{\mu} \int_{s \in \Xi} h(\theta, s) \mu(ds \times \Xi) \\
 &\text{s.t.} \quad \mu(\Xi \times \Xi^i) = 1/m, \quad i \in [m] \\
 &\quad \mu(\Xi \times \Xi^{m+1}) = 0, \\
 293 \quad (\text{CLP}) \quad &\sum_{i \in [m]} \int_{s \in \Xi} d(s, s^i) \mu(ds \times \Xi^i) \leq r_0, \\
 &\mu \succeq 0,
 \end{aligned}$$

294 where $\mu \succeq 0$ denotes that μ is a positive measure.

295 *Proof.* The inner problem in (WRO) can be written as:

$$296 \quad (10) \quad \begin{aligned} & \sup_{P \in \mathcal{M}(\Xi, \mathcal{F})} \int_{s \in \Xi} h(\theta, s) P(ds) \\ & \text{s.t.} \quad \mathcal{W}(P, P_0) \leq r_0. \end{aligned}$$

297 Let $\text{val}(\text{eqn}\#)$ denote the optimal value of a problem given by (eqn#). By Assump-
298 tion 3.1 Ξ is compact, and $h(\theta, \cdot)$ is bounded in Ξ . Hence, the objective $\mathbb{E}_P[h(\theta, \xi)]$ is
299 finite, and therefore $\text{val}(\text{10})$ is finite. We first show that $\text{val}(\text{10}) = \text{val}(\text{CLP})$. For this
300 purpose we consider the following auxiliary problem:

$$301 \quad (11) \quad \begin{aligned} & \sup_{P \in \mathcal{M}(\Xi, \mathcal{F})} \int_{s \in \Xi} h(\theta, s) P(ds) \\ & \text{s.t.} \quad \mathcal{W}(P, P_0) < r_0, \end{aligned}$$

302 whose feasible set is the interior of \mathcal{P} . We will prove that $\text{val}(\text{10}) \geq \text{val}(\text{CLP}) \geq$
303 $\text{val}(\text{11})$, and $\text{val}(\text{10}) = \text{val}(\text{11})$, and hence $\text{val}(\text{10}) = \text{val}(\text{CLP})$.

304 We now show that $\text{val}(\text{CLP}) \geq \text{val}(\text{11})$. Let \tilde{P} be a feasible solution of (11). Since
305 we have $\mathcal{W}(\tilde{P}, P_0) < r_0$, by the definition of the Wasserstein metric in (2), there exists
306 a $\tilde{K} \in \mathcal{S}(\tilde{P}, P_0)$ satisfying:

$$307 \quad \int_{\Xi \times \Xi} d(s_1, s_2) \tilde{K}(ds_1 \times ds_2) = \sum_{i \in [m]} \int_{s \in \Xi} d(s, s^i) \tilde{K}(ds \times \Xi^i) \leq r_0.$$

308 Therefore, \tilde{K} is a feasible solution of (CLP) with the objective value $\int_{s \in \Xi} h(\theta, s) \tilde{K}(ds \times$
309 $\Xi)$. Now observe that $\int_{s \in \Xi} h(\theta, s) \tilde{K}(ds \times \Xi) = \int_{s \in \Xi} h(\theta, s) \tilde{P}(ds)$, because $\tilde{K} \in$
310 $\mathcal{S}(\tilde{P}, P_0)$. Consequently, for any sequence $\{\tilde{P}_k\}_1^\infty$ such that $\int_{s \in \Xi} h(\theta, s) \tilde{P}_k(ds) \rightarrow$
311 $\text{val}(\text{11})$, there exist $\{\tilde{K}_k\}_1^\infty$ satisfying $\int_{s \in \Xi} h(\theta, s) \tilde{K}_k(ds \times \Xi) = \int_{s \in \Xi} h(\theta, s) \tilde{P}_k(ds)$,
312 hence $\int_{s \in \Xi} h(\theta, s) \tilde{K}_k(ds \times \Xi) \rightarrow \text{val}(\text{11})$. It follows that $\text{val}(\text{CLP}) \geq \text{val}(\text{11})$.

313 Next we show that $\text{val}(\text{10}) \geq \text{val}(\text{CLP})$. Suppose \hat{K} is a feasible solution of (CLP).
314 Let $\hat{P} \in \mathcal{M}(\Xi, \mathcal{F})$ be the marginal distribution of \hat{K} such that $\hat{P}(A) := \hat{K}(A \times$
315 $\Xi)$, $\forall A \in \mathcal{F}$. Due to the constraints of (CLP), \hat{P} satisfies $\mathcal{W}(\hat{P}, P_0) \leq r_0$; and hence
316 \hat{P} is a feasible solution of (10). Because $\hat{K} \in \mathcal{S}(\hat{P}, P_0)$, we have $\int_{s \in \Xi} h(\theta, s) \hat{P}(ds) =$
317 $\int_{s \in \Xi} h(\theta, s) \hat{K}(ds \times \Xi)$. Consequently, for any sequence $\{\hat{K}_k\}_1^\infty$ such that
318 $\int_{s \in \Xi} h(\theta, s) \hat{K}_k(ds \times \Xi) \rightarrow \text{val}(\text{CLP})$, there exist a sequence $\{\hat{P}_k\}_1^\infty$ satisfying
319 $\int_{s \in \Xi} h(\theta, s) \hat{P}_k(ds) = \int_{s \in \Xi} h(\theta, s) \hat{K}_k(ds \times \Xi)$. It follows that $\int_{s \in \Xi} h(\theta, s) \hat{P}_k(ds) \rightarrow$
320 $\text{val}(\text{CLP})$, and hence $\text{val}(\text{10}) \geq \text{val}(\text{CLP})$.

321 We now show that $\text{val}(\text{10}) = \text{val}(\text{11})$. Since $\text{val}(\text{10})$ is finite, there exists a sequence
322 of probability measures $\{P_m\}_{m=1}^\infty \subseteq \mathcal{P}$ such that

$$323 \quad (12) \quad \lim_{m \rightarrow \infty} \int_{s \in \Xi} h(\theta, s) P_m(ds) = \text{val}(\text{10}),$$

324 where $\mathcal{P} := \{P \in \mathcal{M}(\Xi, \mathcal{F}) : \mathcal{W}(P, P_0) \leq r_0\}$. Let $\mathcal{P}' := \{P \in \mathcal{M}(\Xi, \mathcal{F}) : \mathcal{W}(P, P_0) <$
325 $r_0\}$ be the feasible set of (11). By Lemma 3.5, for any P_m ($m \geq 1$), there exists a se-
326 quence $\{P_m^n\}_{n=1}^\infty \subseteq \mathcal{P}'$ such that $\lim_{n \rightarrow \infty} \mathcal{W}(P_m^n, P_m) = 0$. Let $f(\theta, P) := \mathbb{E}_P[h(\theta, \xi)]$,
327 then we have that $\text{val}(\text{10}) \geq \text{val}(\text{11}) \geq f(\theta, P_m^n)$, $\forall m, n \geq 1$. Moreover, we also have

328 that

$$\begin{aligned}
 & \lim_{m \rightarrow \infty} \lim_{n \rightarrow \infty} |f(\theta, P_m^n) - \text{val}(\mathbf{10})| \\
 329 \quad & \leq \lim_{m \rightarrow \infty} \lim_{n \rightarrow \infty} \left(|f(\theta, P_m^n) - f(\theta, P_m)| + |f(\theta, P_m) - \text{val}(\mathbf{10})| \right) \\
 & = 0, \quad [\text{using Proposition 3.3 and (12)}]
 \end{aligned}$$

330 which implies that $\text{val}(\mathbf{11}) = \text{val}(\mathbf{10})$. \square

331 Note that the infimum used in (CLP) is over the joint measure μ . Moreover, in (CLP)
 332 both the objective and constraints are in linear form of the measure μ , which is a
 333 ‘new decision variable’ in the inner problem. Theorem 3.7 provides a dual of (CLP)
 334 as a semi-infinite linear program in the linear space of signed measures using conic
 335 programming duality. This theorem shows that the dual problem has a formulation
 336 which decomposes in ξ^i , $i \in [m]$.

337 **THEOREM 3.7.** *Let Assumption 3.1 hold. The dual of (CLP) can be written as*
 338 *the following semi-infinite program:*

$$\begin{aligned}
 339 \quad (\text{CLP-D}) \quad & \min_{v \in \mathbb{R}^{m+1}} \quad \frac{1}{m} \sum_{i=1}^m v_i + r_0 \cdot v_{m+1} \\
 & \text{s.t.} \quad h(\theta, s) - v_i - v_{m+1} \cdot d(s, \xi^i) \leq 0, \quad s \in \Xi, i \in [m] \\
 & \quad v_1, \dots, v_m \in \mathbb{R}, \quad v_{m+1} \geq 0.
 \end{aligned}$$

340 Furthermore, strong duality holds, i.e., $\text{val}(\text{CLP}) = \text{val}(\text{CLP-D})$. Additionally,
 341 for $r_0 > 0$ the optimum solution of (CLP-D) can be bounded by the following polytope:

$$\begin{aligned}
 342 \quad (13) \quad \mathcal{H} := & \left\{ v \in \mathbb{R}^{m+1} : C_1 \leq v_i \leq (m+1)C_2 - mC_1, \text{ for } i \in [m], \right. \\
 & \left. 0 \leq v_{m+1} \leq (C_2 - C_1)/r_0 \right\},
 \end{aligned}$$

343 where C_1 and C_2 are lower and upper bounds of $h(\cdot, \cdot)$ on $\Theta \times \Xi$, respectively.

