

Dynamic optimization with side information

Dimitris Bertsimas, Christopher McCord, Bradley Sturt

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

We develop a tractable and flexible approach for incorporating side information into dynamic optimization under uncertainty. The proposed framework uses predictive machine learning methods (such as k -nearest neighbors, kernel regression, and random forests) to weight the relative importance of various data-driven uncertainty sets in a robust optimization formulation. Through a novel measure concentration result for a class of machine learning methods, we prove that the proposed approach is asymptotically optimal for multi-period stochastic programming with side information. We also describe a general-purpose approximation for these optimization problems, based on overlapping linear decision rules, which is computationally tractable and produces high-quality solutions for dynamic problems with many stages. Across a variety of examples in inventory management, finance, and shipment planning, our method achieves improvements of up to 15% over alternatives and requires less than one minute of computation time on problems with twelve stages.

Key words: Distributionally robust optimization; machine learning; dynamic optimization.

History: This paper was first submitted in May 2019. A revision was submitted in May 2020.

1. Introduction

Dynamic decision making under uncertainty forms the foundation for numerous fundamental problems in operations research and management science. In these problems, a decision maker attempts to minimize an uncertain objective over time, as information incrementally becomes available. For example, consider a retailer with the goal of managing the inventory of a new short life cycle product. Each week, the retailer must decide an ordering quantity to replenish its inventory. Future demand for the product is unknown, but the retailer can base its ordering decisions on the remaining inventory level, which depends on the realized demands in previous weeks. A risk-averse investor faces a similar problem when constructing and adjusting a portfolio of assets in order to achieve a desirable risk-return tradeoff over a horizon of many months. Additional examples abound in energy planning, airline routing, and ride sharing, as well as in many other areas.

To make high quality decisions in dynamic environments, the decision maker must accurately model future uncertainty. Often, practitioners have access to *side information* or *auxiliary covariates*, which can help predict that uncertainty. For a retailer, although the future demand for a newly introduced clothing item is unknown, data on the brand, style, and color of the item, as well

as data on market trends and social media, can help predict it. For a risk-averse investor, while the returns of the assets in future stages are uncertain, recent asset returns and prices of relevant options can provide crucial insight into upcoming volatility. Consequently, organizations across many industries are continuing to prioritize the use of predictive analytics in order to leverage vast quantities of data to understand future uncertainty and make better operational decisions.

In this paper, we address these applications by studying the following class of multi-period stochastic decision problems. Specifically, we consider problems faced by organizations in which decisions $\mathbf{x}_1 \in \mathcal{X}_1 \subseteq \mathbb{R}^{d_x^1}, \dots, \mathbf{x}_T \in \mathcal{X}_T \subseteq \mathbb{R}^{d_x^T}$ are chosen sequentially, as random vectors $\boldsymbol{\xi}_1 \in \Xi_1 \subseteq \mathbb{R}^{d_\xi^1}, \dots, \boldsymbol{\xi}_T \in \Xi_T \subseteq \mathbb{R}^{d_\xi^T}$ become incrementally available at each temporal period. Before selecting any decisions, we observe side information, $\boldsymbol{\gamma} \in \Gamma \subseteq \mathbb{R}^{d_\gamma}$, which may be predictive of the uncertain quantities observed in the subsequent periods. The goal is to choose a decision rule (policy) which minimizes the conditional expected cost over the entire problem horizon:

$$v^*(\bar{\boldsymbol{\gamma}}) \triangleq \underset{\mathbf{x}_t: \Xi_1 \times \dots \times \Xi_{t-1} \rightarrow \mathcal{X}_t}{\text{minimize}} \quad \mathbb{E} \left[c(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T, \mathbf{x}_1, \mathbf{x}_2(\boldsymbol{\xi}_1), \dots, \mathbf{x}_T(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{T-1})) \mid \boldsymbol{\gamma} = \bar{\boldsymbol{\gamma}} \right]. \quad (1)$$

However, the only insight into the joint probability distribution $(\boldsymbol{\gamma}, \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T)$ comes from historical data, $(\boldsymbol{\gamma}^1, \boldsymbol{\xi}_1^1, \dots, \boldsymbol{\xi}_T^1), \dots, (\boldsymbol{\gamma}^N, \boldsymbol{\xi}_1^N, \dots, \boldsymbol{\xi}_T^N)$, which are assumed to be independent and identically distributed (i.i.d.) realizations of the underlying joint distribution. Throughout the paper, we do not impose any parametric structure on the correlations across $(\boldsymbol{\gamma}, \boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T)$, and presume that the structure of optimal decision rules to (1) is unknown. The aim of the present paper is to develop general-purpose approaches to harness this data to approximately solve the stochastic problem (1).

Such dynamic optimization problems with an initial observation of side information arise in many operational contexts. For example, fashion retailers have access to data on the brand, style, and color of a new clothing item prior to any sales, which are predictive of demand for the product in each week of its lifecycle. Similarly, in finance, important economic data (such as the consumer price index CPI and key numbers from the US Bureau of Labor Statistics report) are released monthly on a fixed schedule, and this data serves as side information for a fund manager who seeks to balance the risk of a portfolio in each day of the ensuing month. Consequently, from a modeling perspective, (1) encompasses the variety of decision problems faced by organizations in which side information does not change over time (*e.g.*, the fashion retailer) or varies on a much longer time scale than the length of the decision horizon (*e.g.*, the fund manager).

A recent body of work has aimed to leverage predictive analytics to address (1) in the particular case of single-period problems ($T = 1$). For example, [Hannah et al. \(2010\)](#), [Ban and Rudin \(2018\)](#), [Bertsimas and Kallus \(2020\)](#), [Ho and Hanasusanto \(2019\)](#) investigate prescriptive approaches, based on sample average approximation, that use local machine learning to assign weights to the historical

data based on side information. Bertsimas and Van Parys (2017) propose adding robustness to those weights to achieve optimal asymptotic budget guarantees. Elmachtoub and Grigas (2017) develop an approach for linear optimization problems in which a machine learning model is trained to minimize the decision cost. Unfortunately, prescriptive approaches designed for single-period problems do not generally extend to (1), as illustrated by the following example.

EXAMPLE 1. Suppose a decision maker attempted to approximate (1) by solving

$$\underset{\mathbf{x}_t: \Xi_1 \times \dots \times \Xi_{t-1} \rightarrow \mathcal{X}_t}{\text{minimize}} \quad \sum_{i=1}^N w_N^i(\bar{\gamma}) c(\boldsymbol{\xi}_1^i, \dots, \boldsymbol{\xi}_T^i, \mathbf{x}_1, \mathbf{x}_2(\boldsymbol{\xi}_1^i), \dots, \mathbf{x}_T(\boldsymbol{\xi}_1^i, \dots, \boldsymbol{\xi}_{T-1}^i)), \quad (2)$$

where $w_N^i(\cdot)$ are weight functions (satisfying $\sum_{i=1}^N w_N^i(\bar{\gamma}) = 1$) derived from machine learning methods applied to historical data. When $T = 1$ and the weight functions are constructed using a suitable class of machine learning methods, Bertsimas and Kallus (2020) show under certain conditions that the above optimization problem is asymptotically optimal, and will thus provide a near-optimal approximation of (1) in big data settings. However, it is readily observed that approaches such as (2) will result in a poor approximation of the underlying multi-period stochastic decision problem with side information when $T \geq 2$, as the optimal decision rules produced by (2) will generally be “anticipative” with respect to the historical data.¹ Such anticipativity (a form of *overfitting*) is ultimately of practical importance, as it implies that the (2) can provide an unsuitable approximation of (1) even in the presence of big data.

To circumvent overfitting in the context of multi-period problems with side information, recent literature have aimed to address (1) by constructing scenario trees. Scenario trees have been long studied in the stochastic programming literature, and essentially address overfitting by encoding the various ways that uncertainty can unfold across time. For a class of multi-period inventory management problems with side information, Ban et al. (2019) propose fitting historical data and side information to a parametric regression model, and establish asymptotic optimality when the model is correctly specified. Bertsimas and McCord (2019) propose a different approach based on dynamic programming that uses nonparametric machine learning methods to handle auxiliary side information. These papers also extend to problems where side information is observed at multiple periods. However, these dynamic approaches require scenario tree enumeration and suffer from the curse of dimensionality. As a result, and despite their asymptotic optimality, the existing approaches for addressing (1) can require hours or days to obtain high-quality solutions for problems with ten or fewer time periods.

1.1. Contributions

The aim of the present paper, in a nutshell, is to develop a machine learning-based approach for addressing (1) which remains computationally tractable for operational problems with many

periods. To this end, we develop a new approach to (1) by a natural combination of prescriptive analytics (2) with recent techniques from robust optimization to avoid overfitting (Bertsimas et al. 2018a), and the present paper unifies our understanding of these disparate models through a novel asymptotic theory.

Our proposed combination of two streams of literature (prescriptive analytics and robust optimization) is ultimately viewed as attractive from a practical standpoint. Across multi-period and single-period problems from several applications (shipment planning, inventory management, and finance), the proposed approach produces solutions with up to 15% improvement in average out-of-sample cost compared to alternatives. In particular, the approach does not require a scenario tree, and as a result, is significantly more tractable compared to existing approaches for dynamic optimization with side information. To the best of our knowledge, this is the first approach to address (1) which offers asymptotic optimality guarantees while remaining practically tractable for problems with many periods, thus offering organizations a general-purpose tool for better decision making with predictive analytics.

In greater detail, the key results of this paper are the following:

- (a) We propose addressing (1) by combining the prescriptive analytics approach (2) with a technique of Bertsimas et al. (2018a) to avoid overfitting in multi-period problems.
- (b) We prove under mild conditions that this combination of machine learning and robust optimization is asymptotically optimal for (1) for general spaces of decision rules (Theorem 1).
- (c) To establish the above guarantee, we show for the first time that an *empirical conditional probability distribution* that is constructed from machine learning methods will, as more data is obtained, converge to the underlying *conditional probability distribution* with respect to the type-1 Wasserstein distance (Theorem 2).
- (d) As a byproduct of the new measure concentration result, we show how side information and machine learning can be tractably incorporated into (single-period) Wasserstein-based distributionally robust optimization problems while maintaining its attractive asymptotic optimality.
- (e) To find high quality solutions for problems with many stages in practical computation times, we develop a tractable approximation algorithm for these robust optimization problems by extending an approach of Bertsimas et al. (2019), Chen et al. (2020) to multi-period problems.
- (f) Across multi-period and single-period problems from several applications (shipment planning, inventory management, and finance), we show that the proposed combination of machine

learning and robust optimization outperforms alternatives with up to 15% improvement in average out-of-sample cost. In particular, the proposed approach is practical and scalable, requiring less than one minute on examples with up to twelve stages.

