

# A data-driven approach to multi-stage stochastic linear optimization

Dimitris Bertsimas

Operations Research Center, Massachusetts Institute of Technology, dbertsim@mit.edu

Shimrit Shtern

The William Davidson Faculty of Industrial Engineering & Management, Technion - Israel Institute of Technology, shimrit@technion.ac.il

Bradley Sturt

Operations Research Center, Massachusetts Institute of Technology, bsturt@mit.edu

We propose a new data-driven approach for addressing multi-stage stochastic linear optimization problems with unknown distributions. The approach consists of solving a robust optimization problem that is constructed from sample paths of the underlying stochastic process. As more sample paths are obtained, we prove that the optimal cost of the robust problem converges to that of the underlying stochastic problem. To the best of our knowledge, this is the first data-driven approach for multi-stage stochastic linear optimization which is asymptotically optimal when uncertainty is arbitrarily correlated across time. Finally, we develop approximation algorithms for the proposed approach by extending techniques from the robust optimization literature, and demonstrate their practical value through numerical experiments on stylized data-driven inventory management problems.

*Key words:* Stochastic programming. Robust optimization. Sample-path approximations.

*History:* This paper was first submitted in November 2018. A revision was submitted in March 2020.

---

## 1. Introduction

In the traditional formulation of linear optimization, one makes a decision which minimizes a known objective function and satisfies a known set of constraints. Linear optimization has, by all measures, succeeded as a framework for modeling and solving numerous real world problems. However, in many practical applications, the objective function and constraints are unknown at the time of decision making. To incorporate uncertainty into the linear optimization framework, [Dantzig \(1955\)](#) proposed partitioning the decision variables across multiple stages, which are made sequentially as more uncertain parameters are revealed. This formulation is known today as multi-stage stochastic linear optimization, which has become an integral modeling paradigm in many applications (*e.g.*, supply chain management, energy planning, finance) and remains a focus of the stochastic optimization community ([Birge and Louveaux 2011](#), [Shapiro et al. 2009](#)).

In practice, decision makers increasingly have access to historical data which can provide valuable insight into future uncertainty. For example, consider a manufacturer which sells short lifecycle products. The manufacturer does not know a joint probability distribution of the demand over a new product’s lifecycle, but has access to historical demand trajectories over the lifecycle of similar products. Another example is energy planning, where operators must coordinate and commit to production levels throughout a day, the output of wind turbines is subject to uncertain weather conditions, and data on historical daily wind patterns is increasingly available. Other examples include portfolio management, where historical asset returns over time are available to investors, and transportation planning, where data comes in the form of historical ride usage of transit and ride sharing systems over the course of a day. Such historical data provides significant potential for operators to better understand how uncertainty unfolds through time, which can in turn be used for better planning.

When the underlying probability distribution is unknown, data-driven approaches to multi-stage stochastic linear optimization traditionally follow a two-step procedure. The historical data is first fit to a parametric model (*e.g.*, an autoregressive moving average process), and decisions are then obtained by solving a multi-stage stochastic linear optimization problem using the estimated distribution. The estimation step is considered essential, as techniques for solving multi-stage stochastic linear optimization (*e.g.*, scenario tree discretization) generally require knowledge of the correlation structure of uncertainty across time; see [Shapiro et al. \(2009, Section 5.8\)](#). A fundamental difficulty in this approach is choosing a parametric model which will accurately estimate the underlying correlation structure and lead to good decisions.

Nonparametric data-driven approaches to multi-stage stochastic linear optimization where uncertainty is correlated across time are surprisingly scarce. [Pflug and Pichler \(2016\)](#) propose a nonparametric estimate-then-optimize approach based on applying a kernel density estimator to the historical data, which enjoys asymptotic optimality guarantees under a variety of strong technical conditions. [Hanasusanto and Kuhn \(2013\)](#) present another nonparametric approach wherein the conditional distributions in stochastic dynamic programming are estimated using kernel regression. [Krokhmal and Uryasev \(2007\)](#) discuss nonparametric path-grouping heuristics for constructing scenario trees from historical data. In the case of multi-stage stochastic linear optimization, to the best of our knowledge, there are no previous nonparametric data-driven approaches which are asymptotically optimal in the presence of time-dependent correlations. Moreover, in the absence of additional assumptions on the estimated distribution or on the problem setting, multi-stage stochastic linear optimization problems are notorious for being computationally demanding.

The main contribution of this paper is a new data-driven approach for multi-stage stochastic linear optimization that is *asymptotically optimal*, even when uncertainty is arbitrarily correlated across time. In other words, we propose a data-driven approach for addressing multi-stage stochastic linear optimization with unknown distributions that (i) does not require any parametric modeling assumptions on the correlation structure of the underlying probability distribution, and (ii) converges to the underlying multi-stage stochastic linear optimization problem as the size of the dataset tends to infinity. Such an asymptotic optimality guarantee is of practical importance, as it ensures that the approach offers a near-optimal approximation of the underlying stochastic problem in the presence of big data.

Our approach for multi-stage stochastic linear optimization is based on robust optimization. Specifically, given sample paths of the underlying stochastic process, the proposed approach consists of constructing and solving a multi-stage robust linear optimization problem with multiple uncertainty sets. The main result of this paper (Theorem 1) establishes, under certain assumptions, that the optimal cost of this robust optimization problem converges nearly to that of the stochastic problem as the number of sample paths tends to infinity. While this robust optimization problem is computationally demanding to solve exactly, we provide evidence that it can be tractably approximated to reasonable accuracy by leveraging approximation techniques from the robust optimization literature. To the best of our knowledge, there was no similar work in the literature which addresses multi-stage stochastic linear optimization by solving a sequence of robust optimization problems.

The paper is organized as follows. Section 2 introduces multi-stage stochastic linear optimization in a data-driven setting. Section 3 presents the new data-driven approach to multi-stage stochastic linear optimization. Section 4 states the main asymptotic optimality guarantees. Section 5 presents two examples of approximation algorithms by leveraging techniques from robust optimization. Section 6 discusses implications of our asymptotic optimality guarantees in the context of Wasserstein-based distributionally robust optimization. Section 7 demonstrates the out-of-sample performance and computational tractability of the proposed methodologies in computational experiments. Section 8 offers concluding thoughts. All technical proofs are relegated to the attached appendices.

### 1.1. Related literature

Originating with Soyster (1973) and Ben-Tal and Nemirovski (1999), robust optimization has been widely studied as a general framework for decision-making under uncertainty, in which “optimal” decisions are those which perform best under the worst-case parameter realization from an “uncertainty set”. Beginning with the seminal work of Ben-Tal et al. (2004), robust optimization has

been viewed with particular success as a computationally tractable framework for addressing multi-stage problems. Indeed, by restricting the space of decision rules, a stream of literature showed that multi-stage robust linear optimization problems can be solved in polynomial time by using duality-based reformulations or cutting-plane methods. For a modern overview of decision-rule approximations, we refer the reader to [Delage and Iancu \(2015\)](#), [Georghiou et al. \(2018\)](#), [Ben-Tal et al. \(2009\)](#), [Bertsimas et al. \(2011a\)](#). A variety of non-decision rule approaches to solving multi-stage robust optimization have been proposed as well, such as [Zeng and Zhao \(2013\)](#), [Zhen et al. \(2018\)](#), [Xu and Burer \(2018\)](#), [Georghiou et al. \(2019\)](#).

Despite its computational tractability for multi-stage problems, a central critique of traditional robust optimization is that it does not aspire to find solutions which perform well on average. Several works have aimed to quantify the quality of solutions from multi-stage robust linear optimization from the perspective of multi-stage stochastic linear optimization ([Chen et al. 2007](#), [Bertsimas and Goyal 2010](#), [Bertsimas et al. 2011b](#)). By and large, it is fair to say that multi-stage robust linear optimization is viewed today as a distinct framework from multi-stage stochastic linear optimization, aiming to find solutions with good worst-case as opposed to good average performance.

Providing a potential tradeoff between the stochastic and robust frameworks, distributionally robust optimization has recently received significant attention. First proposed by [Scarf \(1958\)](#), distributionally robust optimization models the uncertain parameters with a probability distribution, but the distribution is presumed to be unknown and contained in an ambiguity set of distributions. Even though single-stage stochastic optimization is generally intractable, the introduction of ambiguity can surprisingly emit tractable reformulations ([Delage and Ye 2010](#), [Wiesemann et al. 2014](#)). Consequently, the extension of distributionally robust optimization to multi-stage decision making is an active area of research, including [Bertsimas et al. \(2019b\)](#) for multi-stage distributionally robust linear optimization with moment-based ambiguity sets.

There has been a proliferation of data-driven constructions of ambiguity sets which offer various probabilistic performance guarantees, including those based on the  $p$ -Wasserstein distance for  $p \in [1, \infty)$  ([Pflug and Wozabal 2007](#), [Mohajerin Esfahani and Kuhn 2018](#)), phi-divergences ([Ben-Tal et al. 2013](#), [Bayraksan and Love 2015](#), [Van Parys et al. 2017](#)), and statistical hypothesis tests ([Bertsimas et al. 2018](#)). Many of these data-driven approaches have since been applied to the particular case of two-stage distributionally robust linear optimization, including [Jiang and Guan \(2018\)](#) for phi-divergence and [Hanasusanto and Kuhn \(2018\)](#) for  $p$ -Wasserstein ambiguity sets when  $p \in [1, \infty)$ . To the best of our knowledge, no previous work has demonstrated whether such

distributionally robust approaches, if extended to solve multi-stage stochastic linear optimization (with three or more stages) directly from data, retain their asymptotic optimality guarantees.

In contrast to the above literature, our motivation for robust optimization in this paper is not to find solutions which perform well on the worst-case realization in an uncertainty set, are risk-averse, or have finite-sample probabilistic guarantees. Rather, our proposed approach to multi-stage stochastic linear optimization adds robustness to the historical data as a tool to avoid overfitting as the number of data points tends to infinity. In this spirit, our work is perhaps closest related to several papers in the context of machine learning (Xu et al. 2012, Shafieezadeh-Abadeh et al. 2019), which showed that adding robustness to historical data can be used to develop machine learning methods which have nonparametric performance guarantees when the solution space (of classification or regression models) is not finite-dimensional. To the best of our knowledge, this paper is the first to apply this use of robust optimization in the context of multi-stage stochastic linear optimization to achieve asymptotic optimality without restricting the space of decision rules.

As far as we are aware, our data-driven approach of averaging over multiple uncertainty sets is novel in the context of multi-stage stochastic linear optimization, and its asymptotic optimality guarantees do not follow from existing literature. Xu et al. (2012) considered averaging over multiple uncertainty sets to establish convergence guarantees for predictive machine learning methods, drawing connections with distributionally robust optimization and kernel density estimation. Their convergence results require that the objective function is continuous, the underlying distribution is continuous, and there are no constraints on the support. Absent strong assumptions on the problem setting and on the space of decision rules (which in general can be discontinuous), these properties do not hold in multi-stage problems. Erdoğan and Iyengar (2006) provide feasibility guarantees on robust constraints over unions of uncertainty sets with the goal of approximating ambiguous chance constraints using the Prohorov metric. Their probabilistic guarantees require that the constraint functions have a finite VC-dimension (Erdoğan and Iyengar 2006, Theorem 5), an assumption which does not hold in general for two- or multi-stage problems (Erdoğan and Iyengar 2007). In this paper, we instead establish general asymptotic optimality guarantees for the proposed data-driven approach for multi-stage stochastic linear optimization by developing new bounds for distributionally robust optimization with the 1-Wasserstein ambiguity set and connections with nonparametric support estimation (Devroye and Wise 1980).

Under a particular construction of the uncertainty sets, we show that the proposed data-driven approach to multi-stage stochastic linear optimization can also be interpreted as distributionally robust optimization using the  $\infty$ -Wasserstein ambiguity set (see Section 6). However, the asymptotic optimality guarantees in our paper do not make use of this interpretation, as there were surprisingly

few previous convergence results for this ambiguity set, even in single-stage settings. Indeed, when an underlying distribution is unbounded, the  $\infty$ -Wasserstein distance between an empirical distribution and true distribution is always infinite (Givens and Shortt 1984) and thus does not converge to zero as more data is obtained. Therefore, it is not possible to develop measure concentration guarantees for the  $\infty$ -Wasserstein distance (akin to those of Fournier and Guillin (2015)) which hold in general for light-tailed but unbounded probability distributions. Consequently, the proof techniques used by Mohajerin Esfahani and Kuhn (2018, Theorem 3.6) to establish convergence guarantees for the 1-Wasserstein ambiguity set do not appear to extend to the  $\infty$ -Wasserstein ambiguity set. As a byproduct of the results in this paper, we obtain asymptotic optimality guarantees for distributionally robust optimization with the  $\infty$ -Wasserstein ambiguity set under the same mild probabilistic assumptions as Mohajerin Esfahani and Kuhn (2018) for the first time.

## 1.2. Notation

We denote the real numbers by  $\mathbb{R}$ , the nonnegative real numbers by  $\mathbb{R}_+$ , and the integers by  $\mathbb{Z}$ . Lowercase and uppercase bold letters refer to vectors and matrices. We assume throughout that  $\|\cdot\|$  refers to an  $\ell_p$ -norm in  $\mathbb{R}^d$ , such as  $\|\mathbf{v}\|_1 = \sum_{i=1}^d |v_i|$  or  $\|\mathbf{v}\|_\infty = \max_{i \in [d]} |v_i|$ . We let  $\emptyset$  denote the empty set,  $\text{int}(\cdot)$  be the interior of a set, and  $[K]$  be shorthand for the set of consecutive integers  $\{1, \dots, K\}$ . Throughout the paper, we let  $\boldsymbol{\xi} := (\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T) \in \mathbb{R}^d$  denote a stochastic process with a joint probability distribution  $\mathbb{P}$ , and assume that  $\hat{\boldsymbol{\xi}}^1, \dots, \hat{\boldsymbol{\xi}}^N$  are independent and identically distributed (i.i.d.) samples from that distribution. Let  $\mathbb{P}^N := \mathbb{P} \times \dots \times \mathbb{P}$  denote the  $N$ -fold probability distribution over the historical data. We let  $S \subseteq \mathbb{R}^d$  denote the support of  $\mathbb{P}$ , that is, the smallest closed set where  $\mathbb{P}(\boldsymbol{\xi} \in S) = 1$ . The extended real numbers are defined as  $\bar{\mathbb{R}} := \mathbb{R} \cup \{-\infty, \infty\}$ , and we adopt the convention that  $\infty - \infty = \infty$ . The expectation of a measurable function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$  applied to the stochastic process is denoted by  $\mathbb{E}[f(\boldsymbol{\xi})] = \mathbb{E}_{\mathbb{P}}[f(\boldsymbol{\xi})] = \mathbb{E}_{\mathbb{P}}[\max\{f(\boldsymbol{\xi}), 0\}] - \mathbb{E}_{\mathbb{P}}[\max\{-f(\boldsymbol{\xi}), 0\}]$ . Finally, for any set  $\mathcal{Z} \subseteq \mathbb{R}^d$ , we let  $\mathcal{P}(\mathcal{Z})$  denote the set of all probability distributions on  $\mathbb{R}^d$  which satisfy  $\mathbb{Q}(\boldsymbol{\xi} \in \mathcal{Z}) \equiv \mathbb{E}_{\mathbb{Q}}[\mathbb{I}\{\boldsymbol{\xi} \in \mathcal{Z}\}] = 1$ .

## 2. Problem Setting

We consider multi-stage stochastic linear optimization problems with  $T \geq 1$  stages. The uncertain parameters observed over the time horizon are represented by a stochastic process  $\boldsymbol{\xi} := (\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T) \in \mathbb{R}^d$  with an underlying joint probability distribution, where  $\boldsymbol{\xi}_t \in \mathbb{R}^{d_t}$  is a random variable that is observed immediately after the decision in stage  $t$  is selected. We assume throughout that the random variables  $\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T$  may be correlated. A decision rule  $\mathbf{x} := (\mathbf{x}_1, \dots, \mathbf{x}_T)$  is a collection of policies which specify what decision to make in each stage based on the information observed

up to that point. More precisely, a policy in each stage is a measurable function of the form  $\mathbf{x}_t: \mathbb{R}^{d_1} \times \dots \times \mathbb{R}^{d_{t-1}} \rightarrow \mathbb{R}^{n_t - p_t} \times \mathbb{Z}^{p_t}$ . We use the shorthand notation  $\mathbf{x} \in \mathcal{X}$  to denote such decision rules.

In multi-stage stochastic linear optimization, our goal is to find a decision rule which minimizes a linear cost function in expectation while satisfying a system of linear inequalities almost surely. These problems are represented by

$$\begin{aligned} & \underset{\mathbf{x} \in \mathcal{X}}{\text{minimize}} && \mathbb{E} \left[ \sum_{t=1}^T \mathbf{c}_t(\boldsymbol{\xi}) \cdot \mathbf{x}_t(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}) \right] \\ & \text{subject to} && \sum_{t=1}^T \mathbf{A}_t(\boldsymbol{\xi}) \mathbf{x}_t(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}) \leq \mathbf{b}(\boldsymbol{\xi}) \quad \text{a.s.} \end{aligned} \tag{1}$$

Following standard convention, we assume that the problem parameters  $\mathbf{c}_1(\boldsymbol{\xi}) \in \mathbb{R}^{n_1}, \dots, \mathbf{c}_T(\boldsymbol{\xi}) \in \mathbb{R}^{n_T}$ ,  $\mathbf{A}_1(\boldsymbol{\xi}) \in \mathbb{R}^{m \times n_1}, \dots, \mathbf{A}_T(\boldsymbol{\xi}) \in \mathbb{R}^{m \times n_T}$ , and  $\mathbf{b}(\boldsymbol{\xi}) \in \mathbb{R}^m$  are affine functions of the stochastic process.

In this paper, we assume that the underlying joint probability distribution of the stochastic process is unknown. Instead, our information comes from historical data of the form

$$\hat{\boldsymbol{\xi}}^j \equiv (\hat{\boldsymbol{\xi}}_1^j, \dots, \hat{\boldsymbol{\xi}}_T^j), \quad j = 1, \dots, N.$$

We refer to each of these trajectories as a sample path of the stochastic process. This setting corresponds to many real-life applications. For example, consider managing the inventory of a new short lifecycle product, in which production decisions must be made over the product's lifecycle. In this case, each sample path represents the historical sales data observed over the lifecycle of a comparable product. Further examples are readily found in energy planning and finance, among many others. We assume that  $\hat{\boldsymbol{\xi}}^1, \dots, \hat{\boldsymbol{\xi}}^N$  are independent and identically distributed (i.i.d.) realizations of the stochastic process  $\boldsymbol{\xi} \equiv (\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T)$ . Our goal in this paper is a general-purpose, nonparametric sample-path approach for solving Problem (1) in practical computation times.

We will also assume that the support of the stochastic process is unknown. For example, in inventory management, an upper bound on the demand, if one exists, is generally unknown. On the other hand, we often have partial knowledge on the underlying support. For example, when the stochastic process captures the demand for a new product or the energy produced by a wind turbine, it is often the case that the uncertainty will be nonnegative. To allow any partial knowledge on the support to be incorporated, we assume knowledge of a convex superset  $\Xi \subseteq \mathbb{R}^d$  of the support of the underlying joint distribution, that is,  $\mathbb{P}(\boldsymbol{\xi} \in \Xi) = 1$ .

### 3. A Robust Approach to Multi-Stage Stochastic Linear Optimization

We now present the proposed data-driven approach, based on robust optimization, for solving multi-stage stochastic linear optimization. First, we construct an uncertainty set  $\mathcal{U}_N^j \subseteq \Xi$  around each sample path, consisting of realizations  $\zeta \equiv (\zeta_1, \dots, \zeta_T)$  which are slight perturbations of  $\hat{\xi}^j \equiv (\hat{\xi}_1^j, \dots, \hat{\xi}_T^j)$ . Then, we optimize for decision rules by averaging over the worst-case costs from each uncertainty set, and require that the decision rule is feasible for all realizations in all of the uncertainty sets. Formally, the proposed approach is the following:

$$\begin{aligned} & \underset{\mathbf{x} \in \mathcal{X}}{\text{minimize}} && \frac{1}{N} \sum_{j=1}^N \sup_{\zeta \in \mathcal{U}_N^j} \sum_{t=1}^T \mathbf{c}_t(\zeta) \cdot \mathbf{x}_t(\zeta_1, \dots, \zeta_{t-1}) \\ & \text{subject to} && \sum_{t=1}^T \mathbf{A}_t(\zeta) \mathbf{x}_t(\zeta_1, \dots, \zeta_{t-1}) \leq \mathbf{b}(\zeta) \quad \forall \zeta \in \cup_{j=1}^N \mathcal{U}_N^j. \end{aligned} \tag{2}$$

In contrast to traditional robust optimization, Problem (2) involves averaging over multiple uncertainty sets. Thus, the explicit goal here is to obtain decision rules which perform well on average while simultaneously not overfitting the historical data. We note that Problem (2) only requires that the decision rules are feasible for the realizations in the uncertainty sets. These feasibility requirements are justified when the overlapping uncertainty sets encompass the variability of future realizations of the uncertainty; see Section 4.

Out of the various possible constructions of the uncertainty sets, our investigation shall henceforth be focused on uncertainty sets constructed as balls of the form

$$\mathcal{U}_N^j := \left\{ \zeta \equiv (\zeta_1, \dots, \zeta_T) \in \Xi : \|\zeta - \hat{\xi}^j\| \leq \epsilon_N \right\},$$

where  $\epsilon_N \geq 0$  is a parameter which controls the size of the uncertainty sets. The parameter is indexed by  $N$  to allow for the size of the uncertainty sets to change as more data is obtained. The rationale for this particular uncertainty set is three-fold. First, it is conceptually simple, requiring only a single parameter to both estimate the expectation in the objective and the support of the distribution in the constraints. Second, under appropriate choice of the robustness parameter, we will show that Problem (2) with these uncertainty sets provides a near-optimal approximation of Problem (1) in the presence of big data (see Section 4). Finally, the uncertainty sets are of similar structure, which can be exploited to obtain tractable reformulations (see Section 5).

Our approach, in a nutshell, uses robust optimization as a tool for solving multi-stage stochastic linear optimization directly from data. More specifically, we obtain decision rules and estimate the optimal cost of Problem (1) by solving Problem (2). We refer the proposed data-driven approach for solving multi-stage stochastic linear optimization problems as *sample* or *sample-path* robust



optimization. As mentioned previously, the purpose of robustness is to ensure that resulting decision rules do not overfit the historical sample paths. To illustrate this role performed by robustness, we consider the following example.

EXAMPLE 1. Consider a supplier which aims to satisfy uncertain demand over two phases at minimal cost. The supplier selects an initial production quantity at \$1 per unit after observing preorders, and produces additional units at \$2 per unit after the regular orders are received. To determine the optimal production levels, we wish to solve

$$\begin{aligned} & \text{minimize} && \mathbb{E}[x_2(\xi_1) + 2x_3(\xi_1, \xi_2)] \\ & \text{subject to} && x_2(\xi_1) + x_3(\xi_1, \xi_2) \geq \xi_1 + \xi_2 \quad \text{a.s.} \\ & && x_2(\xi_1), x_3(\xi_1, \xi_2) \geq 0 \quad \text{a.s.} \end{aligned} \tag{3}$$

The output of the optimization problem are decision rules,  $x_2 : \mathbb{R} \rightarrow \mathbb{R}$  and  $x_3 : \mathbb{R}^2 \rightarrow \mathbb{R}$ , which specify what production levels to choose as a function of the demands observed up to that point. The joint probability distribution of the demand process  $(\xi_1, \xi_2) \in \mathbb{R}^2$  is unknown, and the supplier's knowledge comes from historical demand realizations of past products, denoted by  $(\hat{\xi}_1^1, \hat{\xi}_2^1), \dots, (\hat{\xi}_1^N, \hat{\xi}_2^N)$ . For the sake of illustration, suppose we attempted to approximate Problem (3) by choosing the decision rules which perform best when averaging over the historical data without any robustness. Such a sample average approach amounts to solving

$$\begin{aligned} & \text{minimize} && \frac{1}{N} \sum_{j=1}^N \left( x_2(\hat{\xi}_1^j) + 2x_3(\hat{\xi}_1^j, \hat{\xi}_2^j) \right) \\ & \text{subject to} && x_2(\hat{\xi}_1^j) + x_3(\hat{\xi}_1^j, \hat{\xi}_2^j) \geq \hat{\xi}_1^j + \hat{\xi}_2^j \quad \forall j \in [N] \\ & && x_2(\hat{\xi}_1^j), x_3(\hat{\xi}_1^j, \hat{\xi}_2^j) \geq 0 \quad \forall j \in [N]. \end{aligned}$$

Suppose that the random variable  $\xi_1$  for preorders has a continuous distribution. In that case, it immediately follows that  $\hat{\xi}_1^1 \neq \dots \neq \hat{\xi}_1^N$  almost surely, and thus an optimal decision rule for the above optimization problem is

$$x_2(\xi_1) = \begin{cases} \hat{\xi}_1^j + \hat{\xi}_2^j, & \text{if } \xi_1 = \hat{\xi}_1^j \text{ for } j \in [N], \\ 0, & \text{otherwise;} \end{cases} \quad x_3(\xi_1, \xi_2) = 0.$$

Unfortunately, these decision rules are *nonsensical* with respect to Problem (3). Indeed, the decision rules will not result in feasible decisions for the true stochastic problem with probability one. Moreover, the optimal cost of the above optimization problem will converge almost surely to  $\mathbb{E}[\xi_1 + \xi_2]$  as the number of sample paths  $N$  tends to infinity, which can in general be far from that of the stochastic problem. Clearly, such a sample average approach results in overfitting, even in big data settings, and thus provides an unsuitable approximation of Problem (3).  $\square$

The key takeaway from this paper is that this overfitting phenomenon in the above example is eliminated by adding robustness to the historical data. In particular, we will show in the following section that when the robustness parameter  $\epsilon_N$  is chosen appropriately, Problem (2) converges to a near-optimal approximation of Problem (1) as more data is obtained, without requiring any parametric assumptions on the stochastic process nor restrictions on the space of decision rules.

## 4. Asymptotic Optimality

In this section, we present theoretical guarantees showing that Problem (2) provides a near-optimal approximation of Problem (1) in the presence of big data. In Section 4.1, we describe the assumptions used in the subsequent convergence results. In Section 4.2, we present the main result of this paper (Theorem 1), which establishes asymptotic optimality of the proposed data-driven approach. In Section 4.3, we interpret Theorem 1 through several examples. In Section 4.4, we present asymptotic feasibility guarantees.

### 4.1. Assumptions

We begin by introducing our assumptions which will be used for establishing asymptotic optimality guarantees. First, we will assume that the joint probability distribution of the stochastic process satisfies the following light-tail assumption:

ASSUMPTION 1. *There exists a constant  $a > 1$  such that  $b := \mathbb{E}[\exp(\|\xi\|^a)] < \infty$ .*

For example, this assumption is satisfied if the stochastic process has a multivariate Gaussian distribution, and is not satisfied if the stochastic process has a multivariate exponential distribution. Importantly, Assumption 1 does not require any parametric assumptions on the correlation structure of the random variables across stages, and we do not assume that the coefficient  $a > 1$  is known.

Second, we will assume that the robustness parameter  $\epsilon_N$  is chosen to be strictly positive and decreases to zero as more data is obtained at the following rate:

ASSUMPTION 2. *There exists a constant  $\kappa > 0$  such that  $\epsilon_N := \kappa N^{-\frac{1}{\max\{3, d+1\}}}$ .*

In a nutshell, Assumption 2 provides a theoretical requirement on how to choose the robustness parameter to ensure that Problem (2) will not overfit the historical data (see Example 1 from Section 3). The rate also provides practical guidance on how the robustness parameter can be updated as more data is obtained. We note that, for many of the following results, the robustness parameter can decrease to zero at a faster rate; nonetheless, we shall impose Assumption 2 for all our results for simplicity.

Finally, our convergence guarantees for Problem (2) do not require any restrictions on the space of decision rules. Our analysis will only require the following assumption on the problem structure.

ASSUMPTION 3. *There exists a  $L \geq 0$  such that, for all  $N \in \mathbb{N}$ , the optimal cost of Problem (2) would not change if we added the following constraints:*

$$\sup_{\zeta \in \cup_{j=1}^N \mathcal{U}_N^j} \|\mathbf{x}_t(\zeta_1, \dots, \zeta_{t-1})\| \leq \sup_{\zeta \in \cup_{j=1}^N \mathcal{U}_N^j} L(1 + \|\zeta\|) \quad \forall t \in [T].$$

This assumption says that there always exists a near-optimal decision rule to Problem (2) where the decisions which result from realizations in uncertainty sets are bounded by the largest realization in the uncertainty sets. Moreover, this is a mild assumption that we find can be easily verified in many practical examples. In Appendix A, we show that every example presented in this paper satisfies this assumption.

## 4.2. Main result

We now present the main result of this paper (Theorem 1), which shows that the optimal cost of Problem (2) nearly converges to the optimal cost of Problem (1) as  $N \rightarrow \infty$ . For notational convenience, let  $J^*$  be the optimal cost of Problem (1),  $\hat{J}_N$  be the optimal cost of Problem (2), and  $S \subseteq \Xi$  be the support of the underlying joint probability distribution of the stochastic process.

Our main result presents tight asymptotic lower and upper bounds on the optimal cost  $\hat{J}_N$  of Problem (2). First, let  $\underline{J}$  be defined as the maximal optimal cost of any chance-constrained variant of the multi-stage stochastic linear optimization problem:

$$\begin{aligned} \underline{J} := & \lim_{\rho \downarrow 0} \underset{\mathbf{x} \in \mathcal{X}, \tilde{S} \subseteq \Xi}{\text{minimize}} \quad \mathbb{E} \left[ \sum_{t=1}^T \mathbf{c}_t(\boldsymbol{\xi}) \cdot \mathbf{x}_t(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}) \mathbb{I} \left\{ \boldsymbol{\xi} \in \tilde{S} \right\} \right] \\ & \text{subject to} \quad \sum_{t=1}^T \mathbf{A}_t(\zeta) \mathbf{x}_t(\zeta_1, \dots, \zeta_{t-1}) \leq \mathbf{b}(\zeta) \quad \forall \zeta \in \tilde{S} \\ & \quad \quad \quad \mathbb{P}(\boldsymbol{\xi} \in \tilde{S}) \geq 1 - \rho. \end{aligned}$$

We observe that the above limit must exist, as the optimal cost of the chance-constrained optimization problem is monotone in  $\rho$ . We also observe that  $\underline{J}$  is always a lower bound on  $J^*$ , since for every  $\rho > 0$ , adding the constraint  $\mathbb{P}(\boldsymbol{\xi} \in \tilde{S}) = 1$  to the above chance-constrained optimization problem would increase its optimal cost to  $J^*$ .<sup>1</sup>

<sup>1</sup>The definition does not preclude the possibility that  $\underline{J}$  is equal to  $-\infty$  or  $\infty$ . However, we do not expect either of those values to occur outside of pathological cases; see Section 4.3. The same remark applies to the upper bound  $\bar{J}$ .

Second, let  $\bar{J}$  be the optimal cost of the multi-stage stochastic linear optimization problem with an additional restriction that the decision rules are feasible on an expanded support:

$$\begin{aligned} \bar{J} := & \lim_{\rho \downarrow 0} \underset{\mathbf{x} \in \mathcal{X}}{\text{minimize}} \quad \bar{\mathbb{E}} \left[ \sum_{t=1}^T \mathbf{c}_t(\boldsymbol{\xi}) \cdot \mathbf{x}_t(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}) \right] \\ & \text{subject to} \quad \sum_{t=1}^T \mathbf{A}_t(\boldsymbol{\zeta}) \mathbf{x}_t(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_{t-1}) \leq \mathbf{b}(\boldsymbol{\zeta}) \quad \forall \boldsymbol{\zeta} \in \Xi: \text{dist}(\boldsymbol{\zeta}, S) \leq \rho. \end{aligned}$$

We remark that the limit as  $\rho$  tends down to zero must exist as well, since the optimal cost of the above optimization problem with expanded support is monotone in  $\rho$ . Note also that the expectation in the objective function has been replaced with  $\bar{\mathbb{E}}[\cdot]$ , which we define here as the local upper semicontinuous envelope of an expectation, *i.e.*,  $\bar{\mathbb{E}}[f(\boldsymbol{\xi})] := \lim_{\epsilon \rightarrow 0} \mathbb{E}[\sup_{\boldsymbol{\zeta} \in \Xi: \|\boldsymbol{\zeta} - \boldsymbol{\xi}\| \leq \epsilon} f(\boldsymbol{\xi})]$ . We similarly observe that  $\bar{J}$  is an upper bound on  $J^*$ , since the above optimization problem involves additional constraints and an upper envelope of the objective function.