344 *Proof.* We note that (CLP) can be rewritten as:

$$\begin{aligned}
 345 \quad (14) \quad & \sup_{\mu} \int_{(s_1 \times s_2) \in \Xi \times \Xi} h(\theta, s_1) \mu(ds_1 \times ds_2) \\
 & \text{s.t.} \quad \int_{(s_1 \times s_2) \in \Xi \times \Xi} \mathbf{1}_{\Xi^i}(s_2) \mu(s_1 \times s_2) = 1/m, \quad i \in [m] \\
 & \int_{(s_1 \times s_2) \in \Xi \times \Xi} \mathbf{1}_{\Xi^{m+1}}(s_2) \mu(s_1 \times s_2) = 0, \\
 & \int_{(s_1 \times s_2) \in \Xi \times \Xi} \left[\sum_{i \in [m]} d(s_1, s_2) \cdot \mathbf{1}_{\Xi^i}(s_2) \right] \mu(ds_1 \times ds_2) \leq r_0, \\
 & \mu \succeq 0,
 \end{aligned}$$

346 where $\mathbf{1}_{\Xi^i}(s)$ is an indicator function. Note that the second constraint in (14) is
 347 defined on the set $\Xi^{m+1} = \Xi \setminus (\cup_{i \in [m]} \Xi^i)$. For a given θ , we define functions $\{\psi_j\}_{j=0}^{m+2}$
 348 as follows:

$$349 \quad (15) \quad \psi_j(s_1, s_2) = \begin{cases} h(\theta, s_1) & j = 0 \\ \mathbf{1}_{\Xi^j}(s_2) & j \in [m+1] \\ \sum_{i=1}^m d(s_1, s_2) \cdot \mathbf{1}_{\Xi^i}(s_2) & j = m+2. \end{cases}$$

350 Clearly $\{\psi_j\}_{j=0}^{m+2}$ are bounded $\mathcal{F} \times \mathcal{F}$ -measurable functions on $\Xi \times \Xi$. It follows
 351 that $\{\psi_j\}_{j=0}^{m+2}$ are μ -integrable for any $\mu \succeq 0$. We will first put (14) in a standard
 352 form of conic linear program. Let \mathcal{X} be the linear space of *finite signed measures*,
 353 and $\mathcal{X}^+ := \{\mu \in \mathcal{X} : \mu \succeq 0\}$ be the set of non-negative measures which is a
 354 convex cone in \mathcal{X} . Let \mathcal{X}' be the set of functions that are μ -integrable for all $\mu \in$
 355 \mathcal{X}^+ , i.e., $\mathcal{X}' := \{f : \Xi \times \Xi \rightarrow \mathbb{R} \mid f \in L^1(\Xi \times \Xi, \mu), \forall \mu \in \mathcal{X}^+\}$. Define $\langle \mu, f \rangle :=$
 356 $\int_{(s_1 \times s_2) \in \Xi \times \Xi} f(s_1, s_2) \mu(ds_1 \times ds_2)$, $\forall \mu \in \mathcal{X}, f \in \mathcal{X}'$. Define the linear operator $\mathcal{A} :$
 357 $\mathcal{X} \rightarrow \mathbb{R}^{m+2}$ as

$$358 \quad (16) \quad \mathcal{A}\mu := [\langle \mu, \psi_1 \rangle, \dots, \langle \mu, \psi_{m+2} \rangle]^T.$$

359 Define a vector b as $b = \underbrace{[1/m, \dots, 1/m, 0, r_0]}_m^T \in \mathbb{R}^{m+2}$, and a convex cone $\mathcal{K} :=$
 360 $\mathbf{0}^{m+1} \times (-\infty, 0]$ in \mathbb{R}^{m+2} . Using the above notations, (14) can be rewritten as:

$$361 \quad (\text{CLP1}) \quad \begin{array}{l} \sup_{\mu} \quad \langle \mu, \psi_0 \rangle \\ \text{s.t.} \quad \mathcal{A}\mu - b \in \mathcal{K} \\ \mu \in \mathcal{X}^+. \end{array}$$

362 Using Lemma A.1(a) in Appendix A from [55], we have the following dual of (CLP1):

$$363 \quad (17) \quad \begin{array}{l} \inf_w \quad -b^T \cdot w \\ \text{s.t.} \quad \mathcal{A}^* w + \psi_0 \in -(\mathcal{X}^+)^* \\ w \in \mathcal{K}^*, \end{array}$$

364 where $\mathcal{K}^* = \mathbb{R}^{m+1} \times (-\infty, 0]$ is the dual cone of \mathcal{K} , and $-(\mathcal{X}^+)^* = \{f \in \mathcal{X}' : \langle \mu, f \rangle \leq$
 365 $0, \forall \mu \succeq 0\}$ is the polar cone of \mathcal{X}^+ . The adjoint operator \mathcal{A}^* acts on \mathbb{R}^{m+2} as:
 366 $\mathcal{A}^* v = \sum_{i=1}^{m+2} v_i \psi_j$, $\forall v \in \mathbb{R}^{m+2}$. Therefore, the dual problem (17) can be expressed
 367 as follows:

$$368 \quad (18) \quad \begin{array}{l} \min_w \quad -\frac{1}{m} \sum_{i=1}^m w_i - r_0 \cdot w_{m+2} \\ \text{s.t.} \quad \langle \mu, \psi_0 \rangle + \sum_{i=1}^{m+2} w_i \langle \mu, \psi_i \rangle \leq 0, \quad \forall \mu \in \mathcal{X}^+, \\ w_{m+2} \leq 0. \end{array}$$

369 We have

$$370 \quad \sum_{i=1}^{m+2} w_i \langle \mu, \psi_i \rangle + \langle \mu, \psi_0 \rangle = \sum_{i=1}^m \int_{\Xi \times \Xi^i} [w_i + w_{m+2} \cdot d(s_1, \xi^i) + h(\theta, s_1)] \mu(ds_1 \times ds_2) \\ 371 \quad (19) \quad + \int_{\Xi \times \Xi^{m+1}} [w_{m+1} + h(\theta, s_1)] \mu(ds_1 \times ds_2). \\ 372$$

373 Since \mathcal{F} contains all singletons in Ξ , all integrands on the RHS of (19) must be non-
 374 positive to ensure the constraint $\sum_{i=1}^{m+2} w_i \langle \mu, \psi_i \rangle + \langle \mu, \psi_0 \rangle \leq 0$, $\forall \mu \in \mathcal{X}^+$. Otherwise,
 375 we concentrate the measure μ on the points $(\hat{s} \times \xi^i)$, $i \in [m]$, or the set (\hat{s}, Ξ^{m+1})
 376 at which one of the integrands is positive to achieve a contradiction. In other words,
 377 $\sum_{i=1}^{m+2} w_i \langle \mu, \psi_i \rangle + \langle \mu, \psi_0 \rangle \leq 0$, $\forall \mu \in \mathcal{X}^+$ are equivalent to the following constraints:

$$378 \quad (20) \quad w_i + w_{m+2} \cdot d(s, \xi^i) + h(\theta, s) \leq 0, \quad \forall s \in \Xi \quad i \in [m],$$

$$379 \quad (21) \quad w_{m+1} + h(\theta, s) \leq 0, \quad \forall s \in \Xi.$$

381 Since w_{m+1} does not appear in the objective function of (18), and constraint (21) can
 382 be satisfied by choosing w_{m+1} sufficiently small, we can remove w_{m+1} and constraint
 383 (21) from the dual problem. By rewriting variables w_i as $-v_i$ for $i \in [m]$ and w_{m+2}
 384 as $-v_{m+1}$, the dual problem (18) can be rewritten as (CLP-D).

385 We now show that strong duality holds. We only need to verify that $\text{val}(\text{CLP})$ is
 386 finite and show that $\text{Sol}(\text{CLP-D})$ is nonempty and bounded [55] (see Lemma A.1(b),
 387 in Appendix A). First, since the objective in (CLP) is an expectation of a bounded
 388 function, $\text{val}(\text{CLP})$ is finite. Second, we note that the feasible set of (CLP-D), denoted
 389 as $\text{Fsb}(\text{CLP-D})$, is a closed convex subset in \mathbb{R}^{m+1} . Since $C_1 \leq h \leq C_2$ on $\Theta \times \Xi$,
 390 $v = [\underbrace{C_2, \dots, C_2}_m, 0]^T$ is a feasible solution of (CLP-D), we have $\text{val}(\text{CLP-D}) \leq C_2$.

391 Moreover, substituting $s = \xi^i$ in the i^{th} constraint of (CLP-D), we have $v_i \geq h(\theta, \xi^i) \geq$
 392 C_1 , $i \in [m]$, which means that C_1 is a lower bound on v_i , $i \in [m]$. To show that
 393 $\text{Sol}(\text{CLP-D})$ is bounded, consider any $v^* \in \text{Sol}(\text{CLP-D})$. Since $\text{val}(\text{CLP-D}) \leq C_2$, we
 394 have $\frac{1}{m} \sum_{i=1}^m v_i^* + r_0 \cdot v_{m+1}^* \leq C_2$. It follows that

$$395 \quad (22) \quad v_{m+1}^* \leq \frac{1}{r_0} \left(C_2 - \frac{1}{m} \sum_{i=1}^m v_i^* \right) \leq \frac{1}{r_0} (C_2 - C_1),$$

$$v_i^* \leq m(C_2 + r_0 \cdot v_{m+1}^*) - \sum_{j=1}^{i-1} v_j^* - \sum_{j=i+1}^m v_j^* \leq (m+1)C_2 - mC_1, \quad i \in [m],$$

396 which provide upper bounds on v_i^* , $i \in [m+1]$. Therefore, $\text{Sol}(\text{CLP-D})$ is bounded
 397 by the following compact set:

$$398 \quad (23) \quad \mathcal{H} := \left\{ v \in \mathbb{R}^{m+1} : C_1 \leq v_i \leq (m+1)C_2 - mC_1, \text{ for } i \in [m], \right. \\ \left. 0 \leq v_{m+1} \leq (C_2 - C_1)/r_0 \right\}.$$

399 Now $\text{Fsb}(\text{CLP-D}) \cap \mathcal{H}$ is a non-empty compact set in \mathbb{R}^{m+1} , and $\text{Sol}(\text{CLP-D})$ is a
 400 subset of $\text{Fsb}(\text{CLP-D}) \cap \mathcal{H}$. By Weierstrass Theorem [8], there exists an optimal
 401 solution to (CLP-D). Therefore, $\text{Sol}(\text{CLP-D})$ is nonempty, bounded, and attains its
 402 optimum at a solution in \mathcal{H} . \square

403 **COROLLARY 3.8.** *Let Assumption 3.1 hold. The Wasserstein-robust optimization*
 404 *problem (WRO) is equivalent to the following semi-infinite program:*

$$405 \quad (\text{WRO-D}) \quad \min_{\theta, v} \quad \frac{1}{m} \sum_{i=1}^m v_i + r_0 \cdot v_{m+1}$$

$$\text{s.t.} \quad h(\theta, s) - v_i - v_{m+1} \cdot d(s, \xi^i) \leq 0, \quad s \in \Xi, i \in [m]$$

$$\theta \in \Theta, v \in \mathcal{H}.$$

406 *Proof.* From Theorem 3.7, the inner problem of (WRO) is reformulated as a
 407 minimization problem with semi-infinite constraints. We can now combine (CLP-D)
 408 with the outer problem of (WRO), and have an equivalent combined formulation
 409 as (WRO-D). Since from Theorem 3.7 the optimal solution of the inner problem is
 410 bounded in the polytope \mathcal{H} , these additional constraints are added to (WRO-D). \square

411 We note that the constraints in (WRO-D) decompose in the scenarios ξ^i , $i \in [m]$; and
 412 for a given s , $d(s, \xi^i)$ is a constant.