The paper is organized as follows. Section 2 introduces the problem setting and notation. Section 3 proposes the new framework for incorporating machine learning into dynamic optimization. Section 4 develops theoretical guarantees on the proposed approach. Section 5 discusses implications of these results in the context of single-period distributionally robust optimization with the type-1 Wasserstein ambiguity set. Section 6 presents the general multi-policy approximation scheme for dynamic optimization with side information. Section 7 presents a detailed investigation and computational simulations of the proposed methodology in shipment planning, inventory management, and finance. We conclude in Section 8.

1.2. Comparison to related work

This paper follows a recent body of literature on data-driven optimization under uncertainty in operations research and management science. Much of this work has focused on the paradigm of distributionally robust optimization, in which the optimal solution is that which performs best in expectation over a worst-case probability distribution from an ambiguity set. Motivated by probabilistic guarantees, distributionally robust optimization has found particular applicability in data-driven settings in which the ambiguity set is constructed using historical data, such as Delage and Ye (2010), Xu et al. (2012), Mohajerin Esfahani and Kuhn (2018), Van Parys et al. (2017). In particular, the final steps in our convergence result (Section 4.4) draw heavily from similar techniques from Mohajerin Esfahani and Kuhn (2018) and Bertsimas et al. (2018a). In contrast to previous work, this paper develops a new measure concentration result for the empirical conditional probability distribution (Section 4.3) which enables machine learning and side information to be incorporated into sample robust optimization and Wasserstein-based distributionally robust optimization for the first time.

To the best of our knowledge, the proposed combination of machine learning and robust optimization for addressing (1) is novel and its theoretical justification does not follow from the existing literature. With respect to prescriptive analytics, Bertsimas and Kallus (2020) establish asymptotic optimality guarantees for problems of the form (2) in the case of $T = 1$. Their result requires that the cost function is equicontinuous. Their proof relies on results from the machine learning literature (Walk 2010), which show that an appropriately constructed *empirical conditional probability distribution* (with weights $\{w_N^i(\bar{\gamma})\}$ assigned to each historical observation ξ^i) weakly converges

to the *true conditional probability distribution* of ξ given $\gamma = \bar{\gamma}$, under certain assumptions. However, the asymptotic optimality and proof techniques do not apply to (2) when $T \geq 2$, since the cost function resulting from decision rules is not equicontinuous in general. For problems without side information, Bertsimas et al. (2018a) circumvent the requirement of equicontinuity by adding robustness to the historical data. To establish asymptotic optimality, they use the fact that the empirical probability distribution of the uncertainties concentrates around the true distribution with respect to the type-1 Wasserstein distance. In the present paper, we unify these proof techniques by developing a new measure concentration result for machine learning which shows that the empirical conditional probability distribution produced by appropriate weight functions concentrates around the true conditional probability distribution with respect to the type-1 Wasserstein distance. This establishes the asymptotic optimality of our robustification of (2) for multi-stage stochastic decision problems with side information.

Several recent papers have focused on tractable approximations of two- and multi-stage *distributionally* and *sample* robust optimization. Many approaches are based around policy approximation schemes, including lifted linear decision rules (Bertsimas et al. 2018b), K -adaptivity (Hanasusanto et al. 2016), and finite adaptability (Bertsimas et al. 2018a). Alternative approaches include tractable approximations of copositive formulations (Natarajan et al. 2011, Hanasusanto and Kuhn 2018). Closest related to the approximation scheme in this paper are Chen et al. (2020) and Bertsimas et al. (2019), which address two-stage problems via overlapping decision rules. Chen et al. (2020) propose a *scenario-wise* modeling approach that leads to novel approximations of various distributionally robust applications, including two-stage distributionally robust optimization using Wasserstein ambiguity sets and expectations of piecewise convex objective functions in single-stage problems. Independently, Bertsimas et al. (2019) investigate a *multi-policy approximation* of two-stage sample robust optimization by optimizing a separate linear decision rule for each uncertainty set and prove that this approximation gap converges to zero as the amount of data goes to infinity. In Section 6 of this paper, we show how to extend similar techniques to dynamic problems with many stages for the first time.

2. Problem Setting

As described in the introduction, we consider finite-horizon discrete-time stochastic decision problems. The uncertain quantities observed in each stage are denoted by random variables $\xi_1 \in \Xi_1 \subseteq \mathbb{R}^{d_\xi^1}, \dots, \xi_T \in \Xi_T \subseteq \mathbb{R}^{d_\xi^T}$, and the decisions made in each stage are denoted by $\mathbf{x}_1 \in \mathcal{X}_1 \subseteq \mathbb{R}^{d_x^1}, \dots, \mathbf{x}_T \in \mathcal{X}_T \subseteq \mathbb{R}^{d_x^T}$. Given realizations of the uncertain quantities and decisions, we incur a cost of

$$c(\xi_1, \dots, \xi_T, \mathbf{x}_1, \dots, \mathbf{x}_T) \in \mathbb{R}.$$

Let a decision rule $\pi = (\pi_1, \dots, \pi_T)$ denote a collection of measurable functions $\pi_t : \Xi_1 \times \dots \times \Xi_{t-1} \rightarrow \mathcal{X}_t$ which specify what decision to make in stage t based of the information observed up to that point. For notational convenience, let Π denote the space of all measurable non-anticipative decision rules. Given realizations of the uncertain quantities and choice of decision rules, the resulting cost is

$$c^\pi(\xi_1, \dots, \xi_T) \triangleq c(\xi_1, \dots, \xi_T, \pi_1, \pi_2(\xi_1), \dots, \pi_T(\xi_1, \dots, \xi_{T-1})).$$

Before selecting the decision rules, we observe auxiliary side information $\gamma \in \Gamma \subseteq \mathbb{R}^{d_\gamma}$. For example, in the aforementioned fashion setting, the side information may contain information on the brand, style, and color of a new clothing item and the remaining uncertainties representing the demand for the product in each week of the lifecycle.

Given a realization of the side information $\gamma = \bar{\gamma}$, our goal is to find decision rules which minimize the conditional expected cost:

$$v^*(\bar{\gamma}) \triangleq \underset{\pi \in \Pi}{\text{minimize}} \quad \mathbb{E} \left[c^\pi(\xi_1, \dots, \xi_T) \mid \gamma = \bar{\gamma} \right]. \quad (1)$$

We refer to (1) as *dynamic optimization with side information*. The optimization takes place over a collection Π which is any subset of the space of all non-anticipative decision rules. In this paper, we assume that the joint distribution of the side information and uncertain quantities $(\gamma, \xi_1, \dots, \xi_T)$ is unknown, and our knowledge consists of historical data of the form

$$(\gamma^1, \xi_1^1, \dots, \xi_T^1), \dots, (\gamma^N, \xi_1^N, \dots, \xi_T^N),$$

where each of these tuples consists of a realization of the side information and the following realization of the random variables over the stages. For example, in the aforementioned fashion setting, each tuple corresponds to the side information of a past fashion item as well as its demand over its lifecycle. We will not assume any parametric structure on the relationship between the side information and future uncertainty.

The goal of this paper is a general-purpose, computationally tractable, data-driven approach for approximately solving dynamic optimization with side information. In the following sections, we propose and analyze a new framework which leverages nonparametric machine learning, trained from historical data, to predict future uncertainty from side information in a way that leads to near-optimal decision rules to (1).

2.1. Notation

The joint probability distribution of the side information γ and uncertain quantities $\xi = (\xi_1, \dots, \xi_T)$ is denoted by \mathbb{P} . For the purpose of proving theorems, we assume throughout this paper that the historical data are independent and identically distributed (i.i.d.) samples from this distribution \mathbb{P} . In other words, we assume that the historical data satisfies

$$((\gamma^1, \xi^1), \dots, (\gamma^N, \xi^N)) \sim \mathbb{P}^N,$$

where $\mathbb{P}^N \triangleq \mathbb{P} \times \dots \times \mathbb{P}$ is the product measure. The set of all probability distributions supported on $\Xi \triangleq \Xi_1 \times \dots \times \Xi_T \subseteq \mathbb{R}^{d_\xi}$ is denoted by $\mathcal{P}(\Xi)$. For each of the side information $\bar{\gamma} \in \Gamma$, we assume that its conditional probability distribution satisfies $\mathbb{P}_{\bar{\gamma}} \in \mathcal{P}(\Xi)$, where $\mathbb{P}_{\bar{\gamma}}(\cdot)$ is shorthand for $\mathbb{P}(\cdot | \gamma = \bar{\gamma})$. We use “i.o.” as shorthand for “infinitely often”. We sometimes use subscript notation for expectations to specify the underlying probability distribution; for example, the following two expressions are equivalent:

$$\mathbb{E}_{\xi \sim \mathbb{P}_{\bar{\gamma}}} [f(\xi_1, \dots, \xi_T)] \equiv \mathbb{E} [f(\xi_1, \dots, \xi_T) | \gamma = \bar{\gamma}].$$

Finally, we say that the cost function resulting from a policy π is upper semicontinuous if

$$\limsup_{\zeta \rightarrow \bar{\zeta}} c^\pi(\zeta_1, \dots, \zeta_T) \leq c^\pi(\bar{\zeta}_1, \dots, \bar{\zeta}_T)$$

for all $\bar{\zeta} \in \Xi$.

3. Sample Robust Optimization with Side Information

In this section, we present our approach for incorporating machine learning in dynamic optimization. We first review sample robust optimization, and then we introduce the proposed *sample robust optimization with side information* approach to (1).

3.1. Preliminary: sample robust optimization

Consider a stochastic dynamic optimization problem of the form (1) in which there is no side information. The underlying joint distribution of the random variables $\xi \equiv (\xi_1, \dots, \xi_T)$ is unknown, but we have data consisting of sample paths, $\xi^1 \equiv (\xi_1^1, \dots, \xi_T^1), \dots, \xi^N \equiv (\xi_1^N, \dots, \xi_T^N)$. For this setting, sample robust optimization can be used to find approximate solutions in stochastic dynamic optimization. To apply the framework, one constructs an uncertainty set around each sample path in the training data and then chooses the decision rules that optimize the average of the worst-case realizations of the cost. Formally, this framework results in the following robust optimization problem:

$$\underset{\pi \in \Pi}{\text{minimize}} \quad \sum_{i=1}^N \frac{1}{N} \sup_{\zeta \in \mathcal{U}_N^i} c^\pi(\zeta_1, \dots, \zeta_T), \quad (3)$$

where $\mathcal{U}_N^i \subseteq \Xi$ is an uncertainty set around ξ^i . Intuitively speaking, (3) chooses the decision rules by averaging over the historical sample paths which are adversarially perturbed. Under mild probabilistic assumptions on the underlying joint distribution and appropriately constructed uncertainty sets, Bertsimas et al. (2018a) show that sample robust optimization converges asymptotically to the underlying stochastic problem and that (3) is amenable to approximations similar to dynamic robust optimization.