Our main result is the following:

**THEOREM 1.** *Suppose Assumptions 1, 2, and 3 hold. Then,  $\mathbb{P}^\infty$ -almost surely we have*

$$J \leq \liminf_{N \rightarrow \infty} \hat{J}_N \leq \limsup_{N \rightarrow \infty} \hat{J}_N \leq \bar{J}.$$

*Proof.* See Appendix B.  $\square$

The above theorem provides assurance that the proposed data-driven approach becomes a near-optimal approximation of multi-stage stochastic linear optimization in the presence of big data. Note that Theorem 1 holds in very general cases; for example, it does not require boundedness on the decisions or random variables, requires no parametric assumptions on the correlations across stages, and holds when the decisions contain both continuous and integer components. Moreover, these asymptotic bounds for Problem (2) do not necessitate imposing any restrictions on the space of decision rules. To the best of our knowledge, such nonparametric asymptotic optimality guarantees for a sample-path approach to multi-stage stochastic linear optimization are the first of their kind when uncertainty is correlated across time.

Our proof of Theorem 1 is based on a new uniform convergence result (Theorem 2) which establishes a general relationship for arbitrary functions between the in-sample worst-case cost and the expected out-of-sample cost over the uncertainty sets. We state this theorem below due to its independent interest.

**THEOREM 2.** *If Assumptions 1 and 2 hold, then there exists a  $\bar{N} \in \mathbb{N}$ ,  $\mathbb{P}^\infty$ -almost surely, such that*

$$\mathbb{E} [f(\boldsymbol{\xi}) \mathbb{I} \{ \boldsymbol{\xi} \in \cup_{j=1}^N \mathcal{U}_N^j \}] \leq \frac{1}{N} \sum_{j=1}^N \sup_{\boldsymbol{\zeta} \in \mathcal{U}_N^j} f(\boldsymbol{\zeta}) + M_N \sup_{\boldsymbol{\zeta} \in \cup_{j=1}^N \mathcal{U}_N^j} |f(\boldsymbol{\zeta})|$$

for all  $N \geq \bar{N}$  and all measurable functions  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ , where  $M_N := N^{-\frac{1}{(d+1)(d+2)}} \log N$ .

*Proof.* See Appendix C.  $\square$

We note that our proofs of Theorems 1 and 2 also utilize a feasibility guarantee (Theorem 3) which can be found Section 4.4.

### 4.3. Examples where $\bar{J} - \underline{J}$ is zero or strictly positive

In general, we do not expect the gap between the lower and upper bound in Theorem 1 to be large. In fact, we will now show that the lower and upper bounds can be equal, *i.e.*,  $\underline{J} = \bar{J}$ , in which case the optimal cost of Problem (2) provably converges to the optimal cost of Problem (1). We provide such an example by revisiting the stochastic inventory management problem from Example 1. The following proposition, in combination with Theorem 1, shows that adding robustness to the historical data provably overcomes the overfitting phenomenon discussed in Section 3.

**PROPOSITION 1.** *For Problem (3),  $\underline{J} = J^*$ . If there is an optimal  $x_2^*: \mathbb{R} \rightarrow \mathbb{R}$  for Problem (3) which is continuous, then  $\bar{J} = J^*$ .*

The proof of this proposition, found in Appendix D, holds for any underlying probability distribution which satisfies Assumption 1.

While the above example shows that the lower and upper bounds may be equal, unless further restrictions are placed on the space of multi-stage stochastic linear optimization problems, they can have a nonzero gap. In the following, we present three examples that provide intuition on the situations in which this gap may be strictly positive. Our first example presents a problem in which the lower bound  $\underline{J}$  is equal to  $J^*$  but is strictly less than the upper bound  $\bar{J}$ .

**EXAMPLE 2.** Consider the single-stage stochastic problem

$$\begin{aligned} & \underset{x_1 \in \mathbb{Z}}{\text{minimize}} && x_1 \\ & \text{subject to} && x_1 \geq \xi_1 \quad \text{a.s.}, \end{aligned}$$

where the random variable  $\xi_1$  is governed by the probability distribution  $\mathbb{P}(\xi_1 > \alpha) = (1 - \alpha)^k$  for fixed  $k > 0$ , and  $\Xi = [0, 2]$ . We observe that the support of the random variable is  $S = [0, 1]$ , and thus the optimal cost of the stochastic problem is  $J^* = 1$ . We similarly observe that the lower bound is  $\underline{J} = 1$  and the upper bound, due to the integrality of the first stage decision, is  $\bar{J} = 2$ . If  $\epsilon_N = N^{-\frac{1}{3}}$ , then we prove in Appendix E that the bounds in Theorem 1 are tight under different choices of  $k$ :

Range of $k$	Result
$k \in (0, 3)$	$\mathbb{P}^\infty \left( J < \liminf_{N \rightarrow \infty} \widehat{J}_N = \limsup_{N \rightarrow \infty} \widehat{J}_N = \bar{J} \right) = 1$
$k = 3$	$\mathbb{P}^\infty \left( J = \liminf_{N \rightarrow \infty} \widehat{J}_N < \limsup_{N \rightarrow \infty} \widehat{J}_N = \bar{J} \right) = 1$
$k \in (3, \infty)$	$\mathbb{P}^\infty \left( J = \liminf_{N \rightarrow \infty} \widehat{J}_N = \limsup_{N \rightarrow \infty} \widehat{J}_N < \bar{J} \right) = 1$

This example shows that gaps can arise between the lower and upper bounds when mild changes in the support of the underlying probability distribution lead to significant changes in the optimal cost of Problem (1). Moreover, this example illustrates that each of the inequalities in Theorem 1 can hold with equality or strict inequality when the feasibility of decisions depends on random variables that have not yet been realized.  $\square$

Our second example presents a problem in which the upper bound  $\bar{J}$  is equal to  $J^*$  but is strictly greater than the lower bound  $J$ . This example deals with the special case in which any chance constrained version of a stochastic problem leads to a decision which is infeasible for the true stochastic problem.

EXAMPLE 3. Consider the single-stage stochastic problem

$$\begin{aligned} & \underset{x_1 \in \mathbb{R}^2}{\text{minimize}} && x_{12} \\ & \text{subject to} && \xi_1(1 - x_{12}) \leq x_{11} \quad \text{a.s.} \\ & && 0 \leq x_{12} \leq 1, \end{aligned}$$

where  $\xi_1 \sim \text{Gaussian}(0, 1)$  and  $\Xi = \mathbb{R}$ . The constraints are satisfied only if  $x_{12} = 1$ , and so the optimal cost of the stochastic problem is  $J^* = 1$ . Since there is no expectation in the objective and  $\Xi$  equals the true support, we also observe that  $\bar{J} = 1$ . However, we readily observe that there is always a feasible solution to the sample robust optimization problem (Problem (2)) where  $x_{12} = 0$ , and therefore  $J = \widehat{J}_N = 0$  for all  $N \in \mathbb{N}$ .  $\square$

Our third and final example demonstrates the necessity of the upper semicontinuous envelope  $\bar{\mathbb{E}}[\cdot]$  in the definition of the upper bound.

EXAMPLE 4. Consider the two-stage stochastic problem

$$\begin{aligned} & \underset{x_2: \mathbb{R} \rightarrow \mathbb{Z}}{\text{minimize}} && \mathbb{E}[x_2(\xi_1)] \\ & \text{subject to} && x_2(\xi_1) \geq \xi_1 \quad \text{a.s.}, \end{aligned}$$

where  $\theta \sim \text{Bernoulli}(0.5)$  and  $\psi \sim \text{Uniform}(0, 1)$  are independent random variables,  $\xi_1 = \theta\psi$ , and  $\Xi = [0, 1]$ . An optimal decision rule  $x_2^*: \mathbb{R} \rightarrow \mathbb{Z}$  to the stochastic problem is given by  $x_2^*(\xi_1) = 0$  for all

$\xi_1 \leq 0$  and  $x_2^*(\xi_1) = 1$  for all  $\xi_1 > 0$ , which implies that  $J^* = \frac{1}{2}$ . It follows from similar reasoning that  $\underline{J} = \frac{1}{2}$ . Since  $\Xi$  equals the support of the random variable, the only difference between the stochastic problem and the upper bound is that the latter optimizes over the local upper semicontinuous envelope, and we observe that  $\lim_{N \rightarrow \infty} \widehat{J}_N = \bar{J} = \bar{\mathbb{E}}[x_2^*(\xi_1)] = 1$ .  $\square$

In each of the above examples, we observe that the bounds in Theorem 1 are tight, in the sense that the optimal cost of Problem (2) converges either to the lower bound or the upper bound. This provides some indication that the bounds in Theorem 1 offer an accurate depiction of how Problem (2) can behave in the asymptotic regime. On the other hand, the above examples which illustrate a nonzero gap seem to require intricate construction, and future work may identify (sub-classes) of Problem (1) where the equality of the bounds can be ensured.

#### 4.4. Feasibility guarantees

We conclude Section 4 by discussing out-of-sample feasibility guarantees for decision rules obtained from Problem (2). Recall that Problem (2) finds decision rules which are feasible for each realization in the uncertainty sets. However, one cannot guarantee that these decision rules will be feasible for realizations outside of the uncertainty sets. Thus, a pertinent question is whether a decision rule obtained from approximately solving Problem (2) is feasible with high probability. To address the question of feasibility, we leverage classic results from detection theory.

Let  $S_N := \cup_{j=1}^N \mathcal{U}_N^j$  be shorthand for the union of the uncertainty sets. We say that a decision rule is  $S_N$ -feasible if

$$\sum_{t=1}^T \mathbf{A}_t(\zeta) \mathbf{x}_t(\zeta_1, \dots, \zeta_{t-1}) \leq \mathbf{b}(\zeta) \quad \forall \zeta \in S_N.$$

In other words, the set of feasible decision rules to Problem (2) are exactly those which are  $S_N$ -feasible. Our subsequent analysis utilizes the following (seemingly tautological) observation: for any decision rule that is  $S_N$ -feasible,

$$\mathbb{P} \left( \sum_{t=1}^T \mathbf{A}_t(\boldsymbol{\xi}) \mathbf{x}_t(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}) \leq \mathbf{b}(\boldsymbol{\xi}) \right) \geq \mathbb{P}(\boldsymbol{\xi} \in S_N),$$

where  $\mathbb{P}(\boldsymbol{\xi} \in S_N)$  is shorthand for  $\mathbb{P}(\boldsymbol{\xi} \in S_N \mid \hat{\boldsymbol{\xi}}^1, \dots, \hat{\boldsymbol{\xi}}^N)$ . Indeed, this inequality follows from the fact that a decision rule which is  $S_N$ -feasible is definitionally feasible for all realizations  $\zeta \in S_N$ , and thus the probability of feasibility is at least the probability that  $\boldsymbol{\xi} \in S_N$ .

We have thus transformed the analysis of feasible decision rules for Problem (2) to the problem of analyzing the performance of  $S_N$  as an estimate for the support  $S$  of a stochastic process. Interestingly, this nonparametric estimator for the support of a joint probability distribution has

been widely studied in the statistics literature, with perhaps the earliest results coming from Devroye and Wise (1980) in detection theory. Since then, the performance of  $S_N$  as a nonparametric estimate of  $S$  has been studied with applications in cluster analysis and image recognition (Korostelev and Tsybakov 1993, Schölkopf et al. 2001). Leveraging this connection between stochastic optimization and support estimation, we obtain the following guarantee on feasibility.

**THEOREM 3.** *Suppose Assumptions 1 and 2 hold. Then,  $\mathbb{P}^\infty$ -almost surely we have*

$$\lim_{N \rightarrow \infty} \left( \frac{N^{\frac{1}{d+1}}}{(\log N)^{d+1}} \right) \mathbb{P}(\boldsymbol{\xi} \notin S_N) = 0.$$

*Proof.* See Appendix F.  $\square$

Intuitively speaking, Theorem 3 provides a guarantee that *any* feasible decision rule to Problem (2) will be feasible with high probability on future data when the number of sample paths is large. To illustrate why robustness is indeed necessary to achieve such feasibility guarantees, we recall from Example 1 that decision rules may prohibitively overfit the data and be infeasible with probability one if the robustness parameter  $\epsilon_N$  is set to zero.

## 5. Approximation Techniques

In the previous section, we developed theoretical guarantees which demonstrated that Problem (2) provides a good approximation of multi-stage stochastic linear optimization when the number of sample paths is large. In this section, we demonstrate that Problem (2) can be addressed using approximation techniques from the field of robust optimization. Specifically, we show that two decision-rule approximation schemes from robust optimization, *linear decision rules* and *finite adaptability*, can be extended to obtain approximations of Problem (2). In particular, we present a novel duality argument (Theorem 4) which allows the computational cost of these techniques to scale efficiently in the number of sample paths. The computational tractability and out-of-sample performance of these approximation schemes is illustrated via numerical experiments in Section 7.

### 5.1. Linear decision rules

Generally speaking, multi-stage optimization problems are computationally demanding due to optimizing over an unrestricted space of decision rules. To overcome this challenge, a common approximation technique in robust optimization is to restrict the space of decision rules to a space which can more easily be optimized. As described in Section 1.1, the success of robust optimization as a modeling framework for addressing real-world multi-stage problems is often attributed the computational tractability of such decision rule approximations. This section extends one such



decision rule scheme, known as linear decision rules, to approximately solve Problem (2) and illustrates its computational tractability in big data settings.

Specifically, we consider approximating Problem (2) by restricting its decision rules to those of the form

$$\mathbf{x}_t(\zeta_1, \dots, \zeta_{t-1}) = \mathbf{x}_{t,0} + \sum_{s=1}^{t-1} \mathbf{X}_{t,s} \zeta_s.$$

Thus, rather than optimizing over the space of all possible decision rules (functions), we instead optimize over a finite collection of decision variables which parameterize a linear decision rule. For the setting where  $\mathbf{c}_t(\boldsymbol{\xi})$  and  $\mathbf{A}_t(\boldsymbol{\xi})$  do not depend on the uncertain parameters and all decision variables are continuous, the resulting linear decision rule approximation of Problem (2) is given by

$$\begin{aligned} & \text{minimize} && \frac{1}{N} \sum_{j=1}^N \sup_{\zeta \in \mathcal{U}_N^j} \sum_{t=1}^T \mathbf{c}_t \cdot \left( \mathbf{x}_{t,0} + \sum_{s=1}^{t-1} \mathbf{X}_{t,s} \zeta_s \right) \\ & \text{subject to} && \sum_{t=1}^T \mathbf{A}_t \left( \mathbf{x}_{t,0} + \sum_{s=1}^{t-1} \mathbf{X}_{t,s} \zeta_s \right) \leq \mathbf{b}(\zeta) \quad \forall \zeta \in \cup_{j=1}^N \mathcal{U}_N^j, \end{aligned} \tag{4}$$

where the decision variables are  $\mathbf{x}_{t,0} \in \mathbb{R}^{n_t}$  and  $\mathbf{X}_{t,s} \in \mathbb{R}^{n_t \times d_s}$  for all  $1 \leq s < t \leq T$  and the affine function  $\mathbf{b}(\zeta) \in \mathbb{R}^m$  is shorthand for  $\mathbf{b}^0 + \sum_{t=1}^T \mathbf{B}_t \zeta_t$ .

Much like linear decision rules in robust optimization, we observe that Problem (4), when feasible, always produces a feasible decision rule for Problem (2) and an upper bound on its optimal cost. Nonetheless, Problem (4) has semi-infinite constraints, which must be eliminated in order for the optimization problem to be solvable by off-the-shelf solvers. A standard technique from robust optimization for eliminating semi-infinite constraints is to introduce (dual) auxiliary decision variables and constraints for each uncertainty set. Importantly, for Problem (4) to be practically tractable in the presence of big data, the size of an equivalent finite-dimensional optimization problem must scale efficiently in the number of sample paths.

We now show that Problem (4) can be reformulated as a linear optimization with size that scales linearly in the number of sample paths (Theorem 4). The central idea enabling the following reformulation is that the worst-case realizations over the various uncertainty sets are found by optimizing over identical linear functions. Thus, when constructing the robust counterparts for each uncertainty set, we can combine the dual auxiliary decision variables from different uncertainty sets, resulting in a reformulation where the number of auxiliary decision variables is independent of the number of sample paths. To illustrate this reformulation technique, we focus on uncertainty sets which satisfy the following construction:

**ASSUMPTION 4.** *The uncertainty sets have the form  $\mathcal{U}_N^j := \{\zeta \in \mathbb{R}^d : \boldsymbol{\ell}^j \leq \zeta \leq \mathbf{u}^j\}$ .*

For example, Assumption 4 holds if we choose the set  $\Xi$  to be  $\mathbb{R}_+^d$  and use the  $\|\cdot\|_\infty$  norm in the uncertainty sets from Section 3. The following illustrates the novel duality technique described above:

**THEOREM 4.** *If Assumption 4 holds, then Problem (4) can be reformulated as a linear optimization problem with  $O(md)$  auxiliary decision variables and  $O(md + mN)$  linear constraints.*

*Proof.* By introducing epigraph variables  $v_1, \dots, v_N \in \mathbb{R}$ , the constraints in Problem (4) can be rewritten as

$$\begin{aligned} \sum_{t=1}^T \left( \sum_{s=t+1}^T \mathbf{X}_{s,t}^\top \mathbf{c}_s \right) \cdot \boldsymbol{\zeta}_t &\leq v_j - \sum_{t=1}^T \mathbf{c}_t \cdot \mathbf{x}_{t,0} & \forall \boldsymbol{\zeta} \in \mathcal{U}_N^j, j \in \{1, \dots, N\}, \\ \sum_{t=1}^T \left( -\mathbf{B}_t + \sum_{s=t+1}^T \mathbf{A}_s \mathbf{X}_{t,s} \right) \boldsymbol{\zeta}_t &\leq \mathbf{b}^0 - \sum_{t=1}^T \mathbf{A}_t \mathbf{x}_{t,0} & \forall \boldsymbol{\zeta} \in \mathcal{U}_N^j, j \in \{1, \dots, N\}. \end{aligned}$$

We will now reformulate each of these semi-infinite constraints by introducing auxiliary variables. First, we observe that each of the above semi-infinite constraints can be rewritten as

$$\max_{\boldsymbol{\zeta} \in \mathcal{U}_N^j} \sum_{t=1}^T \mathbf{d}_t \cdot \boldsymbol{\zeta}_t \leq \gamma$$

for some vector  $\mathbf{d} := (\mathbf{d}_1, \dots, \mathbf{d}_T) \in \mathbb{R}^d$  and scalar  $\gamma \in \mathbb{R}$ . Moreover, it follows from strong duality for linear optimization that

$$\max_{\boldsymbol{\zeta} \in \mathcal{U}_N^j} \sum_{t=1}^T \mathbf{d}_t \cdot \boldsymbol{\zeta}_t = \begin{cases} \text{minimize} & \sum_{t=1}^T (\mathbf{u}_t^j \cdot \boldsymbol{\mu}_t - \boldsymbol{\ell}_t^j \cdot \boldsymbol{\lambda}_t) \\ \text{subject to} & \boldsymbol{\mu}_t - \boldsymbol{\lambda}_t = \mathbf{d}_t \quad \forall t \in [T], \end{cases}$$

where  $\mathbf{u}^j := (\mathbf{u}_1^j, \dots, \mathbf{u}_T^j) \in \mathbb{R}^d$  and  $\boldsymbol{\ell}^j := (\boldsymbol{\ell}_1^j, \dots, \boldsymbol{\ell}_T^j) \in \mathbb{R}^d$  are the upper and lower bounds which define the uncertainty set. We readily observe that the solutions  $\boldsymbol{\mu}_t = [\mathbf{d}_t]_+$  and  $\boldsymbol{\lambda}_t = [-\mathbf{d}_t]_+$  are optimal for the above optimization problem. Importantly, these optimal solutions to the dual problem are *independent* of the index  $j$ . Thus, the semi-infinite constraints in the epigraph formulation of Problem (4) are satisfied if and only if there exists  $\boldsymbol{\alpha} := (\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_T) \in \mathbb{R}_+^d$  and  $\boldsymbol{\beta} := (\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_T) \in \mathbb{R}_+^d$  which satisfy

$$\begin{aligned} \sum_{t=1}^T (\boldsymbol{\alpha}_t \cdot \mathbf{u}_t^j - \boldsymbol{\beta}_t \cdot \boldsymbol{\ell}_t^j + \mathbf{c}_t \cdot \mathbf{x}_{t,0}) &\leq v_j \quad \forall j \in [N] \\ \boldsymbol{\alpha}_t - \boldsymbol{\beta}_t &= \sum_{s=t+1}^T \mathbf{X}_{s,t}^\top \mathbf{c}_s \quad \forall t \in [T] \end{aligned}$$

and there exists  $\mathbf{M} := (\mathbf{M}_1, \dots, \mathbf{M}_T) \in \mathbb{R}_+^{m \times d}$  and  $\boldsymbol{\Lambda} := (\boldsymbol{\Lambda}_1, \dots, \boldsymbol{\Lambda}_T) \in \mathbb{R}_+^{m \times d}$  which satisfy

$$\begin{aligned} \sum_{t=1}^T (\mathbf{M}_t \mathbf{u}_t^j - \boldsymbol{\Lambda}_t \boldsymbol{\ell}_t^j + \mathbf{A}_t \mathbf{x}_{t,0}) &\leq \mathbf{b}^0 \quad \forall j \in [N] \\ \mathbf{M}_t - \boldsymbol{\Lambda}_t &= -\mathbf{B}_t + \sum_{s=t+1}^T \mathbf{A}_s \mathbf{X}_{t,s} \quad \forall t \in [T] \end{aligned}$$

Removing the epigraph decision variables, the resulting reformulation of Problem (4) is

$$\begin{aligned}
& \text{minimize} && \frac{1}{N} \sum_{j=1}^N \sum_{t=1}^T (\boldsymbol{\alpha}_t \cdot \mathbf{u}_t^j - \boldsymbol{\beta}_t \cdot \boldsymbol{\ell}_t^j + \mathbf{c}_t \cdot \mathbf{x}_{t,0}) \\
& \text{subject to} && \boldsymbol{\alpha}_t - \boldsymbol{\beta}_t = \sum_{s=t+1}^T \mathbf{X}_{s,t}^\top \mathbf{c}_s && t \in [T] \\
& && \sum_{t=1}^T (\mathbf{M}_t \mathbf{u}_t^j - \boldsymbol{\Lambda}_t \boldsymbol{\ell}_t^j + \mathbf{A}_t \mathbf{x}_{t,0}) \leq \mathbf{b}^0 && j \in [N] \\
& && \mathbf{M}_t - \boldsymbol{\Lambda}_t = -\mathbf{B}_t + \sum_{s=t+1}^T \mathbf{A}_s \mathbf{X}_{t,s} && t \in [T]
\end{aligned}$$

where the auxiliary decision variables are  $\boldsymbol{\alpha} \equiv (\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_T), \boldsymbol{\beta} \equiv (\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_T) \in \mathbb{R}_+^d$  and  $\mathbf{M} \equiv (\mathbf{M}_1, \dots, \mathbf{M}_T), \boldsymbol{\Lambda} \equiv (\boldsymbol{\Lambda}_1, \dots, \boldsymbol{\Lambda}_T) \in \mathbb{R}_+^{m \times d}$ . Thus, the reformulation technique allowed us to decrease the number of auxiliary decision variables from  $O(Nmd)$  to  $O(md)$ .  $\square$

While linear decision rules can sometimes provide a near-optimal approximation of Problem (2) (see Section 7.2), we do not expect this to be the case in general. Indeed, we recall from Section 4 that Problem (2) can provide a near-optimal approximation of Problem (1), and it has been known from the early literature that linear decision rules generally provide a poor approximation for multi-stage stochastic linear optimization; see, *e.g.*, Garstka and Wets (1974, Section 6). Nonetheless, we can obtain tighter approximations of Problem (2) by selecting a richer space of decision rules, an abundance of which can be found in the robust optimization literature (see Section 5.2 for an example). Moreover, Problem (2) is also amenable to new approximation schemes which exploit its particular structure; we refer to our companion paper Bertsimas et al. (2019a) for such an approximation algorithm for two-stage problems. In all cases, and as a result of the convergence guarantees from Section 4, Problem (2) offers an opportunity to extend algorithmic advances from robust optimization to obtain approximations of multi-stage stochastic linear optimization.

## 5.2. Finite adaptability

In this section, we show how to extend the decision rule approximation scheme of finite adaptability from robust optimization (Bertsimas and Caramanis 2010) to obtain tighter approximations of Problem (2). Specifically, finite adaptability partitions the set  $\Xi$  into smaller regions, and then optimizes a separate static or linear decision rule in each region. The approach of finite adaptability extends to problems with integer decision variables, and the practitioner can trade off the tightness of their approximations with an increase in computational cost. We show that the duality techniques from the previous section (Theorem 4) readily extend to this richer class of decision rules, and the practical performance of this approach in a stylized example is demonstrated in Section 7.1.

We begin by describing the approximation scheme of finite adaptability from robust optimization. In finite adaptability, one partitions the uncertainty set into different regions, and optimizes a separate linear decision rule for each region. Let  $P^1, \dots, P^K \subseteq \mathbb{R}^d$  be regions which form a partition of  $\Xi \subseteq \mathbb{R}^d$ . For each stage  $t$ , let  $P_t^k \subseteq \mathbb{R}^{d_1 + \dots + d_t}$  be the projection of the region  $P^k$  onto the first  $t$  stages. Then, we consider approximating Problem (2) by restricting its decision rules to those of the form

$$\mathbf{x}_t(\zeta_1, \dots, \zeta_{t-1}) = \begin{cases} \mathbf{x}_{t,0}^1 + \sum_{s=1}^{t-1} \mathbf{X}_{t,s}^1 \zeta_s, & \text{if } (\zeta_1, \dots, \zeta_{t-1}) \in P_{t-1}^1, \\ \vdots \\ \mathbf{x}_{t,0}^K + \sum_{s=1}^{t-1} \mathbf{X}_{t,s}^K \zeta_s, & \text{if } (\zeta_1, \dots, \zeta_{t-1}) \in P_{t-1}^K. \end{cases}$$

In contrast to a single linear decision rule, finite adaptability allows for greater degrees of freedom at a greater computational cost. Indeed, for each region  $P^k$ , we choose a separate linear decision rule which is locally optimal for that region. To accommodate integer decision variables, we restrict the corresponding component of each  $\mathbf{x}_{t,0}^k$  to be integer and restrict the associated rows of each matrix  $\mathbf{X}_{t,s}^k$  to be zero.

A complication of finite adaptability is that we may not have enough information at any intermediary stage to determine which region  $P^k$  will contain the entire trajectory. In other words, at the start of stage  $t$ , a decision must be chosen after only observing the values of  $(\zeta_1, \dots, \zeta_{t-1})$ , and there may be two or more regions of the partition for which their projections  $P_{t-1}^k$  and  $P_{t-1}^{k'}$  are overlapping. Fortunately, the following proposition shows that the aforementioned complication caused by overlapping projections can be resolved by adding constraints of the form  $\mathbf{x}_t^k = \mathbf{x}_t^{k'}$  and  $\mathbf{X}_{t,s}^k = \mathbf{X}_{t,s}^{k'}$  for every  $1 \leq s < t$  when the regions  $P^k$  and  $P^{k'}$  are indistinguishable at stage  $t$ .

**PROPOSITION 2 (Proposition 4, Bertsimas and Dunning (2016)).** *If there exists  $\zeta \equiv (\zeta_1, \dots, \zeta_T) \in P^k$  and  $\zeta' \equiv (\zeta'_1, \dots, \zeta'_T) \in P^{k'}$  such that  $(\zeta_1, \dots, \zeta_{t-1}) = (\zeta'_1, \dots, \zeta'_{t-1})$ , and  $\zeta \in \text{int}(P^{k'})$  or  $\zeta' \in \text{int}(P^k)$  hold, then we must enforce the constraints that  $\mathbf{x}_{t,0}^k = \mathbf{x}_{t,0}^{k'}$  and  $\mathbf{X}_{t,s}^k = \mathbf{X}_{t,s}^{k'}$  for all  $1 \leq s < t$  at stage  $t$  as the two regions cannot be distinguished with the uncertain parameters realized by that stage. Otherwise, we do not need to enforce any constraints at stage  $t$  for this pair.*

For brevity, we let  $\mathcal{T}(P^1, \dots, P^K)$  denote the collection of tuples  $(k, k', t)$  for which  $P^k$  and  $P^{k'}$  cannot be distinguished at stage  $t$ , which we assume can be tractably computed.

We now extend the approach of finite adaptability to Problem (2). Let  $P^1, \dots, P^K$  be a given partition of  $\Xi$ , and let the intersections between regions of the partition and uncertainty sets be denoted by  $\mathcal{K}^j := \{k \in [K] : \mathcal{U}_N^j \cap P^k \neq \emptyset\}$ . For the setting where  $\mathbf{c}_t(\xi)$  and  $\mathbf{A}_t(\xi)$  do not depend on

the uncertain parameters, the resulting linear decision rule approximation of Problem (2) is given by

$$\begin{aligned}
& \text{minimize} && \frac{1}{N} \sum_{j=1}^N \max_{k \in \mathcal{K}^j} \max_{\zeta \in \mathcal{U}_N^j \cap P^k} \sum_{t=1}^T \mathbf{c}_t \cdot \left( \mathbf{x}_{t,0}^k + \sum_{s=1}^{t-1} \mathbf{X}_{t,s}^k \zeta_s \right) \\
& \text{subject to} && \sum_{t=1}^T \mathbf{A}_t \left( \mathbf{x}_{t,0}^k + \sum_{s=1}^{t-1} \mathbf{X}_{t,s}^k \zeta_s \right) \leq \mathbf{b}(\zeta) \quad \forall \zeta \in \cup_{j=1}^N \mathcal{U}_N^j \cap P^k, k \in [K] \\
& && \mathbf{x}_t^k = \mathbf{x}_t^{k'}, \mathbf{X}_{t,s}^k = \mathbf{X}_{t,s}^{k'} \quad \forall (k, k', t) \in \mathcal{T}(P^1, \dots, P^K), 1 \leq s < t.
\end{aligned} \tag{5}$$

where the decision variables are  $\mathbf{x}_{t,0}^k \in \mathbb{R}^{n_t}$  and  $\mathbf{X}_{t,s}^k \in \mathbb{R}^{n_t \times d_s}$  for all  $1 \leq s < t$  and  $k \in [K]$ .

Speaking intuitively, the approximation gap between Problem (2) and Problem (5) depends on the selection and granularity of the partition. By choosing partitions with a greater number of regions, Problem (5) can produce a tighter approximation of Problem (2), although this comes with an increase in problem size. For heuristic algorithms for selecting the partitions, we refer the reader to Postek and den Hertog (2016) and Bertsimas and Dunning (2016). Once the partitions are determined, we obtain a reformulation of Problem (5) by employing the same duality techniques as in Section 5.1. To this end, we will assume that the intersections between the regions of the partition and uncertainty sets take a rectangular form:

**ASSUMPTION 5.** *The intersection between each uncertainty set and region of the partition either has the form  $\mathcal{U}_N^j \cap P^k := \{\zeta \in \mathbb{R}^d : \ell^{jk} \leq \zeta \leq \mathbf{u}^{jk}\}$  or is empty.*

We remark that this assumption can be guaranteed under the same conditions as Assumption 4 when the partition's regions are constructed as hyperrectangles. We now show that Problem (5) can be reformulated as a finite-dimensional linear optimization problem which scales lightly in the number of sample paths  $N$  as well as the number of regions  $K$ .

**COROLLARY 1.** *If Assumption 5 holds, then (5) can be reformulated by adding at most  $O(N + Kmd)$  auxiliary continuous decision variables and  $O(m \sum_{j=1}^N |\mathcal{K}_j| + Kmd)$  linear constraints. The*

reformulation is

$$\begin{aligned}
& \text{minimize} && \frac{1}{N} \sum_{j=1}^N v_j \\
& \text{subject to} && \sum_{t=1}^T (\mathbf{u}_t^{jk} \cdot \boldsymbol{\alpha}_t^k - \boldsymbol{\ell}_t^{jk} \cdot \boldsymbol{\beta}_t^k + \mathbf{c}_t \cdot \mathbf{x}_{t,0}^k) \leq v_j \quad j \in [N], k \in \mathcal{K}_j \\
& && \boldsymbol{\alpha}_t^k - \boldsymbol{\beta}_t^k = \sum_{s=t+1}^T (\mathbf{X}_{s,t}^k)^\top \mathbf{c}_s \quad t \in [T], k \in [K] \\
& && \sum_{t=1}^T (\mathbf{M}_t^k \mathbf{u}_t^{jk} - \boldsymbol{\Lambda}_t^k \boldsymbol{\ell}_t^{jk} + \mathbf{A}_t \mathbf{x}_{t,0}^k) \leq \mathbf{b}^0 \quad j \in [N], k \in \mathcal{K}_j \\
& && \mathbf{M}_t^k - \boldsymbol{\Lambda}_t^k = -\mathbf{B}_t + \sum_{s=t+1}^T \mathbf{A}_s \mathbf{X}_{t,s}^k \quad t \in [T], k \in [K] \\
& && \mathbf{x}_t^k = \mathbf{x}_t^{k'}, \mathbf{X}_{t,s}^k = \mathbf{X}_{t,s}^{k'} \quad (k, k', t) \in \mathcal{T}(P^1, \dots, P^K), 1 \leq s < t,
\end{aligned}$$

where the auxiliary decision variables are  $\mathbf{v} \in \mathbb{R}^N$  as well as  $\boldsymbol{\alpha}^k := (\boldsymbol{\alpha}_1^k, \dots, \boldsymbol{\alpha}_T^k), \boldsymbol{\beta}^k := (\boldsymbol{\beta}_1^k, \dots, \boldsymbol{\beta}_T^k) \in \mathbb{R}_+^d$  and  $\mathbf{M}^k := (\mathbf{M}_1^k, \dots, \mathbf{M}_T^k), \boldsymbol{\Lambda}^k := (\boldsymbol{\Lambda}_1^k, \dots, \boldsymbol{\Lambda}_T^k) \in \mathbb{R}_+^{m \times d}$  for each  $k \in [K]$ . Note that  $\mathbf{b}(\boldsymbol{\zeta}) := \mathbf{b}^0 + \sum_{t=1}^T \mathbf{B}_t \boldsymbol{\zeta}_t \in \mathbb{R}^m$ .