413 **4. Algorithms for the (WRO) Refomulation.** Corollary 3.8 shows that the
 414 Wasserstein-robust optimization problem (WRO) is equivalent to (WRO-D), which
 415 is a semi-infinite program. Any algorithm for solving a general semi-infinite program
 416 can now be applied to solve (WRO-D). For a general continuous function $h(\theta, s)$ in θ
 417 we present a modification of the exchange algorithm [29] in Section 4.1, which ensures
 418 ε -optimality after solving a finite number of finitely constrained master problems. In
 419 the special case where $h(\theta, s)$ is a convex function of θ , (WRO-D) is a convex semi-
 420 infinite program. We adapt the cutting surface algorithm in [43] for (WRO-D) and
 421 use its structure to achieve a global linear rate of convergence.

422 Let $x = [\theta, v]$ be the decision variables, and define the following functions:

$$423 \quad (24) \quad \begin{aligned} f(x) &:= \frac{1}{m} \sum_{i=1}^m v_i + r_0 \cdot v_{m+1}, \\ g_i(x, s) &:= h(\theta, s) - v_i - v_{m+1} \cdot d(s, \xi^i), \quad i \in [m]. \end{aligned}$$

424 The problem (WRO-D) can be rewritten as:

$$425 \quad (\text{SIP}) \quad \begin{aligned} \min_x \quad & f(x) \\ \text{s.t.} \quad & g_i(x, s) \leq 0, \quad s \in \Xi, i \in [m] \\ & x \in X, \end{aligned}$$

426 where $X = \Theta \times \mathcal{H}$. Problem (SIP) is a semi-infinite program. An approach to obtain
 427 a solution of such problems is to solve relaxation problems (master problems) with
 428 a finite number of constraints, and add a violated constraint obtained from solving
 429 a separation problem (defined for $\forall i \in [m]$) to tighten the formulation iteratively.
 430 In particular, the separation problem of (SIP) for identifying a violated constraint at
 431 the solution (\tilde{x}, \tilde{v}) of the current master problem can be written as follows:

$$432 \quad (\text{Sep-}i) \quad \max_{s \in \Xi} g_i(\tilde{x}, s) = h(\tilde{\theta}, s) - \tilde{v}_i - \tilde{v}_{m+1} \cdot d(s, \xi^i), \quad \text{for } i \in [m].$$

433 The inequality generated from solving (Sep- i) is called a feasibility cut.

434 For clarity of notations, we consider the following general form of a semi-infinite
 435 program:

$$436 \quad (\text{gen-SIP}) \quad \begin{aligned} \min_x \quad & f(x) \\ \text{s.t.} \quad & g(x, t) \leq 0, \quad t \in T, \\ & x \in X, \end{aligned}$$

437 where $X \subseteq \mathbb{R}^{k_1}$ and $T \subseteq \mathbb{R}^{k_2} \times \mathbb{Z}^{k_3}$, allowing that T may be defined as a mixed-integer
 438 set.

439 **4.1. A modified exchange algorithm for (WRO) model.** We now assume
 440 that we have an oracle to solve the master and separation problems of (gen-SIP).
 441 The modified exchange algorithm given in Algorithm 1 allows an ε -optimal solution
 442 to (gen-SIP), when compared with the original exchange algorithm. Theorem 4.2
 443 shows that our modified exchange method finds a solution of a desired accuracy in
 444 finitely many iterations.

445 **DEFINITION 4.1.** A point $z_0 \in Z$ is an ε -feasible solution of (gen-SIP) if
 446 $\max_{t \in T} g(z_0, t) \leq \varepsilon$. A point $z_0 \in Z$ is an ε -optimal solution of (gen-SIP) if z_0 is
 447 ε -feasible and $f(z_0) \leq \text{val}(\text{gen-SIP})$.

448 **THEOREM 4.2.** *Let $Z \times T$ be compact, and $g(z, t)$ be continuous on $Z \times T$. Suppose*
 449 *we have an oracle that generates an optimal solution of the problem $\min_{z \in Z} \{f(z) :$*
 450 *$g(z, t) \leq 0, t \in T'\}$ for any finite set $T' \subseteq T$, and an oracle that generates an ε -*
 451 *optimal solution of the problem $\max_{t \in T} g(z, t)$ for any $z \in Z$ and $\varepsilon > 0$. Then Algorithm 1*
 452 *returns an ε -optimal solution of (gen-SIP) in finitely many iterations.*

453 *Proof.* Since $g(z, t)$ is continuous on $Z \times T$, and $Z \times T$ is compact, it follows that
 454 $g(z, t)$ is uniformly continuous on $Z \times T$. Therefore, there exists an $\alpha > 0$ such that

$$455 \quad (25) \quad |g(z', t') - g(z, t)| \leq \frac{\varepsilon}{2}, \quad \text{if } \|z' - z\| + \|t' - t\| \leq \alpha.$$

456 First, we prove by contradiction that the algorithm terminates in finitely many
 457 iterations. Suppose the algorithm generates infinite sequences $\mathcal{Z} = \{z_k\}_0^\infty$ and
 458 $\mathcal{T} = \{t_k\}_0^\infty$ without terminating. We claim that $t_{k+1} \notin \cup_{i=1}^k B(t_i, \alpha)$ for every k ,
 459 where $B(t_i, \alpha)$ is the closed ball of center t_i and radius α in $\mathbb{R}^{k_2+k_3}$. If not, there
 460 exists t_k such that $t_k \in B(t_i, \alpha)$ for some $i < k$. Using (25) and $t_i \in T_{z_{i-1}}$, we have
 461 $g(z_{k-1}, t_k) \leq g(z_{k-1}, t_i) + \frac{\varepsilon}{2} \leq \frac{\varepsilon}{2}$, indicating that the termination criteria is satisfied.
 462 Therefore, the above claim holds. Now consider the compact set: $\mathcal{B} = \cup_{\{t \in T\}} B(t, \alpha)$.
 463 Since $B(t_k, \alpha) \subseteq \mathcal{B}$ for every k , and by the claim we have $B(t_i, \alpha) \cap B(t_j, \alpha) = \emptyset$ for
 464 every $t_i, t_j \in \mathcal{T}$ with $i \neq j$, it follows that $\sum_{i=0}^\infty \text{vol}(B(t_i, \alpha)) \leq \text{vol}(\mathcal{B})$, which leads to
 465 a contradiction.

466 We now prove that once the algorithm terminates, it returns an ε -solution. Sup-
 467 pose it terminates at the end of the n th iteration. Then z_n is an optimal solution
 468 of the problem $\min_{z \in Z} \{f(z) : g(z, t) \leq 0, t \in T_n\}$. It follows that $f(z_n) \leq \text{val}(T_n) \leq$
 469 $\text{val}(\text{gen-SIP})$, where $\text{val}(T_n)$ is the optimal value of the problem (gen-SIP) with con-
 470 straint index set T_n . By the separation oracle and termination criteria, it also holds
 471 that $\max_{t \in T} g(z_k, t) \leq g(z_n, t_{n+1}) + \frac{\varepsilon}{2} \leq \varepsilon$. Hence, z_n is an ε -optimal solution of
 472 (gen-SIP). \square

Algorithm 1 A modified exchange algorithm to solve (gen-SIP).

Prerequisites: Two oracles specified in Theorem 4.2.

Output: An ε -optimal solution of (gen-SIP).

Step 1 Set $T_0 \leftarrow \emptyset, k \leftarrow 0$.

Step 2 Determine an optimal solution z_k of the problem $\min_{z \in Z} \{f(z) : \text{s.t. } g(z, t) \leq 0, t \in T_k\}$.

Step 3 Determine a $\frac{\varepsilon}{2}$ -optimal solution t_{k+1} of the problem $\max_{t \in T} g(z_k, t)$. If $g(z_k, t_{k+1}) \leq \frac{\varepsilon}{2}$, stop and return z_k ; otherwise let $T_{k+1} \leftarrow T_k \cup \{t_{k+1}\}, k \leftarrow k + 1$ and go to Step 2

473 **4.2. A central cutting-surface algorithm for the convex case.** In this
 474 section we make the following additional assumption.

475 **ASSUMPTION 4.3.** *The feasible region Θ is convex, and the function $h(\cdot, s)$ is con-*
 476 *vex for every $s \in \Xi$. Furthermore, there exists a parameter B satisfying the following*
 477 *condition:*

$$478 \quad (26) \quad B > \|\eta\|, \quad \forall \eta \in \partial_\theta h(\theta, s) \quad \forall \theta \in \Theta, s \in \Xi,$$

479 where $\partial_\theta h(\theta, s)$ is the sub differential set of $h(\theta, s)$ at θ .

480 Since the master problem of (SIP) is a convex optimization problem, we assume it
 481 can be solved efficiently up to optimality. We now present a central cutting-surface
 482 algorithm to solve (SIP). A pseudo-code for this algorithm is given in Algorithm 2.
 483 The algorithm is initialized with a solution $x^{(0)} = [\theta^0, \mathbf{0}_{m+1}]$, where θ^0 may be taken
 484 as a solution of the empirical deterministic optimization problem:

$$485 \quad (\text{EDO}) \quad \min_{\theta \in \Theta} \frac{1}{m} \sum_{i=1}^m h(\theta, \xi^i).$$

486 At the k th iteration of this algorithm (Step 1) a master problem with a centering
 487 argumentation is solved. This problem is defined by a set of constraints (with the
 488 index set $Q^{(k-1)} \subseteq \Xi$) inherited from the $(k-1)$ th iteration. The master problem at
 489 the k th iteration is formulated as follows:

$$490 \quad (\text{Master}) \quad \begin{aligned} & \max_{x, w, t} && w \\ & \text{s.t.} && t + w \leq M^{(k-1)}, \\ & && f(x) - t + B \cdot w \leq 0, \\ & && g_i(x, s) + B \cdot w \leq 0, \quad \forall i \in [m], \quad s \in Q^{(k-1)}, \\ & && x \in X. \end{aligned}$$

491 The (Master) problem is a convex optimization problem. For clarity, we drop the
 492 index i in (Master) in Algorithm 2. We assume that there is an oracle to find an
 493 ε -optimal solution to (Sep- i) for any $\varepsilon > 0$. The algorithm terminates if no feasibility
 494 cut is found. Otherwise, the newly found feasibility cut is added to the working set
 495 (Step 4). At the end of each iteration, we may optionally drop certain constraints
 496 that are not binding at the current solution of the master problem (Step 6).

Algorithm 2 A central cutting-surface algorithm from [43] to solve (SIP).