3.2. Incorporating side information into sample robust optimization

We now present our new framework, based on sample robust optimization, for solving dynamic optimization with side information. In the proposed framework, we first train a machine learning algorithm on the historical data to predict future uncertainty (ξ_1, \dots, ξ_T) as a function of the side information. From the trained learner, we obtain weight functions $w_N^i(\bar{\gamma})$, for $i = 1, \dots, N$, each of which captures the relevance of the i th training sample to the new side information, $\bar{\gamma}$. We incorporate the weights into sample robust optimization by multiplying the cost associated with each training example by the corresponding weight function. The resulting *sample robust optimization with side information* framework is as follows:

$$\hat{v}^N(\bar{\gamma}) \triangleq \underset{\pi \in \Pi}{\text{minimize}} \sum_{i=1}^N w_N^i(\bar{\gamma}) \sup_{\zeta \in \mathcal{U}_N^i} c^\pi(\zeta_1, \dots, \zeta_T), \quad (4)$$

where the uncertainty sets are defined

$$\mathcal{U}_N^i \triangleq \{\zeta \in \Xi : \|\zeta - \xi^i\| \leq \epsilon_N\},$$

and $\|\cdot\|$ is some ℓ_p norm with $p \geq 1$.

The above framework provides the flexibility for the practitioner to construct weights from a variety of machine learning algorithms. We focus in this paper on weight functions which come from nonparametric machine learning methods. Examples of viable predictive models include k -nearest neighbors (kNN), kernel regression, classification and regression trees (CART), and random forests (RF). We describe these four classes of weight functions.

DEFINITION 1. The k -nearest neighbor weight functions are given by:

$$w_{N,k\text{NN}}^i(\bar{\gamma}) \triangleq \begin{cases} \frac{1}{k_N}, & \text{if } \gamma^i \text{ is a } k_N\text{-nearest neighbor of } \bar{\gamma}, \\ 0, & \text{otherwise.} \end{cases}$$

Formally, γ^i is a k_N -nearest neighbor of $\bar{\gamma}$ if $|\{j \in \{1, \dots, N\} \setminus i : \|\gamma^j - \bar{\gamma}\| < \|\gamma^i - \bar{\gamma}\|\}| < k_N$. For more technical details, we refer the reader to Biau and Devroye (2015).

DEFINITION 2. The kernel regression weight functions are given by:

$$w_{N,\text{KR}}^i(\bar{\gamma}) \triangleq \frac{K(\|\gamma^i - \bar{\gamma}\|/h_N)}{\sum_{j=1}^N K(\|\gamma^j - \bar{\gamma}\|/h_N)},$$

where $K(\cdot)$ is the kernel function and h_N is the bandwidth parameter. Examples of kernel functions include the Gaussian kernel, $K(u) = \frac{1}{\sqrt{2\pi}}e^{-u^2/2}$, the triangular kernel, $K(u) = (1-u)\mathbb{1}\{u \leq 1\}$, and the Epanechnikov kernel, $K(u) = \frac{3}{4}(1-u^2)\mathbb{1}\{u \leq 1\}$. For more information on kernel regression, see [Friedman et al. \(2001, Chapter 6\)](#).

The next two types of weight functions we present are based on classification and regression trees ([Breiman et al. 1984](#)) and random forests ([Breiman 2001](#)). We refer the reader to [Bertsimas and Kallus \(2020\)](#) for technical implementation details.

DEFINITION 3. The classification and regression tree weight functions are given by:

$$w_{N,\text{CART}}^i(\bar{\gamma}) \triangleq \begin{cases} \frac{1}{|l^N(\bar{\gamma})|}, & i \in l^N(\bar{\gamma}), \\ 0, & \text{otherwise,} \end{cases}$$

where $l^N(\bar{\gamma})$ is the set of indices i such that γ^i is contained in the same leaf of the tree as $\bar{\gamma}$.

DEFINITION 4. The random forest weight functions are given by:

$$w_{N,\text{RF}}^i(\bar{\gamma}) \triangleq \frac{1}{B} \sum_{b=1}^B w_{N,\text{CART}}^{i,b}(\bar{\gamma}),$$

where B is the number of trees in the ensemble, and $w_{N,\text{CART}}^{i,b}(\bar{\gamma})$ refers to the weight function of the b th tree in the ensemble.

All of the above weight functions come from nonparametric machine learning methods. They are highly effective as predictive methods because they can learn complex relationships between the side information and the response variable without requiring the practitioner to state an explicit parametric form. Similarly, as we prove in [Section 4](#), solutions to [\(4\)](#) with these weight functions are asymptotically optimal for [\(1\)](#) without any parametric restrictions on the relationship between γ and ξ . In other words, incorporating side information into sample robust optimization via [\(4\)](#) leads to better decisions asymptotically, even without specific knowledge of how the side information affects the uncertainty.

4. Asymptotic Optimality

In this section, we establish asymptotic optimality guarantees for sample robust optimization with side information. We prove that, under mild conditions, [\(4\)](#) converges to [\(1\)](#) as the number of training samples goes to infinity. Thus, as the amount of data grows, sample robust optimization

with side information becomes an optimal approximation of the underlying stochastic dynamic optimization problem. Crucially, our convergence guarantee does not require parametric restrictions on the space of decision rules (*e.g.*, linearity) or parametric restrictions on the joint distribution of the side information and uncertain quantities.

4.1. Main result

We begin by presenting our main result. The proof of the result depends on some technical assumptions and concepts from distributionally robust optimization. For simplicity, we defer the statement and discussion of technical assumptions regarding the underlying probability distribution and cost until Sections 4.3 and 4.4, and first discuss what is needed to apply the method in practice. The practitioner needs to select a weight function, parameters associated with that weight function, and the radius, ϵ_N , of the uncertainty sets. While these may be selected by cross validation, we show that the method will in general converge if the parameters are selected to satisfy the following:

ASSUMPTION 1. *The weight functions and uncertainty set radius satisfy one of the following:*

1. $\{w_N^i(\cdot)\}$ are k -nearest neighbor weight functions with $k_N = \min(\lceil k_3 N^\delta \rceil, N - 1)$ for constants $k_3 > 0$ and $\delta \in (\frac{1}{2}, 1)$, and $\epsilon_N = \frac{k_1}{N^p}$ for constants $k_1 > 0$ and $0 < p < \min\left(\frac{1-\delta}{d_\gamma}, \frac{2\delta-1}{d_\xi+2}\right)$.
2. $\{w_N^i(\cdot)\}$ are kernel regression weight functions with the Gaussian, triangular, or Epanechnikov kernel function and $h_N = k_4 N^{-\delta}$ for constants $k_4 > 0$ and $\delta \in \left(0, \frac{1}{2d_\gamma}\right)$, and $\epsilon_N = \frac{k_1}{N^p}$ for constants $k_1 > 0$ and $0 < p < \min\left(\delta, \frac{1-\delta d_\gamma}{2+d_\xi}\right)$.

Given Assumption 1, our main result is the following.

THEOREM 1. *Suppose the weight function and uncertainty sets satisfy Assumption 1, the joint probability distribution of (γ, ξ) satisfies Assumptions 2-4 from Section 4.3, and the cost function satisfies Assumptions 5-6 from Section 4.4. Then, for every $\bar{\gamma} \in \Gamma$,*

$$\lim_{N \rightarrow \infty} \hat{v}^N(\bar{\gamma}) = v^*(\bar{\gamma}), \quad \mathbb{P}^\infty\text{-almost surely.}$$

The theorem says that objective value of (4) will converge almost surely to the optimal value of the full-information problem, (1), as N goes to infinity. The assumptions of the theorem require that the joint distribution and the feasible decision rules are well behaved. We will discuss these technical assumptions in more detail in the following sections.

In order to prove the asymptotic optimality of sample robust optimization with side information, we view (4) through the more general lens of Wasserstein-based distributionally robust optimization. We first review some properties of the Wasserstein metric and then prove a key intermediary result, from which our main result follows.

4.2. Review of the Wasserstein metric

The Wasserstein metric provides a distance function between probability distributions. In particular, given two probability distributions $\mathbb{Q}, \mathbb{Q}' \in \mathcal{P}(\Xi)$, the type-1 Wasserstein distance is defined as the optimal objective value of a minimization problem:

$$d_1(\mathbb{Q}, \mathbb{Q}') \triangleq \inf \left\{ \mathbb{E}_{(\xi, \xi') \sim \Pi} \|\xi - \xi'\| : \begin{array}{l} \Pi \text{ is a joint distribution of } \xi \text{ and } \xi' \\ \text{with marginals } \mathbb{Q} \text{ and } \mathbb{Q}', \text{ respectively} \end{array} \right\}.$$

The Wasserstein metric is particularly appealing because a distribution with finite support can have a finite distance to a continuous distribution. This allows us to construct a Wasserstein ball around an empirical distribution that includes continuous distributions, which cannot be done with other popular measures such as the Kullback-Leibler divergence (Kullback and Leibler 1951). We remark that the type-1 Wasserstein metric satisfies the axioms of a metric, including the triangle inequality (Clement and Desch 2008):

$$d_1(\mathbb{Q}_1, \mathbb{Q}_2) \leq d_1(\mathbb{Q}_1, \mathbb{Q}_3) + d_1(\mathbb{Q}_3, \mathbb{Q}_2), \quad \forall \mathbb{Q}_1, \mathbb{Q}_2, \mathbb{Q}_3 \in \mathcal{P}(\Xi).$$

Important to this paper, the type-1 Wasserstein metric admits a dual form, as shown by Kantorovich and Rubinstein (1958),

$$d_1(\mathbb{Q}, \mathbb{Q}') = \sup_{\text{Lip}(h) \leq 1} |\mathbb{E}_{\xi \sim \mathbb{Q}}[h(\xi)] - \mathbb{E}_{\xi \sim \mathbb{Q}'}[h(\xi)]|,$$

where the supremum is taken over all 1-Lipschitz functions. Note that the absolute value is optional in the dual form of the metric, and the space of Lipschitz functions can be restricted to those which satisfy $h(0) = 0$ without loss of generality. Finally, we remark that Fournier and Guillin (2015) prove under a light-tailed assumption that the 1-Wasserstein distance between the empirical distribution and its underlying distribution concentrates around zero with high probability. Theorem 2 in the following section extends this concentration result to the setting with side information.