*Proof.* The proof follows from similar duality techniques as Theorem 4 and is thus omitted.

This result suggests that Problem (2) with finite adaptability is scalable, in the sense that the size of the resulting reformulation for a given partition  $P^1, \dots, P^K$  scales lightly in the number of sample paths  $N$ . Assuming that the partition's regions and uncertainty sets are hyperrectangles, we remark that  $\boldsymbol{\ell}^{jk}, \mathbf{u}^{jk}$ , and  $\mathcal{T}(P^1, \dots, P^K)$  can be obtained efficiently by computing the intersection of each uncertainty set and region of the partition.

## 6. Relationships with Distributionally Robust Optimization

In the previous sections, we discussed the theoretical underpinnings and computational tractability of Problem (2) as a data-driven approach to multi-stage stochastic linear optimization. An attractive aspect of the proposed approach is its simplicity, interpretable as a straightforward robustification of historical sample paths. In this section, we explore connections between our data-driven approach to multi-stage stochastic linear optimization and distributionally robust optimization, and discuss implications of our results to the latter.

Our exposition in this section focuses on the following formulation of multi-stage distributionally robust linear optimization:

$$\begin{aligned}
& \underset{\mathbf{x} \in \mathcal{X}}{\text{minimize}} && \sup_{\mathbb{Q} \in \mathcal{A}_N} \mathbb{E}_{\mathbb{Q}} \left[ \sum_{t=1}^T \mathbf{c}_t(\boldsymbol{\xi}) \cdot \mathbf{x}_t(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}) \right] \\
& \text{subject to} && \sum_{t=1}^T \mathbf{A}_t(\boldsymbol{\xi}) \mathbf{x}_t(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}) \leq \mathbf{b}(\boldsymbol{\xi}) \quad \mathbb{Q}\text{-a.s.}, \forall \mathbb{Q} \in \mathcal{A}_N.
\end{aligned} \tag{6}$$

Intuitively speaking, this framework chooses the decision rules which minimize the expected cost with respect to an adversarially chosen probability distribution from an ambiguity set. The requirement that the constraints hold almost surely for every distribution in the ambiguity set ensures that the objective function will evaluate the cost function on realizations of the stochastic process where the decision rules are feasible. Examples of this formulation in multi-stage and data-driven two-stage problems include [Bertsimas et al. \(2019b\)](#) and [Hanasusanto and Kuhn \(2018\)](#).

Our following discussion focuses on ambiguity sets which are constructed using historical data and Wasserstein-based distances between probability distributions. Given two bounded probability distributions, their  $\infty$ -Wasserstein distance is defined as

$$\mathbf{d}_{\infty}(\mathbb{Q}, \mathbb{Q}') := \inf \left\{ \Pi\text{-ess sup}_{\Xi \times \Xi} \|\boldsymbol{\xi} - \boldsymbol{\xi}'\| : \begin{array}{l} \Pi \text{ is a joint distribution of } \boldsymbol{\xi} \text{ and } \boldsymbol{\xi}' \\ \text{with marginals } \mathbb{Q} \text{ and } \mathbb{Q}', \text{ respectively} \end{array} \right\},$$

where the essential supremum of the joint distribution is given by

$$\Pi\text{-ess sup}_{\Xi \times \Xi} \|\boldsymbol{\xi} - \boldsymbol{\xi}'\| := \inf \{ M : \Pi(\|\boldsymbol{\xi} - \boldsymbol{\xi}'\| > M) = 0 \}.$$

For any  $p \in [1, \infty)$ , the  $p$ -Wasserstein distance between two probability distributions is defined as

$$\mathbf{d}_p(\mathbb{Q}, \mathbb{Q}') = \inf \left\{ \left( \int_{\Xi \times \Xi} \|\boldsymbol{\xi} - \boldsymbol{\xi}'\|^p d\Pi(\boldsymbol{\xi}, \boldsymbol{\xi}') \right)^{\frac{1}{p}} : \begin{array}{l} \Pi \text{ is a joint distribution of } \boldsymbol{\xi} \text{ and } \boldsymbol{\xi}' \\ \text{with marginals } \mathbb{Q} \text{ and } \mathbb{Q}', \text{ respectively} \end{array} \right\}.$$

For technical details on these distances, we refer the reader to [Givens and Shortt \(1984\)](#). For any  $p \in [1, \infty]$ , let the  $p$ -Wasserstein ambiguity set be defined as

$$\mathcal{A}_N = \left\{ \mathbb{Q} \in \mathcal{P}(\Xi) : \mathbf{d}_p(\mathbb{Q}, \widehat{\mathbb{P}}_N) \leq \epsilon_N \right\},$$

where  $\epsilon_N \geq 0$  is a robustness parameter which controls the size of the ambiguity set and  $\widehat{\mathbb{P}}_N$  is the empirical probability distribution which assigns equal weight to each of the historical sample paths  $\hat{\boldsymbol{\xi}}^1, \dots, \hat{\boldsymbol{\xi}}^N$ . We henceforth refer to Problem (6) with the  $p$ -Wasserstein ambiguity set as  $p$ -WDRO.

As discussed at the end of Section 1.1, there are relatively few previous convergence guarantees for distributionally robust optimization with the  $\infty$ -Wasserstein ambiguity set, even for single-stage problems. Indeed, when the underlying distribution is unbounded, the  $\infty$ -Wasserstein ambiguity set will never contain the true distribution, even as  $N$  tends to infinity, since the distance  $\mathbf{d}_{\infty}(\mathbb{P}, \widehat{\mathbb{P}}_N)$

from the true to the empirical distribution will always be infinite. Thus, except under stronger assumptions than Assumption 1, the techniques used by Mohajerin Esfahani and Kuhn (2018, Theorems 3.5 and 3.6) to establish finite-sample and convergence guarantees for the 1-Wasserstein ambiguity set do not extend to the  $\infty$ -Wasserstein ambiguity set. Nonetheless, distributionally robust optimization with the  $\infty$ -Wasserstein ambiguity set has recently received interest in the context of regularization and adversarial training in machine learning (Gao et al. 2017, Staib and Jegelka 2017).

The following proposition shows that Problem (2), under the particular construction of uncertainty sets from Section 3, can also be interpreted as Problem (6) with the  $\infty$ -Wasserstein ambiguity set.

PROPOSITION 3. *Problem (2) with uncertainty sets of the form*

$$\mathcal{U}_N^j := \left\{ \zeta \equiv (\zeta_1, \dots, \zeta_T) \in \Xi : \|\zeta - \hat{\xi}^j\| \leq \epsilon_N \right\}$$

*is equivalent to  $\infty$ -WDRO.*

*Proof.* See Appendix G.  $\square$

Therefore, as a byproduct of Theorem 1 from Section 4, we have obtained general convergence guarantees for distributionally robust optimization using the  $\infty$ -Wasserstein ambiguity set under mild probabilistic assumptions.

For comparison, we now show that similar asymptotic optimality guarantees for multi-stage stochastic linear optimization are not obtained by  $p$ -WDRO for any  $p \in [1, \infty)$ . Indeed, the following proposition shows that the constraints induced by such an approach are overly conservative in general.

PROPOSITION 4. *If  $p \in [1, \infty)$  and  $\epsilon_N > 0$ , then a decision rule is feasible for  $p$ -WDRO only if*

$$\sum_{t=1}^T \mathbf{A}_t(\zeta) \mathbf{x}_t(\zeta_1, \dots, \zeta_{t-1}) \leq \mathbf{b}(\zeta) \quad \forall \zeta \in \Xi.$$

*Proof.* See Appendix H.  $\square$

As discussed in Section 2, the set  $\Xi$  is not necessarily a tight approximation of the true (unknown) support of the stochastic process, and may be strictly and significantly larger. Thus, the constraints induced from  $p$ -WDRO with  $p \in [1, \infty)$  may eliminate optimal or high-quality decision rules for Problem (1). Consequently,  $p$ -WDRO with  $p \in [1, \infty)$  is not asymptotically optimal for multi-stage stochastic linear optimization in general. We conclude this section with two further remarks.



REMARK 1. If we relaxed the constraints of  $p$ -WDRO with  $p \in [1, \infty)$  in an attempt to decrease its conservatism, then the resulting decision rules are not guaranteed to be feasible for the stochastic problem. Thus, the finite-sample guarantees provided by [Mohajerin Esfahani and Kuhn \(2018, Equation 2\)](#), which served as one of the principle justifications for using  $p$ -WDRO, would no longer provide meaningful insight into the true out-of-sample performance of this decision rule.  $\square$

REMARK 2. The conservatism of  $p$ -WDRO can lead to suboptimal decisions, even for problems where uncertainty does not impact feasibility, if the true support is not known exactly. For example, consider the problem

$$\begin{aligned} & \underset{x_2: \mathbb{R} \rightarrow \mathbb{R}, x_3: \mathbb{R}^2 \rightarrow \mathbb{R}}{\text{minimize}} && \mathbb{E}[x_2(\xi_1) + 2x_3(\xi_1, \xi_2)] \\ & \text{subject to} && x_2(\xi_1) + x_3(\xi_1, \xi_2) \geq \xi_1 + \xi_2 \quad \text{a.s.} \\ & && x_2(\xi_1) + x_3(\xi_1, \xi_2) \geq \xi_1 - \xi_2 \quad \text{a.s.} \end{aligned}$$

We observe that  $x_2(\xi_1) = \xi_1$  and  $x_3(\xi_1, \xi_2) = |\xi_2|$  are feasible decision rules, regardless of the underlying probability distribution. Suppose that the probability distribution and support of  $(\xi_1, \xi_2)$  is unknown, and our only information comes from historical data. If we approximate this stochastic problem using  $p$ -WDRO for any  $p \in [1, \infty)$  and linear decision rules, we are tasked with solving

$$\begin{aligned} & \underset{x_{2,0}, x_{2,1}, x_{3,0}, x_{3,1}, x_{3,2} \in \mathbb{R}}{\text{minimize}} && \sup_{\mathbb{Q} \in \mathcal{A}_N} \mathbb{E}_{\mathbb{Q}}[(x_{2,0} + x_{2,1}\zeta_1) + 2(x_{3,0} + x_{3,1}\zeta_1 + x_{3,2}\zeta_2)] \\ & \text{subject to} && (x_{2,0} + x_{2,1}\zeta_1) + (x_{3,0} + x_{3,1}\zeta_1 + x_{3,2}\zeta_2) \geq \zeta_1 + \zeta_2 \quad \forall \zeta \in \mathbb{R}^2 \\ & && (x_{2,0} + x_{2,1}\zeta_1) + (x_{3,0} + x_{3,1}\zeta_1 + x_{3,2}\zeta_2) \geq \zeta_1 - \zeta_2 \quad \forall \zeta \in \mathbb{R}^2. \end{aligned}$$

It follows from identical reasoning as [Bertsimas et al. \(2019b, Section 3\)](#) that there are no linear decision rules which are feasible for the above optimization problem. In particular, the above optimization problem will remain infeasible even if the true support of the random variable happens to be bounded but the bound is unknown. In contrast, the sample robust optimization approach (Problem (2)) to this example will always have a feasible linear decision rule. A similar example is found in [Section 7.2](#).  $\square$

## 7. Experimental Results

In this section, we assess the empirical performance of Problem (2) with the approximation algorithms from [Section 5](#) in the context of data-driven inventory management problems. In [Section 7.1](#), we revisit [Example 1](#) to demonstrate the impact of the robustness parameter when Problem (2) is approximated with a rich space of decision rules. In [Section 7.2](#), we compare the performance of Problem (2) with linear decision rules to alternative approaches in a practically-motivated inventory control problem.

### 7.1. Three-stage stochastic inventory management

In the first experiment, we revisit the stochastic inventory management problem from Example 1 in Section 3. Our aim is to provide a simple experiment that demonstrates the impact of the robustness parameter when Problem (2) is approximated by a rich restricted space of decision rules.

**7.1.1. Problem description.** We consider the three-stage stochastic inventory problem

$$\begin{aligned} & \underset{x_2: \mathbb{R} \rightarrow \mathbb{R}, x_3: \mathbb{R}^2 \rightarrow \mathbb{R}}{\text{minimize}} && \mathbb{E}[x_2(\xi_1) + 2x_3(\xi_1, \xi_2)] \\ & \text{subject to} && x_2(\xi_1) + x_3(\xi_1, \xi_2) \geq \xi_1 + \xi_2 \quad \text{a.s.} \\ & && x_2(\xi_1), x_3(\xi_1, \xi_2) \geq 0 \quad \text{a.s.} \end{aligned} \quad (3)$$

We assume that the joint probability distribution of the preorder and regular demand  $(\xi_1, \xi_2) \in \mathbb{R}^2$  is unknown. Instead, we have access to historical data consisting of preorder and regular demands of past products  $(\hat{\xi}_1^1, \hat{\xi}_2^1), \dots, (\hat{\xi}_1^N, \hat{\xi}_2^N)$ , which are independent and identically distributed sample paths of the underlying stochastic process, and knowledge that the stochastic process  $(\xi_1, \xi_2)$  will be contained in  $\Xi = \mathbb{R}_+^2$  almost surely.

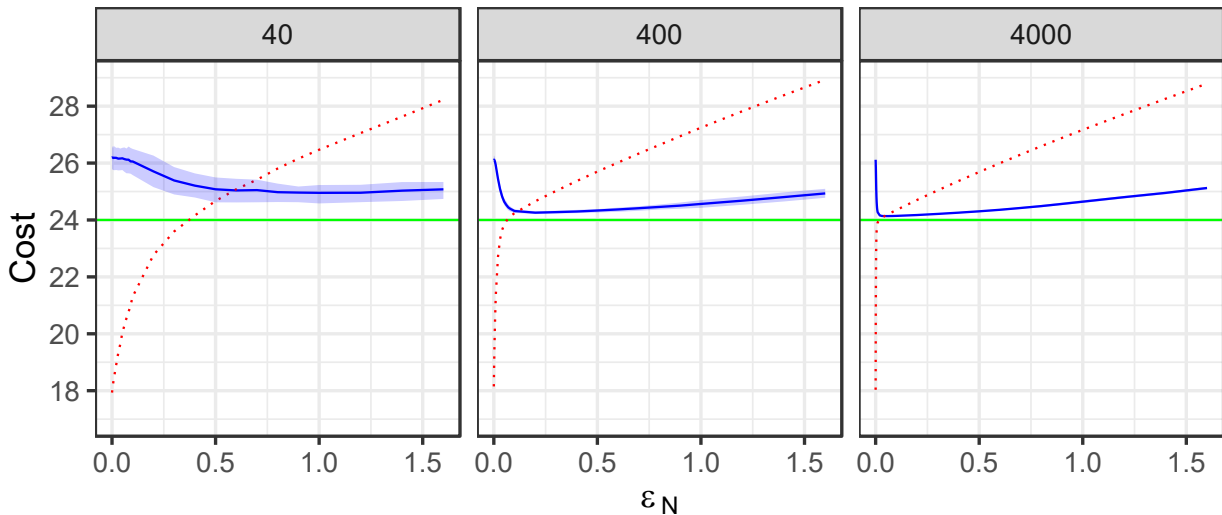
**7.1.2. Experiments.** In the following experiment, we obtain decision rules for Problem (3) by approximating Problem (2) with finite adaptability. Specifically, given a historical dataset, we first construct the uncertainty sets as described in Section 3 using the  $\ell_\infty$ -norm. Assume that the preorder demand is a continuous random variable and, without loss of generality, let the historical data be sorted such that  $\hat{\xi}_1^1 < \dots < \hat{\xi}_1^N$ . Then, we approximate Problem (2) by restricting the decision rule  $x_2: \mathbb{R} \rightarrow \mathbb{R}$  to be a piecewise static function of the form

$$x_2(\zeta_1) = \begin{cases} x_2^1, & \text{if } (\zeta_1, \zeta_2) \in P^1, \\ \vdots \\ x_2^N, & \text{if } (\zeta_1, \zeta_2) \in P^N; \end{cases} \quad P^k = \begin{cases} \left(-\infty, \frac{\hat{\xi}_1^1 + \hat{\xi}_1^2}{2}\right] \times \mathbb{R}, & \text{if } k = 1, \\ \left[\frac{\hat{\xi}_1^{k-1} + \hat{\xi}_1^k}{2}, \frac{\hat{\xi}_1^k + \hat{\xi}_1^{k+1}}{2}\right) \times \mathbb{R}, & \text{if } k \in \{2, \dots, N-1\}, \\ \left[\frac{\hat{\xi}_1^{N-1} + \hat{\xi}_1^N}{2}, \infty\right) \times \mathbb{R}, & \text{if } k = N. \end{cases}$$

These piecewise static decision rules are defined over the partition  $P^1, \dots, P^N$ , which is constructed from the midpoints of consecutive historical preorder demands. Note that the partition is constructed such that each historical sample path lies in its own region, *i.e.*,  $(\hat{\xi}_1^k, \hat{\xi}_2^k) \in P^k$  for each  $k \in [N]$ . The best piecewise static decision rule for Problem (2) over the partition  $P^1, \dots, P^N$  is obtained by solving

$$\begin{aligned} & \underset{\substack{x_2^1, \dots, x_2^N \in \mathbb{R}, \\ x_3: \mathbb{R}^2 \rightarrow \mathbb{R}}}{\text{minimize}} && \frac{1}{N} \sum_{j=1}^N \max_{k \in \mathcal{K}^j} \sup_{\zeta \in \mathcal{U}_N^j \cap P^k} \{x_2^k + 2x_3(\zeta_1, \zeta_2)\} \\ & \text{subject to} && x_2^k + x_3(\zeta_1, \zeta_2) \geq \zeta_1 + \zeta_2 && \forall \zeta \in \cup_{j=1}^N \mathcal{U}_N^j \cap P^k, k \in [N] \\ & && x_2^k, x_3(\zeta_1, \zeta_2) \geq 0 && \forall \zeta \in \cup_{j=1}^N \mathcal{U}_N^j \cap P^k, k \in [N], \end{aligned} \quad (7)$$

Figure 1 Three-stage stochastic inventory management: impact of robustness parameter.



*Note.* The solid blue line is the average out-of-sample cost of decision rules produced by Problem (7), and the shaded blue region is the 20th and 80th percentiles over the 100 training datasets. The dotted red line is the average in-sample cost of Problem (7), and the solid green line is the optimal cost of Problem (3). Results are shown for  $N \in \{40, 400, 4000\}$ .

where  $\mathcal{K}_j := \{k \in [N] : \mathcal{U}_N^j \cap P^k \neq \emptyset\}$  denotes the indices of regions  $P^1, \dots, P^N$  that intersect the uncertainty set  $\mathcal{U}_N^j$ . By following techniques from Section 5.2 and exploiting problem structure, Problem (7) can be reformulated as a linear optimization problem with  $O(N)$  decision variables and  $O(\sum_{j=1}^N |\mathcal{K}_j|)$  constraints; see Appendix I.

We perform experiments where the preorder and regular demands have a joint probability distribution given by  $\xi_1 \sim \text{Unif}[0, 12]$  and  $\xi_2 \sim \text{Unif}[0, \xi_1^2/2]$ .<sup>2</sup> For various choices of  $N$ , we generate training datasets of size  $N$  and evaluate the out-of-sample cost of decision rules obtained by Problem (7) with robustness parameter  $\epsilon_N$ . The out-of-sample cost of the obtained decision rules,  $\mathbb{E}[x_2(\xi_1) + 2x_3(\xi_1, \xi_2)]$ , is approximated by the average cost of the decision rules on a common testing dataset of 10000 sample paths. All experiments are repeated over 100 training datasets, and all sample paths in the training and testing datasets are independently generated from the true joint probability distribution.

**7.1.3. Results.** In Figure 1, we show the impact of the robustness parameter on the in-sample and out-of-sample cost of the piecewise static decision rules obtained by Problem (7). The results demonstrate that a *strictly positive* choice of the robustness parameter is essential in order to obtain the best out-of-sample cost for each  $N$ . This is due to the fact that Problem (2) has

<sup>2</sup> Under this distribution, an optimal decision rule for Problem (3) is  $x_2^*(\xi_1) = \xi_1 + \frac{\xi_1^2}{4}$  and  $x_3^*(\xi_1, \xi_2) = \max\{0, \xi_2 - \frac{\xi_1^2}{4}\}$ . This can be readily derived by applying the newsvendor model to the conditional distribution of  $\xi_2 | \xi_1$ .

been approximated with a rich space of decision rules, and in particular, a space of decision rules that becomes increasingly flexible as more historical sample paths are obtained. In contrast, when the robustness parameter is set to zero, Figure 1 shows that the in-sample cost of Problem (7) converges to the suboptimal  $\mathbb{E}[\xi_1 + \xi_2] = 18$ ; these findings are consistent with the discussion in Example 1, and show that approximating Problem (2) with a rich restricted space of decision rules will asymptotically overfit the historical data when the robustness parameter is set to zero.

These results highlight a practical strength of Problem (2). Indeed, the approximation of Problem (2) using piecewise static decision rules (Problem (7)) did not require nor utilize any information about the structure of optimal decision rules for the underlying stochastic problem. At the same time, despite searching over a rich space of decision rules, Problem (7) with an appropriate robustness parameter produces decision rules which do overfit the historical data. This shows that Problem (2) provides a opportunity to find high-quality decision rules for multi-stage stochastic linear optimization problems, even when (i) the only information on the underlying distribution comes from limited data, and (ii) the structure of optimal decision rules for the stochastic problem is complex or unknown.

## 7.2. Multi-stage stochastic inventory management

In the second experiment, we consider a multi-stage stochastic inventory management problem with (unknown) autoregressive demand in the data-driven setting. Our goal is to assess the performance of Problem (2) with linear decision rules in comparison to alternative data-driven approaches.

**7.2.1. Problem description.** We consider an inventory management problem of a single product over a finite planning horizon. At the beginning of each time period  $t \in [T]$ , we start with  $I_t \in \mathbb{R}$  units of product in inventory. We then select a production quantity of  $x_t \in [0, \bar{x}_t]$  with zero lead time at a cost of  $c_t$  per unit. The product demand  $\xi_t \geq 0$  is then revealed, the inventory is updated to  $I_{t+1} = I_t + x_t - \xi_t$ , and we incur a holding cost of  $h_t \max\{I_{t+1}, 0\}$  and a backorder cost of  $b_t \max\{-I_{t+1}, 0\}$ . We begin with zero units of inventory in the first period. Our goal is to dynamically select the production quantities to minimize the expected total cost over the planning horizon, captured by

$$\begin{aligned}
& \underset{\mathbf{x}, \mathbf{I}, \mathbf{y}}{\text{minimize}} && \mathbb{E} \left[ \sum_{t=1}^T (c_t x_t(\xi_1, \dots, \xi_{t-1}) + y_{t+1}(\xi_1, \dots, \xi_t)) \right] \\
& \text{subject to} && I_{t+1}(\xi_1, \dots, \xi_t) = I_t(\xi_1, \dots, \xi_{t-1}) + x_t(\xi_1, \dots, \xi_{t-1}) - \xi_t && \text{a.s., } \forall t \in [T] \\
& && y_{t+1}(\xi_1, \dots, \xi_t) \geq h_t I_{t+1}(\xi_1, \dots, \xi_t) && \text{a.s., } \forall t \in [T] \\
& && y_{t+1}(\xi_1, \dots, \xi_t) \geq -b_t I_{t+1}(\xi_1, \dots, \xi_t) && \text{a.s., } \forall t \in [T] \\
& && 0 \leq x_t(\xi_1, \dots, \xi_{t-1}) \leq \bar{x}_t && \text{a.s., } \forall t \in [T].
\end{aligned} \tag{8}$$

We consider the setting where the joint probability distribution of the stochastic process  $(\xi_1, \dots, \xi_T) \in \mathbb{R}^T$  is unknown. Our only information on the distribution comes from historical data consisting of demand realizations for past products  $(\hat{\xi}_1^1, \dots, \hat{\xi}_T^1), \dots, (\hat{\xi}_1^N, \dots, \hat{\xi}_T^N)$ , which are independent and identically distributed sample paths of the underlying stochastic process, and knowledge that the stochastic process will be contained in  $\Xi = \mathbb{R}_+^T$  almost surely.

**7.2.2. Experiments.** We perform computational experiments on the following data-driven approaches for obtaining decision rules for Problem (8):

- *SRO-LDR*: This is the proposed data-driven approach for multi-stage stochastic linear optimization (Problem (2)), where the uncertainty sets are constructed as described in Section 3 with the  $\ell_\infty$ -norm. The approach is approximated using linear decision rules (see Section 5.1) and solved using the reformulation developed in Theorem 4. We choose the robustness parameter for each training dataset using 5-fold cross validation, where the range of possible values considered in the cross-validation procedure is  $\epsilon_N \in \{b \cdot 10^a : a \in \{-2, -1, 0, 1\}, b \in \{1, \dots, 9\}\}$ .
- *SAA-LDR*: This is the same approach as SRO-LDR, except the robustness parameter is set to zero.
- *Approx PCM*: This is a data-driven extension of the approach developed in Bertsimas et al. (2019b). In this approach, decision rules are obtained by solving a multi-stage distributionally robust optimization problem (Problem (6)) in which  $\mathcal{A}_N$  is the set of joint probability distributions with the same mean and covariance as those estimated from the historical data. This distributionally robust optimization problem is solved approximately by restricting to lifted linear decision rules, as described in Bertsimas et al. (2019b, Section 3).
- *DDP* and *RDDP*: This is the robust data-driven dynamic programming approach proposed by Hanasusanto and Kuhn (2013). The approach estimates cost-to-go functions by applying kernel regression to the historical sample paths. Decisions are obtained from optimizing over the cost-to-go functions, which are evaluated approximately using the algorithm described in Hanasusanto and Kuhn (2013, Section 4). Since the algorithm requires both input sample paths and initial state paths, we use half of the training dataset as the input sample paths, and the other half to generate the state paths via the lifted linear decision rules obtained by Approx PCM. The approach also requires a robustness parameter  $\gamma$ , which we choose to be either  $\gamma = 0$  (DDP) or  $\gamma = 10$  (RDDP).
- *WDRO-LDR*: Described in Section 6, this approach obtains decision rules by solving a multi-stage distributionally robust optimization problem (Problem (6)) in which  $\mathcal{A}_N$  is chosen to be

the 1-Wasserstein ambiguity set with the  $\ell_1$ -norm. Similarly as SRO-LDR, the distributionally robust optimization problem is approximated using linear decision rules, which is solved using a duality-based reformulation provided in Appendix J. The robustness parameter is chosen using the same procedure as SRO-LDR.

We perform computational simulations using the same parameters and data generation as See and Sim (2010). Specifically, the demand is a nonstationary autoregressive stochastic process of the form  $\xi_t = \varsigma_t + \alpha\varsigma_{t-1} + \dots + \alpha\varsigma_1 + \mu$ , where  $\varsigma_1, \dots, \varsigma_T$  are independent random variables distributed uniformly over  $[-\bar{\varsigma}, \bar{\varsigma}]$ . The parameters of the stochastic process are  $\mu = 200$  and  $\bar{\varsigma} = 40$  when  $T = 5$ , and  $\mu = 200$  and  $\bar{\varsigma} = 20$  when  $T = 10$ . The capacities and costs are  $\bar{x}_t = 260$ ,  $c_t = 0.1$ ,  $h_t = 0.02$  for all  $t \in [T]$ ,  $b_t = 0.2$  for all  $t \in [T - 1]$ , and  $b_T = 2$ .

To compare the above data-driven approaches, we take the following steps. For various choices of  $N$ , we generate 100 training datasets of size  $N$  and obtain decision rules by applying the above data-driven approaches to each training dataset. The out-of-sample costs of the obtained decision rules are approximated using a common testing dataset of 10000 sample paths.<sup>3</sup> Specifically, for each sample path  $(\xi_1^i, \dots, \xi_T^i) \in \mathbb{R}^T$  in the testing dataset, we calculate production quantities  $(x_1^{A,i,\ell}, \dots, x_T^{A,i,\ell}) \in \mathbb{R}^T$  by applying the decision rule obtained from approach  $\mathcal{A}$  on the  $\ell$ -th training dataset. The out-of-sample cost of the decision rule is then approximated as

$$\frac{1}{10000} \sum_{i=1}^{10000} \sum_{t=1}^T (c_t x_t^{A,i,\ell} + \max \{h_t I_{t+1}^{A,i,\ell}, -b_t I_{t+1}^{A,i,\ell}\}),$$

where the inventory levels  $(I_1^{A,i,\ell}, \dots, I_T^{A,i,\ell})$  are computed from the production quantities  $(x_1^{A,i,\ell}, \dots, x_T^{A,i,\ell}) \in \mathbb{R}^T$  and the test sample path  $(\xi_1^i, \dots, \xi_T^i) \in \mathbb{R}^T$ . All sample paths in the training and testing datasets are drawn independently from the true joint probability distribution.

As discussed earlier, Problem (2) is not guaranteed to find decision rules which are feasible for all realizations in  $\Xi$ . Therefore, the linear decision rules obtained by SRO-LDR and SAA-LDR, when applied to sample paths in the testing dataset, may result in production quantities which exceed  $\bar{x}_1, \dots, \bar{x}_T$  or are negative. Thus, before computing the out-of-sample costs, we first project each production quantity  $x_t^{A,i,\ell}$  onto the interval  $[0, \bar{x}_t]$  to ensure it is feasible. We discuss the impact of this projection procedure at the end of the results section.

**7.2.3. Results.** In Table 1 and Figure 2, we report the out-of-sample costs and computation times of the various approaches. SRO-LDR produces an out-of-sample cost which outperforms the other approaches, most notably when the size of the training dataset is small, and requires less

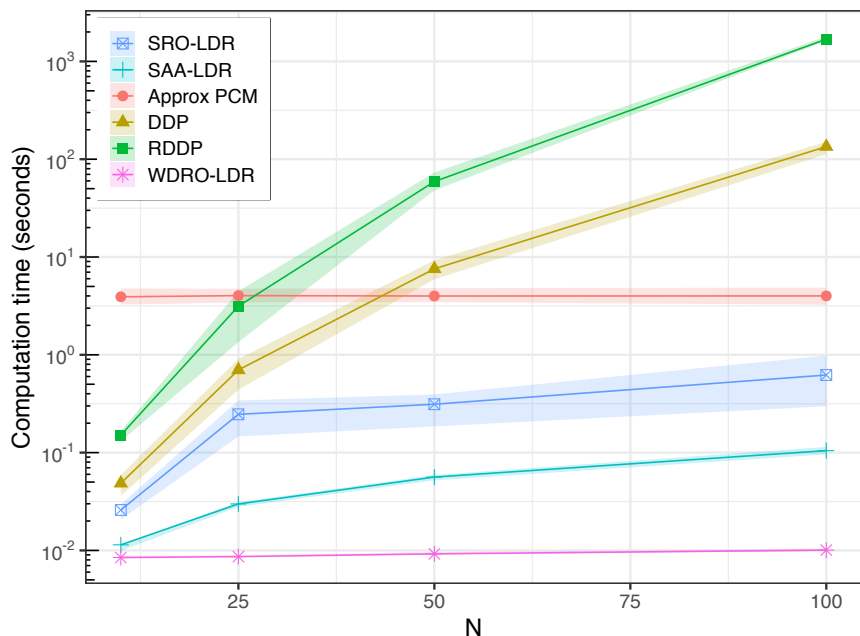
<sup>3</sup> For DDP and RDDP, we only evaluated on the first 1000 sample paths in the testing dataset, due to the computational cost of optimizing over the cost-to-go functions for each testing sample path.