Prerequisites: Assumptions 3.1 and 4.3 hold. There exists an oracle to find an ε -optimal solution to (Sep- i).

Input: A strict upper bound U of the objective function h , a centering parameter $B > 0$ which satisfies (26), a tolerance error ε , an arbitrary $\alpha > 1$ specifying how aggressively cuts are dropped.

Output: An ε -optimal solution to (SIP).

Step 1 (Initialization). Set $k \leftarrow 1$, $M^{(0)} \leftarrow U$, $v^{(0)} \leftarrow \mathbf{0}_{m+1}$, $x^{(0)} \leftarrow [\theta^0, \mathbf{0}_{m+1}]$,
 $\tilde{x}^{(0)} \leftarrow x^{(0)}$, where θ^0 is a solution to (EDO). Let $Q^{(0)} \leftarrow \{s^{(0)}\}$, where $s^{(0)} = \xi$.

Step 2 (Solve a master problem). Determine the optimal solution $(x^{(k)}, w^{(k)})$ to (Master).

Step 3 (Optimal solution?). If $w^{(k)} = 0$, stop and return $\tilde{x}^{(k-1)}$.

Step 4 (Feasible solution?). Find an ε -optimal solution denoted by $s^{(k)}$ to (Sep- i) using the oracle. If $g(x^{(k)}, s^{(k)}) > 0$ and go to Step 5, otherwise go to Step 6.

Step 5 (Add a cut). Set $Q^{(k)} \leftarrow Q^{(k-1)} \cup \{s^{(k)}\}$, $\tilde{x}^{(k)} \leftarrow \tilde{x}^{(k-1)}$ and $M^{(k)} \leftarrow M^{(k-1)}$.
 Go to Step 6.

Step 6 (Update best know ε -feasible solution). Set $Q^{(k)} \leftarrow Q^{(k-1)}$, $\tilde{x}^{(k)} \leftarrow x^{(k)}$
 $M^{(k)} \leftarrow f(x^{(k)})$.

Step 7 (Optionally drop cuts). Let $D = \{s^{(l)} \in Q^{(k)} : l \in \{0\} \cup [k] \mid w^{(l)} \geq \alpha w^{(k)} \text{ and } g(x^{(k)}, s^{(l)}) + B \cdot w^{(k)} < 0\}$ and set $Q^{(k)} \leftarrow Q^{(k)} \setminus D$.
 Set $k \leftarrow k + 1$ and go to Step 2.

497 The convergence of Algorithm 2 is given by Theorem 4.4. This theorem is a
 498 refinement of the linear rate of convergence result proved in Theorem 8 of [43] (See
 499 Appendix A) for the central cutting surface algorithm when specialized to (WRO-D).
 500

501 **THEOREM 4.4.** *Convergence of the central cutting-surface algorithm.*

- 502 (a) *Algorithm 2 either finds an ε -optimal solution of (SIP) in finitely many it-*
 503 *erations or generates $\{\tilde{x}^{(k)}\}_{k=1}^{\infty}$ that each accumulated point is an ε -optimal*
 504 *solution of (SIP).*
 505 (b) *Algorithm 2 converges linearly in objective function value between consecutive*
 506 *feasibility cuts. The rate of convergence satisfies:*

$$507 \quad (27) \quad \rho \leq 1 - \frac{1}{2B' + 1},$$

508 where $B' := \max\{\sqrt{r_0^2 + 1 + (1/m)}, \sqrt{B^2 + L^2 + 1}\}$ and $L = \max\{d(s, \xi^i) :$
 509 $s \in \Xi, i \in [m]\}$.

510 *Proof.* Note that (SIP) is equivalent to the following reformulation:

$$511 \quad (28) \quad \begin{array}{ll} \min_{x,t} & t \\ \text{s.t.} & f(x) - t \leq 0, \\ & g_i(x, s) \leq 0, \quad \forall s \in \Xi, \forall i \in [m] \\ & x \in X. \end{array}$$

512 Treating t as x_0 in (38), we now verify that (28) satisfies Assumption A.2 in Ap-
 513 pendix A. Since $\Theta \times \Xi$ is bounded and $h(\cdot, \cdot)$ is continuous, $\exists C_1, C_2$ such that $h(\cdot, \cdot) \in$
 514 $[C_1, C_2]$ on $\Theta \times \Xi$, and the objective value of (WRO) is in $[C_1, C_2]$. Also, since
 515 there is no duality gap, we can set the dual objective $C_1 \leq f(x) \leq C_2$ for all
 516 $x \in X$, and set $t \in [C_1, C_2]$, without affecting the optimal solution and the op-
 517 timal value. Hence, Assumption A.2 (1) is satisfied. To verify Assumption A.2
 518 (2), we note that for any $\eta > 0$, we can verify that $[\bar{t}, \bar{x}]$ is a Slater point of (28),
 519 where $\bar{t} = C_2 + 2\eta$ and $\bar{x} = [\theta^0, (C_2 + \eta)\mathbf{1}_m, 0]$ ($\mathbf{1}_m$ is the m -dimensional vector
 520 with all entries being 1). Now we show that Assumption A.2 (3) is also satis-
 521 fied. Now let us take a subgradient of $f(x) - t$ and $g_i(x, s)$ ($\forall i \in [m]$) with re-
 522 spect to the decision variables $[t, x]$, and use Assumption 4.3, we can set the cen-
 523 terality parameter B' to be $B' := \max\{\sqrt{r_0^2 + 1 + (1/m)}, \sqrt{B^2 + L^2 + 1}\}$, where
 524 $L = \max\{d(s, \xi^i) : s \in \Xi, i \in [m]\}$. The oracle is assumed to be given based on the
 525 prerequisites of Algorithm 2, hence Assumption A.2 (4) is satisfied.

526 We now apply Theorem A.3 on the semi-infinite program (28) and the correspond-
 527 ing master problem (Master). Part (a) directly follows from Theorem A.3 (a)-(c). By
 528 Theorem A.3 (d), for $k \geq \hat{k}$ iterations, where $w^{(\hat{k})} < \eta/B$, we have

$$529 \quad (29) \quad \rho^{(k)} \leq 1 - \frac{\eta - B'w^{(k)}}{\eta + B'(\bar{t} - f^*)},$$

530 where f^* is the optimal value of (28). Since (29) holds for every $\eta > 0$, we can
 531 select an η to maximize $\frac{\eta - B'w^{(k)}}{\eta + B'(\bar{t} - f^*)}$, hence minimizing the upper bound of $\rho^{(k)}$ in
 532 Theorem A.3. Let $F(\eta) := \frac{\eta - B'w^{(k)}}{\eta + B'(\bar{t} - f^*)}$, and substitute $\bar{t} = C_2 + 2\eta$ in $F(\eta)$, we have

$$533 \quad (30) \quad F(\eta) = \frac{\eta - B'w^{(k)}}{(2B' + 1)\eta + B'(C_2 - f^*)}, \quad F'(\eta) = \frac{B'w^{(k)}(2B' + 1) + B'(C_2 - f^*)}{[(2B' + 1)\eta + B'(C_2 - f^*)]^2}.$$

534 Since $w^{(k)} > 0$, $f^* \leq C_2$, we have $F'(\eta) > 0$ for all $\eta > 0$. Therefore, the maximum
 535 value of $F(\cdot)$ is:

$$536 \quad (31) \quad \max_{\eta > 0} F(\eta) = F(\infty) = \frac{1}{2B' + 1}.$$

537 It follows that the uniform rate of convergence satisfies: $\rho \leq 1 - \frac{1}{2B'+1}$. \square

538 **4.3. Computational tractability of the separation problem.** We now dis-
 539 cuss the computational difficulty of solving (Sep- i), which depends on the function
 540 form of $h(\theta, s)$ in s for a given θ , the metric d , and the uncertainty set Ξ . Since
 541 $\tilde{v}_{m+1} \geq 0$ and in most applications $d(\cdot, \cdot)$ is chosen to be a vector norm, the term
 542 $-\tilde{v}_i - \tilde{v}_{m+1}d(s, \xi^i)$ in (Sep- i) is concave in s . Therefore, in the case where $h(\theta, s)$
 543 is concave in s , and Ξ is a convex set, (Sep- i) becomes a convex optimization prob-
 544 lem. For a very general case where $h(\theta, \cdot)$ and $d(\xi^i, \cdot)$ are continuously differentiable
 545 over the compact (not necessarily convex) uncertainty set Ξ , drawing a sufficiently
 546 many independent uniform samples $s^t \in \Xi$ and verifying if objective value of (Sep- i)
 547 is greater than ε can either identify a violated constraint (not necessarily the most
 548 violated constraint), or conclude that the solution of the current master problem is
 549 ε -optimal with high probability (See [43] Section 5.2). Furthermore, for cases where
 550 $h(\theta, \cdot)$ is a polynomial function, $d(\cdot, \cdot)$ is an Euclidean norm, and Ξ is specified by
 551 polynomial inequalities (e.g., an ellipsoid), (Sep- i) falls in to the category of poly-
 552 nomial optimization. The global optimal solution in such cases can be obtained by
 553 solving a sequence of SDP relaxations (a primal approach) [35] or a sequence of SOS
 554 relaxations (a dual approach) [48].

555 For some important models from statistical learning such as linear regression,
 556 linear support vector machine, logistic regression, etc., the loss function h has the
 557 form $h(\theta_0 + \theta^T x, y)$, where x is the feature vector and y is the response value. For
 558 this case, (Sep- i) can be solved efficiently using a branch-and-bound scheme based
 559 on piece-wise linear approximations of $h(\theta_0 + \theta^T x, y)$. For more details about this
 560 approach, see [42].