4.3. Concentration of the empirical conditional probability distribution

Given a local predictive method, let the corresponding empirical conditional measure be defined as

$$\hat{\mathbb{P}}_{\bar{\gamma}}^N := \sum_{i=1}^N w_N^i(\bar{\gamma}) \delta_{\xi^i},$$

where δ_{ξ} denotes the Dirac probability distribution which places point mass at ξ . In this section, we prove under mild assumptions that the empirical conditional measure $\hat{\mathbb{P}}_{\bar{\gamma}}^N$ concentrates quickly to $\mathbb{P}_{\bar{\gamma}}$ with respect to the 1-Wasserstein metric. We introduce the following assumptions on the underlying joint probability distribution:

ASSUMPTION 2 (**Conditional Subgaussianity**). *There exists a parameter $\sigma > 0$ such that*

$$\mathbb{P}(\|\boldsymbol{\xi}\| - \mathbb{E}[\|\boldsymbol{\xi}\| \mid \gamma = \bar{\gamma}] > t \mid \gamma = \bar{\gamma}) \leq \exp\left(-\frac{t^2}{2\sigma^2}\right) \quad \forall t > 0, \bar{\gamma} \in \Gamma.$$

ASSUMPTION 3 (**Lipschitz Continuity**). *There exists $0 < L < \infty$ such that*

$$d_1(\mathbb{P}_{\bar{\gamma}}, \mathbb{P}_{\bar{\gamma}'}) \leq L\|\bar{\gamma} - \bar{\gamma}'\|, \quad \forall \bar{\gamma}, \bar{\gamma}' \in \Gamma.$$

ASSUMPTION 4 (**Smoothness of Side Information**). *The set Γ is compact, and there exists $g > 0$ such that*

$$\mathbb{P}(\|\gamma - \bar{\gamma}\| \leq \epsilon) \geq g\epsilon^{d_\gamma}, \quad \forall \epsilon > 0, \bar{\gamma} \in \Gamma.$$

Let us reflect on the conditions on the underlying joint distribution. Assumption 2 requires that the distribution of the uncertainty is not heavy-tailed, conditional on the side information. This is satisfied, for example, if $\boldsymbol{\xi}$ has bounded support or follows a Gaussian distribution, conditional on $\bar{\gamma}$. Assumption 3 requires that the conditional distribution of $\boldsymbol{\xi}$ is a smooth function of γ . This ensures we can actually learn about the conditional distribution $\mathbb{P}_{\bar{\gamma}}$ from historical data with side information that are similar (but not identical) to $\bar{\gamma}$. Assumption 4 ensures the side information are distributed in such a way that every possible $\bar{\gamma} \in \Gamma$ has nearby observations in the historical data, as $N \rightarrow \infty$.

With these assumptions, we are ready to prove the concentration result, which is proved using a novel technique that relies on the dual form of the Wasserstein metric and a discrete approximation of the space of 1-Lipschitz functions.

THEOREM 2. *Suppose the weight function and uncertainty sets satisfy Assumption 1 and the joint probability distribution of $(\gamma, \boldsymbol{\xi})$ satisfies Assumptions 2-4. Then, for every $\bar{\gamma} \in \Gamma$,*

$$\mathbb{P}^\infty \left(\left\{ d_1(\mathbb{P}_{\bar{\gamma}}, \hat{\mathbb{P}}_{\bar{\gamma}}^N) > \epsilon_N \right\} \text{ i.o. } \right) = 0.$$

Proof. Without loss of generality, we assume throughout the proof that all norms $\|\cdot\|$ refer to the ℓ_∞ norm.² Fix any $\bar{\gamma} \in \Gamma$. It follows from Assumption 1 that

$$\{w_N^i(\bar{\gamma})\} \text{ are not functions of } \boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^N; \tag{5}$$

$$\sum_{i=1}^N w_N^i(\bar{\gamma}) = 1 \text{ and } w_N^1(\bar{\gamma}), \dots, w_N^N(\bar{\gamma}) \geq 0, \quad \forall N \in \mathbb{N}; \tag{6}$$

$$\epsilon_N = \frac{k_1}{N^p}, \quad \forall N \in \mathbb{N}, \tag{7}$$

for constants $k_1, p > 0$. Moreover, Assumption 1 also implies that there exist constants $k_2 > 0$ and $\eta > p(2 + d_\xi)$ such that

$$\lim_{N \rightarrow \infty} \frac{1}{\epsilon_N} \sum_{i=1}^N w_N^i(\bar{\gamma}) \|\gamma^i - \bar{\gamma}\| = 0, \quad \mathbb{P}^\infty\text{-almost surely}; \quad (8)$$

$$\mathbb{E}_{\mathbb{P}^N} \left[\exp \left(\frac{-\theta}{\sum_{i=1}^N w_N^i(\bar{\gamma})^2} \right) \right] \leq \exp(-k_2 \theta N^\eta), \quad \forall \theta \in (0, 1), N \in \mathbb{N}. \quad (9)$$

The proof of the above statements under Assumption 1 is found in Appendix EC.1. Now, choose any fixed $q \in (0, \eta/(2 + d_\xi) - p)$, and let

$$b_N \triangleq N^q, \quad B_N \triangleq \{\zeta \in \mathbb{R}^{d_\xi} : \|\zeta\| \leq b_N\}, \quad I_N \triangleq \mathbf{1}\{\xi^1, \dots, \xi^N \in B_N\}.$$

Finally, we define the following intermediary probability distributions:

$$\hat{\mathbb{Q}}_{\bar{\gamma}}^N \triangleq \sum_{i=1}^N w_N^i(\bar{\gamma}) \mathbb{P}_{\gamma^i}, \quad \hat{\mathbb{Q}}_{\bar{\gamma}|B_N}^N \triangleq \sum_{i=1}^N w_N^i(\bar{\gamma}) \mathbb{P}_{\gamma^i|B_N},$$

where $\mathbb{P}_{\gamma^i|B_N}(\cdot)$ is shorthand for $\mathbb{P}(\cdot | \gamma = \gamma^i, \xi \in B_N)$.

Applying the triangle inequality for the 1-Wasserstein metric and the union bound,

$$\begin{aligned} \mathbb{P}^\infty \left(\{d_1(\mathbb{P}_{\bar{\gamma}}, \hat{\mathbb{P}}_{\bar{\gamma}}^N) > \epsilon_N\} \text{ i.o.} \right) &\leq \mathbb{P}^\infty \left(\left\{ d_1(\mathbb{P}_{\bar{\gamma}}, \hat{\mathbb{Q}}_{\bar{\gamma}}^N) > \frac{\epsilon_N}{3} \right\} \text{ i.o.} \right) \\ &\quad + \mathbb{P}^\infty \left(\left\{ d_1(\hat{\mathbb{Q}}_{\bar{\gamma}}^N, \hat{\mathbb{Q}}_{\bar{\gamma}|B_N}^N) > \frac{\epsilon_N}{3} \right\} \text{ i.o.} \right) \\ &\quad + \mathbb{P}^\infty \left(\left\{ d_1(\hat{\mathbb{Q}}_{\bar{\gamma}|B_N}^N, \hat{\mathbb{P}}_{\bar{\gamma}}^N) > \frac{\epsilon_N}{3} \right\} \text{ i.o.} \right). \end{aligned}$$

We now proceed to bound each of the above terms.

Term 1: $d_1(\mathbb{P}_{\bar{\gamma}}, \hat{\mathbb{Q}}_{\bar{\gamma}}^N)$: By the dual form of the 1-Wasserstein metric,

$$d_1(\mathbb{P}_{\bar{\gamma}}, \hat{\mathbb{Q}}_{\bar{\gamma}}^N) = \sup_{\text{Lip}(h) \leq 1} \left| \mathbb{E}[h(\xi) | \gamma = \bar{\gamma}] - \sum_{i=1}^N w_N^i(\bar{\gamma}) \mathbb{E}[h(\xi) | \gamma = \gamma^i] \right|,$$

where the supremum is taken over all 1-Lipschitz functions. By (6) and Jensen's inequality, we can upper bound this by

$$\begin{aligned} d_1(\mathbb{P}_{\bar{\gamma}}, \hat{\mathbb{Q}}_{\bar{\gamma}}^N) &\leq \sum_{i=1}^N w_N^i(\bar{\gamma}) \left(\sup_{\text{Lip}(h) \leq 1} |\mathbb{E}[h(\xi) | \gamma = \bar{\gamma}] - \mathbb{E}[h(\xi) | \gamma = \gamma^i]| \right) \\ &= \sum_{i=1}^N w_N^i(\bar{\gamma}) d_1(\mathbb{P}_{\bar{\gamma}}, \mathbb{P}_{\gamma^i}) \\ &\leq L \sum_{i=1}^N w_N^i(\bar{\gamma}) \|\bar{\gamma} - \gamma^i\|, \end{aligned}$$

where the final inequality follows from Assumption 3. Therefore, it follows from (8) that

$$\mathbb{P}^\infty \left(\left\{ d_1(\mathbb{P}_{\bar{\gamma}}, \hat{\mathbb{Q}}_{\bar{\gamma}}^N) > \frac{\epsilon_N}{3} \right\} \text{ i.o.} \right) = 0. \quad (10)$$

Term 2: $d_1(\hat{\mathbb{Q}}_{\bar{\gamma}}^N, \hat{\mathbb{Q}}_{\bar{\gamma}|B_N}^N)$: Consider any Lipschitz function $\text{Lip}(h) \leq 1$ for which $h(0) = 0$, and let $\bar{N} \in \mathbb{N}$ satisfy $b_N \geq \sigma + \sup_{\bar{\gamma} \in \Gamma} \mathbb{E}[\|\boldsymbol{\xi}\| | \gamma = \bar{\gamma}]$ (which is finite because of Assumption 4). Then, for all $N \geq \bar{N}$, and all $\bar{\gamma}' \in \Gamma$,

$$\begin{aligned}
& \mathbb{E}[h(\boldsymbol{\xi}) | \gamma = \bar{\gamma}'] - \mathbb{E}[h(\boldsymbol{\xi}) | \gamma = \bar{\gamma}', \boldsymbol{\xi} \in B_N] \\
&= \mathbb{E}[h(\boldsymbol{\xi}) \mathbb{1}\{\boldsymbol{\xi} \notin B_N\} | \gamma = \bar{\gamma}'] + \mathbb{E}[h(\boldsymbol{\xi}) \mathbb{1}\{\boldsymbol{\xi} \in B_N\} | \gamma = \bar{\gamma}'] - \mathbb{E}[h(\boldsymbol{\xi}) | \gamma = \bar{\gamma}', \boldsymbol{\xi} \in B_N] \\
&= \mathbb{E}[h(\boldsymbol{\xi}) \mathbb{1}\{\boldsymbol{\xi} \notin B_N\} | \gamma = \bar{\gamma}'] + \mathbb{E}[h(\boldsymbol{\xi}) | \gamma = \bar{\gamma}', \boldsymbol{\xi} \in B_N] \mathbb{P}(\boldsymbol{\xi} \in B_N | \gamma = \bar{\gamma}') - \mathbb{E}[h(\boldsymbol{\xi}) | \gamma = \bar{\gamma}', \boldsymbol{\xi} \in B_N] \\
&= \mathbb{E}[h(\boldsymbol{\xi}) \mathbb{1}\{\boldsymbol{\xi} \notin B_N\} | \gamma = \bar{\gamma}'] - \mathbb{E}[h(\boldsymbol{\xi}) | \gamma = \bar{\gamma}', \boldsymbol{\xi} \in B_N] \mathbb{P}(\boldsymbol{\xi} \notin B_N | \gamma = \bar{\gamma}') \\
&\leq \mathbb{E}[\|\boldsymbol{\xi}\| \mathbb{1}\{\boldsymbol{\xi} \notin B_N\} | \gamma = \bar{\gamma}'] + b_N \mathbb{P}(\boldsymbol{\xi} \notin B_N | \gamma = \bar{\gamma}') \\
&= \int_{b_N}^{\infty} \mathbb{P}(\|\boldsymbol{\xi}\| > t | \gamma = \bar{\gamma}') dt + b_N \mathbb{P}(\|\boldsymbol{\xi}\| \geq b_N | \gamma = \bar{\gamma}') \\
&\leq (\sigma + b_N) \exp\left(-\frac{1}{2\sigma^2} \left(b_N - \sup_{\bar{\gamma}' \in \Gamma} \mathbb{E}[\|\boldsymbol{\xi}\| | \gamma = \bar{\gamma}']\right)^2\right).
\end{aligned}$$