**Table 1** Multi-stage stochastic inventory management: Average out-of-sample cost.

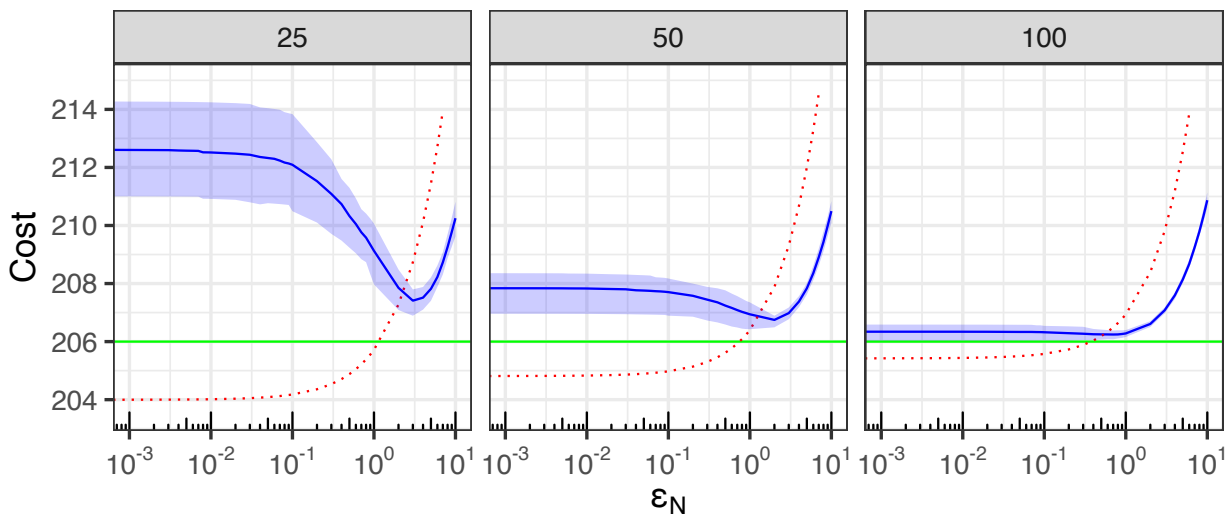
$T$	$\alpha$	Approach	Size of training dataset (N)				DP
			10	25	50	100	
5	0	SRO-LDR	<b>111.8(7.2)</b>	<b>109.5(1.3)</b>	<b>108.5(0.7)</b>	108.0(0.3)	108
		SAA-LDR	127.3(12.6)	111.7(3.1)	108.7(1.0)	<b>107.9(0.3)</b>	
		Approx PCM	118.5(2.2)	117.4(1.0)	117.1(0.7)	117.0(0.5)	
		DDP	2262.7(363.3)	1189.6(854.2)	525.4(510.8)	205.1(201.5)	
		RDDP	2255.5(393.2)	1175.8(856.2)	515.6(506.0)	202.4(195.3)	
		WDRO-LDR	2400.3(0.0)	2400.3(0.0)	2400.3(0.0)	2400.3(0.0)	
	0.25	SRO-LDR	<b>113.0(4.2)</b>	<b>110.0(1.8)</b>	108.7(0.8)	108.0(0.3)	107
		SAA-LDR	127.6(13.0)	111.7(3.1)	<b>108.6(1.0)</b>	<b>108.0(0.2)</b>	
		Approx PCM	126.8(3.8)	125.3(1.5)	124.8(0.9)	124.6(0.8)	
		DDP	2251.5(488.8)	1393.7(897.8)	679.3(656.0)	236.9(240.5)	
		RDDP	2222.0(556.5)	1386.2(900.7)	670.1(654.9)	236.2(237.2)	
		WDRO-LDR	2400.7(0.0)	2400.7(0.0)	2400.7(0.0)	2400.7(0.0)	
	0.5	SRO-LDR	<b>115.5(5.3)</b>	<b>112.0(4.0)</b>	110.8(2.7)	111.7(2.6)	108
		SAA-LDR	129.5(13.3)	113.1(6.8)	<b>110.7(2.9)</b>	<b>111.6(2.6)</b>	
		Approx PCM	136.0(4.8)	134.0(1.9)	133.4(1.2)	133.2(1.0)	
		DDP	2263.8(480.1)	1563.7(917.6)	777.5(787.7)	364.2(488.1)	
		RDDP	2253.8(515.8)	1532.6(940.9)	716.8(758.9)	334.6(477.1)	
		WDRO-LDR	2401.2(0.0)	2401.2(0.0)	2401.2(0.0)	2401.2(0.0)	
10	0	SRO-LDR	<b>208.9(1.0)</b>	<b>207.5(0.6)</b>	<b>206.8(0.5)</b>	<b>206.2(0.2)</b>	206
		SAA-LDR	293.9(70.1)	212.6(2.1)	207.8(1.1)	206.3(0.4)	
		Approx PCM	215.3(2.1)	214.5(0.6)	214.1(0.6)	214.1(0.4)	
		DDP	5211.4(1131.1)	2827.9(1757.5)	1335.6(1206.4)	497.5(550.2)	
		RDDP	5210.1(1133.4)	2820.3(1758.7)	1327.6(1206.0)	500.1(552.7)	
		WDRO-LDR	5800.3(0.0)	5800.3(0.0)	5800.3(0.0)	5800.3(0.0)	
	0.25	SRO-LDR	<b>210.3(2.9)</b>	<b>207.8(1.1)</b>	<b>206.9(0.5)</b>	<b>206.3(0.2)</b>	206
		SAA-LDR	295.1(70.4)	212.7(2.1)	207.8(1.1)	206.3(0.4)	
		Approx PCM	228.7(4.5)	226.2(1.8)	225.7(1.1)	225.5(0.9)	
		DDP	5215.6(1350.1)	3214.7(1984.9)	1598.0(1566.5)	440.2(417.1)	
		RDDP	5202.0(1368.7)	3185.3(1977.1)	1593.6(1566.6)	437.0(418.0)	
		WDRO-LDR	5800.2(0.0)	5800.5(0.0)	5800.5(0.0)	5800.5(0.0)	
	0.5	SRO-LDR	<b>211.1(3.9)</b>	<b>207.9(1.0)</b>	<b>206.9(0.6)</b>	<b>206.3(0.2)</b>	206
		SAA-LDR	297.8(70.7)	213.0(2.3)	207.9(1.1)	206.4(0.4)	
		Approx PCM	245.3(7.0)	242.1(2.7)	241.6(2.0)	240.9(1.6)	
		DDP	5374.8(1052.7)	3676.4(2159.1)	1960.9(1878.3)	644.1(914.6)	
		RDDP	5313.0(1173.0)	3630.9(2161.4)	1949.2(1863.9)	644.0(913.3)	
		WDRO-LDR	5800.7(0.0)	5800.7(0.0)	5800.3(0.0)	5800.7(0.0)	

Mean (standard deviation) for the out-of-sample cost of decision rules obtained by various data-driven approaches for Problem (8). The robustness parameters in SRO-LDR and WDRO-LDR are chosen using cross validation. The column DP presents the dynamic programming approximations of the optimal cost of Problem (8) from See and Sim (2010, Tables EC.1 and EC.2), which have an accuracy of  $\pm 1\%$ .

than one second of computation time. We note that the out-of-sample cost of SRO-LDR roughly converges to the dynamic programming (DP) estimate of the optimal cost of Problem (8), which suggests that linear decision rules provide a good approximation of the optimal production decision rules for this particular stochastic problem. The relationship between the robustness parameter and the in-sample and out-of-sample cost of SRO-LDR is shown in Figure 3.

**Figure 2** Multi-stage stochastic inventory management: Computation times for  $T = 10$ ,  $\alpha = 0.25$ .

*Note.* Computation times for data-driven approaches to the multi-stage stochastic inventory management problem. Results are shown for  $T = 10$  and  $\alpha = 0.25$ , and similar computation times were observed for other choices of  $\alpha$ . The graph shows the mean value of the computation times over 100 training datasets for each value of  $N$ .

**Figure 3** Multi-stage stochastic inventory management: Impact of robustness parameter on SRO-LDR for  $T = 10$ ,  $\alpha = 0.25$ .

*Note.* The solid blue line is the average out-of-sample cost of decision rules produced by SRO-LDR, and the shaded blue region is the 20th and 80th percentiles over the 100 training datasets. The dotted red line is the average in-sample cost of SRO-LDR, and the solid green line is a dynamic programming approximation of the optimal cost of Problem (8) from See and Sim (2010, Table EC.2).



We briefly reflect on some notable differences between SRO-LDR and the other approaches. First, the results demonstrate that a strictly positive choice of the robustness parameter is not necessary to avoid asymptotic overfitting when Problem (2) is approximated with a fixed, finite-dimensional space of decision rules; indeed, Table 1 and Figure 3 show that SAA-LDR can provide an out-of-sample cost which is similar to that of SRO-LDR for moderate to large training datasets. However, SRO-LDR produces an out-of-sample cost which significantly outperforms SAA-LDR when  $N$  is small ( $N \in \{10, 25\}$ ). More generally, this shows that there exist regimes in which a positive choice of the robustness parameter can still provide significant value even when Problem (2) is approximated using linear decision rules. Second, we note that WDRO-LDR consistently produces decision rules with large average out-of-sample cost; this is due to the fact that this approach requires the linear decision rules to satisfy  $0 \leq x_{t,0} + \sum_{s=1}^{t-1} x_{t,s} \zeta_s \leq \bar{x}_t$  for all  $(\zeta_1, \dots, \zeta_T) \in \mathbb{R}_+^T$ , which reduces to a static decision rule for the production quantity in each stage. Finally, we remark that the average out-of-sample cost of DDP and RDDP improved significantly with the size of the training dataset, but produced high variability across training datasets and required long computation time.

We recall that Table 1 reports the out-of-sample costs of SRO-LDR and SAA-LDR after their production quantities are projected onto the feasible region (see Section 7.2.2). In Appendix K, we discuss the impact of this projection procedure on the out-of-sample cost. Specifically, we show across the above experiments that (i) SRO-LDR produces feasible production quantities for more than 93% of the sample paths in the testing dataset, and (ii) the average  $\ell_1$ -distance between the production quantities  $(x_1^{A,i,\ell}, \dots, x_T^{A,i,\ell})$  and the feasible region  $[0, \bar{x}_1] \times \dots \times [0, \bar{x}_T]$  is less than 2 units. This shows that SRO-LDR consistently produces feasible or nearly-feasible decisions, and thus the out-of-sample costs of SRO-LDR are unlikely to be an artifact of this projection procedure.

## 8. Conclusion

In this work, we presented a new data-driven approach, based on robust optimization, for solving multi-stage stochastic linear optimization problems where uncertainty is correlated across time. We showed that the proposed approach is asymptotically optimal, providing assurance that the approach offers a near-optimal approximation of the underlying stochastic problem in the presence of big data. At the same time, the optimization problem resulting from the proposed approach can be addressed by approximation algorithms and reformulation techniques which have underpinned the success of multi-stage robust optimization. The practical value of the proposed approach was illustrated by computational examples inspired by real-world applications, demonstrating that the proposed data-driven approach can produce high-quality decisions in reasonable computation times. Through these contributions, this work provides a step towards helping organizations across domains leverage historical data to make better operational decisions in dynamic environments.

## References

- Amparo Baïllo, Antonio Cuevas, and Ana Justel. Set estimation and nonparametric detection. *Canadian Journal of Statistics*, 28(4):765–782, 2000.
- Güzin Bayraksan and David K Love. Data-driven stochastic programming using phi-divergences. In *The Operations Research Revolution*, chapter 1, pages 1–19. INFORMS, 2015.
- Aharon Ben-Tal and Arkadi Nemirovski. Robust solutions of uncertain linear programs. *Operations Research Letters*, 25(1):1–13, 1999.
- Aharon Ben-Tal, Alexander Goryashko, Elana Guslitzer, and Arkadi Nemirovski. Adjustable robust solutions of uncertain linear programs. *Mathematical Programming*, 99(2):351–376, 2004.
- Aharon Ben-Tal, Laurent El Ghaoui, and Arkadi Nemirovski. *Robust Optimization*. Princeton University Press, 2009.
- Aharon Ben-Tal, Dick den Hertog, Anja De Waegenaere, Bertrand Melenberg, and Gijs Rennen. Robust solutions of optimization problems affected by uncertain probabilities. *Management Science*, 59(2):341–357, 2013.
- Dimitris Bertsimas and Constantine Caramanis. Finite adaptability in multistage linear optimization. *IEEE Transactions on Automatic Control*, 55(12):2751–2766, 2010.
- Dimitris Bertsimas and Iain Dunning. Multistage robust mixed-integer optimization with adaptive partitions. *Operations Research*, 64(4):980–998, 2016.
- Dimitris Bertsimas and Vineet Goyal. On the power of robust solutions in two-stage stochastic and adaptive optimization problems. *Mathematics of Operations Research*, 35(2):284–305, 2010.
- Dimitris Bertsimas, David B Brown, and Constantine Caramanis. Theory and applications of robust optimization. *SIAM Review*, 53(3):464–501, 2011a.
- Dimitris Bertsimas, Vineet Goyal, and Xu Andy Sun. A geometric characterization of the power of finite adaptability in multistage stochastic and adaptive optimization. *Mathematics of Operations Research*, 36(1):24–54, 2011b.
- Dimitris Bertsimas, Vishal Gupta, and Nathan Kallus. Robust sample average approximation. *Mathematical Programming*, 171(1):217–282, 2018.
- Dimitris Bertsimas, Shimrit Shtern, and Bradley Sturt. Two-stage sample robust optimization. *arXiv preprint arXiv:1907.07142*, 2019a.
- Dimitris Bertsimas, Melvyn Sim, and Meilin Zhang. Adaptive distributionally robust optimization. *Management Science*, 65(2):604–618, 2019b.
- John R Birge and Francois Louveaux. *Introduction to Stochastic Programming*. Springer Science & Business Media, 2011.
- Leo Breiman. *Probability*. SIAM, 1992.
- Xin Chen, Melvyn Sim, and Peng Sun. A robust optimization perspective on stochastic programming. *Operations Research*, 55(6):1058–1071, 2007.
- George B Dantzig. Linear programming under uncertainty. *Management Science*, 1(3-4):197–206, 1955.
- Erick Delage and Dan A Iancu. Robust multistage decision making. In *The Operations Research Revolution*, pages 20–46. INFORMS, 2015.
- Erick Delage and Yinyu Ye. Distributionally robust optimization under moment uncertainty with application to data-driven problems. *Operations Research*, 58(3):595–612, 2010.
- Luc Devroye and Gary L Wise. Detection of abnormal behavior via nonparametric estimation of the support. *SIAM Journal on Applied Mathematics*, 38(3):480–488, 1980.
- E. Erdoğan and Garud Iyengar. Ambiguous chance constrained problems and robust optimization. *Mathematical Programming*, 107(1):37–61, Jun 2006.
- E. Erdoğan and Garud Iyengar. On two-stage convex chance constrained problems. *Mathematical Methods of Operations Research*, 65(1):115–140, 2007.

- K Bruce Erickson. The strong law of large numbers when the mean is undefined. *Transactions of the American Mathematical Society*, 185:371–381, 1973.
- Nicolas Fournier and Arnaud Guillin. On the rate of convergence in wasserstein distance of the empirical measure. *Probability Theory and Related Fields*, 162(3-4):707–738, 2015.
- Rui Gao, Xi Chen, and Anton J Kleywegt. Wasserstein distributional robustness and regularization in statistical learning. *arXiv preprint arXiv:1712.06050v2*, 2017.
- Stanley J Garstka and Roger J-B Wets. On decision rules in stochastic programming. *Mathematical Programming*, 7(1):117–143, 1974.
- Angelos Georghiou, Daniel Kuhn, and Wolfram Wiesemann. The decision rule approach to optimization under uncertainty: methodology and applications. *Computational Management Science*, pages 1–32, 2018.
- Angelos Georghiou, Angelos Tsoukalas, and Wolfram Wiesemann. Robust dual dynamic programming. *Operations Research*, 67(3):813–830, 2019.
- Clark R Givens and Rae Michael Shortt. A class of wasserstein metrics for probability distributions. *The Michigan Mathematical Journal*, 31(2):231–240, 1984.
- Grani A Hanasusanto and Daniel Kuhn. Robust data-driven dynamic programming. In *Advances in Neural Information Processing Systems*, pages 827–835, 2013.
- Grani A Hanasusanto and Daniel Kuhn. Conic programming reformulations of two-stage distributionally robust linear programs over wasserstein balls. *Operations Research*, 66(3):849–869, 2018.
- Ruiwei Jiang and Yongpei Guan. Risk-averse two-stage stochastic program with distributional ambiguity. *Operations Research*, 66(5):1390–1405, 2018.
- Aleksandr Petrovich Korostelev and Alexandre B Tsybakov. *Minimax Theory of Image Reconstruction*, volume 82. Springer-Verlag, New York, 1993.
- Pavlo Krokmal and Stanislav Uryasev. A sample-path approach to optimal position liquidation. *Annals of Operations Research*, 152(1):193–225, 2007.
- Peyman Mohajerin Esfahani and Daniel Kuhn. Data-driven distributionally robust optimization using the wasserstein metric: Performance guarantees and tractable reformulations. *Mathematical Programming*, 171(1):115–166, 2018.
- Georg Ch Pflug and Alois Pichler. From empirical observations to tree models for stochastic optimization: Convergence properties. *SIAM Journal on Optimization*, 26(3):1715–1740, 2016.
- Georg Ch Pflug and David Wozabal. Ambiguity in portfolio selection. *Quantitative Finance*, 7(4):435–442, 2007.
- Krzysztof Postek and Dick den Hertog. Multistage adjustable robust mixed-integer optimization via iterative splitting of the uncertainty set. *INFORMS Journal on Computing*, 28(3):553–574, 2016.
- Herbert E Scarf. A min-max solution of an inventory problem. In *Studies in the Mathematical Theory of Inventory and Production*, pages 201–209. Stanford University Press, 1958.
- Bernhard Schölkopf, John C Platt, John Shawe-Taylor, Alex J Smola, and Robert C Williamson. Estimating the support of a high-dimensional distribution. *Neural Computation*, 13(7):1443–1471, 2001.
- Chuen-Teck See and Melvyn Sim. Robust approximation to multiperiod inventory management. *Operations Research*, 58(3):583–594, 2010.
- Soroosh Shafieezadeh-Abadeh, Daniel Kuhn, and Peyman Mohajerin Esfahani. Regularization via mass transportation. *Journal of Machine Learning Research*, 20(103):1–68, 2019.
- Alexander Shapiro, Darinka Dentcheva, and Andrzej Ruszczyński. *Lectures on Stochastic Programming: Modeling and Theory*. SIAM, 2009.
- Allen L Soyster. Convex programming with set-inclusive constraints and applications to inexact linear programming. *Operations Research*, 21(5):1154–1157, 1973.
- Matthew Staib and Stefanie Jegelka. Distributionally robust deep learning as a generalization of adversarial training. In *NIPS Machine Learning and Computer Security Workshop*, 2017.

- Bart PG Van Parys, Peyman Mohajerin Esfahani, and Daniel Kuhn. From data to decisions: Distributionally robust optimization is optimal. *arXiv preprint arXiv:1704.04118*, 2017.
- Jean-Louis Verger-Gaugry. Covering a ball with smaller equal balls in  $\mathbb{R}^n$ . *Discrete & Computational Geometry*, 33(1):143–155, 2005.
- Hong Wang and RJ Tomkins. A zero-one law for large order statistics. *Canadian Journal of Statistics*, 20(3):323–334, 1992.
- Wolfram Wiesemann, Daniel Kuhn, and Melvyn Sim. Distributionally robust convex optimization. *Operations Research*, 62(6):1358–1376, 2014.
- Guanglin Xu and Samuel Burer. A copositive approach for two-stage adjustable robust optimization with uncertain right-hand sides. *Computational Optimization and Applications*, 70(1):33–59, 2018.
- Huan Xu, Constantine Caramanis, and Shie Mannor. A distributional interpretation of robust optimization. *Mathematics of Operations Research*, 37(1):95–110, 2012.
- Bo Zeng and Long Zhao. Solving two-stage robust optimization problems using a column-and-constraint generation method. *Operations Research Letters*, 41(5):457–461, 2013.
- Jianzhe Zhen, Dick den Hertog, and Melvyn Sim. Adjustable robust optimization via fourier–motzkin elimination. *Operations Research*, 66(4):1086–1100, 2018.

## Appendix A: Verifying Assumption 3 in Examples

In this appendix, we show that each multi-stage stochastic linear optimization problem considered in this paper satisfies Assumption 3.

### A.1. Example 1 from Section 3

Consider the sample robust optimization problem

$$\begin{aligned} & \underset{x_2: \mathbb{R} \rightarrow \mathbb{R}, x_3: \mathbb{R}^2 \rightarrow \mathbb{R}}{\text{minimize}} && \frac{1}{N} \sum_{j=1}^N \sup_{\zeta \in \mathcal{U}_N^j} \{x_2(\zeta_1) + 2x_3(\zeta_1, \zeta_2)\} \\ & \text{subject to} && x_2(\zeta_1) + x_3(\zeta_1, \zeta_2) \geq \zeta_1 + \zeta_2 \quad \forall \zeta \in \cup_{j=1}^N \mathcal{U}_N^j \\ & && x_2(\zeta_1), x_3(\zeta_1, \zeta_2) \geq 0 \quad \forall \zeta \in \cup_{j=1}^N \mathcal{U}_N^j. \end{aligned}$$

We observe that the decisions must be nonnegative for every realization in the uncertainty sets. Moreover, the following constraints can be added to the above problem without affecting its optimal cost:

$$\begin{aligned} x_2(\zeta_1) &\leq \sup_{\zeta' \in \cup_{j=1}^N \mathcal{U}_N^j} \{\zeta'_1 + \zeta'_2\} \quad \forall \zeta \in \cup_{j=1}^N \mathcal{U}_N^j, \\ x_3(\zeta_1, \zeta_2) &\leq \sup_{\zeta' \in \cup_{j=1}^N \mathcal{U}_N^j} \{\zeta'_1 + \zeta'_2\} \quad \forall \zeta \in \cup_{j=1}^N \mathcal{U}_N^j. \end{aligned}$$

Indeed, the above constraints ensure that we are never purchasing inventory which exceeds the maximal  $\zeta_1 + \zeta_2$  which can be realized in the uncertainty sets. Thus, we have shown that Assumption 3 holds.

### A.2. Example 2 from Section 4.3

Consider the sample robust optimization problem

$$\begin{aligned} & \underset{x_1 \in \mathbb{Z}}{\text{minimize}} && x_1 \\ & \text{subject to} && x_1 \geq \zeta_1 \quad \forall \zeta_1 \in \cup_{j=1}^N \mathcal{U}_N^j. \end{aligned}$$

We observe that an optimal solution to this problem is  $x_1 = \max_{\zeta_1 \in \cup_{j=1}^N \mathcal{U}_N^j} \lceil \zeta_1 \rceil$ , and thus the constraint

$$x_1 \leq \max_{\zeta_1 \in \cup_{j=1}^N \mathcal{U}_N^j} \zeta_1 + 1$$

can be added to the above problem without affecting its optimal cost. We conclude that Assumption 3 holds.

### A.3. Example 3 from Section 4.3

Consider the sample robust optimization problem

$$\begin{aligned} & \underset{x_1 \in \mathbb{R}^2}{\text{minimize}} && x_{12} \\ & \text{subject to} && \zeta_1(1 - x_{12}) \leq x_{11} \quad \forall \zeta_1 \in \cup_{j=1}^N \mathcal{U}_N^j \\ & && 0 \leq x_{12} \leq 1. \end{aligned}$$

We observe that an optimal solution to this problem is given by  $x_{11} = \max_{\zeta_1 \in \cup_{j=1}^N \mathcal{U}_N^j} \zeta_1$  and  $x_{12} = 0$ . Thus, the constraint

$$x_{11} \leq \max_{\zeta_1 \in \cup_{j=1}^N \mathcal{U}_N^j} \zeta_1$$

can be added to the above problem without affecting its optimal cost. We conclude that Assumption 3 holds.

#### A.4. Example 4 from Section 4.3

Consider the sample robust optimization problem

$$\begin{aligned} & \underset{x_2: \mathbb{R} \rightarrow \mathbb{Z}}{\text{minimize}} && \frac{1}{N} \sum_{j=1}^N \sup_{\zeta_1 \in \mathcal{U}_N^j} x_2(\zeta_1) \\ & \text{subject to} && x_2(\zeta_1) \geq \zeta_1 \quad \forall \zeta_1 \in \bigcup_{j=1}^N \mathcal{U}_N^j. \end{aligned}$$

Since  $\Xi = [0, 1]$ , we observe that the constraint

$$x_2(\zeta_1) \leq 1 \quad \forall \zeta_1 \in \bigcup_{j=1}^N \mathcal{U}_N^j$$

can be added to the above problem without affecting its optimal cost. We conclude that Assumption 3 holds.

#### A.5. Inventory example from Section 7.2

Consider the sample robust optimization problem

$$\begin{aligned} & \underset{\mathbf{x}, \mathbf{I}, \mathbf{y}}{\text{minimize}} && \frac{1}{N} \sum_{j=1}^T \sup_{\zeta \in \mathcal{U}_N^j} \sum_{t=1}^T (c_t x_t(\zeta_1, \dots, \zeta_{t-1}) + y_{t+1}(\zeta_1, \dots, \zeta_t)) \\ & \text{subject to} && I_{t+1}(\zeta_1, \dots, \zeta_t) = I_t(\zeta_1, \dots, \zeta_{t-1}) + x_t(\zeta_1, \dots, \zeta_{t-1}) - \zeta_t \quad \forall \zeta \in \bigcup_{j=1}^N \mathcal{U}_N^j, \forall t \in [T] \\ & && y_{t+1}(\zeta_1, \dots, \zeta_t) \geq h_t I_{t+1}(\zeta_1, \dots, \zeta_t) \quad \forall \zeta \in \bigcup_{j=1}^N \mathcal{U}_N^j, \forall t \in [T] \\ & && y_{t+1}(\zeta_1, \dots, \zeta_t) \geq -b_t I_{t+1}(\zeta_1, \dots, \zeta_t) \quad \forall \zeta \in \bigcup_{j=1}^N \mathcal{U}_N^j, \forall t \in [T] \\ & && 0 \leq x_t(\zeta_1, \dots, \zeta_{t-1}) \leq \bar{x}_t \quad \forall \zeta \in \bigcup_{j=1}^N \mathcal{U}_N^j, \forall t \in [T], \end{aligned}$$

where  $I_1 = 0$  and  $\Xi = \mathbb{R}_+^T$ . For any feasible decision rule to the above problem and for each stage  $t$ , we observe that the following constraint is satisfied:

$$- \sup_{\zeta' \in \bigcup_{j=1}^N \mathcal{U}_N^j} \sum_{s=1}^T \zeta'_s \leq I_{t+1}(\zeta_1, \dots, \zeta_t) \leq \sum_{s=1}^T \bar{x}_s \quad \forall \zeta \in \bigcup_{j=1}^N \mathcal{U}_N^j.$$

Moreover, we can without loss of generality impose the constraint that

$$0 \leq y_{t+1}(\zeta_1, \dots, \zeta_t) = \max \{h_t I_{t+1}(\zeta_1, \dots, \zeta_t), -b_t I_{t+1}(\zeta_1, \dots, \zeta_t)\} \quad \forall \zeta \in \bigcup_{j=1}^N \mathcal{U}_N^j.$$

Applying the aforementioned bounds on  $I_{t+1}(\zeta_1, \dots, \zeta_t)$  over the uncertainty sets, we conclude that Assumption 3 holds.

## Appendix B: Proof of Theorem 1 from Section 4.2

In this appendix, we present the proof of Theorem 1. The theorem consists of asymptotic lower and upper bounds on the optimal cost of Problem (2), and we will address the proofs of the two bounds separately.

We first present the proof of the lower bound, which utilizes Theorem 2 from Section 4.2 and Theorem 3 from Section 4.4.

**THEOREM 1A.** *Suppose Assumptions 1, 2, and 3 hold. Then,  $\mathbb{P}^\infty$ -almost surely we have*

$$J \leq \liminf_{N \rightarrow \infty} \widehat{J}_N.$$

*Proof.* Recall from Assumption 1 that  $b := \mathbb{E}[\exp(\|\xi\|^a)] < \infty$  for some  $a > 1$ , and let  $L \geq 0$  be the constant from Assumption 3. Then,

$$\sum_{N=1}^{\infty} \mathbb{P}^N \left( \sup_{\zeta \in \bigcup_{j=1}^N \mathcal{U}_N^j} L(1 + \|\zeta\|) > \log N \right) = \sum_{N=1}^{\infty} \mathbb{P}^N \left( \max_{j \in [N]} \left\{ L(1 + \|\hat{\xi}^j\| + \epsilon_N) \right\} > \log N \right) \quad (\text{EC.1})$$

$$\leq \sum_{N=1}^{\infty} \mathbb{N}\mathbb{P}(L(1 + \|\boldsymbol{\xi}\| + \epsilon_N) > \log N) \quad (\text{EC.2})$$

$$= \sum_{N=1}^{\infty} \mathbb{N}\mathbb{P}\left(\|\boldsymbol{\xi}\| > \frac{\log N}{L} - 1 - \epsilon_N\right)$$

$$= \sum_{N=1}^{\infty} \mathbb{N}\mathbb{P}\left(\exp(\|\boldsymbol{\xi}\|^a) > \exp\left(\left(\frac{\log N}{L} - 1 - \epsilon_N\right)^a\right)\right)$$

$$\leq \sum_{N=1}^{\infty} \frac{Nb}{\exp\left(\left(\frac{\log N}{L} - 1 - \epsilon_N\right)^a\right)} \quad (\text{EC.3})$$

$$< \infty, \quad (\text{EC.4})$$

where (EC.1) follows from the definition of the uncertainty sets, (EC.2) follows from the union bound, (EC.3) follows from Markov's inequality, and (EC.4) follows from  $a > 1$  and  $\epsilon_N \rightarrow 0$ . Therefore, the Borel-Cantelli lemma and Assumption 3 imply that the following equality holds for all sufficiently large  $N \in \mathbb{N}$ ,  $\mathbb{P}^\infty$ -almost surely:

$$\begin{aligned} \widehat{J}_N = \quad & \underset{\mathbf{x} \in \mathcal{X}}{\text{minimize}} \quad \frac{1}{N} \sum_{j=1}^N \sup_{\boldsymbol{\zeta} \in \mathcal{U}_N^j} \sum_{t=1}^T \mathbf{c}_t(\boldsymbol{\zeta}) \cdot \mathbf{x}_t(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_{t-1}) \\ & \text{subject to} \quad \sum_{t=1}^T \mathbf{A}_t(\boldsymbol{\zeta}) \mathbf{x}_t(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_{t-1}) \leq \mathbf{b}(\boldsymbol{\zeta}) \quad \forall \boldsymbol{\zeta} \in \cup_{j=1}^N \mathcal{U}_N^j \\ & \quad \quad \quad \|\mathbf{x}_t(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_{t-1})\| \leq \log N \quad \forall \boldsymbol{\zeta} \in \cup_{j=1}^N \mathcal{U}_N^j, t. \end{aligned} \quad (\text{EC.5})$$

Moreover, since  $\mathbf{c}_1(\boldsymbol{\zeta}) \in \mathbb{R}^{n_1}, \dots, \mathbf{c}_T(\boldsymbol{\zeta}) \in \mathbb{R}^{n_T}$  are affine functions of the stochastic process, it follows from identical reasoning as (EC.1)-(EC.4) and the equivalence of  $\ell_p$ -norms in finite-dimensional spaces that  $\sup_{\boldsymbol{\zeta} \in \cup_{j=1}^N \mathcal{U}_N^j} \|\mathbf{c}_t(\boldsymbol{\zeta})\|_* \leq \log N$  for all sufficiently large  $N \in \mathbb{N}$ ,  $\mathbb{P}^\infty$ -almost surely.