561 **5. Numerical Experiments.** We investigate the following Wasserstein-robust
 562 logistic regression (WRLR) model as a numerical example to illustrate our algorithmic
 563 ideas and the performance of (WRO):

$$564 \quad (\text{WRLR}) \quad \min_{[\theta_0, \theta] \in \Theta} \max_{P \in \mathcal{P}} \mathbb{E}_P \left[\log(1 + \exp[-y(\theta_0 + \theta^T x)]) \right],$$

565 where $x \in \mathbb{R}^n$ is the feature vector, and $y \in \{0, 1\}$ is the label. The semi-infinite
 566 reformulation of (WRLR) is written as follows:

$$567 \quad (32) \quad \begin{aligned} & \min_{\theta_0, \theta, v} \frac{1}{m} \sum_{i=1}^m v_i + r_0 \cdot v_{m+1} \\ & \text{s.t.} \quad \log(1 + \exp[-y(\theta_0 + \theta^T x)]) - v_i - v_{m+1} \cdot d(s, \xi^i) \leq 0, \quad s \in \Xi, i \in [m], \\ & \quad [\theta_0, \theta] \in \Theta, v \in \mathcal{H}, \end{aligned}$$

568 where $s = [x, y]$, and $\xi^i = [x^i, y^i]$ is the i th ($i \in [m]$) observed sample. We assume
 569 that only the feature vector x has uncertainty but not the label y . Therefore, the
 570 uncertainty set Ξ can be written as $\Xi = \Xi_0 \cup \Xi_1$, where Ξ_0, Ξ_1 are the uncertainty sets
 571 for feature vectors with the label 0, 1 respectively. The uncertainty set Ξ_0 is defined as

Table 1: Summary of data sets from UCI Machine Learning Repository

Data set	Area	No. Attrib.	No. Observ.
BA	Finance	4	1372
VC	Health care	6	310
PID	Health care	8	768
BCW	Health care	9	699
ST-H	Health care	13	270
EES	Health care	14	14980
SPT-H	Health care	22	267
ION	Aerospace	34	351
SPTF-H	Health care	44	267
SPAM	Computer	57	4601
CB	Aerospace	60	208

572 an n -dimensional hyper-rectangle such that each dimension corresponds to an interval
 573 $[\bar{x}_j \pm \sigma_j]$ for the feature x_j , where \bar{x}_j is the sample mean of observations with label 0,
 574 and σ_j is the sample standard deviation. The uncertainty set Ξ_1 is defined similarly.
 575 We used the l_1 norm to define the metric d on Ξ_0 (Ξ_1), i.e., $d(x, x') = \|x - x'\|_1$.

576 **5.1. Numerical setup.** We present computational effort in solving the semi-
 577 infinite dual problem (**WRO-D**) of (**WRLR**) using the cutting-surface algorithm (See
 578 Section 4.2). We also present the out of sample predictive performance of the (**WRLR**)
 579 model compared with the ordinary logistic regression model (**LR**). The algorithm
 580 for solving (**WRO-D**) were implemented in C++, and the computational tests were
 581 performed on an Intel Xeon CPU with 4 GB of RAM. The cutting-surface algorithm
 582 for (**WRO-D**) consists of iteratively solving the master problem (**Master**) and the
 583 separation problem (**Sep- i**). The convex master problem (**Master**) is solved using Ipopt
 584 3.12.4 [59] which implements a primal-dual interior point method. The separation
 585 problem (**Sep- i**) is solved using the branch-and-bound scheme based on sequentially
 586 piece-wise linear approximating $h_\theta(u) := \log(1 + e^{-u})$, and each convex optimization
 587 subproblem is solved using CPLEX 12.6.3. We used 11 data sets (those with less
 588 than 60 features) from the UCI machine learning repository in our computational
 589 testing, which are: Banknote authentication (BA), Vertebral column (VC), Pima
 590 Indians diabetes (PID), Breast cancer Wisconsin (BCW), Statlog heart (ST-H), EEG
 591 eye state (EES), SPECT heart (SPT-H), Ionosphere (ION), SPECTF heart (SPTF-
 592 H), Spambase (SPAM), Connectionist bench (CB). A summary of these datasets is
 593 given in Table 1. We now describe the data generation for our test problems. We
 594 chose m ($m = 50, 75, 100, 150$) observations from each of the UCI data sets. We kept
 595 the class labels of the chosen observations unchanged.

596 **5.2. Computational effort in solving WRLR.** The computational effort in
 597 solving the semi-infinite dual problem (**WRO-D**) of (**WRLR**) for 11 data sets is given in
 598 Table 2 for different choices of m . The numbers are averaged over the 100 experiments.
 599 Columns 4-8 give the number of outer iterations in Algorithm 2, total number of
 600 constraints (cuts) added to the master problem at termination, the CPU time for
 601 solving a problem instance, and the percentage of time spend in solving the master
 602 and separation problems, respectively. We also provide the CPU time for solving
 603 (**LR**) in Column 3. The results show that the number of calls to the master problem

604 is approximately $4 \sim 40$ when solving (WRLR). Over these calls approximately $15m \sim$
 605 $50m$ cutting surfaces are added. In other words, approximately $15m \sim 50m$ artificial
 606 samples are identified. For data sets with more features, the program spends a larger
 607 fraction of time on solving master problems since their scale becomes larger. The
 608 computational time of solving (LR) models is less than 1 second for data sets with
 609 feature size less than 20, and for data sets with feature size between $30 \sim 60$ the
 610 computational time is less than 20 seconds. The average time of solving (WRLR)
 611 models is $\lesssim 100$ times that of solving the (LR) models.

612 **5.3. Predictive performance of the WRLR model.** We now compare the
 613 predictive performance of the (WRLR) model with the (LR) model. For each data set,
 614 we randomly select m samples out of all the observed samples to train both models.
 615 For each combination of data set and the training sample size m , we performed 100
 616 experiments. Both trained (LR) and (WRLR) models are used to predict the remain-
 617 ing observations from the data set, and the corresponding AUC values (area under the
 618 ROC curve) are recorded. AUC value is the most popular metric used for evaluating
 619 the performance of a model used in medical literature. In each experiment, training
 620 samples are selected randomly and independently. The mean AUC values ($\overline{\text{AUC}}$)
 621 over 100 experiments and the standard errors are listed in Table 2. The p-values
 622 in this table are based on the hypothesis test: $H_0 : \overline{\text{AUC}}_{\text{WRLR}}^{\text{OOS}} \leq \overline{\text{AUC}}_{\text{LR}}^{\text{OOS}}$ versus
 623 $H_1 : \overline{\text{AUC}}_{\text{WRLR}}^{\text{OOS}} > \overline{\text{AUC}}_{\text{LR}}^{\text{OOS}}$, where $\overline{\text{AUC}}_{\text{WRLR}}^{\text{OOS}}$ and $\overline{\text{AUC}}_{\text{LR}}^{\text{OOS}}$ denote the out of sample
 624 mean AUC values corresponding to (WRLR) and (LR), respectively. Statistically,
 625 the smaller the p-value, the more likely H_1 is true.

626 To train the (WRLR) model, one needs to specify the radius r_0 of the Wasser-
 627 stein ball. One way to determine this radius is to use the concentration inequality
 628 $\Pr(\mathcal{W}(P_{\text{true}}, P_0) \leq r_0) \geq 1 - \gamma$. The theoretical bounds in [22] are of limited value.
 629 Instead we used six candidate empirical Wasserstein radii: $\{0, 0.01, 0.05, 0.1, 0.5, 1\}$ ¹
 630 and a cross-validation approach to select the best r_0 from these. Specifically, we used
 631 the following 4-fold cross-validation approach: we divided m training samples into
 632 4 subsets and used any three of them to train WRLR with every candidate value
 633 of r_0 and tested the model on the remaining subset. We finally picked the r_0 value
 634 corresponding to the best mean AUC value over 4 folds. Once this r_0 is selected, it is
 635 used for out of sample testing on the remaining observations that is not part of the
 636 chosen m samples.

637 The comparison shows that the quantity $(\overline{\text{AUC}}_{\text{WRLR}}^{\text{OOS}} - \overline{\text{AUC}}_{\text{LR}}^{\text{OOS}})$ ranges from -
 638 .0462 to .1005 and the relative difference ranges from -.3791 to .7122, in the studied
 639 cases. We observe that $\overline{\text{AUC}}_{\text{WRLR}}^{\text{OOS}}$ is greater than $\overline{\text{AUC}}_{\text{LR}}^{\text{OOS}}$ in 34 (77%) cases out of
 640 all 44 cases. The standard errors associated with (WRLR) are smaller than that of
 641 (LR) in 31 (72%) cases. This suggests that not only the distributionally-robust model
 642 is better, its performance is also more stable. It is seen from the p-values at the
 643 significance level $\alpha = 0.05$, (WRLR) outperforms (LR) in 24 (55%) cases which
 644 are from seven data sets: BA, BCW, ST-H, SPT-H, ION, SPTF-H and SPAM. For 7
 645 (16%) cases which are from data sets: VC, PID and EES, $\overline{\text{AUC}}_{\text{WRLR}}^{\text{OOS}}$ is significantly
 646 smaller than $\overline{\text{AUC}}_{\text{LR}}^{\text{OOS}}$, indicating (WRLR) is not as good as (LR) for these three data
 647 sets. For the remaining 13 (29%) cases that are from data sets: PID, ST-H, EES,
 648 SPT-H, SPAM and CB, the difference in mean AUC is not statistically significant.

¹Note that with $r_0 = 0$, the (WRLR) reduces to the (LR) model.

Table 2: Computational performance of solving WRLR and predictive performance of WRLR compared with LR for 11 data sets. Listed values are average of 100 experiments.