The first inequality follows because $|h(\boldsymbol{\xi})| \leq b_N$ for all $\boldsymbol{\xi} \in B_N$ and $|h(\boldsymbol{\xi})| \leq \|\boldsymbol{\xi}\|$ otherwise. For the second inequality, we used the Gaussian tail inequality $\int_x^{\infty} e^{-t^2/2} dt \leq e^{-x^2/2}$ for $x \geq 1$ (Vershynin 2018) along with Assumption 2. Because this bound holds uniformly over all h , and all $\bar{\gamma}' \in \Gamma$, it follows that

$$\begin{aligned}
d_1(\hat{\mathbb{Q}}_{\bar{\gamma}}^N, \hat{\mathbb{Q}}_{\bar{\gamma}|B_N}^N) &= \sup_{\text{Lip}(h) \leq 1, h(0)=0} \left| \sum_{i=1}^N w_N^i(\bar{\gamma}) (\mathbb{E}[h(\boldsymbol{\xi}) | \gamma = \gamma^i] - \mathbb{E}[h(\boldsymbol{\xi}) | \gamma = \gamma^i, \boldsymbol{\xi} \in B_N]) \right| \\
&\leq \sum_{i=1}^N w_N^i(\bar{\gamma}) \sup_{\text{Lip}(h) \leq 1, h(0)=0} |\mathbb{E}[h(\boldsymbol{\xi}) | \gamma = \gamma^i] - \mathbb{E}[h(\boldsymbol{\xi}) | \gamma = \gamma^i, \boldsymbol{\xi} \in B_N]| \\
&\leq \sup_{\bar{\gamma}' \in \Gamma} \sup_{\text{Lip}(h) \leq 1, h(0)=0} |\mathbb{E}[h(\boldsymbol{\xi}) | \gamma = \bar{\gamma}'] - \mathbb{E}[h(\boldsymbol{\xi}) | \gamma = \bar{\gamma}', \boldsymbol{\xi} \in B_N]| \\
&\leq (\sigma + b_N) \exp\left(-\frac{1}{2\sigma^2} \left(b_N - \sup_{\bar{\gamma}' \in \Gamma} \mathbb{E}[\|\boldsymbol{\xi}\| | \gamma = \bar{\gamma}']\right)^2\right),
\end{aligned}$$

for all $N \geq \bar{N}$. It is easy to see that the right hand side above divided by $\epsilon_N/3$ goes to 0 as N goes to infinity, so

$$\mathbb{P}^{\infty} \left(\left\{ d_1(\hat{\mathbb{Q}}_{\bar{\gamma}}^N, \hat{\mathbb{Q}}_{\bar{\gamma}|B_N}^N) > \frac{\epsilon_N}{3} \right\} \text{ i.o.} \right) = 0.$$

Term 3: $d_1(\hat{\mathbb{Q}}_{\bar{\gamma}|B_N}^N, \hat{\mathbb{P}}_{\bar{\gamma}}^N)$: By the law of total probability,

$$\mathbb{P}^N \left(d_1(\hat{\mathbb{Q}}_{\bar{\gamma}|B_N}^N, \hat{\mathbb{P}}_{\bar{\gamma}}^N) > \frac{\epsilon_N}{3} \right) \leq \mathbb{P}^N(I_N = 0) + \mathbb{P}^N \left(d_1(\hat{\mathbb{Q}}_{\bar{\gamma}|B_N}^N, \hat{\mathbb{P}}_{\bar{\gamma}}^N) > \frac{\epsilon_N}{3} \mid I_N = 1 \right).$$

We now show that each of the above terms have finite summations. First,

$$\sum_{N=1}^{\infty} \mathbb{P}^N(I_N = 0) \leq \sum_{N=1}^{\infty} N \sup_{\bar{\gamma}' \in \Gamma} \mathbb{P}(\boldsymbol{\xi} \notin B_N | \gamma = \bar{\gamma}') \leq \sum_{N=1}^{\infty} N \sup_{\bar{\gamma}' \in \Gamma} \exp\left(-\frac{(b_N - \mathbb{E}[\|\boldsymbol{\xi}\| | \gamma = \bar{\gamma}'])^2}{2\sigma^2}\right) < \infty.$$

The first inequality follows from the union bound, the second inequality follows from Assumption 2, and the final inequality follows because $\sup_{\bar{\gamma}' \in \Gamma} \mathbb{E}[\|\xi\| | \gamma = \bar{\gamma}'] < \infty$ and the definition of b_N .

Second, for each $l \in \mathbb{N}$, we define several quantities. Let \mathcal{P}_l be the partitioning of $B_N = [-b_N, b_N]^{d_\xi}$ into 2^{ld_ξ} translations of $(-b_N 2^{-l}, b_N 2^{-l})^{d_\xi}$. Let \mathcal{H}_l be the set of piecewise constant functions which are constant on each region of the partition \mathcal{P}_l , taking values on $\{kb_N 2^{-l} : k \in \{0, \pm 1, \pm 2, \pm 3, \dots, \pm 2^l\}\}$. Note that $|\mathcal{H}_l| = (2^{l+1} + 1)^{2^{ld_\xi}}$. Then, we observe that for all Lipschitz functions $\text{Lip}(h) \leq 1$ which satisfy $h(0) = 0$, there exists a $\hat{h} \in \mathcal{H}_l$ such that

$$\sup_{\zeta \in B_N} |h(\zeta) - \hat{h}(\zeta)| \leq b_N 2^{-l+1}.$$

Indeed, within each region of the partition, h can vary by no more than $b_N 2^{-l+1}$. The possible function values for \hat{h} are separated by $b_N 2^{-l}$. Because h is bounded by $\pm b_N$, this implies the existence of $\hat{h} \in \mathcal{H}_l$ such that \hat{h} has a value within $b_N 2^{-l+1}$ of h everywhere within that region. The identical reasoning holds for all other regions of the partition.

Therefore, for every $l \in \mathbb{N}$,

$$\begin{aligned} & \mathbb{P}^N \left(\mathbf{d}_1(\hat{\mathbb{Q}}_{\bar{\gamma}|B_N}^N, \hat{\mathbb{P}}_{\bar{\gamma}}^N) > \frac{\epsilon_N}{3} \mid I_N = 1 \right) \\ &= \mathbb{P}^N \left(\sup_{\substack{\text{Lip}(h) \leq 1 \\ h(0) = 0}} \sum_{i=1}^N w_N^i(\bar{\gamma}) (h(\xi^i) - \mathbb{E}[h(\xi) \mid \gamma = \gamma^i, \xi \in B_N]) > \frac{\epsilon_N}{3} \mid I_N = 1 \right) \\ &\leq \mathbb{P}^N \left(\sup_{\hat{h} \in \mathcal{H}_l} \sum_{i=1}^N w_N^i(\bar{\gamma}) (\hat{h}(\xi^i) - \mathbb{E}[\hat{h}(\xi) \mid \gamma = \gamma^i, \xi \in B_N]) > \frac{\epsilon_N}{3} - 2 \cdot b_N 2^{-l+1} \mid I_N = 1 \right) \\ &\leq |\mathcal{H}_l| \sup_{\hat{h} \in \mathcal{H}_l} \mathbb{P}^N \left(\sum_{i=1}^N w_N^i(\bar{\gamma}) (\hat{h}(\xi^i) - \mathbb{E}[\hat{h}(\xi) \mid \gamma = \gamma^i, \xi \in B_N]) > \frac{\epsilon_N}{3} - b_N 2^{-l+2} \mid I_N = 1 \right), \end{aligned}$$

where the final inequality follows from the union bound. We choose $l = \left\lceil 2 + \log_2 \frac{6b_N}{\epsilon_N} \right\rceil$, in which case

$$\frac{\epsilon_N}{3} - b_N 2^{-l+2} \geq \frac{\epsilon_N}{6}.$$

Furthermore, for all sufficiently large N ,

$$|\mathcal{H}_l| = (2^{l+1} + 1)^{2^{ld_\xi}} \leq \left(96 \frac{b_N}{\epsilon_N} \right)^{24^{d_\xi} (b_N/\epsilon_N)^{d_\xi}} = \exp \left(24^{d_\xi} \left(\frac{b_N}{\epsilon_N} \right)^{d_\xi} \log \frac{96b_N}{\epsilon_N} \right).$$

Applying Hoeffding's inequality, and noting $|\hat{h}(\xi^i)|$ is bounded by b_N when $\xi^i \in B_N$, we have the following for all $\hat{h} \in \mathcal{H}_l$:

$$\mathbb{P}^N \left(\sum_{i=1}^N w_N^i(\bar{\gamma}) (\hat{h}(\xi^i) - \mathbb{E}[\hat{h}(\xi) \mid \xi \in B_N, \gamma = \gamma^i]) > \frac{\epsilon_N}{6} \mid I_N = 1 \right)$$

$$\begin{aligned}
&= \mathbb{E} \left[\mathbb{P}^N \left(\sum_{i=1}^N w_N^i(\bar{\gamma}) \left(\hat{h}(\xi^i) - \mathbb{E}[\hat{h}(\xi) | \xi \in B_N, \gamma = \gamma^i] \right) > \frac{\epsilon_N}{6} \middle| I_N = 1, \gamma^1, \dots, \gamma^N \right) \middle| I_N = 1 \right] \\
&\leq \mathbb{E} \left[\exp \left(-\frac{\epsilon_N^2}{72 \sum_{i=1}^N (w_N^i(\bar{\gamma}))^2 b_N^2} \right) \middle| I_N = 1 \right] \\
&= \mathbb{E} \left[\exp \left(-\frac{\epsilon_N^2}{72 \sum_{i=1}^N (w_N^i(\bar{\gamma}))^2 b_N^2} \right) I_N \right] \left(\frac{1}{\mathbb{P}^N(I_N = 1)} \right) \\
&\leq 2 \mathbb{E} \left[\exp \left(-\frac{\epsilon_N^2}{72 \sum_{i=1}^N (w_N^i(\bar{\gamma}))^2 b_N^2} \right) \right] \\
&\leq 2 \exp \left(-\frac{k_2 N^\eta \epsilon_N^2}{72 b_N^2} \right),
\end{aligned}$$

for N sufficiently large that $\mathbb{P}(I_N = 1) \geq 1/2$ and $\epsilon_N^2/72b_N^2 < 1$. Note that (9) was used for the final inequality. Combining these results, we have

$$\mathbb{P}^N \left(d_1(\hat{\mathbb{P}}_{\bar{\gamma}}^N, \hat{\mathbb{Q}}_{\bar{\gamma}|B_N}^N) > \epsilon_N/3 \middle| I_N = 1 \right) \leq 2 \exp \left(24^{d_\xi} \left(\frac{b_N}{\epsilon_N} \right)^{d_\xi} \log \frac{96b_N}{\epsilon_N} - \frac{k_2 \epsilon_N^2 N^\eta}{72 N b_N^2} \right),$$

for N sufficiently large. For some constants $c_1, c_2 > 0$, and sufficiently large N , this is upper bounded by

$$2 \exp \left(-c_1 N^{\eta-2(p+q)} + c_2 N^{d_\xi(q+p)} \log N \right).$$

Since $0 < d_\xi(p+q) < \eta - 2(p+q)$, we can conduct a limit comparison test with $1/N^2$ to see that this term has a finite sum over N , which completes the proof. \square

4.4. Proof of main result

Theorem 2 provides the key ingredient for the proof of the main consistency result. We state our final two assumptions on the dynamic optimization problem to establish our main result.