We now apply Theorem 2 to obtain an asymptotic lower bound on the optimization problem in (EC.5). Indeed, let  $M_N$  be shorthand for  $N^{-\frac{1}{(d+1)(d+2)}} \log N$ . Then, for all sufficiently large  $N \in \mathbb{N}$ ,  $\mathbb{P}^\infty$ -almost surely, and for any decision rule  $\mathbf{x} \in \mathcal{X}$  which is feasible for the optimization problem in (EC.5),

$$\begin{aligned} & \mathbb{E} \left[ \sum_{t=1}^T \mathbf{c}_t(\boldsymbol{\xi}) \cdot \mathbf{x}_t(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}) \mathbb{I} \{ \boldsymbol{\xi} \in \cup_{j=1}^N \mathcal{U}_N^j \} \right] \\ & \leq \frac{1}{N} \sum_{j=1}^N \sup_{\boldsymbol{\zeta} \in \mathcal{U}_N^j} \sum_{t=1}^T \mathbf{c}_t(\boldsymbol{\zeta}) \cdot \mathbf{x}_t(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_{t-1}) + M_N \sup_{\boldsymbol{\zeta} \in \cup_{j=1}^N \mathcal{U}_N^j} \left| \sum_{t=1}^T \mathbf{c}_t(\boldsymbol{\zeta}) \cdot \mathbf{x}_t(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_{t-1}) \right| \\ & \leq \frac{1}{N} \sum_{j=1}^N \sup_{\boldsymbol{\zeta} \in \mathcal{U}_N^j} \sum_{t=1}^T \mathbf{c}_t(\boldsymbol{\zeta}) \cdot \mathbf{x}_t(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_{t-1}) + M_N \sum_{t=1}^T \sup_{\boldsymbol{\zeta} \in \cup_{j=1}^N \mathcal{U}_N^j} \|\mathbf{c}_t(\boldsymbol{\zeta})\|_* \|\mathbf{x}_t(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_{t-1})\| \\ & \leq \frac{1}{N} \sum_{j=1}^N \sup_{\boldsymbol{\zeta} \in \mathcal{U}_N^j} \sum_{t=1}^T \mathbf{c}_t(\boldsymbol{\zeta}) \cdot \mathbf{x}_t(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_{t-1}) + TM_N (\log N)^2, \end{aligned}$$

where the first inequality follows from Theorem 2, the second inequality follows from the triangle inequality and the Cauchy-Schwartz inequality, and the third inequality follows because  $\|\mathbf{c}_t(\boldsymbol{\zeta})\|_* \leq \log N$  and  $\|\mathbf{x}_t(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_T)\| \leq \log N$  for all sufficiently large  $N \in \mathbb{N}$  and all realizations in the uncertainty sets. We remark that the above bound holds uniformly for all decision rules which are feasible for the optimization problem in (EC.5). Therefore, we have shown that the following inequality holds for all sufficiently large  $N \in \mathbb{N}$ ,  $\mathbb{P}^\infty$ -almost surely:

$$\begin{aligned} \widehat{J}_N + TM_N (\log N)^2 \geq \quad & \underset{\mathbf{x} \in \mathcal{X}}{\text{minimize}} \quad \mathbb{E} \left[ \sum_{t=1}^T \mathbf{c}_t(\boldsymbol{\xi}) \cdot \mathbf{x}_t(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}) \mathbb{I} \{ \boldsymbol{\xi} \in \cup_{j=1}^N \mathcal{U}_N^j \} \right] \\ & \text{subject to} \quad \sum_{t=1}^T \mathbf{A}_t(\boldsymbol{\zeta}) \mathbf{x}_t(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_{t-1}) \leq \mathbf{b}(\boldsymbol{\zeta}) \quad \forall \boldsymbol{\zeta} \in \cup_{j=1}^N \mathcal{U}_N^j \\ & \quad \quad \quad \|\mathbf{x}_t(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_{t-1})\| \leq \log N \quad \forall \boldsymbol{\zeta} \in \cup_{j=1}^N \mathcal{U}_N^j, t. \end{aligned}$$

Next, we obtain a lower bound on the right side of the above inequality by removing the last row of constraints and relaxing  $\cup_{j=1}^N \mathcal{U}_N^j$  to any set which contains the stochastic process with sufficiently high probability:

$$\begin{aligned} & \underset{\mathbf{x} \in \mathcal{X}, \tilde{S} \subseteq \Xi}{\text{minimize}} && \mathbb{E} \left[ \sum_{t=1}^T \mathbf{c}_t(\boldsymbol{\xi}) \cdot \mathbf{x}_t(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}) \mathbb{I} \left\{ \boldsymbol{\xi} \in \tilde{S} \right\} \right] \\ & \text{subject to} && \sum_{t=1}^T \mathbf{A}_t(\boldsymbol{\zeta}) \mathbf{x}_t(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_{t-1}) \leq \mathbf{b}(\boldsymbol{\zeta}) \quad \forall \boldsymbol{\zeta} \in \tilde{S} \\ & && \mathbb{P} \left( \boldsymbol{\xi} \in \tilde{S} \right) \geq \mathbb{P} \left( \boldsymbol{\xi} \in \cup_{j=1}^N \mathcal{U}_N^j \right). \end{aligned} \tag{EC.6}$$

Finally, for any fixed  $\rho \in (0, 1)$ , it follows from Theorem 3 that  $\mathbb{P}(\boldsymbol{\xi} \in \cup_{j=1}^N \mathcal{U}_N^j \cap S) \geq 1 - \rho$  for all sufficiently large  $N \in \mathbb{N}$ ,  $\mathbb{P}^\infty$ -almost surely.<sup>4</sup> Furthermore, we observe that  $TM_N(\log N)^2$  converges to zero as the number of sample paths  $N$  tends to infinity. Therefore, we have shown that the following inequality holds,  $\mathbb{P}^\infty$ -almost surely:

$$\begin{aligned} \liminf_{N \rightarrow \infty} \widehat{J}_N & \geq \underset{\mathbf{x} \in \mathcal{X}, \tilde{S} \subseteq \Xi}{\text{minimize}} && \mathbb{E} \left[ \sum_{t=1}^T \mathbf{c}_t(\boldsymbol{\xi}) \cdot \mathbf{x}_t(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}) \mathbb{I} \left\{ \boldsymbol{\xi} \in \tilde{S} \right\} \right] \\ & \text{subject to} && \sum_{t=1}^T \mathbf{A}_t(\boldsymbol{\zeta}) \mathbf{x}_t(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_{t-1}) \leq \mathbf{b}(\boldsymbol{\zeta}) \quad \forall \boldsymbol{\zeta} \in \tilde{S} \\ & && \mathbb{P} \left( \boldsymbol{\xi} \in \tilde{S} \right) \geq 1 - \rho. \end{aligned}$$

Since the inequality holds true for every  $\rho \in (0, 1)$ , and since the optimal cost of the above optimization problem is monotone in  $\rho$ , we obtain the desired result.  $\square$

We now conclude the proof of Theorem 1 by establishing its upper bound.

**THEOREM 1B.** *Suppose Assumption 2 holds. Then,  $\mathbb{P}^\infty$ -almost surely we have*

$$\limsup_{N \rightarrow \infty} \widehat{J}_N \leq \bar{J}.$$

*Proof.* Consider any  $\rho > 0$  such that there is a decision rule  $\mathbf{x} \in \mathcal{X}$  which satisfies

$$\bar{\mathbb{E}} \left[ \sum_{t=1}^T \mathbf{c}_t(\boldsymbol{\xi}) \cdot \mathbf{x}_t(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}) \right] < \infty, \tag{EC.7}$$

$$\sum_{t=1}^T \mathbf{A}_t(\boldsymbol{\zeta}) \mathbf{x}_t(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_{t-1}) \leq \mathbf{b}(\boldsymbol{\zeta}) \quad \forall \boldsymbol{\zeta} \in \Xi: \text{dist}(\boldsymbol{\zeta}, S) \leq \rho. \tag{EC.8}$$

Indeed, if no such  $\rho > 0$  and  $\mathbf{x} \in \mathcal{X}$  existed, then  $\bar{J} = \infty$  and the desired result follows immediately. We recall from Assumption 2 that  $\epsilon_N \rightarrow 0$  as  $N \rightarrow \infty$ . Therefore,

$$\begin{aligned} \limsup_{N \rightarrow \infty} \widehat{J}_N & \leq \limsup_{N \rightarrow \infty} \frac{1}{N} \sum_{j=1}^N \sup_{\boldsymbol{\zeta} \in \Xi: \|\boldsymbol{\zeta} - \hat{\boldsymbol{\xi}}^j\| \leq \epsilon_N} \sum_{t=1}^T \mathbf{c}_t(\boldsymbol{\zeta}) \cdot \mathbf{x}_t(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_{t-1}) \\ & \leq \lim_{k \rightarrow \infty} \limsup_{N \rightarrow \infty} \frac{1}{N} \sum_{j=1}^N \sup_{\boldsymbol{\zeta} \in \Xi: \|\boldsymbol{\zeta} - \hat{\boldsymbol{\xi}}^j\| \leq \epsilon_k} \sum_{t=1}^T \mathbf{c}_t(\boldsymbol{\zeta}) \cdot \mathbf{x}_t(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_{t-1}) \\ & = \lim_{k \rightarrow \infty} \mathbb{E} \left[ \sup_{\boldsymbol{\zeta} \in \Xi: \|\boldsymbol{\zeta} - \boldsymbol{\xi}\| \leq \epsilon_k} \sum_{t=1}^T \mathbf{c}_t(\boldsymbol{\zeta}) \cdot \mathbf{x}_t(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_{t-1}) \right] \quad \mathbb{P}^\infty\text{-almost surely} \\ & = \bar{\mathbb{E}} \left[ \sum_{t=1}^T \mathbf{c}_t(\boldsymbol{\xi}) \cdot \mathbf{x}_t(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}) \right]. \end{aligned} \tag{EC.9}$$

<sup>4</sup> We remark that [Devroye and Wise \(1980, Theorem 2\)](#) could be used here in lieu of Theorem 3. Our primary use of Theorem 3 is in the proof of Theorem 2.



The first inequality holds because the decision rule is feasible but possibly suboptimal for Problem (2) for all  $N \geq \min\{\bar{N} : \epsilon_{\bar{N}} \leq \rho\}$ . The second inequality follows because  $\epsilon_k \rightarrow 0$  monotonically as  $k \rightarrow \infty$ . The first equality follows from the strong law of large numbers (Erickson 1973), which holds since (EC.7) ensures that

$$\mathbb{E} \left[ \max \left\{ \sup_{\zeta \in \Xi: \|\zeta - \hat{\xi}\| \leq \epsilon_k} \sum_{t=1}^T \mathbf{c}_t(\zeta) \cdot \mathbf{x}_t(\zeta_1, \dots, \zeta_{t-1}), 0 \right\} \right] < \infty$$

for all sufficiently large  $k$ . The final equality follows the definition of the local upper semicontinuous envelope. Since the set of decision rules which satisfy (EC.8) does not get smaller as  $\rho \downarrow 0$ , we conclude that the following holds  $\mathbb{P}^\infty$ -almost surely:

$$\begin{aligned} \limsup_{N \rightarrow \infty} \widehat{J}_N \leq & \lim_{\rho \downarrow 0} \underset{\mathbf{x} \in \mathcal{X}}{\text{minimize}} && \mathbb{E} \left[ \sum_{t=1}^T \mathbf{c}_t(\boldsymbol{\xi}) \cdot \mathbf{x}_t(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}) \right] \\ & \text{subject to} && \sum_{t=1}^T \mathbf{A}_t(\zeta) \mathbf{x}_t(\zeta_1, \dots, \zeta_{t-1}) \leq \mathbf{b}(\zeta) \quad \forall \zeta \in \Xi: \text{dist}(\zeta, S) \leq \rho. \end{aligned}$$

This concludes the proof.  $\square$

## Appendix C: Proof of Theorem 2 from Section 4.2

In this appendix, we present the proof of Theorem 2. The proof is organized as follows. In Appendix C.1, we first develop a helpful intermediary bound (Lemma EC.2). In Appendix C.2, we use that bound to prove Theorem 2. In Appendix C.3, we provide for completeness the proofs of some miscellaneous and rather technical results that were used in Appendix C.2.

### C.1. An intermediary result

Our proof of Theorem 2 relies on an intermediary result (Lemma EC.2), which establishes a relationship between sample robust optimization and distributionally robust optimization with the 1-Wasserstein ambiguity set. We begin by establishing the relationship for the specific case where there is a single data point.

LEMMA EC.1. *Let  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  be measurable,  $\mathcal{Z} \subseteq \mathbb{R}^d$ , and  $\hat{\xi} \in \mathcal{Z}$ . If  $\theta_2 \geq 2\theta_1 \geq 0$ , then*

$$\sup_{\mathbb{Q} \in \mathcal{P}(\mathcal{Z}): \mathbb{E}_{\mathbb{Q}}[\|\boldsymbol{\xi} - \hat{\xi}\|] \leq \theta_1} \mathbb{E}_{\mathbb{Q}}[f(\boldsymbol{\xi})] \leq \sup_{\zeta \in \mathcal{Z}: \|\zeta - \hat{\xi}\| \leq \theta_2} f(\zeta) + \frac{2\theta_1}{\theta_2} \sup_{\zeta \in \mathcal{Z}} |f(\zeta)|. \quad (\text{EC.10})$$

*Proof.* We first apply the Richter-Rogonsinski Theorem (see Theorem 7.32 and Proposition 6.40 of Shapiro et al. (2009)), which says that a distributionally robust optimization problem with  $m$  moment constraints is equivalent to optimizing a weighted average of  $m+1$  points. Thus,

$$\begin{aligned} \sup_{\mathbb{Q} \in \mathcal{P}(\mathcal{Z}): \mathbb{E}_{\mathbb{Q}}[\|\boldsymbol{\xi} - \hat{\xi}\|] \leq \theta_1} \mathbb{E}_{\mathbb{Q}}[f(\boldsymbol{\xi})] &= \begin{cases} \sup_{\zeta^1, \zeta^2 \in \mathcal{Z}, \lambda \in [0,1]} \lambda f(\zeta^1) + (1-\lambda)f(\zeta^2) \\ \text{subject to} & \lambda \|\zeta^1 - \hat{\xi}\| + (1-\lambda) \|\zeta^2 - \hat{\xi}\| \leq \theta_1 \end{cases} \\ &\leq \begin{cases} \sup_{\zeta^1, \zeta^2 \in \mathcal{Z}, \lambda \in [0,1]} \lambda f(\zeta^1) + (1-\lambda)f(\zeta^2) \\ \text{subject to} & \lambda \|\zeta^1 - \hat{\xi}\| \leq \theta_1, (1-\lambda) \|\zeta^2 - \hat{\xi}\| \leq \theta_1, \end{cases} \end{aligned} \quad (\text{EC.11})$$

where the inequality follows from relaxing the constraints on  $\zeta^1$  and  $\zeta^2$ . Let us assume from this point onward that  $\sup_{\zeta \in \mathcal{Z}} |f(\zeta)| < \infty$ ; indeed, if  $\sup_{\zeta \in \mathcal{Z}} |f(\zeta)| = \infty$ , then the inequality in (EC.10) would trivially hold since the right-hand side would equal infinity. Then, it follows from (EC.11) that

$$\sup_{\mathbb{Q} \in \mathcal{P}(\mathcal{Z}): \mathbb{E}_{\mathbb{Q}}[\|\boldsymbol{\xi} - \hat{\xi}\|] \leq \theta_1} \mathbb{E}_{\mathbb{Q}}[f(\boldsymbol{\xi})] \leq \sup_{0 \leq \lambda \leq 1} \left\{ \lambda \left( \sup_{\zeta \in \mathcal{Z}: \|\zeta - \hat{\xi}\| \leq \frac{\theta_1}{\lambda}} f(\zeta) \right) + (1-\lambda) \left( \sup_{\zeta \in \mathcal{Z}: \|\zeta - \hat{\xi}\| \leq \frac{\theta_1}{1-\lambda}} f(\zeta) \right) \right\}. \quad (\text{EC.12})$$

We observe that the supremum over  $0 \leq \lambda \leq 1$  is symmetric with respect to  $\lambda$ , in the sense that  $\lambda$  can be restricted to  $[0, \frac{1}{2}]$  or  $[\frac{1}{2}, 1]$  without loss of generality. Moreover, under the assumption that  $\theta_2 \geq 2\theta_1$ , the interval  $[0, 1 - \frac{\theta_1}{\theta_2}]$  is a superset of the interval  $[0, \frac{1}{2}]$ . Combining these arguments, we conclude that the right side of (EC.12) is equal to

$$\sup_{0 \leq \lambda \leq 1 - \frac{\theta_1}{\theta_2}} \left\{ \lambda \left( \sup_{\zeta \in \mathcal{Z}: \|\zeta - \hat{\xi}\| \leq \frac{\theta_1}{\lambda}} f(\zeta) \right) + (1 - \lambda) \left( \sup_{\zeta \in \mathcal{Z}: \|\zeta - \hat{\xi}\| \leq \frac{\theta_1}{1-\lambda}} f(\zeta) \right) \right\}. \quad (\text{EC.13})$$

Next, we observe that  $\frac{\theta_1}{1-\lambda} \leq \theta_2$  for every feasible  $\lambda$  for the above optimization problem. Using this inequality, we obtain the following upper bound:

$$\begin{aligned} & \sup_{\mathbb{Q} \in \mathcal{P}(\mathcal{Z}): \mathbb{E}_{\mathbb{Q}}[\|\xi - \hat{\xi}\|] \leq \theta_1} \mathbb{E}_{\mathbb{Q}}[f(\xi)] \\ & \leq \sup_{0 \leq \lambda \leq 1 - \frac{\theta_1}{\theta_2}} \left\{ \lambda \left( \sup_{\zeta \in \mathcal{Z}: \|\zeta - \hat{\xi}\| \leq \frac{\theta_1}{\lambda}} f(\zeta) \right) + (1 - \lambda) \left( \sup_{\zeta \in \mathcal{Z}: \|\zeta - \hat{\xi}\| \leq \theta_2} f(\zeta) \right) \right\} \\ & = \sup_{\zeta \in \mathcal{Z}: \|\zeta - \hat{\xi}\| \leq \theta_2} f(\zeta) + \sup_{0 \leq \lambda \leq 1 - \frac{\theta_1}{\theta_2}} \left\{ \lambda \left( \sup_{\zeta \in \mathcal{Z}: \|\zeta - \hat{\xi}\| \leq \frac{\theta_1}{\lambda}} f(\zeta) - \sup_{\zeta \in \mathcal{Z}: \|\zeta - \hat{\xi}\| \leq \theta_2} f(\zeta) \right) \right\}, \end{aligned} \quad (\text{EC.14})$$

where the above equality comes from rearranging terms. For every  $\frac{\theta_1}{\theta_2} \leq \lambda \leq 1 - \frac{\theta_1}{\theta_2}$ , it immediately follows from  $\frac{\theta_1}{\lambda} \leq \theta_2$  that

$$\sup_{\zeta \in \mathcal{Z}: \|\zeta - \hat{\xi}\| \leq \frac{\theta_1}{\lambda}} f(\zeta) - \sup_{\zeta \in \mathcal{Z}: \|\zeta - \hat{\xi}\| \leq \theta_2} f(\zeta) \leq 0,$$

and the above holds at equality when  $\lambda = \frac{\theta_1}{\theta_2}$ . Therefore,

$$\begin{aligned} & \sup_{0 \leq \lambda \leq 1 - \frac{\theta_1}{\theta_2}} \left\{ \lambda \left( \sup_{\zeta \in \mathcal{Z}: \|\zeta - \hat{\xi}\| \leq \frac{\theta_1}{\lambda}} f(\zeta) - \sup_{\zeta \in \mathcal{Z}: \|\zeta - \hat{\xi}\| \leq \theta_2} f(\zeta) \right) \right\} \\ & = \sup_{0 \leq \lambda \leq \frac{\theta_1}{\theta_2}} \left\{ \lambda \left( \sup_{\zeta \in \mathcal{Z}: \|\zeta - \hat{\xi}\| \leq \frac{\theta_1}{\lambda}} f(\zeta) - \sup_{\zeta \in \mathcal{Z}: \|\zeta - \hat{\xi}\| \leq \theta_2} f(\zeta) \right) \right\} \end{aligned} \quad (\text{EC.15})$$

$$\leq \sup_{0 \leq \lambda \leq \frac{\theta_1}{\theta_2}} \left\{ \lambda \left( \sup_{\zeta \in \mathcal{Z}} f(\zeta) - \inf_{\zeta \in \mathcal{Z}} f(\zeta) \right) \right\} \quad (\text{EC.16})$$

$$\leq \frac{2\theta_1}{\theta_2} \sup_{\zeta \in \mathcal{Z}} |f(\zeta)|. \quad (\text{EC.17})$$

Line (EC.15) follows because we can without loss of generality restrict  $\lambda$  to the interval  $[0, \frac{\theta_1}{\theta_2}]$ . Line (EC.16) is obtained by applying the global lower and upper bounds on  $f(\zeta)$ . Finally, we obtain (EC.17) since

$$0 \leq \sup_{\zeta \in \mathcal{Z}} f(\zeta) - \inf_{\zeta \in \mathcal{Z}} f(\zeta) \leq 2 \sup_{\zeta \in \mathcal{Z}} |f(\zeta)|.$$

Combining (EC.14) and (EC.17), we obtain the desired result.  $\square$

We now extend the previous lemma to the general case with more than one data point. In the following, we let  $\hat{\mathbb{P}}_N$  denote the empirical distribution of historical data  $\hat{\xi}^1, \dots, \hat{\xi}^N \in \mathbb{R}^d$ ,  $\mathcal{Z} \subseteq \mathbb{R}^d$  be any set that contains the historical data, and  $d_1(\cdot, \cdot)$  denote the 1-Wasserstein distance between two probability distributions (see Section 6).

LEMMA EC.2. *Let  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  be measurable,  $\mathcal{Z} \subseteq \mathbb{R}^d$ , and  $\hat{\xi}^1, \dots, \hat{\xi}^N \in \mathcal{Z}$ . If  $\theta_2 \geq 2\theta_1 \geq 0$ , then*

$$\sup_{\mathbb{Q} \in \mathcal{P}(\mathcal{Z}): d_1(\mathbb{Q}, \hat{\mathbb{P}}_N) \leq \theta_1} \mathbb{E}_{\mathbb{Q}}[f(\xi)] \leq \frac{1}{N} \sum_{j=1}^N \sup_{\zeta \in \mathcal{Z}: \|\zeta - \hat{\xi}^j\| \leq \theta_2} f(\zeta) + \frac{4\theta_1}{\theta_2} \sup_{\zeta \in \mathcal{Z}} |f(\zeta)|.$$

*Proof.* We recall from the proof of [Mohajerin Esfahani and Kuhn \(2018, Theorem 4.2\)](#) that

$$\left\{ \mathbb{Q} \in \mathcal{P}(\mathcal{Z}) : d_1(\mathbb{Q}, \hat{\mathbb{P}}_N) \leq \theta_1 \right\} = \left\{ \frac{1}{N} \sum_{j=1}^N \mathbb{Q}_j : \frac{1}{N} \sum_{j=1}^N \mathbb{E}_{\mathbb{Q}_j} [\|\boldsymbol{\xi} - \hat{\boldsymbol{\xi}}^j\|] \leq \theta_1, \mathbb{Q}_1, \dots, \mathbb{Q}_N \in \mathcal{P}(\mathcal{Z}) \right\}.$$

Therefore,

$$\sup_{\mathbb{Q} \in \mathcal{P}(\mathcal{Z}) : d_1(\mathbb{Q}, \hat{\mathbb{P}}_N) \leq \theta_1} \mathbb{E}_{\mathbb{Q}} [f(\boldsymbol{\xi})] = \sup_{\boldsymbol{\gamma} \in \mathbb{R}_+^N} \left\{ \frac{1}{N} \sum_{j=1}^N \sup_{\mathbb{Q} \in \mathcal{P}(\mathcal{Z}) : \mathbb{E}_{\mathbb{Q}} [\|\boldsymbol{\xi} - \hat{\boldsymbol{\xi}}^j\|] \leq \gamma_j} \mathbb{E}_{\mathbb{Q}} [f(\boldsymbol{\xi})] : \frac{1}{N} \sum_{j=1}^N \gamma_j \leq \theta_1 \right\}. \quad (\text{EC.18})$$

For any choice of  $\boldsymbol{\gamma} \in \mathbb{R}_+^N$ , we can partition the components  $\gamma_j$  into those that satisfy  $2\gamma_j \leq \theta_2$  and  $2\gamma_j > \theta_2$ . Thus,

$$\begin{aligned} & \sup_{\mathbb{Q} \in \mathcal{P}(\mathcal{Z}) : d_1(\mathbb{Q}, \hat{\mathbb{P}}_N) \leq \theta_1} \mathbb{E}_{\mathbb{Q}} [f(\boldsymbol{\xi})] \\ & \leq \sup_{\boldsymbol{\gamma} \in \mathbb{R}_+^N} \left\{ \frac{1}{N} \sum_{j \in [N] : 2\gamma_j \leq \theta_2} \sup_{\mathbb{Q} \in \mathcal{P}(\mathcal{Z}) : \mathbb{E}_{\mathbb{Q}} [\|\boldsymbol{\xi} - \hat{\boldsymbol{\xi}}^j\|] \leq \gamma_j} \mathbb{E}_{\mathbb{Q}} [f(\boldsymbol{\xi})] + \frac{1}{N} \sum_{j \in [N] : 2\gamma_j > \theta_2} \sup_{\boldsymbol{\zeta} \in \mathcal{Z}} |f(\boldsymbol{\zeta})| : \frac{1}{N} \sum_{j=1}^N \gamma_j \leq \theta_1 \right\}, \quad (\text{EC.19}) \end{aligned}$$

where the inequality follows from upper bounding each of the inner distributionally robust optimization problems for which  $2\gamma_j > \theta_2$  by  $\sup_{\boldsymbol{\zeta} \in \mathcal{Z}} |f(\boldsymbol{\zeta})|$ . Due to the constraints on  $\boldsymbol{\gamma}$ , there can be at most  $\frac{2N\theta_1}{\theta_2}$  components which satisfy  $2\gamma_j > \theta_2$ . It thus follows from [\(EC.19\)](#) that

$$\begin{aligned} & \sup_{\mathbb{Q} \in \mathcal{P}(\mathcal{Z}) : d_1(\mathbb{Q}, \hat{\mathbb{P}}_N) \leq \theta_1} \mathbb{E}_{\mathbb{Q}} [f(\boldsymbol{\xi})] \\ & \leq \sup_{\boldsymbol{\gamma} \in \mathbb{R}_+^N} \left\{ \frac{1}{N} \sum_{j \in [N] : 2\gamma_j \leq \theta_2} \sup_{\mathbb{Q} \in \mathcal{P}(\mathcal{Z}) : \mathbb{E}_{\mathbb{Q}} [\|\boldsymbol{\xi} - \hat{\boldsymbol{\xi}}^j\|] \leq \gamma_j} \mathbb{E}_{\mathbb{Q}} [f(\boldsymbol{\xi})] : \frac{1}{N} \sum_{j=1}^N \gamma_j \leq \theta_1 \right\} + \frac{2\theta_1}{\theta_2} \sup_{\boldsymbol{\zeta} \in \mathcal{Z}} |f(\boldsymbol{\zeta})|. \quad (\text{EC.20}) \end{aligned}$$

To conclude the proof, we apply [Lemma EC.1](#) to each distributionally robust optimization problem in [\(EC.20\)](#) to obtain the following upper bounds:

$$\begin{aligned} & \sup_{\mathbb{Q} \in \mathcal{P}(\mathcal{Z}) : d_1(\mathbb{Q}, \hat{\mathbb{P}}_N) \leq \theta_1} \mathbb{E}_{\mathbb{Q}} [f(\boldsymbol{\xi})] \\ & \leq \sup_{\boldsymbol{\gamma} \in \mathbb{R}_+^N} \left\{ \frac{1}{N} \sum_{j \in [N] : 2\gamma_j \leq \theta_2} \left( \sup_{\boldsymbol{\zeta} \in \mathcal{Z} : \|\boldsymbol{\zeta} - \hat{\boldsymbol{\xi}}^j\| \leq \theta_2} f(\boldsymbol{\zeta}) + \frac{2\gamma_j}{\theta_2} \sup_{\boldsymbol{\zeta} \in \mathcal{Z}} |f(\boldsymbol{\zeta})| \right) : \frac{1}{N} \sum_{j=1}^N \gamma_j \leq \theta_1 \right\} + \frac{2\theta_1}{\theta_2} \sup_{\boldsymbol{\zeta} \in \mathcal{Z}} |f(\boldsymbol{\zeta})| \quad (\text{EC.21}) \end{aligned}$$

$$\leq \sup_{\boldsymbol{\gamma} \in \mathbb{R}_+^N} \left\{ \frac{1}{N} \sum_{j=1}^N \left( \sup_{\boldsymbol{\zeta} \in \mathcal{Z} : \|\boldsymbol{\zeta} - \hat{\boldsymbol{\xi}}^j\| \leq \theta_2} f(\boldsymbol{\zeta}) + \frac{2\gamma_j}{\theta_2} \sup_{\boldsymbol{\zeta} \in \mathcal{Z}} |f(\boldsymbol{\zeta})| \right) : \frac{1}{N} \sum_{j=1}^N \gamma_j \leq \theta_1 \right\} + \frac{2\theta_1}{\theta_2} \sup_{\boldsymbol{\zeta} \in \mathcal{Z}} |f(\boldsymbol{\zeta})| \quad (\text{EC.22})$$

$$= \frac{1}{N} \sum_{j=1}^N \sup_{\boldsymbol{\zeta} \in \mathcal{Z} : \|\boldsymbol{\zeta} - \hat{\boldsymbol{\xi}}^j\| \leq \theta_2} f(\boldsymbol{\zeta}) + \frac{4\theta_1}{\theta_2} \sup_{\boldsymbol{\zeta} \in \mathcal{Z}} |f(\boldsymbol{\zeta})|. \quad (\text{EC.23})$$

Line [\(EC.21\)](#) follows from applying [Lemma EC.1](#) to [\(EC.20\)](#). Line [\(EC.22\)](#) follows because

$$\sup_{\boldsymbol{\zeta} \in \mathcal{Z} : \|\boldsymbol{\zeta} - \hat{\boldsymbol{\xi}}^j\| \leq \theta_2} f(\boldsymbol{\zeta}) + \frac{2\gamma_j}{\theta_2} \sup_{\boldsymbol{\zeta} \in \mathcal{Z}} |f(\boldsymbol{\zeta})| \geq 0$$

for each component that satisfies  $2\gamma_j > \theta_2$ , and thus adding these quantities to [\(EC.21\)](#) results in an upper bound. Finally, [\(EC.23\)](#) follows from the constraint  $\frac{1}{N} \sum_{j=1}^N \gamma_j \leq \theta_1$ . This concludes the proof.  $\square$

## C.2. Proof of [Theorem 2](#)

We have established above a deterministic bound ([Lemma EC.2](#)) between sample robust optimization and distributionally robust optimization with the 1-Wasserstein ambiguity set. We will now combine that bound with a concentration inequality of [Fournier and Guillin \(2015\)](#) to prove [Theorem 2](#). We remark that the following proof will employ [Theorem 3](#) from [Section 4.4](#) as well as notation from [Section 6](#). For clarity of exposition, some intermediary and rather technical details of the following proof have been relegated to [Appendix C.3](#).

**THEOREM 2.** *If Assumptions 1 and 2 hold, then there exists a  $\bar{N} \in \mathbb{N}$ ,  $\mathbb{P}^\infty$ -almost surely, such that*

$$\mathbb{E} [f(\boldsymbol{\xi}) \mathbb{I} \{ \boldsymbol{\xi} \in \cup_{j=1}^N \mathcal{U}_N^j \}] \leq \frac{1}{N} \sum_{j=1}^N \sup_{\boldsymbol{\zeta} \in \mathcal{U}_N^j} f(\boldsymbol{\zeta}) + M_N \sup_{\boldsymbol{\zeta} \in \cup_{j=1}^N \mathcal{U}_N^j} |f(\boldsymbol{\zeta})|$$

for all  $N \geq \bar{N}$  and all measurable functions  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , where  $M_N := N^{-\frac{1}{(d+1)(d+2)}}$   $\log N$ .