Data set	m	Computational Performance										Predictive Performance					
		LR		WRLR				LR				WRLR					
		CPU [sec]	No. main iters. ¹	No. cuts	CPU [sec]	Master ² (%)	Sep. ³ (%)	Mean AUC	S.E.	Diff.	Rel. Diff. ⁴	p-value	Mean AUC	S.E.	Diff.	Rel. Diff. ⁴	p-value
BA	50	0.022	3.8	66.9	1.21	13.74	86.26	.9985	.0002	.9991	.0001	.0006	.4202	.0062	.0000	.0000	.0000
	75	0.025	4.3	90.8	0.86	17.96	82.04	.9985	.0001	.9994	.0000	.0009	.5775	.0000	.0000	.0000	.0000
	100	0.032	3.9	116.7	1.83	13.89	86.11	.9993	.0000	.9995	.0000	.0002	.2724	.0001	.0000	.0000	.0000
VC	50	0.042	4.6	157.5	2.30	14.42	85.58	.9996	.0000	.9997	.0000	.0000	.1179	.0000	.0000	.0000	.0000
	75	0.021	7.2	195	1.65	30.51	69.49	.8782	.0025	.8320	.0034	-.0462	-.3791	1.0000	.0000	.0000	.0000
	100	0.032	7.1	271.3	1.82	29.15	70.85	.8867	.0022	.8441	.0029	-.0426	-.3755	1.0000	.0000	.0000	.0000
PID	50	0.047	6.6	323.3	3.73	32.44	67.56	.8887	.0023	.8504	.0025	-.0383	-.3442	1.0000	.0000	.0000	.0000
	75	0.064	6.5	541.5	5.19	30.63	69.37	.8951	.0023	.8601	.0029	-.0350	-.3337	1.0000	.0000	.0000	.0000
	100	0.032	7.5	203.7	4.19	22.20	77.80	.7542	.0058	.7564	.0037	.0022	.0089	.3756	.0000	.0000	.0000
BCW	50	0.044	6.2	244	5.15	19.78	80.22	.8009	.0020	.8000	.0017	-.0009	-.0043	.6281	.0000	.0000	.0000
	75	0.076	9.4	284	7.88	27.48	72.52	.9773	.0016	.9886	.0010	.0112	.4954	.0000	.0000	.0000	.0000
	100	0.097	8.9	501.4	12.84	34.28	65.72	.9790	.0025	.9940	.0001	.0150	.7122	.0000	.0000	.0000	.0000
ST-H	50	0.124	9.6	786.1	27.31	41.79	58.21	.9889	.0007	.9945	.0001	.0056	.5049	.0000	.0000	.0000	.0000
	75	0.069	11.7	242.7	12.18	30.56	69.44	.8317	.0035	.8808	.0029	.0490	.2914	.0000	.0000	.0000	.0000
	100	0.159	8.5	321.9	13.02	27.58	72.42	.8504	.0036	.8903	.0018	.0398	.2664	.0000	.0000	.0000	.0000
EES	50	0.147	9.6	381.5	10.50	24.15	75.85	.8945	.0022	.9064	.0011	.0120	.1133	.0000	.0000	.0000	.0000
	75	0.193	8.1	472.5	21.14	29.02	70.98	.8986	.0017	.8990	.0017	.0004	.0042	.4319	.0000	.0000	.0000
	100	0.041	15.7	251.9	7.80	37.72	62.28	.5874	.0037	.5902	.0035	.0029	.0070	.2877	.0000	.0000	.0000
SPT-H	50	0.079	12.1	435	11.11	37.37	62.63	.6095	.0038	.5996	.0041	-.0099	-.0252	.9602	.0000	.0000	.0000
	75	0.322	11.9	577.8	17.78	36.55	63.45	.6221	.0033	.6175	.0032	-.0047	-.0124	.8459	.0000	.0000	.0000
	100	0.383	8.5	737.8	31.56	41.78	58.22	.6245	.0023	.6139	.0026	-.0106	-.0282	.9987	.0000	.0000	.0000
ION	50	0.241	21.5	938.8	38.87	88.91	11.09	.8126	.0016	.8176	.0018	.0050	.0269	.0203	.0000	.0000	.0000
	75	0.343	24.5	1031.2	53.43	83.40	16.60	.8221	.0021	.8311	.0030	.0089	.0501	.0078	.0000	.0000	.0000
	100	0.54	19.4	1122.5	63.08	83.44	16.56	.8311	.0022	.8370	.0037	.0058	.0346	.0866	.0000	.0000	.0000
50	0.987	13.4	1384	49.70	77.91	22.09	.8301	.0029	.8567	.0033	.0267	.1569	.0000	.0000	.0000	.0000	
50	0.912	34.7	740.6	173.37	63.56	36.44	.8429	.0024	.8708	.0021	.0279	.1775	.0000	.0000	.0000	.0000	.0000

Continued on the next page

Table 2 – continued from previous page

Data set	m	Computational Performance										Predictive Performance					
		LR					WRLR					LR			WRLR		
		CPU [sec]	No. main iters. ¹	No. cuts	CPU [sec]	Master ² (%)	Sep. ³ (%)	Mean AUC	S.E.	Diff.	Rel. Diff. ⁴	p-value	Mean AUC	S.E.	Diff.	Rel. Diff. ⁴	p-value
	75	2.415	28.2	1208.5	229.49	62.95	37.05	.8582	.0031	.0338	.2381	.8919	.0018	.0338	.2381	.0000	
	100	1.51	23.1	1502.5	441.58	68.51	31.49	.8606	.0029	.0360	.2584	.8967	.0018	.0360	.2584	.0000	
	150	2.774	17.8	1718.3	386.63	61.30	38.70	.8715	.0024	.0291	.2264	.9006	.0021	.0291	.2264	.0000	
SPTF-H	50	2.144	38.9	438.1	62.68	70.38	29.62	.6818	.0048	.0700	.2198	.7517	.0056	.0700	.2198	.0000	
	75	3.718	31.6	791.4	112.66	74.41	25.59	.7255	.0045	.1005	.3661	.8260	.0031	.1005	.3661	.0000	
	100	3.453	29	983.7	144.66	75.50	24.50	.7638	.0033	.0730	.3091	.8368	.0027	.0730	.3091	.0000	
	150	5.694	23.7	1327.6	221.44	77.87	22.13	.8142	.0039	.0490	.2638	.8632	.0037	.0490	.2638	.0000	
SPAM	50	3.655	50.5	538.2	266.40	66.10	33.90	.8969	.0012	.0015	.0142	.8984	.0010	.0015	.0142	.1796	
	75	4.76499	46.9	918.8	422.79	60.15	39.85	.9219	.0014	.0030	.0386	.9249	.0019	.0030	.0386	.0998	
	100	11.065	40.5	1368	409.94	65.66	34.34	.9164	.0024	.0166	.1992	.9331	.0008	.0166	.1992	.0000	
	150	15.76	36.8	1989.7	1414.10	75.31	24.69	.9346	.0010	.0041	.0630	.9387	.0007	.0041	.0630	.0005	
CB	50	5.59802	44.6	420.4	120.13	42.30	57.70	.8081	.0043	.0033	.0286	.8136	.0033	.0033	.0286	.1549	
	75	9.23099	45.7	601.2	228.41	53.83	46.17	.8105	.0027	.0015	.0080	.8120	.0016	.0015	.0080	.3138	
	100	10.67	43.7	763.5	203.12	51.79	48.21	.8218	.0037	.0010	.0056	.8228	.0043	.0010	.0056	.4298	
	150	19.887	39.4	1362	570.95	60.77	39.23	.8405	.0061	.0062	.0352	.8461	.0062	.0056	.0352	.2597	

¹ Number of outer iterations described in Algorithm 2.

² Fraction of time spent on solving master problems described in Algorithm 2.

³ Fraction of time spent on solving separation problems described in Algorithm 2.

⁴ Defined as $(AUC_{WRLR} - AUC_{LR}) / (1 - AUC_{LR})$.

649 **6. Concluding Remarks.** The computational results presented in this paper
 650 used a Wasserstein-robust framework for the logistic regression model, where all the
 651 decision variables are continuous, and the feasible set is assumed to be convex. The
 652 modified exchange algorithm of this paper is applicable to a broad class of nonlin-
 653 ear optimization problems. These algorithms can be implemented further within a
 654 branch-and-bound framework. Since the decomposition framework for solving the
 655 dual of (WRO) is suitable for possibly mixed-integer decision variables and model pa-
 656 rameters, this framework can also be adapted to model and solve (WRO) application
 657 problems such as those arising in scheduling, logistics and supply chain management.

658 **Acknowledgement.** We are grateful to Changhyeok Lee and Liwei Zeng for a
 659 valuable discussion. We also acknowledge NSF grants CMMI-1362003 and CMMI-
 660 1100868 that were used to support this research.

661 REFERENCES

- 662 [1] E. D. BARRIO, E. GINE, AND C. MATRAN, *Central limit theorems for the wasserstein distance*
 663 *between the empirical and the true distributions*, The Annals of Probability, 27 (1999),
 664 pp. 1009–1071.
 665 [2] A. BEN-TAL, L. E. GHAOUI, AND A. NEMIROVSKI, *Robust Optimization*, Princeton Series in
 666 Applied Mathematics, Princeton University Press, August 2009.
 667 [3] A. BEN-TAL, D. HERTOOG, A. D. WAEGENAERE, B. MELENBERG, AND G. RENNEN, *Robust*
 668 *solutions of optimization problems affected by uncertain probabilities*, Management Science,
 669 59 (2013), pp. 341–357.
 670 [4] A. BEN-TAL AND A. NEMIROVSKI, *Robust convex optimization*, Mathematics of Operations
 671 Research, 23 (1998), pp. 769–805.
 672 [5] A. BEN-TAL AND A. NEMIROVSKI, *Robust solutions to uncertain linear programs*, Operations
 673 Research Letters, 25 (1999), pp. 1–13.
 674 [6] A. BEN-TAL AND A. NEMIROVSKI, *Robust optimization - methodology and applications*, Math-
 675 ematical Programming, 92 (2002), pp. 453–480.
 676 [7] A. BEN-TAL, A. NEMIROVSKI, AND C. ROOS, *Robust solutions of uncertain quadratic and conic-*
 677 *quadratic problems*, SIAM J. OPTIM., 13 (2002), pp. 535–560.
 678 [8] D. P. BERTSEKAS, A. NEDIĆ, AND A. E. OZDAGLAR, *Convex Analysis and Optimization*, Athena
 679 Scientific Optimization and Computation, Athena Scientific, April 2003.
 680 [9] D. BERTSIMAS, X. V. DOAN, K. NATARAJAN, AND C. TEO, *Models for minimax stochastic linear*
 681 *optimization problems with risk aversion*, Math. Oper. Res., 35 (2010), pp. 580–602.
 682 [10] D. BERTSIMAS, V. GUPTA, AND N. KALLUS, *Data-driven robust optimization*, 2013. <https://arxiv.org/pdf/1401.0212v2.pdf>.
 683 [11] D. BERTSIMAS AND I. POPESCU, *Optimal inequalities in probability theory: A convex optimiza-*
 684 *tion approach*, SIAM J. OPTIM., 15 (2005), pp. 780–804.
 685 [12] J. R. BIRGE AND R. J.-B. WETS, *Computing bounds for stochastic programming problems*
 686 *by means of a generalized moment problem*, Mathematics of Operations Research, 12(1)
 687 (1987), pp. 149–162.
 688 [13] G. C. CALAFIORE, *Ambiguous risk measures and optimal robust portfolios*, SIAM J. OPTIM.,
 689 18 (2007), pp. 853–877.
 690 [14] G. C. CALAFIORE AND L. E. GHAOUI, *On distributionally robust chance-constrained linear*
 691 *programs*, J. Optim. Theory Appl., 130 (2006), pp. 1–22.
 692 [15] A. CHARNES, W. W. COOPER, AND K. O. KORTANEK, *Duality, haar programs and finite se-*
 693 *quence spaces*, Proceedings of the National Academy of Science, 48 (1962), pp. 783–786.
 694 [16] A. CHARNES, W. W. COOPER, AND K. O. KORTANEK, *Duality in semi-infinite programs and*
 695 *some works of haar and carathéodory*, Management Science, 9 (1963), pp. 209–228.
 696 [17] A. CHARNES, W. W. COOPER, AND K. O. KORTANEK, *On the theory of semi-infinite program-*
 697 *ming and a generalization of the kuhn-tucker saddle point theorem for arbitrary convex*
 698 *functions*, Naval Research Logistics Quarterly, 16 (1969), pp. 41–51.
 699 [18] E. DELAGE AND Y. YE, *Distributional robust optimization under moment uncertainty with*
 700 *application to data-driven problems*, Operations Research, 58 (2010).
 701 [19] J. DUPAČOVÁ, *The minimax approach to stochastic programming and an illustrative applica-*
 702 *tion*, Stochastics, 20 (1987), pp. 73–88.
 703 [20] E. ERDOĞAN AND G. IYENGAR, *Ambiguous chance constrained problems and robust optimiza-*
 704