ASSUMPTION 5 (Regularity of robust problem). *For all $\bar{\gamma} \in \Gamma$, there exists $M \geq 0$ such that the objective value of (4) would not change if we restricted its optimization to the decision rules $\pi \in \Pi$ which satisfy*

$$|c^\pi(\zeta_1, \dots, \zeta_T)| \leq M \left(1 + \max \left\{ \|\zeta\|, \sup_{\zeta' \in \cup_{i=1}^N \mathcal{U}_N^i} \|\zeta'\| \right\} \right), \quad \forall \zeta \in \Xi.$$

ASSUMPTION 6 (Regularity of stochastic problem). *For all $\bar{\gamma} \in \Gamma$, there exists $M' \geq 0$ such that the objective value of (1) would not change if we restricted the optimization to the decision rules $\pi \in \Pi$ for which $c^\pi(\zeta_1, \dots, \zeta_T)$ is upper semicontinuous and $|c^\pi(\zeta_1, \dots, \zeta_T)| \leq M'(1 + \|\zeta\|)$ for all $\zeta \in \Xi$.*

Assumption 5 is a minor modification to Bertsimas et al. (2018a, Assumption 3), and can be verified by decision makers through performing a static analysis; see Bertsimas et al. (2018a, Appendix A). Assumption 6 is a condition on structure of optimal decision rules of the stochastic problem, which is nearly identical to the assumptions of Mohajerin Esfahani and Kuhn (2018, Theorem 3.6(i)) which are used to establish asymptotic optimality for distributionally robust optimization with the type-1 Wasserstein ambiguity set.

Under these assumptions, the proof of Theorem 1 follows from Theorem 2 via arguments similar to those used by Mohajerin Esfahani and Kuhn (2018) and Bertsimas et al. (2018a). We state the proof fully in Appendix EC.2.

5. Implications for Single-Period Distributionally Robust Optimization

Beyond its utility in the context of multi-period problems, the measure concentration result of the previous section (Theorem 2) has potentially valuable implications for distributionally robust optimization with the type-1 Wasserstein ambiguity set. Indeed, consider a single-period optimization problem of the form

$$\underset{\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^{d_x}}{\text{minimize}} \quad \mathbb{E}_{\mathbb{P}} [c(\mathbf{x}, \boldsymbol{\xi})], \quad (11)$$

where $\boldsymbol{\xi} \in \Xi \subseteq \mathbb{R}^{d_\xi}$ is a random vector with a probability distribution \mathbb{P} . When the distribution is unknown and observable only through limited historical data $(\boldsymbol{\xi}^1, \dots, \boldsymbol{\xi}^N) \sim \mathbb{P}^N$, there has been recent interest in approximating the above problems by distributionally robust optimization with the type-1 Wasserstein ambiguity set:

$$\underset{\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^{d_x}}{\text{minimize}} \quad \sup_{\mathbb{Q} \in \mathcal{P}(\Xi): d_1(\mathbb{Q}, \hat{\mathbb{P}}^N) \leq \epsilon_N} \mathbb{E}_{\mathbb{Q}} [c(\mathbf{x}, \boldsymbol{\xi})]. \quad (12)$$

Due to several attractive properties, (12) and its relatives have received considerable recent interest in a variety of single-period operational and statistical applications. Indeed, when the robustness parameter is chosen appropriately and other mild assumptions hold, (12) is guaranteed to be asymptotically optimal (Mohajerin Esfahani and Kuhn 2018, Theorem 3.6) and the worst-case cost can often be reformulated as a tractable optimization problem (Mohajerin Esfahani and Kuhn 2018, Blanchet and Murthy 2019, Gao and Kleywegt 2016). Moreover, there is growing empirical evidence that (12) with a positive choice of the robustness parameter ($\epsilon_N > 0$) can find solutions with significantly better average out-of-sample cost compared to those obtained by the sample average approximation ($\epsilon_N = 0$), particularly when the number of data points is small; see, for example, Mohajerin Esfahani and Kuhn (2018, Section 7.2) and Hanasusanto and Kuhn (2018,

Section 4.2). Theoretical results which aim to explain this improved average out-of-sample cost, both for (12) as well as related robust approaches, are found in [Gotoh et al. \(2018\)](#) and [Anderson and Philpott \(2019\)](#).

In the remainder of this section, using the results from Section 4.3, we now show how side information and machine learning can be easily incorporated into any problem of the form (12), without foregoing its asymptotic optimality or computational tractability. Indeed, consider a single-period optimization problem of the form

$$v^*(\bar{\gamma}) \triangleq \underset{\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^{d_x}}{\text{minimize}} \quad \mathbb{E}_{\mathbb{P}} [c(\mathbf{x}, \boldsymbol{\xi}) \mid \gamma = \bar{\gamma}], \quad (13)$$

where $\boldsymbol{\xi} \in \Xi \subseteq \mathbb{R}^{d_\xi}$ and $\gamma \in \Gamma \subseteq \mathbb{R}^{d_\gamma}$ are random vectors with a joint probability distribution \mathbb{P} . Assume that the distribution is unknown and observable only through limited historical data $((\gamma^1, \boldsymbol{\xi}^1), \dots, (\gamma^N, \boldsymbol{\xi}^N)) \sim \mathbb{P}^N$. We address these problems by a modification of distributionally robust optimization with the type-1 Wasserstein ambiguity set, wherein the empirical probability distribution $\hat{\mathbb{P}}^N$ is replaced with an empirical conditional probability distribution $\hat{\mathbb{P}}_{\bar{\gamma}}^N$ (see Section 4.3):

$$v^N(\bar{\gamma}) \triangleq \underset{\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^{d_x}}{\text{minimize}} \quad \sup_{\mathbb{Q} \in \mathcal{P}(\Xi): d_1(\mathbb{Q}, \hat{\mathbb{P}}_{\bar{\gamma}}^N) \leq \epsilon_N} \mathbb{E}_{\mathbb{Q}} [c(\mathbf{x}, \boldsymbol{\xi})]. \quad (14)$$

As discussed previously, the empirical conditional probability distribution can be constructed using a variety of machine learning methods, such as k -nearest neighbor regression or kernel regression.

For this modification, we obtain the following asymptotic optimality guarantee which is analogous to (12) developed by [Mohajerin Esfahani and Kuhn \(2018\)](#).

THEOREM 3. *Suppose the weight function and uncertainty sets satisfy Assumption 1 and the joint probability distribution of $(\gamma, \boldsymbol{\xi})$ satisfies Assumptions 2-4. Assume that $\hat{\mathbf{x}}_N$ represents an optimizer of (14). Then the following hold for every $\bar{\gamma} \in \Gamma$:*

- (i) *If $c(\mathbf{x}, \boldsymbol{\xi})$ is upper semicontinuous in $\boldsymbol{\xi}$ and there exists $L \geq 0$ with $|c(\mathbf{x}, \boldsymbol{\xi})| \leq L(1 + \|\boldsymbol{\xi}\|)$ for all $\mathbf{x} \in \mathcal{X}$ and $\boldsymbol{\xi} \in \Xi$, then \mathbb{P}^∞ -almost surely we have $\hat{v}^N(\bar{\gamma}) \downarrow v^*(\bar{\gamma})$ as $N \rightarrow \infty$.*
- (ii) *If the assumptions of assertion (i) hold, \mathcal{X} is closed, and $c(\mathbf{x}, \boldsymbol{\xi})$ is lower semicontinuous in \mathbf{x} for every $\boldsymbol{\xi} \in \Xi$, then any accumulation point of $\{\hat{\mathbf{x}}_N\}_{N \in \mathbb{N}}$ is \mathbb{P}^∞ -almost surely an optimal solution for (13).*

Proof. The proof follows from identical reasoning as [Mohajerin Esfahani and Kuhn \(2018, Theorem 3.6\)](#), in which the measure concentration result of [Fournier and Guillin \(2015, Theorem 2\)](#) is replaced by Theorem 2 of the present paper. \square

From the perspective of computational tractability, it is readily observed that (14) retains an identical computational tractability as (12), except where terms of the form $\frac{1}{N}$ are replaced with $w_N^i(\bar{\gamma})$; see, for example, [Mohajerin Esfahani and Kuhn \(2018, Theorem 4.2\)](#). As a result of [Theorem 3](#), we conclude that side information can be tractably incorporated into the variety of operational applications that utilize (single-period) Wasserstein-based distributionally robust optimization.

6. Tractable Approximations

In the previous sections, we presented the new framework of sample robust optimization with side information and established its asymptotic optimality in the context of (1) without any significant structural restrictions on the space of decision rules. In this section, we focus on tractable methods for approximately solving the robust optimization problems that result from this proposed framework. Specifically, we develop a formulation which uses auxiliary decision rules to approximate the cost function. In combination with linear decision rules, this approach enables us to find high-quality decisions for real-world problems with more than ten stages in less than one minute, as we demonstrate in [Section 7](#).