*Proof.* Let  $\kappa > 0$  be the coefficient from Assumption 2, and define  $\bar{\kappa} = \kappa/8$ . For each  $N \in \mathbb{N}$ , define

$$\delta_N := \begin{cases} \bar{\kappa} N^{-\frac{1}{2}} \log N, & \text{if } d = 1, \\ \bar{\kappa} N^{-\frac{1}{d}} (\log N)^2, & \text{if } d \geq 2. \end{cases}$$

It follows from [Fournier and Guillin \(2015\)](#) and Assumption 1 that  $d_1(\mathbb{P}, \widehat{\mathbb{P}}_N) \leq \delta_N$  for all sufficiently large  $N \in \mathbb{N}$ ,  $\mathbb{P}^\infty$ -almost surely (see Lemma EC.3 in Appendix C.3). Therefore, for every measurable function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ ,

$$\begin{aligned} & \mathbb{E} [f(\boldsymbol{\xi}) \mathbb{I} \{ \boldsymbol{\xi} \in \cup_{j=1}^N \mathcal{U}_N^j \}] \\ &= \mathbb{E} \left[ \left( f(\boldsymbol{\xi}) + \sup_{\boldsymbol{\zeta} \in \cup_{j=1}^N \mathcal{U}_N^j} |f(\boldsymbol{\zeta})| \right) \mathbb{I} \{ \boldsymbol{\xi} \in \cup_{j=1}^N \mathcal{U}_N^j \} \right] - \left( \sup_{\boldsymbol{\zeta} \in \cup_{j=1}^N \mathcal{U}_N^j} |f(\boldsymbol{\zeta})| \right) \mathbb{P}(\boldsymbol{\xi} \in \cup_{j=1}^N \mathcal{U}_N^j) \\ &\leq \sup_{\mathbb{Q} \in \mathcal{P}(\Xi) : d_1(\mathbb{Q}, \widehat{\mathbb{P}}_N) \leq \delta_N} \underbrace{\mathbb{E}_{\mathbb{Q}} \left[ \left( f(\boldsymbol{\xi}) + \sup_{\boldsymbol{\zeta} \in \cup_{j=1}^N \mathcal{U}_N^j} |f(\boldsymbol{\zeta})| \right) \mathbb{I} \{ \boldsymbol{\xi} \in \cup_{j=1}^N \mathcal{U}_N^j \} \right]}_{g(\boldsymbol{\xi})} - \left( \sup_{\boldsymbol{\zeta} \in \cup_{j=1}^N \mathcal{U}_N^j} |f(\boldsymbol{\zeta})| \right) \mathbb{P}(\boldsymbol{\xi} \in \cup_{j=1}^N \mathcal{U}_N^j), \end{aligned} \tag{EC.24}$$

where the inequality holds for all sufficiently large  $N \in \mathbb{N}$ ,  $\mathbb{P}^\infty$ -almost surely. Next, we observe that  $g(\boldsymbol{\xi})$  equals zero when  $\boldsymbol{\xi}$  is not an element of  $\cup_{j=1}^N \mathcal{U}_N^j$  and is nonnegative otherwise. Therefore, without loss of generality, we can restrict the supremum over the expectation of  $g(\boldsymbol{\xi})$  to distributions with support contained in  $\cup_{j=1}^N \mathcal{U}_N^j$  (see Lemma EC.4 in Appendix C.3). Therefore, (EC.24) is equal to

$$\begin{aligned} & \sup_{\mathbb{Q} \in \mathcal{P}(\cup_{j=1}^N \mathcal{U}_N^j) : d_1(\mathbb{Q}, \widehat{\mathbb{P}}_N) \leq \delta_N} \mathbb{E}_{\mathbb{Q}} \left[ \left( f(\boldsymbol{\xi}) + \sup_{\boldsymbol{\zeta} \in \cup_{j=1}^N \mathcal{U}_N^j} |f(\boldsymbol{\zeta})| \right) \mathbb{I} \{ \boldsymbol{\xi} \in \cup_{j=1}^N \mathcal{U}_N^j \} \right] - \left( \sup_{\boldsymbol{\zeta} \in \cup_{j=1}^N \mathcal{U}_N^j} |f(\boldsymbol{\zeta})| \right) \mathbb{P}(\boldsymbol{\xi} \in \cup_{j=1}^N \mathcal{U}_N^j) \\ &= \sup_{\mathbb{Q} \in \mathcal{P}(\cup_{j=1}^N \mathcal{U}_N^j) : d_1(\mathbb{Q}, \widehat{\mathbb{P}}_N) \leq \delta_N} \mathbb{E}_{\mathbb{Q}} \left[ \left( f(\boldsymbol{\xi}) + \sup_{\boldsymbol{\zeta} \in \cup_{j=1}^N \mathcal{U}_N^j} |f(\boldsymbol{\zeta})| \right) \right] - \left( \sup_{\boldsymbol{\zeta} \in \cup_{j=1}^N \mathcal{U}_N^j} |f(\boldsymbol{\zeta})| \right) \mathbb{P}(\boldsymbol{\xi} \in \cup_{j=1}^N \mathcal{U}_N^j) \\ &= \sup_{\mathbb{Q} \in \mathcal{P}(\cup_{j=1}^N \mathcal{U}_N^j) : d_1(\mathbb{Q}, \widehat{\mathbb{P}}_N) \leq \delta_N} \mathbb{E}_{\mathbb{Q}} [f(\boldsymbol{\xi})] + \left( \sup_{\boldsymbol{\zeta} \in \cup_{j=1}^N \mathcal{U}_N^j} |f(\boldsymbol{\zeta})| \right) \mathbb{P}(\boldsymbol{\xi} \notin \cup_{j=1}^N \mathcal{U}_N^j), \end{aligned} \tag{EC.25}$$

where the first equality follows because the support of probability distributions in the outer-most supremum is restricted to those which assign measure only on  $\cup_{j=1}^N \mathcal{U}_N^j$ , and the second equality follows because  $\sup_{\boldsymbol{\zeta} \in \cup_{j=1}^N \mathcal{U}_N^j} |f(\boldsymbol{\zeta})|$  is independent of  $\mathbb{Q}$ . By Assumption 2 and the construction of  $\delta_N$ , we have that  $\epsilon_N \geq 2\delta_N$  for all sufficiently large  $N \in \mathbb{N}$ . Thus, it follows from Lemma EC.2 that (EC.25) is upper bounded by

$$\frac{1}{N} \sum_{j=1}^N \sup_{\boldsymbol{\zeta} \in \mathcal{U}_N^j} f(\boldsymbol{\zeta}) + \frac{4\delta_N}{\epsilon_N} \sup_{\boldsymbol{\zeta} \in \cup_{j=1}^N \mathcal{U}_N^j} |f(\boldsymbol{\zeta})| + \left( \sup_{\boldsymbol{\zeta} \in \cup_{j=1}^N \mathcal{U}_N^j} |f(\boldsymbol{\zeta})| \right) \mathbb{P}(\boldsymbol{\xi} \notin \cup_{j=1}^N \mathcal{U}_N^j). \tag{EC.26}$$

By the definition of  $\delta_N$ , and since  $\epsilon_N = \kappa N^{-\frac{1}{3}}$  when  $d = 1$  and  $\epsilon_N = \kappa N^{-\frac{1}{d+1}}$  when  $d \geq 2$ , we have that  $\frac{4\delta_N}{\epsilon_N} \leq \frac{M_N}{2}$  for all sufficiently large  $N$ . Finally, Theorem 3 implies that  $\mathbb{P}(\boldsymbol{\xi} \notin \cup_{j=1}^N \mathcal{U}_N^j) \leq \frac{M_N}{2}$  for all sufficiently large  $N$ ,  $\mathbb{P}^\infty$ -almost surely. Combining (EC.24), (EC.25), and (EC.26), we obtain the desired result.  $\square$

### C.3. Miscellaneous results

We conclude Appendix C with some intermediary and technical lemmas which were used in the proof of Theorem 2. The following lemma is a corollary of [Fournier and Guillin \(2015, Theorem 2\)](#) and is included for completeness.

LEMMA EC.3. *Suppose Assumption 1 holds, and let*

$$\delta_N := \begin{cases} \bar{\kappa} N^{-\frac{1}{2}} \log N, & \text{if } d = 1, \\ \bar{\kappa} N^{-\frac{1}{d}} (\log N)^2, & \text{if } d \geq 2, \end{cases}$$

for any fixed  $\bar{\kappa} > 0$ . Then,  $\mathbf{d}_1(\mathbb{P}, \widehat{\mathbb{P}}_N) \leq \delta_N$  for all sufficiently large  $N \in \mathbb{N}$ ,  $\mathbb{P}^\infty$ -almost surely.

*Proof.* Let  $\bar{N} \in \mathbb{N}$  be any index such that  $\delta_N \leq 1$  for all  $N \geq \bar{N}$ . It follows from Assumption 1 that there exists an  $a > 1$  such that  $b := \mathbb{E}[\exp(\|\xi\|^a)] < \infty$ . Thus, it follows from [Fournier and Guillin \(2015, Theorem 2\)](#) that there exist constants  $c_1, c_2 > 0$  (which depend only  $a, b$ , and  $d$ ) such that for all  $N \geq \bar{N}$ ,

$$\mathbb{P}^N \left( \mathbf{d}_1(\mathbb{P}, \widehat{\mathbb{P}}_N) > \delta_N \right) \leq \begin{cases} c_1 \exp(-c_2 N \delta_N^2), & \text{if } d = 1, \\ c_1 \exp\left(-\frac{c_2 N \delta_N^2}{(\log(2+1/\delta_N))^2}\right), & \text{if } d = 2, \\ c_1 \exp(-c_2 N \delta_N^d), & \text{if } d \geq 3. \end{cases} \quad (\text{EC.27})$$

First, suppose  $d = 1$  and  $N \geq \bar{N}$ . Then, it follows from the definition of  $\delta_N = \bar{\kappa} N^{-\frac{1}{2}} \log N$  and (EC.27) that

$$\mathbb{P}^N \left( \mathbf{d}_1(\mathbb{P}, \widehat{\mathbb{P}}_N) > \delta_N \right) \leq c_1 \exp(-c_2 N \delta_N^2) = c_1 \exp(-c_2 \bar{\kappa}^2 (\log N)^2).$$

Second, suppose  $d = 2$  and  $N \geq \bar{N}$ . Then, it follows from the definition of  $\delta_N = \bar{\kappa} N^{-\frac{1}{2}} (\log N)^2$  and (EC.27) that there exists some constant  $\bar{c} > 0$  (which depends only on  $\bar{\kappa}$  and  $c_2$ ) such that

$$\begin{aligned} \mathbb{P}^N \left( \mathbf{d}_1(\mathbb{P}, \widehat{\mathbb{P}}_N) > \delta_N \right) &\leq c_1 \exp\left(-\frac{c_2 N \delta_N^2}{\log(2+1/\delta_N)^2}\right) \\ &= c_1 \exp\left(-\frac{c_2 \bar{\kappa}^2 (\log N)^4}{\log(2+\bar{\kappa}^{-2} N^{\frac{1}{2}} (\log N)^{-2})^2}\right) \\ &\leq c_1 \exp\left(-\frac{c_2 \bar{\kappa}^2 (\log N)^4}{\log(2+\bar{\kappa}^{-2} N^{\frac{1}{2}})^2}\right) \\ &\leq c_1 \exp(-\bar{c} (\log N)^2). \end{aligned}$$

Third, suppose  $d \geq 3$  and  $N \geq \bar{N}$ . Then, it follows from the definition of  $\delta_N = \bar{\kappa} N^{-\frac{1}{d}} (\log N)^2$  and (EC.27) that

$$\mathbb{P}^N \left( \mathbf{d}_1(\mathbb{P}, \widehat{\mathbb{P}}_N) > \delta_N \right) \leq c_1 \exp(-c_2 N \delta_N^d) = c_1 \exp(-c_2 (\log N)^{2d}).$$

Therefore, for any  $d \geq 1$ , we have shown that

$$\sum_{N=1}^{\infty} \mathbb{P}^N \left( \mathbf{d}_1(\mathbb{P}, \widehat{\mathbb{P}}_N) > \delta_N \right) < \infty,$$

and thus the desired result follows from the Borel-Cantelli lemma.  $\square$

The second lemma (Lemma EC.4) demonstrates that restrictions may be placed on the support of an ambiguity set in distributionally robust optimization without loss of generality when the objective function is nonnegative.

LEMMA EC.4. *Suppose  $\Xi \subseteq \mathbb{R}^d$  and  $\hat{\xi}^1, \dots, \hat{\xi}^N \in \mathcal{Z} \subseteq \Xi$ . Let  $g: \Xi \rightarrow \mathbb{R}$  be any measurable function where  $g(\zeta) \geq 0$  for all  $\zeta \in \mathcal{Z}$ . Then, for all  $\theta \geq 0$ ,*

$$\sup_{\mathbb{Q} \in \mathcal{P}(\Xi): \mathbf{d}_1(\mathbb{Q}, \widehat{\mathbb{P}}_N) \leq \theta} \mathbb{E}_{\mathbb{Q}} [g(\xi) \mathbb{I}\{\xi \in \mathcal{Z}\}] = \sup_{\mathbb{Q} \in \mathcal{P}(\mathcal{Z}): \mathbf{d}_1(\mathbb{Q}, \widehat{\mathbb{P}}_N) \leq \theta} \mathbb{E}_{\mathbb{Q}} [g(\xi)].$$

*Proof.* For notational convenience, let  $\bar{g}(\zeta) := g(\zeta)\mathbb{I}\{\zeta \in \mathcal{Z}\}$  for all  $\zeta \in \Xi$ . It readily follows from  $\mathcal{Z} \subseteq \Xi$  that

$$\sup_{\mathbb{Q} \in \mathcal{P}(\Xi): d_1(\mathbb{Q}, \hat{\mathbb{P}}_N) \leq \theta} \mathbb{E}_{\mathbb{Q}}[\bar{g}(\boldsymbol{\xi})] \geq \sup_{\mathbb{Q} \in \mathcal{P}(\mathcal{Z}): d_1(\mathbb{Q}, \hat{\mathbb{P}}_N) \leq \theta} \mathbb{E}_{\mathbb{Q}}[\bar{g}(\boldsymbol{\xi})] = \sup_{\mathbb{Q} \in \mathcal{P}(\mathcal{Z}): d_1(\mathbb{Q}, \hat{\mathbb{P}}_N) \leq \theta} \mathbb{E}_{\mathbb{Q}}[g(\boldsymbol{\xi})],$$

where the equality holds since  $\bar{g}(\zeta) = g(\zeta)$  for all  $\zeta \in \mathcal{Z}$ .

It remains to show the other direction. By the Richter-Rogonsinski Theorem (see Theorem 7.32 and Proposition 6.40 of [Shapiro et al. \(2009\)](#)),

$$\sup_{\mathbb{Q} \in \mathcal{P}(\Xi): d_1(\mathbb{Q}, \hat{\mathbb{P}}_N) \leq \theta} \mathbb{E}_{\mathbb{Q}}[\bar{g}(\boldsymbol{\xi})] = \begin{cases} \sup_{\zeta^{j1}, \zeta^{j2} \in \Xi, \lambda^j \in [0,1]} & \frac{1}{N} \sum_{j=1}^N (\lambda^j \bar{g}(\zeta^{j1}) + (1 - \lambda^j) \bar{g}(\zeta^{j2})) \\ \text{subject to} & \frac{1}{N} \sum_{j=1}^N (\lambda^j \|\zeta^{j1} - \hat{\boldsymbol{\xi}}^j\| + (1 - \lambda^j) \|\zeta^{j2} - \hat{\boldsymbol{\xi}}^j\|) \leq \theta. \end{cases}$$

For any arbitrary  $\eta > 0$ , let  $(\bar{\zeta}^{j1}, \bar{\zeta}^{j2}, \bar{\lambda}^j)_{j \in [N]}$  be an  $\eta$ -optimal solution to the above optimization problem. We now perform a transformation on this solution. For each  $j \in [N]$ , define  $\check{\lambda}^j = \bar{\lambda}^j$ , and for each  $*$  in  $\{1, 2\}$ , define  $\check{\zeta}^{j*} = \bar{\zeta}^{j*}$  if  $\bar{\zeta}^{j*} \in \mathcal{Z}$  and  $\check{\zeta}^{j*} = \hat{\boldsymbol{\xi}}^j$  otherwise. Since  $\bar{g}(\zeta) \geq 0$  for all  $\zeta \in \Xi$  and  $\bar{g}(\zeta) = 0$  for all  $\zeta \notin \mathcal{Z}$ , it follows that  $\bar{g}(\check{\zeta}^{j*}) \geq \bar{g}(\bar{\zeta}^{j*})$ . By construction,  $(\check{\zeta}^{j1}, \check{\zeta}^{j2}, \check{\lambda}^j)_{j \in [N]}$  is a feasible solution to the above optimization problem, and is also feasible for

$$\begin{aligned} & \sup_{\zeta^{j1}, \zeta^{j2} \in \mathcal{Z}, \lambda^j \in [0,1]} \frac{1}{N} \sum_{j=1}^N (\lambda^j \bar{g}(\zeta^{j1}) + (1 - \lambda^j) \bar{g}(\zeta^{j2})) \\ & \text{subject to} \quad \frac{1}{N} \sum_{j=1}^N (\lambda^j \|\zeta^{j1} - \hat{\boldsymbol{\xi}}^j\| + (1 - \lambda^j) \|\zeta^{j2} - \hat{\boldsymbol{\xi}}^j\|) \leq \theta, \end{aligned}$$

where we replaced the domain of  $\zeta^{j1}$  and  $\zeta^{j2}$  by  $\mathcal{Z}$ . We have thus shown that

$$\begin{aligned} \sup_{\mathbb{Q} \in \mathcal{P}(\Xi): d_1(\mathbb{Q}, \hat{\mathbb{P}}_N) \leq \theta} \mathbb{E}_{\mathbb{Q}}[\bar{g}(\boldsymbol{\xi})] & \leq \frac{1}{N} \sum_{j=1}^N (\bar{\lambda}^j \bar{g}(\bar{\zeta}^{j1}) + (1 - \bar{\lambda}^j) \bar{g}(\bar{\zeta}^{j2})) + \eta \\ & \leq \frac{1}{N} \sum_{j=1}^N (\check{\lambda}^j \bar{g}(\check{\zeta}^{j1}) + (1 - \check{\lambda}^j) \bar{g}(\check{\zeta}^{j2})) + \eta \leq \sup_{\mathbb{Q} \in \mathcal{P}(\mathcal{Z}): d_1(\mathbb{Q}, \hat{\mathbb{P}}_N) \leq \theta} \mathbb{E}_{\mathbb{Q}}[\bar{g}(\boldsymbol{\xi})] + \eta. \end{aligned}$$

Since  $\eta > 0$  was chosen arbitrarily, and by the equivalence of  $\bar{g}$  and  $g$  on  $\mathcal{Z}$ , we have shown the other direction. This concludes the proof.  $\square$

## Appendix D: Proof of Proposition 1 from Section 4.3

In Section 4.2, we introduced a lower bound  $J$  and upper bound  $\bar{J}$  on the optimal cost  $J^*$  of Problem (1). In Theorem 1, we showed under mild assumptions that these quantities also provide an asymptotic lower and upper bound on the optimal cost  $\hat{J}_N$  of Problem (2). In this appendix, we will demonstrate the practical value of these bounds by revisiting the stochastic inventory problem from Example 1 in Section 3.

We recall that this stochastic inventory problem from Example 1 is

$$\begin{aligned} J^* & = \underset{x_2: \mathbb{R} \rightarrow \mathbb{R}, x_3: \mathbb{R}^2 \rightarrow \mathbb{R}}{\text{minimize}} && \mathbb{E}[x_2(\xi_1) + 2x_3(\xi_1, \xi_2)] \\ & \text{subject to} && x_2(\xi_1) + x_3(\xi_1, \xi_2) \geq \xi_1 + \xi_2 \quad \text{a.s.} \\ & && x_2(\xi_1), x_3(\xi_1, \xi_2) \geq 0 \quad \text{a.s.} \end{aligned} \tag{3}$$

where the random variables  $\boldsymbol{\xi} = (\xi_1, \xi_2) \in \mathbb{R}^2$  denote the preorder and regular demand of a new product. We assume throughout this appendix that this stochastic process satisfies Assumption 1 and is contained in  $\Xi := \mathbb{R}_+^2$ . In Appendix A, we showed that the above stochastic inventory problem satisfies Assumption 3.

We will now prove Proposition 1 which, in combination with Theorem 1, shows that adding robustness to the historical data overcomes the overfitting phenomenon discussed in Section 3. For clarity of exposition, we split the proof into two parts.

PROPOSITION 1A. For Problem (3),  $J^* = \underline{J}$ .

*Proof.* Following the notation from Section 4.2, the lower bound for this problem is given by  $\underline{J} = \lim_{\rho \downarrow 0} J_\rho$ , where

$$\begin{aligned} J_\rho := & \underset{\mathbf{x}, \tilde{S} \subseteq \Xi}{\text{minimize}} && \mathbb{E} \left[ (x_2(\xi_1) + 2x_3(\xi_1, \xi_2)) \mathbb{I} \left\{ \boldsymbol{\xi} \in \tilde{S} \right\} \right] \\ & \text{subject to} && x_2(\zeta_1) + x_3(\zeta_1, \zeta_2) \geq \zeta_1 + \zeta_2 && \forall \zeta \in \tilde{S} \\ & && x_2(\zeta_1), x_3(\zeta_1, \zeta_2) \geq 0 && \forall \zeta \in \tilde{S} \\ & && \mathbb{P} \left( \boldsymbol{\xi} \in \tilde{S} \right) \geq 1 - \rho. \end{aligned} \quad (\text{EC.28})$$

Our proof is split into several steps. First, we show that additional constraints can be added to (EC.28) without affecting its optimal cost. Indeed, we observe that the cost of procuring inventory is \$1 per unit immediately after the preorder demand  $\xi_1$  is observed, which is cheaper than producing inventory at \$2 per unit after the regular demand  $\xi_2$  is observed. Since all the demand will eventually need to be satisfied, we can without loss of generality restrict the decision rules in (EC.28) to those which satisfy  $x_2(\zeta_1) \geq \zeta_1$  and  $x_3(\zeta_1, \zeta_2) = \max\{0, \zeta_1 + \zeta_2 - x_2(\zeta_1)\}$  for every realization  $\zeta \in \tilde{S}$ . Therefore,

$$\begin{aligned} J_\rho = & \begin{cases} \underset{x_2: \mathbb{R} \rightarrow \mathbb{R}, \tilde{S} \subseteq \Xi}{\text{minimize}} & \mathbb{E} \left[ (x_2(\xi_1) + 2 \max\{0, \xi_1 + \xi_2 - x_2(\xi_1)\}) \mathbb{I} \left\{ \boldsymbol{\xi} \in \tilde{S} \right\} \right] \\ \text{subject to} & x_2(\zeta_1) \geq \zeta_1 \quad \forall \zeta \in \tilde{S} \\ & \mathbb{P} \left( \boldsymbol{\xi} \in \tilde{S} \right) \geq 1 - \rho \end{cases} \\ = & \begin{cases} \underset{\bar{x}_2: \mathbb{R} \rightarrow \mathbb{R}, \tilde{S} \subseteq \Xi}{\text{minimize}} & \mathbb{E} \left[ (\bar{x}_2(\xi_1) + \xi_1 + 2 \max\{0, \xi_2 - \bar{x}_2(\xi_1)\}) \mathbb{I} \left\{ \boldsymbol{\xi} \in \tilde{S} \right\} \right] \\ \text{subject to} & \bar{x}_2(\zeta_1) \geq 0 \quad \forall \zeta \in \tilde{S} \\ & \mathbb{P} \left( \boldsymbol{\xi} \in \tilde{S} \right) \geq 1 - \rho, \end{cases} \end{aligned} \quad (\text{EC.29})$$

where the second equality follows from substituting  $\bar{x}_2(\zeta_1) = x_2(\zeta_1) - \zeta_1$ . Notice that the constraints and objective function of (EC.29) only apply for realizations in  $\tilde{S}$ . In particular, we may add the constraint that  $\bar{x}_2(\zeta_1) \geq 0$  for all  $\zeta \notin \tilde{S}$  to the above optimization problem without affecting its objective function. Thus, it follows from the aforementioned reasoning that

$$\begin{aligned} J_\rho = & \underset{\bar{x}_2: \mathbb{R} \rightarrow \mathbb{R}, \tilde{S} \subseteq \Xi}{\text{minimize}} && \mathbb{E} \left[ (\bar{x}_2(\xi_1) + \xi_1 + 2 \max\{0, \xi_2 - \bar{x}_2(\xi_1)\}) \mathbb{I} \left\{ \boldsymbol{\xi} \in \tilde{S} \right\} \right] \\ & \text{subject to} && \bar{x}_2(\zeta_1) \geq 0 \quad \forall \zeta \in \Xi \\ & && \mathbb{P} \left( \boldsymbol{\xi} \in \tilde{S} \right) \geq 1 - \rho. \end{aligned} \quad (\text{EC.30})$$

Next, we will develop a lower bound on the optimization problem in (EC.30). Indeed, for any feasible solution to the optimization problem in (EC.30), we observe that the function  $\bar{x}_2(\xi_1) + \xi_1 + 2 \max\{0, \xi_2 - \bar{x}_2(\xi_1)\}$  is nonnegative almost surely. Therefore, for any arbitrary  $r \geq 0$ , a lower bound on  $J_\rho$  is given by

$$\begin{aligned} J_{\rho,r} := & \underset{\bar{x}_2: \mathbb{R} \rightarrow \mathbb{R}, \tilde{S} \subseteq \Xi}{\text{minimize}} && \mathbb{E} \left[ (\bar{x}_2(\xi_1) + \xi_1 + 2 \max\{0, \xi_2 - \bar{x}_2(\xi_1)\}) \mathbb{I} \left\{ \boldsymbol{\xi} \in \tilde{S}, \|\boldsymbol{\xi}\| \leq r \right\} \right] \\ & \text{subject to} && \bar{x}_2(\zeta_1) \geq 0 \quad \forall \zeta \in \Xi \\ & && \mathbb{P} \left( \boldsymbol{\xi} \in \tilde{S} \right) \geq 1 - \rho. \end{aligned} \quad (\text{EC.31})$$

Furthermore, the objective function in the above problem only considers realizations which satisfy  $\xi_2 \leq r$ . Thus, we can impose the constraint that  $\bar{x}_2(\zeta_1) \leq r$  for all  $\zeta \in \Xi$  into (EC.31) without changing its optimal cost:

$$\begin{aligned} J_{\rho,r} = & \underset{\bar{x}_2: \mathbb{R} \rightarrow \mathbb{R}, \tilde{S} \subseteq \Xi}{\text{minimize}} && \mathbb{E} \left[ (\bar{x}_2(\xi_1) + \xi_1 + 2 \max\{0, \xi_2 - \bar{x}_2(\xi_1)\}) \mathbb{I} \left\{ \boldsymbol{\xi} \in \tilde{S}, \|\boldsymbol{\xi}\| \leq r \right\} \right] \\ & \text{subject to} && 0 \leq \bar{x}_2(\zeta_1) \leq r \quad \forall \zeta \in \Xi \\ & && \mathbb{P} \left( \boldsymbol{\xi} \in \tilde{S} \right) \geq 1 - \rho. \end{aligned} \quad (\text{EC.32})$$

We now use Assumption 1 to obtain a lower bound on (EC.32). Indeed, Assumption 1 says that there exists an  $a > 1$  such that  $b := \mathbb{E}[\exp(\|\xi\|^a)] < \infty$ . Therefore, for any feasible solution to the optimization problem in (EC.32),

$$\begin{aligned} & \mathbb{E} \left[ (\bar{x}_2(\xi_1) + \xi_1 + 2 \max\{0, \xi_2 - \bar{x}_2(\xi_1)\}) \mathbb{I} \left\{ \xi \notin \tilde{S} \text{ or } \|\xi\| > r \right\} \right] \\ & \leq \mathbb{E} \left[ (r + \xi_1 + 2\xi_2) \mathbb{I} \left\{ \xi \notin \tilde{S} \text{ or } \|\xi\| > r \right\} \right] \end{aligned} \quad (\text{EC.33})$$

$$\leq \mathbb{E} \left[ (r + \xi_1 + 2\xi_2) \mathbb{I} \left\{ \xi \notin \tilde{S} \right\} \right] + \mathbb{E} \left[ (r + \xi_1 + 2\xi_2) \mathbb{I} \left\{ \|\xi\| > r \right\} \right] \quad (\text{EC.34})$$

$$\leq \sqrt{\mathbb{E} \left[ (r + \xi_1 + 2\xi_2)^2 \right]} \rho + \sqrt{\mathbb{E} \left[ (r + \xi_1 + 2\xi_2)^2 \right]} \mathbb{P}(\|\xi\| > r) \quad (\text{EC.35})$$

$$\leq \underbrace{\sqrt{\mathbb{E} \left[ (r + \xi_1 + 2\xi_2)^2 \right]} \rho + \sqrt{\mathbb{E} \left[ (r + \xi_1 + 2\xi_2)^2 \right]} \frac{b}{\exp(r^a)}}_{h(\rho, r)}. \quad (\text{EC.36})$$

Line (EC.33) follows because  $0 \leq \bar{x}_2(\xi_1) \leq r$  for any feasible solution to the optimization problem in (EC.32), (EC.34) follows from the union bound, (EC.35) follows from  $\mathbb{P}(\xi \in \tilde{S}) \geq 1 - \rho$  and the Cauchy-Schwartz inequality, and (EC.36) follows from Markov's inequality. Therefore,

$$\begin{aligned} \underline{J}_{\rho, r} & \geq -h(\rho, r) + \begin{cases} \text{minimize} & \mathbb{E}[(\bar{x}_2(\xi_1) + \xi_1 + 2 \max\{0, \xi_2 - \bar{x}_2(\xi_1)\})] \\ \text{subject to} & 0 \leq \bar{x}_2(\zeta_1) \leq r \quad \forall \zeta \in \Xi \\ & \mathbb{P}(\xi \in \tilde{S}) \geq 1 - \rho. \end{cases} \\ & \geq -h(\rho, r) + \begin{cases} \text{minimize} & \mathbb{E}[\bar{x}_2(\xi_1) + \xi_1 + 2 \max\{0, \xi_2 - \bar{x}_2(\xi_1)\}] \\ \text{subject to} & \bar{x}_2(\zeta_1) \geq 0 \quad \forall \zeta \in \Xi \end{cases} \\ & = -h(\rho, r) + J^*. \end{aligned} \quad (\text{EC.37})$$

The first inequality follows from the definition of  $\underline{J}_{\rho, r}$  in (EC.32), the definition of  $h(\rho, r)$ , and the law of iterated expectation. The second inequality follows from removing constraints, and the final equality follows from the definition of  $J^*$ .

We now combine the above results to prove the main result. Indeed,

$$\underline{J} = \lim_{\rho \downarrow 0} \underline{J}_\rho \geq \lim_{r \rightarrow \infty} \lim_{\rho \downarrow 0} \underline{J}_{\rho, r} \geq \lim_{r \rightarrow \infty} \lim_{\rho \downarrow 0} -h(\rho, r) + J^* = J^*.$$

The first inequality follows because  $\underline{J}_\rho \geq \underline{J}_{\rho, r}$  for any arbitrary  $r \geq 0$  and the quantity  $\lim_{\rho \downarrow 0} \underline{J}_{\rho, r}$  is monotonically increasing in  $r$ . The second inequality follows from (EC.37). The final equality follows from the definition of  $h(\rho, r)$  and Assumption 1. Since the inequality  $\underline{J} \leq J^*$  always holds, our proof is complete.  $\square$

**PROPOSITION 1B.** *If there is an optimal  $x_2^* : \mathbb{R} \rightarrow \mathbb{R}$  for Problem (3) which is continuous, then  $\bar{J} = J^*$ .*

*Proof.* Let  $x_2^* : \mathbb{R} \rightarrow \mathbb{R}$  denote an optimal second-stage decision rule to Problem (3) which is continuous, and consider any arbitrary  $M \geq 0$ . We define the following new decision rules:

$$x_2^M(\zeta_1) := \max\{0, \min\{x_2^*(\zeta_1), M\}\}, \quad x_3^M(\zeta_1, \zeta_2) := \max\{0, \zeta_1 + \zeta_2 - x_2^M(\zeta_1)\}. \quad (\text{EC.38})$$

It follows from the above construction that

$$\begin{aligned} \bar{\mathbb{E}}[x_2^M(\xi_1) + 2x_3^M(\xi_1, \xi_2)] & = \lim_{\epsilon \downarrow 0} \mathbb{E} \left[ \sup_{\zeta \in \Xi: \|\zeta - \xi\| \leq \epsilon} \{x_2^M(\zeta_1) + 2x_3^M(\zeta_1, \zeta_2)\} \right] \\ & = \mathbb{E} \left[ \lim_{\epsilon \downarrow 0} \sup_{\zeta \in \Xi: \|\zeta - \xi\| \leq \epsilon} \{x_2^M(\zeta_1) + 2x_3^M(\zeta_1, \zeta_2)\} \right] \\ & = \mathbb{E} [x_2^M(\xi_1) + 2x_3^M(\xi_1, \xi_2)]. \end{aligned} \quad (\text{EC.39})$$



The first equality is the definition of the local upper semicontinuous envelope. The second equality follows from the dominated convergence theorem.<sup>5</sup> The third equality follows because  $x_2^M(\zeta_1) + 2x_3^M(\zeta_1, \zeta_2)$  is a continuous function for all  $(\zeta_1, \zeta_2) \in \mathbb{R}^2$ .