- 705 tion, *Mathematical Programming*, 107 (2006), pp. 37–61.
- 706 [21] P. M. ESFAHAN AND D. KUHN, *Data-driven distributionally robust optimization using the*
707 *Wasserstein metric: performance guarantees and tractable reformulations*, 2015. <https://arxiv.org/pdf/1505.05116.pdf>.
- 708 [22] N. FOURNIER AND A. GUILLIN, *On the rate of convergence in wasserstein distance of the*
709 *empirical measure*, *Probab. Theory Relat. Fields*, 162 (2015), pp. 707–738.
- 710 [23] C. R. GIVENS AND R. M. SHORT, *A class of wasserstein metrics for probability distributions*,
711 *Michigan Math. J.*, 31 (1984), pp. 231–240.
- 712 [24] J. GOH AND M. SIM, *Distributionally robust optimization and its tractable approximations*,
713 *Oper. Res.*, 58 (2010), pp. 902–917.
- 714 [25] A. HAAR, *Über linear ungleichungen*, *Acta Mathematica Szeged*, 2 (1924).
- 715 [26] G. HANASUSANTO, V. ROITCH, D. KUHN, AND W. WIESEMANN, *A distributionally robust per-*
716 *spective on uncertainty quantification and chance constrained programming*, *Mathematical*
717 *Programming*, 151 (2015), pp. 35–62.
- 718 [27] R. HETTICH AND H. T. JONGEN, *On first and second order conditions for local optima for*
719 *optimization problems in finite dimensions*, *Methods Oper. Res.*, 23 (1977), pp. 82–97.
- 720 [28] R. HETTICH AND H. T. JONGEN, *Semi-infinite programming: conditions of optimality and ap-*
721 *plications*, in *Optimization Techniques*, *Lecture Notes in Control and Inform. Sci.*, J. Stoer,
722 ed., Springer-Verlag, Berlin, Heidelberg, New York, 1977, ch. 7, pp. 82–97.
- 723 [29] R. HETTICH AND K. O. KORTANEK, *Semi-infinite programming: theory, methods, and applica-*
724 *tions*, *SIAM REVIEW*, 35 (1993), pp. 380–429.
- 725 [30] R. HETTICH AND G. STILL, *Second order optimality conditions for generalized semi-infinite*
726 *programming problems*, *Optimization*, 34 (1995), pp. 195–211.
- 727 [31] P. JAILLET, J. QI, AND M. SIM, *Routing optimization under uncertainty*, *Oper. Res.*, 64 (2016),
728 pp. 186–200.
- 729 [32] R. JIANG AND Y. GUAN, *Risk-averse two-stage stochastic program with distributional ambiguity*,
730 2015. https://www.optimization-online.org/DB_FILE/2015/05/4908.pdf.
- 731 [33] R. JIANG AND Y. GUAN, *Data-driven chance constrained stochastic program*, *Mathematical*
732 *Programming*, 158 (2016), pp. 291–327.
- 733 [34] K. O. KORTANEK AND H. NO, *A central cutting plane algorithm for convex semi-infinite pro-*
734 *gramming problems*, *SIAM J. OPTM*, 3 (1993), pp. 901–918.
- 735 [35] J. B. LASSERRE, *Global optimization with polynomials and the problem of moments*, *SIAM J.*
736 *OPTIM.*, 11 (2001), pp. 796–817.
- 737 [36] C. LEE AND S. MEHROTRA, *A Distributionally-robust Approach for Finding Support*
738 *Vector Machines*, 2015. eprints for the optimization community, [http://www.](http://www.optimization-online.org/DB_HTML/2015/06/4965.html)
739 [optimization-online.org/DB_HTML/2015/06/4965.html](http://www.optimization-online.org/DB_HTML/2015/06/4965.html).
- 740 [37] S. LIAO, C. DELFT, AND J. P. VIAL, *Distributionally robust workforce scheduling in call centers*
741 *with uncertain arrival rates*, *Optim. Methods Softw.*, 28 (2013), pp. 501–522.
- 742 [38] Q. LIN, R. LOXTON, K. L. TEO, Y. H. WU, AND C. J. YU, *A new exact penalty method for*
743 *semi-infinite programming problems*, *J. Comp. Appl. Math.*, 261 (2014), pp. 271–286.
- 744 [39] D. Z. LONG AND J. QI, *Distributionally robust discrete optimization with entropic value-at-risk*,
745 *Oper. Res. Lett.*, 42 (2014), pp. 532–538.
- 746 [40] M. LÓPEZ AND G. STILL, *Semi-infinite programming*, *European Journal of Operational Re-*
747 *search*, 180 (2007), pp. 491–518.
- 748 [41] D. K. LOVE AND G. BAYRAKSAN, *Phi-divergence constrained ambiguous stochastic programs*
749 *for data-driven optimization*, 2016. [http://www.optimization-online.org/DB_HTML/](http://www.optimization-online.org/DB_HTML/2016/03/5350.html)
750 [2016/03/5350.html](http://www.optimization-online.org/DB_HTML/2016/03/5350.html).
- 751 [42] F. Q. LUO, C. LEE, AND S. MEHROTRA, *Application of distributionally robust optimization*
752 *with wasserstein metric on logistic regression and support vector machine*, tech. report.,
753 Northwestern University, 2017.
- 754 [43] S. MEHROTRA AND D. PAPP, *A cutting surface algorithm for semi-infinite convex programming*
755 *with an application to moment robust optimization*, *SIAM J. OPTM*, 24 (2015), pp. 1670–
756 1697.
- 757 [44] S. MEHROTRA AND H. ZHANG, *Models and algorithms for distributionally robust least squares*
758 *problems*, *Mathematical Programming*, 146 (2014), pp. 123–141.
- 759 [45] K. NATARAJAN, D. PACHAMANOVA, AND M. SIM, *Incorporating asymmetric distributional in-*
760 *formation in robust value-at-risk optimization*, *Manag. Sci.*, 54 (2008), pp. 573–585.
- 761 [46] G. NUERNBERGER, *Global unicity in optimization and approximation*, *Z. Angew. Math. Mech.*,
762 65 (1985), pp. T319–T321.
- 763 [47] G. NUERNBERGER, *Global unicity in semi-infinite optimization*, *Numer. Funct. Anal. Optic.*, 8
764 (1985), pp. 173–191.
- 765 [48] P. A. PARRILO, *Semidefinite programming relaxations for semialgebraic problems*, *Math. Pro-*
766

767 gram. Ser. B, 96 (2003), pp. 293–320.
 768 [49] G. C. PFLUG, A. PICHLER, AND D. WOZABAL, *Then 1/n investment strategy is optimal under*
 769 *high model ambiguity*, J. Bank. Financ., 36 (2012), pp. 410–417.
 770 [50] G. PFLUG AND D. WOZABAL, *Ambiguity in portfolio selection*, Quantitative Finance, 7 (2007),
 771 pp. 435–442.
 772 [51] K. POSTEK, D. HERTOOG, AND B. MELENBERG, *Computationally tractable counterparts of dis-*
 773 *tributionally robust constraints on risk measures*, SIAM REV., 58 (2016), pp. 603–650.
 774 [52] A. PRÉKOPA, *Stochastic Programming*, vol. 324 of Mathematics and Its Applications, Springer
 775 Netherlands, 1995.
 776 [53] R. REEMTSEN AND S. GÖRNER, *Numerical methods for semi-infinite programming: A survey*,
 777 in Semi-infinite Programming, Nonconvex Optimization and Its Applications, R. Reemtsen
 778 and J. J. Rückmann, eds., Kluwer Boston, Boston, 1998, ch. 25, pp. 195–275.
 779 [54] S. SHAFIEEZADEH-ABADEH, P. M. ESFAHANI, AND D. KUHN, *Distributional robust logistic re-*
 780 *gression*. arXiv:1509.09259, 2015.
 781 [55] A. SHAPIRO, *On duality theory of conic linear problems*, in Semi-Infinite Programming, vol. 57
 782 of Nonconvex Optimization and Its Applications, Springer US, 2001, pp. 135–165.
 783 [56] A. SHAPIRO AND S. AHMED, *On a class of minimax stochastic programs*, SIAM J. OPTIM., 14
 784 (2004), pp. 1237–1249.
 785 [57] A. SHAPIRO AND A. KLEYWEGT, *Minimax analysis of stochastic problems*, Optimization Meth-
 786 ods and Software, 17 (2002), pp. 523–542.
 787 [58] G. STILL, *Generalized semi-infinite programming: theory and methods*, European Journal of
 788 Operational Research, 119 (1999), pp. 301–313.
 789 [59] A. WÄCHTER AND L. T. BIEGLER, *On the implementation of a primal-dual interior point filter*
 790 *line search algorithm for large-scale nonlinear programming*, Mathematical Programming,
 791 106 (2006), pp. 25–57.
 792 [60] S. WANG AND Y. YUAN, *Feasible method for semi-infinite programs*, SIAM J. OPTM, 25 (2015),
 793 pp. 2537–2560.
 794 [61] Z. WANG, P. W. GLYNN, AND Y. YE, *Likelihood robust optimization for data-driven problems*,
 795 Comput. Manag. Sci., 13 (2016), pp. 241–261.
 796 [62] W. WIESEMANN, D. KUHN, AND M. SIM, *Distributional robust convex optimization*, Operations
 797 Research, 62 (2015), pp. 1358–1376.
 798 [63] D. WOZABAL, *A framework for optimization under ambiguity*, Annals of Operations Research,
 799 193 (2012), pp. 21–47.
 800 [64] H. XU, C. CARAMANIS, AND S. MANNOR, *A distributional interpretation of robust optimization*,
 801 Mathematics of Operations Research, 37 (2012), pp. 95–110.
 802 [65] M. XU, S. Y. WU, AND J. J. YE, *Solving semi-infinite programs by smoothing projected gradient*
 803 *method*, Comp. Optim. Appl., 59 (2014), pp. 591–616.
 804 [66] X. Q. YANG, Z. Y. CHEN, AND J. C. ZHOU, *Optimality conditions for semi-infinite and*
 805 *generalized semi-infinite programs via l_p exact penalty functions*, 2015. <http://www.mypolyuweb.hk/~mayangxq/2015YCZ.pdf>.
 806 [67] İ. YANIKOĞLU AND D. HERTOOG, *Safe approximations of ambiguous chance constraints using*
 807 *historical data*, INFORMS J. Comput., 25 (2013), pp. 666–681.
 808 [68] Y. ZHANG, Z. M. SHEN, AND S. SONG, *Distributionally robust optimization of two-stage lot-*
 809 *sizing problems*, Prod. Oper. Manag., 25 (2016), pp. 2116–2131.
 810 [69] M. H. ZWEIG AND G. CAMPBELL, *Receiver-operating characteristic (roc) plots: a fundamental*
 811 *evaluation tool in clinical medicine*, Clin Chem., 39 (1993), pp. 561–577.
 812 [70] S. ZYMLER, D. KUHN, AND B. RUSTEM, *Distributional robust joint chance constraints with*
 813 *second-order moment information*, Math. Programming, 137 (2013), pp. 167–198.
 814