We focus in this section on dynamic optimization problems with cost functions of the form

$$c(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T, \mathbf{x}_1, \dots, \mathbf{x}_T) = \sum_{t=1}^T \left(\mathbf{f}_t^\top \mathbf{x}_t + \mathbf{g}_t^\top \boldsymbol{\xi}_t + \min_{\mathbf{y}_t \in \mathbb{R}^{d_y^t}} \left\{ \mathbf{h}_t^\top \mathbf{y}_t : \sum_{s=1}^t \mathbf{A}_{t,s} \mathbf{x}_s + \sum_{s=1}^t \mathbf{B}_{t,s} \boldsymbol{\xi}_s + \mathbf{C}_t \mathbf{y}_t \leq \mathbf{d}_t \right\} \right). \quad (15)$$

Such cost functions appear frequently in applications such as inventory management and supply chain networks. Unfortunately, it is well known that these cost functions are convex in the uncertainty $\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T$. Thus, even evaluating the worst-case cost over a convex uncertainty set is computationally demanding in general, as it requires the maximization of a convex function.

As an intermediary step towards developing an approximation scheme for (4) with the above cost function, we consider the following optimization problem:

$$\begin{aligned} \tilde{v}^N(\bar{\gamma}) \triangleq & \quad \underset{\boldsymbol{\pi} \in \Pi, \mathbf{y}_t^i \in \mathcal{R}_t \forall i, t}{\text{minimize}} && \sum_{i=1}^N w_N^i(\bar{\gamma}) \sup_{\boldsymbol{\zeta} \in \mathcal{U}_N^i} \sum_{t=1}^T (\mathbf{f}_t^\top \boldsymbol{\pi}_t(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_{t-1}) + \mathbf{g}_t^\top \boldsymbol{\zeta}_t + \mathbf{h}_t^\top \mathbf{y}_t^i(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_t)) \\ & \text{subject to} && \sum_{s=1}^t \mathbf{A}_{t,s} \boldsymbol{\pi}_s(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_{s-1}) + \sum_{s=1}^t \mathbf{B}_{t,s} \boldsymbol{\zeta}_s + \mathbf{C}_t \mathbf{y}_t^i(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_t) \leq \mathbf{d}_t \\ & && \forall \boldsymbol{\zeta} \in \mathcal{U}_N^i, i \in \{1, \dots, N\}, t \in \{1, \dots, T\}, \end{aligned} \quad (16)$$

where \mathcal{R}_t is the set of all functions $\mathbf{y} : \Xi_1 \times \dots \times \Xi_t \rightarrow \mathbb{R}^{d_y^t}$. In this problem, we have introduced auxiliary decision rules which capture the minimization portion of (15) in each stage. We refer to (16) as a *multi-policy* approach, as it involves different auxiliary decision rules for each uncertainty set. The following theorem shows that (16) is equivalent to (4).

THEOREM 4. For cost functions of the form (15), $\tilde{v}^N(\bar{\gamma}) = \hat{v}^N(\bar{\gamma})$.

Proof. See Appendix EC.3. \square

We observe that (16) involves optimizing over decision rules, and thus is computationally challenging to solve in general. Nonetheless, we can obtain a tractable approximation of (16) by further restricting the space of primary and auxiliary decision rules. For instance, we can restrict all primary and auxiliary decision rules as linear decision rules of the form

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

One can alternatively elect to use a richer class of decision rules, such as lifted linear decision rules (Chen and Zhang 2009, Georghiou et al. 2015). In all cases, feasible approximations that restrict the space of decision rules of (16) provide an upper bound on the cost $\hat{v}^N(\bar{\gamma})$ and produce decision rules that are feasible for (16).

The key benefit of the multi-policy approximation scheme is that it offers many degrees of freedom in approximating the nonlinear cost function. Specifically, in (16), a separate auxiliary decision rule \mathbf{y}_t^i captures the value of the cost function for each uncertainty set in each stage. We approximate each \mathbf{y}_t^i with a linear decision rule, which only needs to be locally accurate, *i.e.*, accurate for realizations in the corresponding uncertainty set. As a result, (16) with linear decision rules results in significantly tighter approximations of (4) compared to using a single linear decision rule, \mathbf{y}_t , for all uncertainty sets in each stage. Moreover, these additional degrees of freedom come with only a mild increase in computation cost, and we substantiate these claims via computational experiments in Section 7.1. In Appendix EC.4, we provide the reformulation of the multi-policy approximation scheme with linear decision rules into a deterministic optimization problem using standard techniques from robust optimization.

7. Computational Experiments

We perform computational experiments to assess the out-of-sample performance and computational tractability of the proposed methodologies across several applications. These examples are dynamic inventory management (Section 7.1), portfolio optimization (Section 7.2), and shipment planning (Section 7.3).

We compare several methods using different machine learning models. These methods include the proposed sample robust optimization with side information, sample average approximation (SAA), the predictions to prescriptions (PtP) approach of Bertsimas and Kallus (2020), and sample robust optimization without side information (SRO). In Table 1, we show that each of the above

Table 1 Relationship of four methods.

ϵ_N	$w_N^i(\bar{\gamma}) = \frac{1}{N}$ for all i	$w_N^i(\bar{\gamma})$ from machine learning
= 0	Sample average approximation	Bertsimas and Kallus (2020)
> 0	Bertsimas et al. (2018a)	This paper

methods are particular instances of (4) from Section 3. The methods in the left column ignore side information by assigning equal weights to each uncertainty set, and the methods in the right column incorporate side information by choosing the weights based on predictive machine learning. The methods in the top row do not incorporate any robustness ($\epsilon_N = 0$), and the methods in the bottom row incorporate robustness via a positive choice of the robustness parameter ($\epsilon_N > 0$) in the uncertainty sets. In addition, for the dynamic inventory management example, we also implement and compare to the residual tree algorithm described in Ban et al. (2019). In each experiment, the relevant methods are applied to the same training datasets, and their solutions are evaluated against a common testing dataset. Further details are provided in each of the following sections.

7.1. Dynamic inventory management

We first consider a dynamic inventory control problem over the first $T = 12$ weeks of a new product. In each week, a retailer observes demand for the product and can replenish inventory via procurement orders to different suppliers with lead times. Our problem setting closely follows Ban et al. (2019), motivated by the fashion industry in which retailers have access to auxiliary side information on the new product (color, brand) which are predictive of how demand unfolds over time.

Problem Description. In each stage $t \in \{1, \dots, T\}$, the retailer procures inventory from multiple suppliers to satisfy demand for a single product. The demands for the product across stages are denoted by $\xi_1, \dots, \xi_T \geq 0$. In each stage t , and before the demand ξ_t is observed, the retailer places procurement orders at various suppliers indexed by $\mathcal{J} = \{1, \dots, |\mathcal{J}|\}$. Each supplier $j \in \mathcal{J}$ has per-unit order cost of $c_{tj} \geq 0$ and a lead time of ℓ_j stages. At the end of each stage, the firm incurs a per-unit holding cost of h_t and a backorder cost of b_t . Inventory is fully backlogged and the firm starts with zero initial inventory. The cost incurred by the firm over the time horizon is captured

by

$$\begin{aligned}
c(\xi_1, \dots, \xi_T, \mathbf{x}_1, \dots, \mathbf{x}_T) &= \sum_{t=1}^T \sum_{j \in \mathcal{J}} c_{tj} x_{tj} + \sum_{t=1}^T \text{minimize}_{y_t \in \mathbb{R}} y_t \\
&\text{subject to } y_t \geq h_t \left(\sum_{j \in \mathcal{J}} \sum_{s=1}^{t-\ell_j} x_{sj} - \sum_{s=1}^t \xi_s \right) \\
& y_t \geq -b_t \left(\sum_{j \in \mathcal{J}} \sum_{s=1}^{t-\ell_j} x_{sj} - \sum_{s=1}^t \xi_s \right).
\end{aligned}$$

Experiments. The parameters of the procurement problem were chosen based on Ban et al. (2019). Specifically, we consider the case of two suppliers where $c_{t1} = 1.0$, $c_{t2} = 0.5$, $h_t = 0.25$, and $b_t = 11$ for each stage. The first supplier has no lead time and the second supplier has a lead time of one stage. We generate training and test data from the same distribution as a shipment planning problem of Bertsimas and Kallus (2020, Section EC.6), with the exception that we generate the side information as i.i.d. samples as opposed to an ARMA process (but with the same marginal distribution). In this case, the demands produced by this data generating process are interpreted as the demands over the $T = 12$ stages. We perform computational experiments comparing the proposed sample robust optimization with side information and the residual tree algorithm proposed by Ban et al. (2019). In particular, we compare sample robust optimization with side information *with* the multi-policy approximation as well as *without* the multi-policy approximation (in which we use a single auxiliary linear decision rule for y_t for all uncertainty sets in each stage). The uncertainty sets from Section 3 are defined with the ℓ_2 norm and $\Xi = \mathbb{R}_+^{12}$. The out-of-sample cost resulting from the decision rules were averaged over 100 training sets of size $N = 40$ and 100 testing points, and sample robust optimization with side information used k -nearest neighbors with varying choices of k and radius $\epsilon \geq 0$ of the uncertainty sets.

Results. In Table 2, we show the average out-of-sample cost resulting from sample robust optimization with side information using linear decision rules, with and without the multi-policy approximation from Section 6. In both settings, we used k -nearest neighbors as the machine learning method and evaluated the out-of-sample performance by applying the linear decision rules for the ordering quantities. The results of these computational experiments in Table 2 demonstrate that significant improvements in average out-of-sample performance are found when combining the multi-policy approximation with side information via k -nearest neighbors. We show in Table 3 that these results are statistically significant. For comparison, we also implemented the residual tree algorithm from Ban et al. (2019). When using their algorithm with a binning of $B = 2$ in each stage, their approach resulted in an average out-of-sample cost of 27142. We were unable to

Table 2 Average out-of-sample cost for dynamic inventory problem.

Method	k	ϵ							
		0	100	200	300	400	500	600	700
Sample robust optimization									
Linear decision rules									
no side information		9669	8783	8590	8789	9150	9604	10102	10614
k-nearest neighbors	26	9600	8566	8411	8642	9030	9494	10001	10528
	20	9640	8544	8375	8603	8996	9464	9974	10505
	13	9862	8561	8365	8573	8960	9433	9943	10473
Linear decision rules with multi-policy									
no side information		8967	7759	7360	7320	7460	7716	8038	8412
k-nearest neighbors	26	11346	8728	7651	7269	7241	7381	7636	7966
	20	13012	9460	7925	7328	7195	7289	7519	7835
	13	16288	10975	8576	7585	7243	7236	7412	7697

Average out-of-sample cost for the dynamic inventory problem using sample robust optimization with $N = 40$. For each uncertainty set radius ϵ and parameter k , average was taken over 100 training sets and 100 test points. Optimal is indicated in bold. The residual tree algorithm with a binning of $B = 2$ in each stage gave an average out-of-sample cost of 27142.

Table 3 Statistical significance for dynamic inventory problem.