Recall that the upper bound is defined as

$$\begin{aligned} \bar{J} &= \lim_{\rho \downarrow 0} \underset{x_2: \mathbb{R} \rightarrow \mathbb{R}, x_3: \mathbb{R}^2 \rightarrow \mathbb{R}}{\text{minimize}} \quad \bar{\mathbb{E}}[x_2(\xi_1) + 2x_3(\xi_1, \xi_2)] \\ &\quad \text{subject to} \quad x_2(\zeta_1) + x_3(\zeta_1, \zeta_2) \geq \zeta_1 + \zeta_2 \quad \forall \zeta \in \Xi: \text{dist}(\zeta, S) \leq \rho \\ &\quad \quad \quad x_2(\zeta_1), x_3(\zeta_1, \zeta_2) \geq 0 \quad \forall \zeta \in \Xi: \text{dist}(\zeta, S) \leq \rho. \end{aligned}$$

We observe that the decision rules from (EC.38) are feasible for the above optimization problem for every  $\rho \geq 0$ . Therefore,

$$\begin{aligned} \bar{J} &\leq \lim_{M \rightarrow \infty} \bar{\mathbb{E}}[x_2^M(\xi_1) + 2x_3^M(\xi_1, \xi_2)] \\ &= \lim_{M \rightarrow \infty} \mathbb{E}[x_2^M(\xi_1) + 2x_3^M(\xi_1, \xi_2)] \end{aligned} \tag{EC.40}$$

$$\leq \lim_{M \rightarrow \infty} \mathbb{E}[x_2^*(\xi_1) + 2x_3^M(\xi_1, \xi_2)] \tag{EC.41}$$

$$\begin{aligned} &= \mathbb{E}[x_2^*(\xi_1)] + \lim_{M \rightarrow \infty} \mathbb{E}[2 \max\{0, \xi_1 + \xi_2 - x_2^*(\xi_1)\} \mathbb{I}\{x_2^*(\xi_1) \leq M\}] \\ &\quad + \lim_{M \rightarrow \infty} \mathbb{E}[2 \max\{0, \xi_1 + \xi_2 - M\} \mathbb{I}\{x_2^*(\xi_1) > M\}] \end{aligned} \tag{EC.42}$$

$$\begin{aligned} &\leq \mathbb{E}[x_2^*(\xi_1)] + \mathbb{E}[2 \max\{0, \xi_1 + \xi_2 - x_2^*(\xi_1)\}] \\ &\quad + \lim_{M \rightarrow \infty} \mathbb{E}[2 \max\{0, \xi_1 + \xi_2 - M\} \mathbb{I}\{x_2^*(\xi_1) > M\}] \end{aligned} \tag{EC.43}$$

$$\begin{aligned} &= J^* + \lim_{M \rightarrow \infty} \mathbb{E}[2 \max\{0, \xi_1 + \xi_2 - M\} \mathbb{I}\{x_2^*(\xi_1) > M\}] \\ &\leq J^* + \lim_{M \rightarrow \infty} \mathbb{E}[2(\xi_1 + \xi_2) \mathbb{I}\{x_2^*(\xi_1) > M\}] \\ &= J^*. \end{aligned} \tag{EC.44}$$

Line (EC.40) follows from (EC.39), (EC.41) holds because  $x_2^*(\xi_1)$  is greater than or equal to  $x_2^M(\xi_1)$  almost surely, (EC.42) follows from the law of total probability, (EC.43) follows from the monotone convergence theorem, and (EC.44) follows from the dominated convergence theorem. Since the inequality  $J^* \leq \bar{J}$  always holds, our proof is complete.  $\square$

## Appendix E: Details for Example 2 from Section 4.3

In this appendix, we provide the omitted technical details of Example 2 from Section 4.3. For convenience, we repeat the example below.

*Consider the single-stage stochastic problem*

$$\begin{aligned} &\underset{x_1 \in \mathbb{Z}}{\text{minimize}} \quad x_1 \\ &\text{subject to} \quad x_1 \geq \xi_1 \quad \text{a.s.}, \end{aligned}$$

where the random variable  $\xi_1$  is governed by the probability distribution  $\mathbb{P}(\xi_1 > \alpha) = (1 - \alpha)^k$  for fixed  $k > 0$ , and  $\Xi = [0, 2]$ . We observe that the support of the random variable is  $S = [0, 1]$ , and thus the optimal cost of the stochastic problem is  $J^* = 1$ . We similarly observe that the lower bound is  $\underline{J} = 1$  and the upper bound,

<sup>5</sup> The dominated convergence theorem can be applied here for two reasons. First, the function

$$h^\epsilon(\xi_1, \xi_2) := \sup_{\zeta \in \Xi: \|\zeta - \xi\| \leq \epsilon} \{x_2^M(\zeta_1) + 2x_3^M(\zeta_1, \zeta_2)\} \geq 0$$

is (pointwise) monotonically decreasing as  $\epsilon \downarrow 0$ . Second, for any  $\epsilon \geq 0$ ,

$$0 \leq \mathbb{E} \left[ \sup_{\zeta \in \Xi: \|\zeta - \xi\| \leq \epsilon} \{x_2^M(\zeta_1) + 2x_3^M(\zeta_1, \zeta_2)\} \right] \leq \mathbb{E}[M + 2(\xi_1 + \xi_2 + 2\epsilon)] < \infty.$$

due to the integrality of the first stage decision, is  $\bar{J} = 2$ . If  $\epsilon_N = N^{-\frac{1}{3}}$ , then we prove in [Appendix E](#) that the bounds in [Theorem 1](#) are tight under different choices of  $k$ :

Range of $k$	Result
$k \in (0, 3)$	$\mathbb{P}^\infty \left( J < \liminf_{N \rightarrow \infty} \hat{J}_N = \limsup_{N \rightarrow \infty} \hat{J}_N = \bar{J} \right) = 1$
$k = 3$	$\mathbb{P}^\infty \left( \underline{J} = \liminf_{N \rightarrow \infty} \hat{J}_N < \limsup_{N \rightarrow \infty} \hat{J}_N = \bar{J} \right) = 1$
$k \in (3, \infty)$	$\mathbb{P}^\infty \left( \underline{J} = \liminf_{N \rightarrow \infty} \hat{J}_N = \limsup_{N \rightarrow \infty} \hat{J}_N < \bar{J} \right) = 1$

We now prove the above bounds. To begin, we recall that  $\mathbb{P}(\xi_1 > \alpha) = (1 - \alpha)^k$ . Thus, for any  $k > 0$ ,

$$\begin{aligned} \underline{J} &= \lim_{\rho \downarrow 0} \min_{x_1 \in \mathbb{Z}} \{x_1 : \mathbb{P}(x_1 \geq \xi_1) \geq 1 - \rho\} = 1, \text{ and} \\ \bar{J} &= \lim_{\rho \downarrow 0} \min_{x_1 \in \mathbb{Z}} \{x_1 : x_1 \geq 1 + \rho\} = 2. \end{aligned}$$

Furthermore, given historical data, the choice of the robustness parameter  $\epsilon_N = N^{-\frac{1}{3}}$ , and  $\Xi = [0, 2]$ ,

$$\hat{J}_N = \min_{x_1 \in \mathbb{Z}} \{x_1 : x_1 \geq \zeta_1, \forall \zeta_1 \in \cup_{j=1}^N \mathcal{U}_N^j\} = \begin{cases} 1, & \text{if } \max_{j \in [N]} \hat{\xi}_1^j \leq 1 - N^{-\frac{1}{3}}, \\ 2, & \text{if } \max_{j \in [N]} \hat{\xi}_1^j > 1 - N^{-\frac{1}{3}}. \end{cases}$$

We first show that

$$\mathbb{P}^\infty \left( \limsup_{N \rightarrow \infty} \hat{J}_N = 1 \right) = \begin{cases} 0, & \text{if } 0 < k \leq 3, \\ 1, & \text{if } k > 3. \end{cases} \quad (\text{Claim 1})$$

Indeed,

$$\begin{aligned} & \mathbb{P}^\infty \left( \limsup_{N \rightarrow \infty} \hat{J}_N = 1 \right) \\ &= \mathbb{P}^\infty \left( \max_{j \in [N]} \hat{\xi}_1^j \leq 1 - N^{-\frac{1}{3}} \text{ for all sufficiently large } N \right) \\ &= \lim_{N \rightarrow \infty} \mathbb{P}^\infty \left( \max_{j \in [n]} \hat{\xi}_1^j \leq 1 - n^{-\frac{1}{3}} \text{ for all } n \geq N \right) \\ &= \lim_{N \rightarrow \infty} \mathbb{P}^\infty \left( \max_{j \in [N]} \hat{\xi}_1^j \leq 1 - N^{-\frac{1}{3}} \text{ and } \max_{j \in [n]} \hat{\xi}_1^j \leq 1 - n^{-\frac{1}{3}} \text{ for all } n \geq N+1 \right) \\ &= \lim_{N \rightarrow \infty} \mathbb{P}^N \left( \max_{j \in [N]} \hat{\xi}_1^j \leq 1 - N^{-\frac{1}{3}} \right) \prod_{n=N+1}^{\infty} \mathbb{P} \left( \max_{j \in [n]} \hat{\xi}_1^j \leq 1 - n^{-\frac{1}{3}} \mid \max_{j \in [n-1]} \hat{\xi}_1^j \leq 1 - (n-1)^{-\frac{1}{3}} \right) \quad (\text{EC.45}) \end{aligned}$$

$$= \lim_{N \rightarrow \infty} \mathbb{P}^\infty \left( \max_{j \in [N]} \hat{\xi}_1^j \leq 1 - N^{-\frac{1}{3}} \right) \prod_{n=N+1}^{\infty} \mathbb{P} \left( \hat{\xi}_1^n \leq 1 - n^{-\frac{1}{3}} \mid \max_{j \in [n-1]} \hat{\xi}_1^j \leq 1 - (n-1)^{-\frac{1}{3}} \right) \quad (\text{EC.46})$$

$$= \lim_{N \rightarrow \infty} \mathbb{P} \left( \xi_1 \leq 1 - N^{-\frac{1}{3}} \right)^N \prod_{n=N+1}^{\infty} \mathbb{P} \left( \xi_1 \leq 1 - n^{-\frac{1}{3}} \right) \quad (\text{EC.47})$$

$$= \lim_{N \rightarrow \infty} \left( 1 - N^{-\frac{k}{3}} \right)^N \prod_{n=N+1}^{\infty} \left( 1 - n^{-\frac{k}{3}} \right). \quad (\text{EC.48})$$

Line [\(EC.45\)](#) follows from the law of total probability. Line [\(EC.46\)](#) follows because, conditional on  $\max_{j \in [n-1]} \hat{\xi}_1^j \leq 1 - (n-1)^{-\frac{1}{3}}$ , we have that  $\hat{\xi}_1^j \leq 1 - n^{-\frac{1}{3}}$  for all  $j \in [n-1]$ . Line [\(EC.47\)](#) follows from the independence of  $\hat{\xi}_1^j$ ,  $j \in \mathbb{N}$ . By evaluating the limit in [\(EC.48\)](#), we conclude the proof of [Claim 1](#).

Next, we show that

$$\mathbb{P}^\infty \left( \liminf_{N \rightarrow \infty} \widehat{J}_N = 1 \right) = 1 \text{ if } k \geq 3. \quad (\text{Claim 2})$$

Indeed,

$$\begin{aligned} \mathbb{P}^\infty \left( \liminf_{N \rightarrow \infty} \widehat{J}_N = 1 \right) &= \mathbb{P}^\infty \left( \max_{j \in [N]} \hat{\xi}_1^j \leq 1 - N^{-\frac{1}{3}} \text{ for infinitely many } N \right) \\ &= \lim_{N \rightarrow \infty} \mathbb{P}^\infty \left( \max_{j \in [n]} \hat{\xi}_1^j \leq 1 - n^{-\frac{1}{3}} \text{ for some } n \geq N \right) \\ &\geq \lim_{N \rightarrow \infty} \mathbb{P}^N \left( \max_{j \in [N]} \hat{\xi}_1^j \leq 1 - N^{-\frac{1}{3}} \right) \\ &= \lim_{N \rightarrow \infty} \mathbb{P} \left( \xi_1 \leq 1 - N^{-\frac{1}{3}} \right)^N \end{aligned} \quad (\text{EC.49})$$

$$= \lim_{N \rightarrow \infty} \left( 1 - N^{-\frac{k}{3}} \right)^N. \quad (\text{EC.50})$$

Line (EC.49) follows from the independence of  $\hat{\xi}^j$ ,  $j \in \mathbb{N}$ . We observe that the limit in (EC.50) is strictly positive when  $k \geq 3$ . It follows from the Hewitt-Savage zero-one law (see, *e.g.*, [Breiman \(1992\)](#), [Wang and Tomkins \(1992\)](#)) that the event  $\{\max_{j \in [N]} \hat{\xi}_1^j \leq 1 - N^{-\frac{1}{3}} \text{ for infinitely many } N\}$  happens with probability zero or one. Thus, (EC.50) implies that the event  $\{\liminf_{N \rightarrow \infty} \widehat{J}_N = 1\}$  must occur with probability one for  $k \geq 3$ .

Finally, we show that

$$\mathbb{P}^\infty \left( \liminf_{N \rightarrow \infty} \widehat{J}_N = 1 \right) = 0 \text{ if } 0 < k < 3. \quad (\text{Claim 3})$$

Indeed, suppose that  $0 < k < 3$ . Then,

$$\sum_{N=1}^{\infty} \mathbb{P}^\infty \left( \widehat{J}_N = 1 \right) = \sum_{N=1}^{\infty} \mathbb{P}^N \left( \max_{j \in [N]} \hat{\xi}_1^j \leq 1 - N^{-\frac{1}{3}} \right) = \sum_{N=1}^{\infty} \left( 1 - N^{-\frac{k}{3}} \right)^N < \infty.$$

Therefore, it follows from the Borel-Cantelli lemma that

$$\mathbb{P}^\infty \left( \liminf_{N \rightarrow \infty} \widehat{J}_N = 1 \right) = \mathbb{P}^\infty \left( \max_{j \in [N]} \hat{\xi}_1^j > 1 - N^{-\frac{1}{3}} \text{ for all sufficiently large } N \right) = 0,$$

when  $0 < k < 3$ , which proves Claim 3.

Combining Claims 1, 2, and 3 with the definitions of  $\underline{J}$  and  $\bar{J}$ , we have shown the desired results.

## Appendix F: Proof of Theorem 3 from Section 4.4

In this appendix, we present the proof of Theorem 3. Our proof techniques follow similar reasoning to [Devroye and Wise \(1980\)](#) and [Baïllo et al. \(2000\)](#) for  $S_N := \cup_{j=1}^N \mathcal{U}_N^j$ , which we adapt to Assumption 1. We remark that the following theorem also provides an intermediary step in the proofs of Theorems 1 and 2, which are found in Appendices B and C.

**THEOREM 3.** *Suppose Assumptions 1 and 2 hold. Then,  $\mathbb{P}^\infty$ -almost surely we have*

$$\lim_{N \rightarrow \infty} \left( \frac{N^{\frac{1}{a+1}}}{(\log N)^{d+1}} \right) \mathbb{P}(\boldsymbol{\xi} \notin S_N) = 0.$$

*Proof.* Choose any arbitrary  $\eta > 0$ , and let  $R_N := N^{\frac{1}{a+1}} (\log N)^{-(d+1)}$ . Moreover, let  $a > 1$  be a fixed constant such that  $b := \mathbb{E}[\exp(\|\boldsymbol{\xi}\|^a)] < \infty$  (the existence of  $a$  and  $b$  follows from Assumption 1). Define

$$A_N := \left\{ \boldsymbol{\zeta} \in \mathbb{R}^d : \|\boldsymbol{\zeta}\| \leq (\log N)^{\frac{a+1}{2a}} \right\}.$$

We begin by showing that  $R_N \mathbb{P}(\boldsymbol{\xi} \notin A_N) \leq \eta$  for all sufficiently large  $N \in \mathbb{N}$ . Indeed,

$$R_N \mathbb{P}(\boldsymbol{\xi} \notin A_N) = R_N \mathbb{P}(\|\boldsymbol{\xi}\| > (\log N)^{\frac{a+1}{2a}}) = R_N \mathbb{P}(\exp(\|\boldsymbol{\xi}\|^a) > \exp((\log N)^{\frac{a+1}{2}})) \leq \frac{bR_N}{\exp((\log N)^{\frac{a+1}{2}})} \leq \eta.$$

The first inequality follows from Markov's inequality and the second inequality holds for all sufficiently large  $N \in \mathbb{N}$  since  $a > 1$ .

Next, define

$$\alpha_N := \frac{\eta}{(\log N)^{\frac{d(a+1)}{2a}} \phi R_N}, \quad B_N := \{\boldsymbol{\zeta} \in \mathbb{R}^d : \mathbb{P}(\|\boldsymbol{\xi} - \boldsymbol{\zeta}\| \leq \epsilon_N) > \alpha_N \epsilon_N^d\},$$

where  $\phi > 0$  is a constant which depends only on  $d$  and will be defined shortly. We now show that  $R_N \mathbb{P}(\boldsymbol{\xi} \notin B_N) \leq 2\eta$  for all sufficiently large  $N$ . Indeed, for all sufficiently large  $N \in \mathbb{N}$ ,

$$\begin{aligned} R_N \mathbb{P}(\boldsymbol{\xi} \notin B_N) &= R_N \mathbb{P}(\boldsymbol{\xi} \in A_N, \boldsymbol{\xi} \notin B_N) + R_N \mathbb{P}(\boldsymbol{\xi} \notin A_N, \boldsymbol{\xi} \notin B_N) \\ &\leq R_N \mathbb{P}(\boldsymbol{\xi} \in A_N, \boldsymbol{\xi} \notin B_N) + R_N \mathbb{P}(\boldsymbol{\xi} \notin A_N) \\ &\leq R_N \mathbb{P}(\boldsymbol{\xi} \in A_N, \boldsymbol{\xi} \notin B_N) + \eta, \end{aligned} \tag{EC.51}$$

where the final inequality follows because  $R_N \mathbb{P}(\boldsymbol{\xi} \notin A_N) \leq \eta$  for all sufficiently large  $N \in \mathbb{N}$ . Now, choose points  $\boldsymbol{\zeta}^1, \dots, \boldsymbol{\zeta}^{K_N} \in A_N$  such that  $\min_{j \in [K_N]} \|\boldsymbol{\zeta} - \boldsymbol{\zeta}^j\| \leq \frac{\epsilon_N}{2}$  for all  $\boldsymbol{\zeta} \in A_N$ . For example, one can place the points on a grid overlaying  $A_N$ . It follows from [Verger-Gaugry \(2005\)](#) that this can be accomplished with a number of points  $K_N$  which satisfies

$$K_N \leq \phi \left( \frac{(\log N)^{\frac{a+1}{2a}}}{\epsilon_N} \right)^d, \tag{EC.52}$$

where  $\phi > 0$  is a constant that depends only on  $d$ . Then, continuing from [\(EC.51\)](#),

$$R_N \mathbb{P}(\boldsymbol{\xi} \notin B_N) \leq R_N \mathbb{P}(\boldsymbol{\xi} \in A_N, \boldsymbol{\xi} \notin B_N) + \eta \leq R_N \sum_{j=1}^{K_N} \mathbb{P}(\|\boldsymbol{\xi} - \boldsymbol{\zeta}^j\| \leq \frac{\epsilon_N}{2}, \boldsymbol{\xi} \notin B_N) + \eta, \tag{EC.53}$$

where the second inequality follows from the union bound. For each  $j \in [K_N]$ , we have two cases to consider. First, suppose there exists a realization  $\boldsymbol{\zeta} \notin B_N$  such that  $\|\boldsymbol{\zeta} - \boldsymbol{\zeta}^j\| \leq \frac{\epsilon_N}{2}$ . Then,

$$\mathbb{P}(\|\boldsymbol{\xi} - \boldsymbol{\zeta}^j\| \leq \frac{\epsilon_N}{2}, \boldsymbol{\xi} \notin B_N) \leq \mathbb{P}(\|\boldsymbol{\xi} - \boldsymbol{\zeta}^j\| \leq \frac{\epsilon_N}{2}) \leq \mathbb{P}(\|\boldsymbol{\xi} - \boldsymbol{\zeta}\| \leq \epsilon_N) \leq \alpha_N \epsilon_N^d,$$

where the second inequality follows because  $\|\boldsymbol{\xi} - \boldsymbol{\zeta}\| \leq \|\boldsymbol{\xi} - \boldsymbol{\zeta}^j\| + \|\boldsymbol{\zeta}^j - \boldsymbol{\zeta}\| \leq \epsilon_N$  whenever  $\|\boldsymbol{\xi} - \boldsymbol{\zeta}^j\| \leq \frac{\epsilon_N}{2}$ , and the third inequality follows from  $\boldsymbol{\zeta} \notin B_N$ . Second, suppose there does not exist a realization  $\boldsymbol{\zeta} \notin B_N$  such that  $\|\boldsymbol{\zeta} - \boldsymbol{\zeta}^j\| \leq \frac{\epsilon_N}{2}$ . Then,

$$\mathbb{P}(\|\boldsymbol{\xi} - \boldsymbol{\zeta}^j\| \leq \frac{\epsilon_N}{2}, \boldsymbol{\xi} \notin B_N) = 0.$$

In each of the two cases, we have shown that

$$\mathbb{P}(\|\boldsymbol{\xi} - \boldsymbol{\zeta}^j\| \leq \frac{\epsilon_N}{2}, \boldsymbol{\xi} \notin B_N) \leq \alpha_N \epsilon_N^d \tag{EC.54}$$

for each  $j \in [K_N]$ . Therefore, we combine [\(EC.53\)](#) and [\(EC.54\)](#) to obtain the following upper bound on  $R_N \mathbb{P}(\boldsymbol{\xi} \notin B_N)$  for all sufficiently large  $N \in \mathbb{N}$ :

$$R_N \mathbb{P}(\boldsymbol{\xi} \notin B_N) \leq R_N K_N \alpha_N \epsilon_N^d + \eta \leq (\log N)^{\frac{d(a+1)}{2a}} \phi R_N \alpha_N + \eta \leq 2\eta. \tag{EC.55}$$

The first inequality follows from [\(EC.53\)](#) and [\(EC.54\)](#), the second inequality follows from [\(EC.52\)](#), and the third inequality follows from the definition of  $\alpha_N$ .

We now prove the main result. Indeed, for all sufficiently large  $N \in \mathbb{N}$ ,

$$R_N \mathbb{P}(\boldsymbol{\xi} \notin S_N) = R_N \mathbb{P}(\boldsymbol{\xi} \notin S_N, \boldsymbol{\xi} \in B_N) + R_N \mathbb{P}(\boldsymbol{\xi} \notin S_N, \boldsymbol{\xi} \notin B_N) \leq R_N \mathbb{P}(\boldsymbol{\xi} \notin S_N, \boldsymbol{\xi} \in B_N) + 2\eta, \tag{EC.56}$$

where the equality follows from the law of total probability and the inequality follows from (EC.55). Let  $\rho := \frac{d(a-1)}{2a} > 0$ . Then, for all sufficiently large  $N \in \mathbb{N}$ :

$$\mathbb{P}^N (R_N \mathbb{P}(\boldsymbol{\xi} \notin S_N) > 3\eta) \leq \mathbb{P}^N (R_N \mathbb{P}(\boldsymbol{\xi} \notin S_N, \boldsymbol{\xi} \in B_N) > \eta) \quad (\text{EC.57})$$

$$\leq \eta^{-1} R_N \mathbb{E}_{\mathbb{P}^N} \left[ \mathbb{P} \left( \boldsymbol{\xi} \notin S_N, \boldsymbol{\xi} \in B_N \mid \hat{\boldsymbol{\xi}}^1, \dots, \hat{\boldsymbol{\xi}}^N \right) \right] \quad (\text{EC.58})$$

$$= \eta^{-1} R_N \mathbb{E} \left[ \mathbb{P}^N (\boldsymbol{\xi} \notin S_N, \boldsymbol{\xi} \in B_N \mid \boldsymbol{\xi}) \right] \quad (\text{EC.59})$$

$$= \eta^{-1} R_N \mathbb{E} \left[ \mathbb{P}^N \left( \|\boldsymbol{\xi} - \hat{\boldsymbol{\xi}}^1\| > \epsilon_N, \dots, \|\boldsymbol{\xi} - \hat{\boldsymbol{\xi}}^N\| > \epsilon_N, \boldsymbol{\xi} \in B_N \mid \boldsymbol{\xi} \right) \right] \quad (\text{EC.60})$$

$$= \eta^{-1} R_N \mathbb{E} \left[ \mathbb{P} \left( \|\boldsymbol{\xi} - \boldsymbol{\xi}'\| > \epsilon_N, \boldsymbol{\xi} \in B_N \mid \boldsymbol{\xi} \right)^N \right] \quad (\text{EC.61})$$

$$\leq \eta^{-1} R_N (1 - \alpha_N \epsilon_N^d)^N \quad (\text{EC.62})$$

$$\leq \eta^{-1} R_N \exp(-N \alpha_N \epsilon_N^d) \quad (\text{EC.63})$$

$$\leq \eta^{-1} R_N \exp\left(-\kappa^d N^{\frac{1}{d+1}} \alpha_N\right) \quad (\text{EC.64})$$

$$= \eta^{-1} R_N \exp\left(-\kappa^d \eta \phi^{-1}(\log N)^{1+\rho}\right). \quad (\text{EC.65})$$

Line (EC.57) follows from (EC.56), (EC.58) follows from Markov's inequality, (EC.59) follows from the law of iterated expectation, and (EC.60) follows from the definition of  $S_N$ . Line (EC.61) follows because, conditional on  $\boldsymbol{\xi}$ , the random variables  $\|\boldsymbol{\xi} - \hat{\boldsymbol{\xi}}^1\|, \dots, \|\boldsymbol{\xi} - \hat{\boldsymbol{\xi}}^N\|$  are independent. Line (EC.62) follows from the definition of  $B_N$ , and (EC.63) follows from the mean value theorem. Line (EC.64) holds since Assumption 2 implies that  $\epsilon_N \geq \kappa N^{-\frac{1}{d+1}}$ . Line (EC.65) follows from the definitions of  $\alpha_N$ ,  $R_N$ , and  $\rho$ . Since  $\rho > 0$ , it follows from (EC.65) and the definition of  $R_N$  that

$$\sum_{N=1}^{\infty} \mathbb{P}^N (R_N \mathbb{P}(\boldsymbol{\xi} \notin S_N) > 3\eta) < \infty, \quad \forall \eta > 0,$$

and thus the Borel-Cantelli lemma implies that  $R_N \mathbb{P}(\boldsymbol{\xi} \notin S_N) \rightarrow 0$  as  $N \rightarrow \infty$ ,  $\mathbb{P}^\infty$ -almost surely.  $\square$

## Appendix G: Proof of Proposition 3 from Section 6

In this appendix, we present the proof of Proposition 3. We begin with the following intermediary lemma.

LEMMA EC.5. *The  $\infty$ -Wasserstein ambiguity set is equivalent to*

$$\left\{ \frac{1}{N} \sum_{j=1}^N \mathbb{Q}_j : \begin{array}{l} \mathbb{Q}_j \left( \|\boldsymbol{\xi} - \hat{\boldsymbol{\xi}}^j\| \leq \epsilon_N \right) = 1 \quad \forall j \in [N], \\ \mathbb{Q}_1, \dots, \mathbb{Q}_N \in \mathcal{P}(\Xi) \end{array} \right\}.$$

*Proof.* By the definition of the  $\infty$ -Wasserstein distance from Section 6,

$$\left\{ \mathbb{Q} \in \mathcal{P}(\Xi) : d_\infty \left( \mathbb{Q}, \hat{\mathbb{P}}_N \right) \leq \epsilon_N \right\} = \left\{ \mathbb{Q} \in \mathcal{P}(\Xi) : \begin{array}{l} \Pi \in \mathcal{P}(\Xi \times \Xi), \\ \Pi \left( \|\boldsymbol{\xi} - \boldsymbol{\xi}'\| \leq \epsilon_N \right) = 1, \text{ and} \\ \Pi \text{ is a joint distribution of } \boldsymbol{\xi} \text{ and } \boldsymbol{\xi}' \\ \text{with marginals } \mathbb{Q} \text{ and } \hat{\mathbb{P}}_N, \text{ respectively} \end{array} \right\}. \quad (\text{EC.66})$$

Let  $\bar{\boldsymbol{\xi}}^1, \dots, \bar{\boldsymbol{\xi}}^L$  be the distinct vectors among  $\hat{\boldsymbol{\xi}}^1, \dots, \hat{\boldsymbol{\xi}}^N$ , and let  $I_1, \dots, I_L$  be index sets defined as

$$I_\ell := \{j \in [N] : \hat{\boldsymbol{\xi}}^j = \bar{\boldsymbol{\xi}}^\ell\}.$$

For any joint distribution  $\Pi$  that satisfies the constraints in the ambiguity set in (EC.66), let  $\mathbb{Q}_\ell$  be the conditional distribution of  $\boldsymbol{\xi}$  given  $\boldsymbol{\xi}' = \bar{\boldsymbol{\xi}}^\ell$ . Then, for every Borel set  $A \subseteq \Xi$ ,

$$\mathbb{Q}(\boldsymbol{\xi} \in A) = \Pi((\boldsymbol{\xi}, \boldsymbol{\xi}') \in A \times \Xi) = \sum_{\ell=1}^L \Pi(\boldsymbol{\xi} \in A \mid \boldsymbol{\xi}' = \bar{\boldsymbol{\xi}}^\ell) \hat{\mathbb{P}}_N(\boldsymbol{\xi}' = \bar{\boldsymbol{\xi}}^\ell) = \sum_{\ell=1}^L \mathbb{Q}_\ell(\boldsymbol{\xi} \in A) \frac{|I_\ell|}{N}.$$

The first equality follows because  $\Pi$  is a joint distribution of  $\boldsymbol{\xi}$  and  $\boldsymbol{\xi}'$  with marginals  $\mathbb{Q}$  and  $\widehat{\mathbb{P}}_N$ , respectively. The second equality follows from the law of total probability. The final equality follows from the definitions of  $\mathbb{Q}_\ell$  and  $\widehat{\mathbb{P}}_N$ . Since the above equalities holds for every Borel set, we have shown that

$$\mathbb{Q} = \sum_{\ell=1}^L \frac{|I_\ell|}{N} \mathbb{Q}_\ell.$$

Furthermore, by using similar reasoning as above, we observe that

$$\Pi(\|\boldsymbol{\xi} - \boldsymbol{\xi}'\| \leq \epsilon_N) = \sum_{\ell=1}^L \Pi(\|\boldsymbol{\xi} - \boldsymbol{\xi}'\| \leq \epsilon_N \mid \boldsymbol{\xi}' = \bar{\boldsymbol{\xi}}^\ell) \widehat{\mathbb{P}}_N(\boldsymbol{\xi}' = \bar{\boldsymbol{\xi}}^\ell) = \sum_{\ell=1}^L \mathbb{Q}_\ell(\|\boldsymbol{\xi} - \bar{\boldsymbol{\xi}}^\ell\| \leq \epsilon_N) \frac{|I_\ell|}{N}.$$

Combining the above results, the ambiguity set from (EC.66) can be rewritten as

$$\begin{aligned} \left\{ \sum_{\ell=1}^L \frac{|I_\ell|}{N} \mathbb{Q}_\ell : \sum_{\ell=1}^L \mathbb{Q}_\ell(\|\boldsymbol{\xi} - \bar{\boldsymbol{\xi}}^\ell\| \leq \epsilon_N) \frac{|I_\ell|}{N} = 1, \right. \\ \left. \mathbb{Q}_1, \dots, \mathbb{Q}_L \in \mathcal{P}(\Xi) \right\} &= \left\{ \sum_{\ell=1}^L \frac{|I_\ell|}{N} \mathbb{Q}_\ell : \mathbb{Q}_\ell(\|\boldsymbol{\xi} - \bar{\boldsymbol{\xi}}^\ell\| \leq \epsilon_N) = 1 \quad \forall \ell \in [L], \right. \\ &\quad \left. \mathbb{Q}_1, \dots, \mathbb{Q}_L \in \mathcal{P}(\Xi) \right\} \\ &= \left\{ \frac{1}{N} \sum_{j=1}^N \mathbb{Q}_j : \mathbb{Q}_j(\|\boldsymbol{\xi} - \hat{\boldsymbol{\xi}}^j\| \leq \epsilon_N) = 1 \quad \forall j \in [N], \right. \\ &\quad \left. \mathbb{Q}_1, \dots, \mathbb{Q}_N \in \mathcal{P}(\Xi) \right\}. \end{aligned}$$

The first equality follows because  $\mathbb{Q}_\ell(\|\boldsymbol{\xi} - \bar{\boldsymbol{\xi}}^\ell\| \leq \epsilon_N) \leq 1$  for each  $\ell \in [L]$ . The second equality follows because  $\mathbb{Q}_\ell(\|\boldsymbol{\xi} - \bar{\boldsymbol{\xi}}^\ell\| \leq \epsilon_N) = 1$  if and only if there exists  $\mathbb{Q}_j \in \mathcal{P}(\Xi)$  for each  $j \in I_\ell$  such that  $\mathbb{Q}_j(\|\boldsymbol{\xi} - \hat{\boldsymbol{\xi}}^j\| \leq \epsilon_N) = 1$  and  $\sum_{j \in I_\ell} \frac{1}{|I_\ell|} \mathbb{Q}_j = \mathbb{Q}_\ell$ . This concludes the proof.  $\square$

We now present the proof of Proposition 3.

**PROPOSITION 3.** *Problem (2) with uncertainty sets of the form*

$$\mathcal{U}_N^j := \left\{ \boldsymbol{\zeta} \equiv (\zeta_1, \dots, \zeta_T) \in \Xi : \|\boldsymbol{\zeta} - \hat{\boldsymbol{\xi}}^j\| \leq \epsilon_N \right\}$$

is equivalent to  $\infty$ -WDRO.