815 **Appendix A. Proofs and supplements for Sections 3 and 4.**

816 **Proof of Proposition 3.3.**

817 *Proof.* By the definition of Wasserstein metric in (2), if two probability measures
 818 P_1 and P_2 satisfying $\mathcal{W}(P_1, P_2) < \varepsilon$, $\varepsilon > 0$, then there exists a $K \in \mathcal{M}(\Xi \times \Xi, \mathcal{F} \times \mathcal{F})$
 819 such that

820 (33)
$$K(A \times \Xi) = P_1(A), \quad K(\Xi \times A) = P_2(A), \quad \forall A \in \mathcal{F}$$

821 (34)
$$\int_{(s_1 \times s_2) \in \Xi \times \Xi} d(s_1, s_2) K(ds_1 \times ds_2) \leq \varepsilon.$$

 822

823 Therefore, we have

$$\begin{aligned}
& |f(P_1) - f(P_2)| = \left| \int_{s_1 \in \Xi} h(\theta, s_1) P_1(ds_1) - \int_{s_2 \in \Xi} h(\theta, s_2) P_2(ds_2) \right| \\
& = \left| \int_{s_1 \in \Xi} h(\theta, s_1) K(ds_1 \times \Xi) - \int_{s_2 \in \Xi} h(\theta, s_2) K(\Xi \times ds_2) \right| \\
824 \quad (35) \quad & = \left| \int_{(s_1 \times s_2) \in \Xi \times \Xi} h(\theta, s_1) K(ds_1 \times ds_2) - \int_{(s_1 \times s_2) \in \Xi \times \Xi} h(\theta, s_2) K(ds_1 \times ds_2) \right| \\
& \leq \int_{(s_1 \times s_2) \in \Xi \times \Xi} |h(\theta, s_1) - h(\theta, s_2)| K(ds_1 \times ds_2) \\
& \leq C(\theta) \int_{(s_1 \times s_2) \in \Xi \times \Xi} d(s_1, s_2) K(ds_1 \times ds_2) \\
& \leq C(\theta) \varepsilon,
\end{aligned}$$

825 which shows that f is continuous in $(\mathcal{M}(\Xi, \mathcal{F}), \mathcal{W})$. \square

826 Proof of Proposition 3.4.

827 *Proof.* We divide the proof into the following two steps: (a) Show that such a joint
828 probability measure K exists. Define a probability measure $K \in \mathcal{M}(\Xi \times \Xi, \mathcal{F} \times \mathcal{F})$ such
829 that $K(A \times B) := P(A \cap B)$ for all $A, B \in \mathcal{F}$, and $K(\cup_{i=1}^{\infty} A_i \times B_i) = \sum_{i=1}^{\infty} K(A_i \times B_i)$
830 for all countable collections $\{A_i \times B_i\}_{i=1}^{\infty}$ of pairwise disjoint sets in $\mathcal{F} \times \mathcal{F}$. It is
831 straightforward to verify that such K is as desired.

832 (b) Show that $\int_{(s_1 \times s_2) \in \Xi \times \Xi} d(s_1, s_2) K(ds_1 \times ds_2) = 0$. We prove by contradiction.
833 Assume $\int_{(s_1 \times s_2) \in \Xi \times \Xi} d(s_1, s_2) K(ds_1 \times ds_2) > \varepsilon$, for some $\varepsilon > 0$, then we have $K(A) >$
834 0 , where $A := \{(s_1 \times s_2) \in \Xi \times \Xi : d(s_1, s_2) > 0\}$. Let $A_n := \{(s_1 \times s_2) \in \Xi \times \Xi :$
835 $d(s_1, s_2) \geq 1/n\}$. Since $A = \cup_{n=1}^{\infty} A_n$, it follows that $\exists m$ such that $K(A_m) > 0$.
836 Since $A_m \in \mathcal{F} \times \mathcal{F}$, it can be expressed as: $A_m = \cup_{i=1}^{\infty} S_i^1 \times S_i^2$, where $S_i^1, S_i^2 \in \mathcal{F}$.
837 This implies that $\exists i$ such that $K(S_i^1 \times S_i^2) > 0$. By the definition of K , we have
838 $K(S_i^1 \times S_i^2) = P(S_i^1 \cap S_i^2) > 0$, hence $S_i^1 \cap S_i^2$ is nonempty, implying that $\exists s \in S_i^1 \cap S_i^2$
839 and hence $(s \times s) \in S_i^1 \times S_i^2 \subseteq A_m$. However, since $d(\cdot, \cdot)$ is a metric, we have $d(s, s) = 0$
840 which contradicts to $d(s, s) \geq 1/m$. \square

841 Duality theorem of conic linear programming.

842 THEOREM A.1 (Proposition 2.8(iii) from [55]). *Let V and W be linear spaces*
843 *(over real numbers), with V' and W' being their dual space respectively. Also let*
844 *$C \subseteq V$ and $K \subseteq W$ be convex cones. Let $A : V \rightarrow W$ be a linear mapping and*
845 *$A^* : W' \rightarrow V'$ be its adjoint mapping. Consider the conic linear optimization problem*
846 *of the following form:*

$$\begin{aligned}
847 \quad (36a) \quad & \text{(P)} \quad \min_{v \in C} \langle c, v \rangle \\
848 \quad (36b) \quad & \text{s.t.} \quad Av + b \in K.
\end{aligned}$$

850 Then the dual of (P) can be formulated as

$$\begin{aligned}
851 \quad (37a) \quad & \text{(D)} \quad \max_{w^* \in -K^*} \langle w^*, b \rangle \\
852 \quad (37b) \quad & \text{s.t.} \quad A^* w^* + c \in C^*,
\end{aligned}$$

854 where C^* and K^* are the dual cones of C and K respectively.

- 855 (a) *The weak duality holds, e.g., $\text{val}(\text{P}) \geq \text{val}(\text{D})$.*
 856 (b) *If $\text{val}(\text{P})$ is finite, Y is a finite dimensional Banach space, and the optimal*
 857 *solution set $\text{Sol}(\text{D})$ of (D) is nonempty and bounded, then $\text{val}(\text{P})=\text{val}(\text{D})$.*

858 **Convergence of the central cutting surface algorithm for convex semi-**
 859 **infinite programs.** We summarize the convergence of the central cutting surface
 860 algorithm from [43] for solving a general semi-infinite convex optimization problem of
 861 the form:

$$\begin{aligned}
 & \text{minimize } x_0 \\
 & \text{subject to } g(x, t) \leq 0 \quad \forall t \in T, \\
 & \quad \quad \quad x \in X,
 \end{aligned}
 \tag{38}$$

863 where the decision variable is x whose first coordinate (and also the objective) is
 864 denoted by x_0 . Assume the following conditions are satisfied (following the notation
 865 from [43]):

- 866 ASSUMPTION A.2. (1) *The set $X \subseteq \mathbb{R}^n$ is compact and convex.*
 867 (2) *There exists a Slater point \bar{x} and $\eta > 0$ satisfying $\bar{x} \in X$ and $g(\bar{x}, t) \leq -\eta$*
 868 *for every $t \in T$.*
 869 (3) *The function $f(\cdot)$ and $g(\cdot, t)$ are convex and sub-differentiable for every $t \in T$;*
 870 *moreover, these sub-differentials are uniformly bounded: there exists a $B > 0$*
 871 *such that for every $x \in X$ and $t \in T$, every subgradient $d \in \partial_x g(x, t)$ satisfies*
 872 *$\|d\| \leq B$.*
 873 (4) *For every point $x \in X$ that is not ε -feasible, there exists an oracle that can*
 874 *find in finite time a $t \in T$ satisfying $g(x, t) > 0$.*

875 The master problem of (38) at the k th iteration has the following form:

$$\begin{aligned}
 & \text{maximize } w \\
 & \text{subject to } x_0 + w \leq M^{(k-1)}, \\
 & \quad \quad \quad g(x, t) + Bw \leq 0 \quad \forall t \in Q^{(k-1)}, \\
 & \quad \quad \quad x \in X.
 \end{aligned}
 \tag{39}$$

877

878 THEOREM A.3. ([43] Theorems 2-4 and Theorem 8) *Let $(x^{(k)}, w^{(k)})$ be the solu-*
 879 *tion to the master problem at the k th iteration, and $\tilde{x}^{(k)}$ be the best know ε -feasible*
 880 *solution at the end of the k th iteration. The following statements hold:*

- 881 (a) *If Algorithm 1 terminates in the k th iteration, then $\tilde{x}^{(k-1)}$ is an ε -optimal*
 882 *solution to (38).*
 883 (b) *If Algorithm 1 does not terminate, then there exists an index \hat{k} such that the*
 884 *sequence $\{\tilde{x}^{(\hat{k}+i)}\}_{i=1}^{\infty}$ consists entirely of ε -feasible solutions.*
 885 (c) *If Algorithm 1 does not terminate, then the sequence $\{\tilde{x}^{(k)}\}_{i=1}^{\infty}$ has an ac-*
 886 *cumulation point, and each accumulation point is an ε -optimal solution to*
 887 *(38).*
 888 (d) *Algorithm 1 converges linearly in objective function value between consecutive*
 889 *feasible cuts after the first \hat{k} iterations, where \hat{k} satisfies $w^{(\hat{k})} < \eta/B$. Denote*
 890 *by x^* an optimal solution of (38), and let $\rho^{(k)} = \frac{\tilde{x}_0^{(k)} - x_0^*}{\tilde{x}_0^{(k-1)} - x_0^*}$, then we have*

$$\rho^{(k)} \leq 1 - \frac{\eta - Bw^{(k)}}{\eta + B(\bar{x}_0 - x_0^*)}
 \tag{40}$$

891