*Proof.* It follows from Lemma EC.5 that the  $\infty$ -Wasserstein ambiguity set can be decomposed into separate distributions, each having a support that is contained in  $\{\boldsymbol{\zeta} \in \Xi : \|\boldsymbol{\zeta} - \hat{\boldsymbol{\xi}}^j\| \leq \epsilon_N\}$  for  $j \in [N]$ . Of course, these sets are exactly equal to the uncertainty sets from Section 3, and thus Lemma EC.5 implies that the  $\infty$ -Wasserstein ambiguity set is equivalent to

$$\left\{ \frac{1}{N} \sum_{j=1}^N \mathbb{Q}_j : \mathbb{Q}_j \in \mathcal{P}(\mathcal{U}_N^j) \text{ for each } j \in [N] \right\}.$$

Therefore, when  $\mathcal{A}_N$  is the  $\infty$ -Wasserstein ambiguity set and each  $\mathcal{U}_N^j$  is a closed balls around  $\hat{\boldsymbol{\xi}}^j$  which is intersected with  $\Xi$ ,

$$\begin{aligned} \sup_{\mathbb{Q} \in \mathcal{A}_N} \mathbb{E}_{\mathbb{Q}} \left[ \sum_{t=1}^T \mathbf{c}_t(\boldsymbol{\xi}) \cdot \mathbf{x}_t(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}) \right] &= \frac{1}{N} \sum_{j=1}^N \sup_{\mathbb{Q} \in \mathcal{P}(\mathcal{U}_N^j)} \mathbb{E}_{\mathbb{Q}} \left[ \sum_{t=1}^T \mathbf{c}_t(\boldsymbol{\xi}) \cdot \mathbf{x}_t(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}) \right] \\ &= \frac{1}{N} \sum_{j=1}^N \sup_{\boldsymbol{\zeta} \in \mathcal{U}_N^j} \sum_{t=1}^T \mathbf{c}_t(\boldsymbol{\zeta}) \cdot \mathbf{x}_t(\zeta_1, \dots, \zeta_{t-1}). \end{aligned}$$

Moreover, it similarly follows from Lemma EC.5 that the following inequalities are equivalent:

$$\begin{aligned} \mathbb{Q} \left( \sum_{t=1}^T \mathbf{A}_t(\boldsymbol{\xi}) \mathbf{x}_t(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}) \leq \mathbf{b}(\boldsymbol{\xi}) \right) &= 1 \quad \forall \mathbb{Q} \in \mathcal{A}_N \\ \frac{1}{N} \sum_{j=1}^N \mathbb{Q}_j \left( \sum_{t=1}^T \mathbf{A}_t(\boldsymbol{\xi}) \mathbf{x}_t(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}) \leq \mathbf{b}(\boldsymbol{\xi}) \right) &= 1 \quad \forall \mathbb{Q}_j \in \mathcal{P}(\mathcal{U}_N^j), j \in [N] \\ \mathbb{Q}_j \left( \sum_{t=1}^T \mathbf{A}_t(\boldsymbol{\xi}) \mathbf{x}_t(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}) \leq \mathbf{b}(\boldsymbol{\xi}) \right) &= 1 \quad \forall \mathbb{Q}_j \in \mathcal{P}(\mathcal{U}_N^j), j \in [N] \\ \sum_{t=1}^T \mathbf{A}_t(\boldsymbol{\zeta}) \mathbf{x}_t(\zeta_1, \dots, \zeta_{t-1}) \leq \mathbf{b}(\boldsymbol{\zeta}) &\quad \forall \boldsymbol{\zeta} \in \mathcal{U}_N^j, j \in [N]. \end{aligned}$$

We have thus shown that Problem (2) and Problem (6) have equivalent objective functions and constraints under the specified constructions of the uncertainty sets and ambiguity set. This concludes the proof.  $\square$

## Appendix H: Proof of Proposition 4 from Section 6

In this appendix, we present the proof of Proposition 4.

**PROPOSITION 4.** *If  $p \in [1, \infty)$  and  $\epsilon_N > 0$ , then a decision rule is feasible for  $p$ -WDRO only if*

$$\sum_{t=1}^T \mathbf{A}_t(\zeta) \mathbf{x}_t(\zeta_1, \dots, \zeta_{t-1}) \leq \mathbf{b}(\zeta) \quad \forall \zeta \in \Xi.$$

*Proof.* Consider any arbitrary  $\bar{\xi} \in \Xi$  such that  $\bar{\xi} \neq \hat{\xi}^j$  for each  $j \in [N]$ . Let  $\delta_{\bar{\xi}}$  denote the Dirac delta distribution which satisfies  $\delta_{\bar{\xi}}(\xi = \bar{\xi}) = 1$ , and let  $\widehat{\mathbb{P}}_N := \frac{1}{N} \sum_{j=1}^N \delta_{\hat{\xi}^j}$  be the empirical distribution of the sample paths. For any  $\lambda \in (0, 1)$ , let the convex combination of the two distributions be given by

$$\mathbb{Q}_{\bar{\xi}}^\lambda := (1 - \lambda) \widehat{\mathbb{P}}_N + \lambda \delta_{\bar{\xi}}.$$

We recall the definition of the  $p$ -Wasserstein distance between  $\widehat{\mathbb{P}}_N$  and  $\mathbb{Q}_{\bar{\xi}}^\lambda$ :

$$d_p(\widehat{\mathbb{P}}_N, \mathbb{Q}_{\bar{\xi}}^\lambda) = \inf \left\{ \left( \int_{\Xi \times \Xi} \|\xi - \xi'\|^p d\Pi(\xi, \xi') \right)^{\frac{1}{p}} : \begin{array}{l} \Pi \text{ is a joint distribution of } \xi \text{ and } \xi' \\ \text{with marginals } \widehat{\mathbb{P}}_N \text{ and } \mathbb{Q}_{\bar{\xi}}^\lambda, \text{ respectively} \end{array} \right\}. \quad (\text{EC.67})$$

Consider a feasible joint distribution  $\bar{\Pi}$  for the above optimization problem in which  $\xi' \sim \mathbb{Q}_{\bar{\xi}}^\lambda$ ,  $\xi'' \sim \widehat{\mathbb{P}}_N$ , and

$$\xi = \begin{cases} \xi', & \text{if } \xi' = \hat{\xi}^j \text{ for some } j \in [N], \\ \xi'', & \text{otherwise.} \end{cases}$$

Indeed, we readily verify that the marginal distributions of  $\xi$  and  $\xi'$  are  $\widehat{\mathbb{P}}_N$  and  $\mathbb{Q}_{\bar{\xi}}^\lambda$ , respectively, and thus this joint distribution is feasible for the optimization problem in (EC.67). Moreover,

$$\begin{aligned} d_p(\widehat{\mathbb{P}}_N, \mathbb{Q}_{\bar{\xi}}^\lambda) &\leq \left( \int_{\Xi \times \Xi} \|\xi - \xi'\|^p d\bar{\Pi}(\xi, \xi') \right)^{\frac{1}{p}} \\ &= \left( \int_{\Xi \times \Xi} \|\xi - \xi'\|^p \mathbb{I}\{\xi' = \bar{\xi}\} d\bar{\Pi}(\xi, \xi') + \underbrace{\int_{\Xi \times \Xi} \|\xi - \xi'\|^p \mathbb{I}\{\xi' \neq \bar{\xi}\} d\bar{\Pi}(\xi, \xi')}_{=0} \right)^{\frac{1}{p}} \\ &= \left( \frac{1}{N} \sum_{j=1}^N \lambda \|\hat{\xi}^j - \bar{\xi}\|^p \right)^{\frac{1}{p}}. \end{aligned}$$

The inequality follows since  $\bar{\Pi}$  is a feasible but possibly suboptimal joint distribution for the optimization problem in (EC.67). The first equality follows from splitting the integral into two cases, and observing that the second case equals zero since  $\xi = \xi'$  whenever  $\xi' \neq \bar{\xi}$ . The final equality follows because  $\xi = \xi''$  whenever  $\xi' = \bar{\xi}$ , and  $\xi''$  is distributed uniformly over the historical sample paths. Thus, for any arbitrary choice of  $\xi \in \Xi$ , we have shown that  $\mathbb{Q}_{\bar{\xi}}^\lambda$  is contained in the  $p$ -Wasserstein ambiguity set whenever  $\lambda \in (0, 1)$  satisfies

$$\begin{aligned} \left( \frac{1}{N} \sum_{j=1}^N \lambda \|\hat{\xi}^j - \bar{\xi}\|^p \right)^{\frac{1}{p}} &\leq \epsilon_N \\ \frac{1}{N} \sum_{j=1}^N \lambda \|\hat{\xi}^j - \bar{\xi}\|^p &\leq \epsilon_N^p \\ \lambda &\leq \frac{\epsilon_N^p}{\frac{1}{N} \sum_{j=1}^N \|\hat{\xi}^j - \bar{\xi}\|^p}. \end{aligned}$$

Now, consider any feasible decision rule for Problem (6), *i.e.*, a decision rule  $\mathbf{x} \in \mathcal{X}$  which satisfies

$$\sum_{t=1}^T \mathbf{A}_t(\boldsymbol{\xi}) \mathbf{x}_t(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}) \leq \mathbf{b}(\boldsymbol{\xi}) \quad \mathbb{Q}\text{-a.s.}, \forall \mathbb{Q} \in \mathcal{A}_N. \quad (\text{EC.68})$$

Let  $\mathcal{A}_N$  be the  $p$ -Wasserstein ambiguity set for  $1 \leq p < \infty$  and  $\epsilon_N > 0$ . Then, for any arbitrary  $\bar{\boldsymbol{\xi}} \in \Xi$ , there exists a  $\lambda \in (0, 1)$  such that  $\mathbb{Q}_\xi^\lambda$  is contained in  $\mathcal{A}_N$ , and so it follows from (EC.68) that the decision rule must satisfy

$$\sum_{t=1}^T \mathbf{A}_t(\bar{\boldsymbol{\xi}}) \mathbf{x}_t(\bar{\boldsymbol{\xi}}_1, \dots, \bar{\boldsymbol{\xi}}_{t-1}) \leq \mathbf{b}(\bar{\boldsymbol{\xi}}).$$

Since  $\bar{\boldsymbol{\xi}} \in \Xi$  was chosen arbitrarily, we conclude that the decision rule must satisfy

$$\sum_{t=1}^T \mathbf{A}_t(\zeta) \mathbf{x}_t(\zeta_1, \dots, \zeta_{t-1}) \leq \mathbf{b}(\zeta) \quad \forall \zeta \in \Xi,$$

which is what we wished to show.  $\square$

## Appendix I: Reformulation of Problem (7) from Section 7.1.2

In this appendix, we develop a reformulation for

$$\begin{aligned} & \underset{\substack{x_2^1, \dots, x_2^N \in \mathbb{R}, \\ x_3: \mathbb{R}^2 \rightarrow \mathbb{R}}}{\text{minimize}} && \frac{1}{N} \sum_{j=1}^N \max_{k \in \mathcal{K}_j} \sup_{\zeta \in \mathcal{U}_N^j \cap P^k} \{x_2^k + 2x_3(\zeta_1, \zeta_2)\} \\ & \text{subject to} && x_2^k + x_3(\zeta_1, \zeta_2) \geq \zeta_1 + \zeta_2 && \forall \zeta \in \cup_{j=1}^N \mathcal{U}_N^j \cap P^k, k \in [N] \\ & && x_2^k, x_3(\zeta_1, \zeta_2) \geq 0 && \forall \zeta \in \cup_{j=1}^N \mathcal{U}_N^j \cap P^k, k \in [N]. \end{aligned} \quad (7)$$

Indeed, for each partition index  $k \in \{1, \dots, N\}$ , we observe that an optimal decision rule for third stage in Problem (7) is given by

$$x_3^*(\zeta_1, \zeta_2) = \max \{0, \zeta_1 + \zeta_2 - x_2^k\} \quad \forall (\zeta_1, \zeta_2) \in \cup_{j=1}^N \mathcal{U}_N^j \cap P^k.$$

Therefore, for each term in the objective of Problem (7),

$$\sup_{\zeta \in \mathcal{U}_N^j \cap P^k} \{x_2^k + 2x_3^*(\zeta_1, \zeta_2)\} = \max_{\zeta \in \mathcal{U}_N^j \cap P^k} \{x_2^k + 2 \max \{0, \zeta_1 + \zeta_2 - x_2^k\}\} = x_2^k + \max \left\{ 0, \max_{\zeta \in \mathcal{U}_N^j \cap P^k} \{\zeta_1 + \zeta_2\} - x_2^k \right\}.$$

Under the specified construction of the partitions, and using the fact that the uncertainty sets are constructed with the  $\ell_\infty$ -norm, we observe that  $M^{jk} := \max_{\zeta \in \mathcal{U}_N^j \cap P^k} \{\zeta_1 + \zeta_2\}$  can be computed in closed form. Therefore, Problem (7) can be reformulated as the following linear optimization problem:

$$\begin{aligned} & \underset{\substack{x_2^1, \dots, x_2^N \in \mathbb{R}, \\ v^1, \dots, v^N \in \mathbb{R}}}{\text{minimize}} && \frac{1}{N} \sum_{j=1}^N v^j \\ & \text{subject to} && v^j \geq x_2^k && \forall j \in [N], k \in \mathcal{K}_j \\ & && v^j \geq x_2^k + 2(M^{jk} - x_2^k) && \forall j \in [N], k \in \mathcal{K}_j \\ & && x_2^k \geq 0 && \forall k \in [N]. \end{aligned}$$

The number of decision variables in the above linear optimization problem is  $2N = O(N)$  and the number of constraints is  $2 \sum_{j=1}^N |\mathcal{K}_j| + N = O(\sum_{j=1}^N |\mathcal{K}_j|)$ .



## Appendix J: Linear Decision Rules for Problem (6) with 1-Wasserstein Ambiguity Sets

In this appendix, we present a reformulation of linear decision rules for Problem (6) using the 1-Wasserstein ambiguity set. The performance of this data-driven approach is illustrated in Section 7.

We first review the necessary notation. Following Section 5, we focus on a specific case of Problem (6) of the form

$$\begin{aligned} & \underset{\mathbf{x} \in \mathcal{X}}{\text{minimize}} && \sup_{\mathbb{Q} \in \mathcal{A}_N} \mathbb{E}_{\mathbb{Q}} \left[ \sum_{t=1}^T \mathbf{c}_t \cdot \mathbf{x}_t(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}) \right] \\ & \text{subject to} && \sum_{t=1}^T \mathbf{A}_t \mathbf{x}_t(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}) \leq \mathbf{b}(\boldsymbol{\xi}) \quad \mathbb{Q}\text{-a.s.}, \forall \mathbb{Q} \in \mathcal{A}_N, \end{aligned} \quad (\text{EC.69})$$

in which  $\mathbf{A}_t(\boldsymbol{\xi})$  and  $\mathbf{c}_t(\boldsymbol{\xi})$  do not depend on the stochastic process. The ambiguity set is constructed as

$$\mathcal{A}_N = \left\{ \mathbb{Q} \in \mathcal{P}(\Xi) : d_1(\mathbb{Q}, \widehat{\mathbb{P}}_N) \leq \epsilon_N \right\},$$

where  $\widehat{\mathbb{P}}_N$  is the empirical distribution of the historical data,  $\epsilon_N \geq 0$  is the robustness parameter, and the 1-Wasserstein distance between two distributions is given by

$$d_1(\mathbb{Q}, \mathbb{Q}') = \inf \left\{ \int_{\Xi \times \Xi} \|\boldsymbol{\xi} - \boldsymbol{\xi}'\| d\Pi(\boldsymbol{\xi}, \boldsymbol{\xi}') : \begin{array}{l} \Pi \text{ is a joint distribution of } \boldsymbol{\xi} \text{ and } \boldsymbol{\xi}' \\ \text{with marginals } \mathbb{Q} \text{ and } \mathbb{Q}', \text{ respectively} \end{array} \right\}.$$

We refer to Section 6 for more details on the 1-Wasserstein ambiguity set. We assume that the robustness parameter satisfies  $\epsilon_N > 0$ , in which case it follows from Proposition 4 in Section 6 that Problem (EC.69) is equivalent to

$$\begin{aligned} & \underset{\mathbf{x} \in \mathcal{X}}{\text{minimize}} && \sup_{\mathbb{Q} \in \mathcal{A}_N} \mathbb{E}_{\mathbb{Q}} \left[ \sum_{t=1}^T \mathbf{c}_t \cdot \mathbf{x}_t(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1}) \right] \\ & \text{subject to} && \sum_{t=1}^T \mathbf{A}_t \mathbf{x}_t(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_{t-1}) \leq \mathbf{b}(\boldsymbol{\zeta}) \quad \forall \boldsymbol{\zeta} \in \Xi. \end{aligned} \quad (\text{EC.70})$$

We next present an extension of the linear decision rule approach to Problem (EC.70), in which we restrict the space of decision rules to those of the form

$$\mathbf{x}_t(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_{t-1}) = \mathbf{x}_{t,0} + \sum_{s=1}^{t-1} \mathbf{X}_{t,s} \boldsymbol{\zeta}_s.$$

The resulting approximation of Problem (EC.70) is given by

$$\begin{aligned} & \underset{\mathbf{x} \in \mathcal{X}}{\text{minimize}} && \sup_{\mathbb{Q} \in \mathcal{A}_N} \mathbb{E}_{\mathbb{Q}} \left[ \sum_{t=1}^T \mathbf{c}_t \cdot \left( \mathbf{x}_{t,0} + \sum_{s=1}^{t-1} \mathbf{X}_{t,s} \boldsymbol{\xi}_s \right) \right] \\ & \text{subject to} && \sum_{t=1}^T \mathbf{A}_t \left( \mathbf{x}_{t,0} + \sum_{s=1}^{t-1} \mathbf{X}_{t,s} \boldsymbol{\zeta}_s \right) \leq \mathbf{b}(\boldsymbol{\zeta}) \quad \forall \boldsymbol{\zeta} \in \Xi, \end{aligned} \quad (\text{EC.71})$$

where the decision variables are  $\mathbf{x}_{t,0} \in \mathbb{R}^{n_t}$  and  $\mathbf{X}_{t,s} \in \mathbb{R}^{n_t \times d_s}$  for all  $1 \leq s < t$ .

In the remainder of this appendix, we develop a tractable reformulation of Problem (EC.71). Our reformulation, which will use similar duality techniques to those presented in Section 5, is presented as Theorem EC.1. Our reformulation requires the following assumption:

**ASSUMPTION EC.1.** *The set  $\Xi \subseteq \mathbb{R}^d$  is a nonempty multi-dimensional box of the form  $[\boldsymbol{\ell}, \mathbf{u}]$ , where any component of  $\boldsymbol{\ell}$  might be  $-\infty$  and any component of  $\mathbf{u}$  may be  $\infty$ . Moreover, the norm in the 1-Wasserstein distance is equal to  $\|\cdot\|_1$ .*

We now present the reformulation of Problem (EC.71) given Assumption EC.1.

**THEOREM EC.1.** *If Assumption EC.1 holds, then Problem (EC.71) can be reformulated by adding at most  $O(md)$  additional continuous decision variables and  $O(md)$  additional linear constraints. The reformulation is given by*

$$\begin{aligned}
& \text{minimize} && \lambda \epsilon_N + \frac{1}{N} \sum_{j=1}^N \left( \sum_{t=1}^T \mathbf{c}_t \cdot \left( \mathbf{x}_{t,0} + \sum_{s=1}^{t-1} \mathbf{X}_{t,s} \hat{\boldsymbol{\xi}}_s^j \right) + \boldsymbol{\alpha} \cdot (\mathbf{u} - \hat{\boldsymbol{\xi}}^j) + \boldsymbol{\beta}(\hat{\boldsymbol{\xi}}^j - \boldsymbol{\ell}) \right) \\
& \text{subject to} && \left\| \sum_{s=t+1}^T (\mathbf{X}_{s,t})^\top \mathbf{c}_s - \boldsymbol{\alpha}_t + \boldsymbol{\beta}_t \right\|_\infty \leq \lambda && t \in [T] \\
& && \mathbf{M}_t - \boldsymbol{\Lambda}_t = -\mathbf{B}_t + \sum_{s=t+1}^T \mathbf{A}_s \mathbf{X}_{s,t} && t \in [T] \\
& && \sum_{t=1}^T (\mathbf{M}_t \mathbf{u}_t - \boldsymbol{\Lambda}_t \boldsymbol{\ell}_t + \mathbf{A}_t \mathbf{x}_{t,0}) \leq \mathbf{b}^0,
\end{aligned}$$

where the auxiliary decision variables are  $\boldsymbol{\alpha} := (\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_T), \boldsymbol{\beta} := (\boldsymbol{\beta}_1, \dots, \boldsymbol{\beta}_T) \in \mathbb{R}_+^d$ , as well as  $\mathbf{M} := (\mathbf{M}_1, \dots, \mathbf{M}_T), \boldsymbol{\Lambda} := (\boldsymbol{\Lambda}_1, \dots, \boldsymbol{\Lambda}_T) \in \mathbb{R}_+^{m \times d}$ .

*Proof.* Following similar reasoning to Theorem 4, the constraints

$$\sum_{t=1}^T \mathbf{A}_t \left( \mathbf{x}_{t,0} + \sum_{s=1}^{t-1} \mathbf{X}_{t,s} \boldsymbol{\zeta}_s \right) \leq \mathbf{b}^0 + \sum_{t=1}^T \mathbf{B}_t \boldsymbol{\zeta}_t \quad \forall \boldsymbol{\zeta} \in \Xi$$

are satisfied if and only if there exist  $\mathbf{M} := (\mathbf{M}_1, \dots, \mathbf{M}_T), \boldsymbol{\Lambda} := (\boldsymbol{\Lambda}_1, \dots, \boldsymbol{\Lambda}_T) \in \mathbb{R}_+^{m \times d}$  which satisfy

$$\begin{aligned}
& \sum_{t=1}^T (\mathbf{M}_t \mathbf{u}_t - \boldsymbol{\Lambda}_t \boldsymbol{\ell}_t + \mathbf{A}_t \mathbf{x}_{t,0}) \leq \mathbf{b}^0, \\
& \mathbf{M}_t - \boldsymbol{\Lambda}_t = \sum_{s=t+1}^T \mathbf{A}_s \mathbf{X}_{s,t} - \mathbf{B}_t, \quad t \in [T].
\end{aligned}$$

The remainder of the proof focuses on the objective function. Note that for any fixed solution to Problem (EC.71) one can define a function  $f: \Xi \rightarrow \mathbb{R}$  as follows

$$f(\boldsymbol{\zeta}) = \sum_{t=1}^T \mathbf{c}_t \cdot \left( \mathbf{x}_{t,0} + \sum_{s=1}^{t-1} \mathbf{X}_{t,s} \boldsymbol{\zeta}_s \right).$$

It follows from Assumption EC.1 that  $\Xi \subseteq \mathbb{R}^d$  is nonempty, convex, and closed, and  $-f(\cdot)$  is proper, convex, and lower semicontinuous on  $\Xi$ , thus satisfying [Mohajerin Esfahani and Kuhn \(2018, Assumption 4.1\)](#). Therefore, we conclude from [Mohajerin Esfahani and Kuhn \(2018, Equation 12b\)](#) that

$$\begin{aligned}
\sup_{\mathbf{Q} \in \mathcal{A}_N} \mathbb{E}_{\mathbf{Q}} \left[ \sum_{t=1}^T \mathbf{c}_t \cdot \left( \mathbf{x}_{t,0} + \sum_{s=1}^{t-1} \mathbf{X}_{t,s} \boldsymbol{\xi}_s \right) \right] &= \sup_{\mathbf{Q} \in \mathcal{A}_N} \mathbb{E}_{\mathbf{Q}} [f(\boldsymbol{\xi})] \\
&= \inf_{\lambda \geq 0} \lambda \epsilon_N + \frac{1}{N} \sum_{j=1}^N \sup_{\boldsymbol{\zeta} \in \Xi} \left\{ f(\boldsymbol{\zeta}) - \lambda \|\boldsymbol{\zeta} - \hat{\boldsymbol{\xi}}^j\|_1 \right\} \\
&= \inf_{\lambda \geq 0} \lambda \epsilon_N + \frac{1}{N} \sum_{j=1}^N \underbrace{\sup_{\boldsymbol{\zeta} \in \Xi} \left\{ \sum_{t=1}^T \mathbf{c}_t \cdot \left( \mathbf{x}_{t,0} + \sum_{s=1}^{t-1} \mathbf{X}_{t,s} \boldsymbol{\zeta}_s \right) - \lambda \|\boldsymbol{\zeta} - \hat{\boldsymbol{\xi}}^j\|_1 \right\}}_{\gamma_j}.
\end{aligned} \tag{EC.72}$$

We now reformulate the expression  $\gamma_j$  for each  $j \in [N]$ . Indeed, it follows from strong duality for linear programming that

$$\begin{aligned} \gamma_j = & \underset{\alpha, \beta \in \mathbb{R}_+^d}{\text{minimize}} && \sum_{t=1}^T \left( \mathbf{c}_t \cdot \mathbf{x}_{t,0} + \alpha_t \cdot (\mathbf{u}_t - \hat{\xi}_t^j) + \beta_t \cdot (\hat{\xi}_t^j - \ell_t) \right) \\ & \text{subject to} && \left\| \sum_{s=t+1}^T (\mathbf{X}_{s,t})^\top \mathbf{c}_s - \alpha_t + \beta_t \right\|_\infty \leq \lambda, \quad t \in [T]. \end{aligned} \quad (\text{EC.73})$$

*Remark:* For any index  $l$  such that  $u_l = \infty$  (alternatively,  $\ell_l = -\infty$ ), the corresponding decision variable  $\alpha_l$  (alternatively,  $\beta_l$ ) should be set to zero and the term  $\alpha_l(u_l - \hat{\xi}_l^j)$  (alternatively,  $\beta_l(\hat{\xi}_l^j - \ell_l)$ ) should be dropped from the objective.

Note that problem (EC.73) is component-wise separable to  $d$  problems of the form

$$\begin{aligned} & \underset{\alpha_l, \beta_l \in \mathbb{R}_+}{\text{minimize}} && \alpha_l(u_l^k - \hat{\xi}_l^j) + \beta_l(\hat{\xi}_l^j - \ell_l) \\ & \text{subject to} && |g_l - \alpha_l + \beta_l| \leq \lambda, \end{aligned} \quad (\text{EC.74})$$

where  $\mathbf{g} := (\sum_{s=2}^T (\mathbf{X}_{s,1})^\top \mathbf{c}_s, \sum_{s=3}^T (\mathbf{X}_{s,2})^\top \mathbf{c}_s, \dots, (\mathbf{X}_{T,T-1})^\top \mathbf{c}_T, 0) \in \mathbb{R}^d$ . Moreover,  $\hat{\xi}^j \in \Xi$  implies that both  $(u_l - \hat{\xi}_l^j)$  and  $(\hat{\xi}_l^j - \ell_l)$  are nonnegative, and so for any fixed  $\lambda$  and  $g_l$ , an optimal solution of (EC.74) is given by  $\alpha_l = \max\{g_l - \lambda, 0\}$  and  $\beta_l = \max\{-g_l - \lambda, 0\}$  (their corresponding minimal values). This solution is independent of the value of  $\hat{\xi}_l^j$ , and therefore, the same variables  $\alpha$  and  $\beta$  can be used in (EC.73) for all values of  $j \in [N]$ . Combining (EC.72) and (EC.73) and plugging the result to the objective function of (EC.71), we obtain the desired formulation.  $\square$

## Appendix K: Supplement to Section 7.2.3

In this appendix, we provide supplemental numerical results for the multi-stage stochastic inventory management problem from Section 7.2. Specifically, the aim of this appendix is to evaluate the impact of the projection procedure, described at the end of Section 7.2.2, on the out-of-sample costs of SRO-LDR and SAA-LDR reported in Table 1.

Following the same notation in Section 7.2.2, let  $\mathbf{x}^{\mathcal{A},i,\ell} = (x_1^{\mathcal{A},i,\ell}, \dots, x_T^{\mathcal{A},i,\ell})$  be the production quantities obtained when the decision rule from approach  $\mathcal{A}$  on training dataset  $\ell$  is applied to the  $i$ th sample path in the testing dataset. For each approach  $\mathcal{A}$  and training dataset  $\ell$ , the probability that the resulting decision rule is *feasible* is approximated by

$$P^{\mathcal{A},\ell} = \frac{1}{10000} \sum_{i=1}^{10000} \mathbb{1}_{\{\mathbf{x}^{\mathcal{A},i,\ell} \in [0, \bar{\mathbf{x}}]\}},$$

and the *infeasibility magnitude* is approximated by

$$C^{\mathcal{A},\ell} = \frac{1}{10000} \sum_{i=1}^{10000} \min_{\mathbf{y} \in [0, \bar{\mathbf{x}}]} \|\mathbf{x}^{\mathcal{A},i,\ell} - \mathbf{y}\|_1.$$

In other words,  $P^{\mathcal{A},\ell}$  tells us how frequently the projection procedure needs to be applied, and  $C^{\mathcal{A},\ell}$  captures the average number of production units which are changed due to the projection procedure.

In Tables EC.1 and EC.2, for each experiment in Section 7.2.3 and for each approach  $\mathcal{A} \in \{\text{SRO-LDR}, \text{SAA-LDR}\}$ , we report the average and standard deviations for  $P^{\mathcal{A},\ell}$  and  $C^{\mathcal{A},\ell}$  over the 100 training datasets. For almost all choices of  $T$ ,  $\alpha$ , and  $N$ , the decision rules produced by SRO-LDR are feasible for over 93% of the sample paths in testing dataset, and their *infeasibility magnitude* is below 2 units. These results imply that the projection procedure does not significantly impact the out-of-sample cost of SRO-LDR reported in Table 1. In contrast, SAA-LDR produces decision rules which have low feasibility and high *infeasibility magnitude* when  $N$  is small. This shows, for small training datasets, that the decision rules obtained by SAA-LDR can be unreliable and require significant corrections to obtain feasible production quantities.

**Table EC.1 Multi-stage stochastic inventory management: out-of-sample feasibility.**

$T$	$\alpha$	Approach	Size of training dataset (N)			
			10	25	50	100
5	0	SRO-LDR	96.3(6.6)	98.9(2.1)	99.7(0.8)	100.0(0.1)
		SAA-LDR	83.0(13.4)	96.5(4.3)	99.5(1.3)	100.0(0.1)
	0.25	SRO-LDR	93.8(7.3)	95.9(3.5)	97.3(2.3)	98.1(1.1)
		SAA-LDR	79.2(12.9)	92.3(5.3)	96.5(2.8)	98.1(1.1)
	0.5	SRO-LDR	89.7(8.6)	91.0(4.9)	91.1(3.7)	94.1(2.4)
		SAA-LDR	73.4(11.3)	85.4(4.9)	89.9(3.5)	94.0(2.3)
10	0	SRO-LDR	99.6(1.0)	100.0(0.1)	100.0(0.0)	100.0(0.0)
		SAA-LDR	61.5(24.6)	99.0(1.6)	100.0(0.1)	100.0(0.0)
	0.25	SRO-LDR	99.4(1.8)	99.9(0.3)	100.0(0.1)	100.0(0.0)
		SAA-LDR	60.2(23.9)	97.8(2.2)	99.8(0.4)	100.0(0.0)
	0.5	SRO-LDR	96.7(2.9)	97.7(1.4)	98.6(0.7)	98.9(0.3)
		SAA-LDR	57.6(22.4)	93.9(3.0)	97.7(1.2)	98.9(0.3)

Mean (standard deviation) of the percentage of the 10,000 sample paths in the testing dataset for which the linear decision rule resulted in feasible production quantities ( $P^{A,i}$ ). In other words, 100% minus the above values indicates the percentage of sample paths in the testing dataset for which the production quantities needed correction. The mean and standard deviation are computed over 100 training datasets for each value of  $N$ ,  $T$ ,  $\alpha$ .

**Table EC.2 Multi-stage stochastic inventory management: infeasibility magnitude.**

$T$	$\alpha$	Approach	Size of training dataset (N)			
			10	25	50	100
5	0	SRO-LDR	0.5(1.6)	0.1(0.3)	0.0(0.0)	0.0(0.0)
		SAA-LDR	4.6(6.5)	0.4(0.9)	0.0(0.1)	0.0(0.0)
	0.25	SRO-LDR	0.8(1.4)	0.4(0.6)	0.2(0.3)	0.1(0.1)
		SAA-LDR	5.7(6.6)	0.9(1.1)	0.2(0.3)	0.1(0.1)
	0.5	SRO-LDR	1.7(2.1)	1.1(0.8)	1.0(0.6)	0.6(0.4)
		SAA-LDR	7.8(7.3)	2.0(1.2)	1.1(0.7)	0.6(0.4)
10	0	SRO-LDR	0.0(0.1)	0.0(0.0)	0.0(0.0)	0.0(0.0)
		SAA-LDR	218.2(1417.0)	0.1(0.2)	0.0(0.0)	0.0(0.0)
	0.25	SRO-LDR	0.0(0.2)	0.0(0.0)	0.0(0.0)	0.0(0.0)
		SAA-LDR	218.8(1417.2)	0.2(0.3)	0.0(0.0)	0.0(0.0)
	0.5	SRO-LDR	0.4(0.4)	0.2(0.2)	0.1(0.1)	0.1(0.0)
		SAA-LDR	220.3(1417.4)	0.7(0.5)	0.2(0.2)	0.1(0.0)

Mean (standard deviation) of the average *infeasibility magnitude* on the testing dataset resulting from applying the projecting procedure on the production quantity produced by the linear decision rule ( $C^{A,i}$ ). The mean and standard deviation are computed over 100 training datasets for each value of  $N$ ,  $T$ ,  $\alpha$ .