Method	k	ϵ							
		0	100	200	300	400	500	600	700
Sample robust optimization									
Linear decision rules									
no side information		*	*	*	*	*	*	*	*
k-nearest neighbors	26	*	*	*	*	*	*	*	*
	20	*	*	*	*	*	*	*	*
	13	*	*	*	*	*	*	*	*
Linear decision rules with multi-policy									
no side information		*	*	*	*	*	*	*	*
k-nearest neighbors	26	*	*	*	*	*	*	*	*
	20	*	*	*	*	-	*	*	*
	13	*	*	*	*	5.8×10^{-3}	1×10^{-3}	*	*

The p -values of the Wilcoxon signed rank test for comparison with sample robust optimization using linear decision rules with multi-policy, $k = 20$, and $\epsilon = 400$. An asterisk denotes that the p -value was less than 10^{-8} . After adjusting for multiple hypothesis testing, each result is significant at the $\alpha = 0.05$ significance level if its p -value is less than $\frac{0.05}{63} \approx 7.9 \times 10^{-4}$.

run with a binning of $B = 3$ in each stage due to time limitations of 10^3 seconds, as the size of the resulting linear optimization problem scales on the order $O(B^T)$. Such results are consistent with the estimations of computation times presented in (Ban et al. 2019, Section 6.3). The running times of the various methods are displayed in Table 4.

7.2. Portfolio optimization

The guarantees developed in this paper (Theorem 1) and the above numerical experiment shows that (4) is practically tractable and performs well in problems where $T \geq 1$. In the current and the following section, we provide numerical evidence that (4) can also outperform existing approaches on single-period problems.

Table 4 Average computation time (seconds) for dynamic inventory problem.

Method	k	ϵ							
		0	100	200	300	400	500	600	700
Sample robust optimization									
Linear decision rules									
no side information		3.86	25.04	24.75	25.82	28.70	35.37	31.13	31.95
k-nearest neighbors	26	4.02	25.43	23.39	25.15	27.88	33.42	30.87	31.60
	20	3.99	25.98	23.56	24.93	27.41	32.67	30.69	31.50
	13	4.19	26.53	24.89	24.99	26.79	31.64	30.23	31.32
Linear decision rules with multi-policy									
no side information		0.16	28.31	30.01	29.05	31.13	36.03	35.57	36.09
k-nearest neighbors	26	0.15	27.74	28.69	27.78	30.54	34.44	35.50	36.15
	20	0.15	27.87	28.51	27.74	30.60	34.36	35.65	36.99
	13	0.14	27.78	28.30	27.27	30.00	33.67	35.91	37.76

Average computation time (seconds) for the dynamic inventory problem using sample robust optimization with $N = 40$. For each choice of uncertainty set radius ϵ and parameter k , average was taken over 100 training sets. The residual tree algorithm of Ban et al. (2019) with a binning of $B = 2$ in each stage had an average computation time of 23.20 seconds. We were unable to run this algorithm with binning of $B = 3$ in each stage.

Specifically, in this section we consider a single-stage portfolio optimization problem in which we wish to find an allocation of a fixed budget to n assets. Our goal is to simultaneously maximize the expected return while minimizing the conditional value at risk (cVaR) of the portfolio. Before selecting our portfolio, we observe auxiliary side information which include general market indicators such as index performance as well as macroeconomic numbers released by the US Bureau of Labor Statistics.

Problem Description. We denote the portfolio allocation among the assets by $\mathbf{x} \in \mathcal{X} \triangleq \{\mathbf{x} \in \mathbb{R}_+^n : \sum_{j=1}^n x_j = 1\}$, and the returns of the assets by the random variables $\boldsymbol{\xi} \in \mathbb{R}^n$. The conditional value at risk at the $\alpha \in (0, 1)$ level measures the expected loss of the portfolio, conditional on losses being above the $1 - \alpha$ quantile of the loss distribution. Rockafellar and Uryasev (2000) showed that the cVaR of a portfolio can be computed as the optimal objective value of a convex minimization problem. Therefore, our portfolio optimization problem can be expressed as a convex optimization problem with an auxiliary decision variable, $\beta \in \mathbb{R}$. Thus, given an observation $\bar{\gamma}$ of the auxiliary side information, our goal is to solve

$$\underset{\mathbf{x} \in \mathcal{X}, \beta \in \mathbb{R}}{\text{minimize}} \quad \mathbb{E} \left[\beta + \frac{1}{\alpha} \max(0, -\mathbf{x}^\top \boldsymbol{\xi} - \beta) - \lambda \mathbf{x}^\top \boldsymbol{\xi} \mid \gamma = \bar{\gamma} \right], \quad (17)$$

where $\lambda \in \mathbb{R}_+$ is a trade-off parameter that balances the risk and return objectives.

Experiment. Our experiments are based on a similar setting from Bertsimas and Van Parys (2017, Section 5.2). Specifically, we perform computational experiments on an instance with parameters $\alpha = 0.05$ and $\lambda = 1$, and the joint distribution of the side information and asset returns are chosen the same as Bertsimas and Van Parys (2017, Section 5.2). In our experiments, we compare sample robust optimization with side information, sample average approximation, sample robust

optimization, and predictions to prescriptions. For the robust approaches (bottom row of Table 1), we construct the uncertainty sets from Section 3 using the ℓ_1 norm. For each training sample size, we compute the out-of-sample objective on a test set of size 1000, and we average the results over 100 instances of training data.

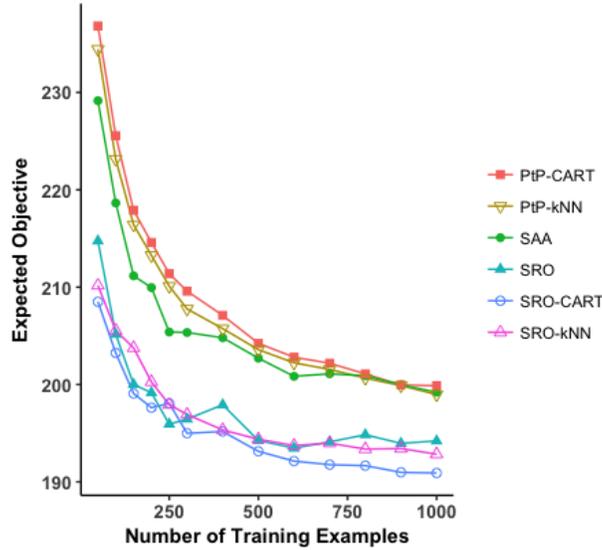
In order to select ϵ_N and other tuning parameters associated with the machine learning weight functions, we first split the data into a training and validation set. We then train the weight functions using the training set, compute decisions for each of the instances in the validation set, and compute the out-of-sample cost on the validation set. We repeat this for a variety of parameter values and select the combination that achieves the best cost on the validation set.

Following a similar reformulation approach as Mohajerin Esfahani and Kuhn (2018), we solve the robust approaches *exactly* by observing that

$$\begin{aligned}
& \underset{\mathbf{x} \in \mathcal{X}, \beta \in \mathbb{R}}{\text{minimize}} && \sum_{i=1}^N w_N^i(\bar{\gamma}) \sup_{\zeta \in \mathcal{U}_N^i} \left\{ \beta + \frac{1}{\alpha} \max\{0, -\mathbf{x}^\top \zeta - \beta\} - \lambda \mathbf{x}^\top \zeta \right\} \\
= & \underset{\mathbf{x} \in \mathcal{X}, \beta \in \mathbb{R}}{\text{minimize}} && \sum_{i=1}^N w_N^i(\bar{\gamma}) \sup_{\zeta \in \mathcal{U}_N^i} \left\{ \max \left\{ \beta - \lambda \mathbf{x}^\top \zeta, \left(\frac{1}{\alpha} + \lambda \right) \mathbf{x}^\top \zeta \right\} \right\} \\
= & \underset{\mathbf{x} \in \mathcal{X}, \beta \in \mathbb{R}}{\text{minimize}} && \sum_{i=1}^N w_N^i(\bar{\gamma}) \max \left\{ \sup_{\zeta \in \mathcal{U}_N^i} \{ \beta - \lambda \mathbf{x}^\top \zeta \}, \sup_{\zeta \in \mathcal{U}_N^i} \left(\frac{1}{\alpha} + \lambda \right) \mathbf{x}^\top \zeta \right\}, \\
= & \underset{\mathbf{x} \in \mathcal{X}, \beta \in \mathbb{R}, \mathbf{v} \in \mathbb{R}^N}{\text{minimize}} && \sum_{i=1}^N w_N^i(\bar{\gamma}) v_i \\
& \text{subject to} && v_i \geq \beta - \lambda \mathbf{x}^\top \zeta \\
& && v_i \geq \left(\frac{1}{\alpha} + \lambda \right) \mathbf{x}^\top \zeta \\
& && \forall \zeta \in \mathcal{U}_N^i, i \in \{1, \dots, N\}.
\end{aligned}$$

The final expression can be reformulated as a deterministic optimization problem by reformulating the robust constraints.

Results. In Figure 1, we show the average out-of-sample objective values using the various methods. Consistent with the computational results of Mohajerin Esfahani and Kuhn (2018) and Bertsimas and Van Parys (2017), the results underscore the importance of robustness in preventing overfitting and achieving good out-of-sample performance in the small data regime. Indeed, we observe that the sample average approximation, which ignores the auxiliary data, outperforms PtP- k NN and PtP-CART when the amount of training data is limited. We believe this is due to the fact the latter methods both throw out training examples, so the methods overfit when the training data is limited, leading to poor out-of-sample performance. In contrast, our methods (SRO- k NN and SRO-CART) typically achieve the strongest out-of-sample performance, even though the amount of training data is limited.

Figure 1 Out-of-sample objective for the portfolio optimization example.

7.3. Shipment planning

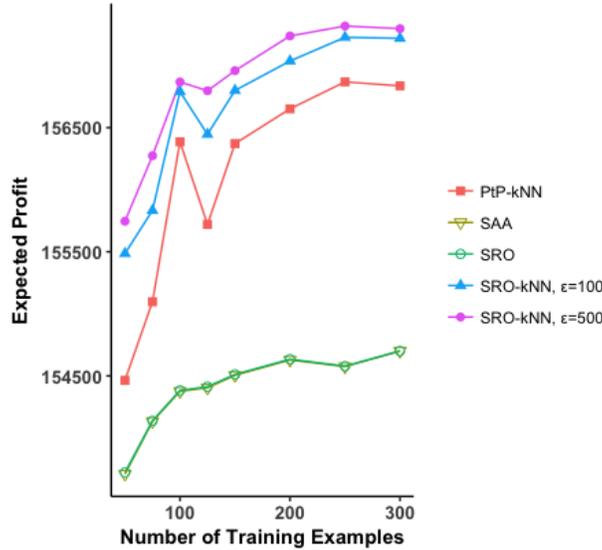
We finally consider a shipment planning problem in which a decision maker seeks to satisfy demand in several locations from several production facilities while minimizing production and transportation costs. Our problem setting closely follows Bertsimas and Kallus (2020), in which the decision maker has access to auxiliary side information (promotions, social media, market trends), which may be predictive of future sales in each retail location.

Problem Description. The decision maker first decides the quantity of inventory $x_f \geq 0$ to produce in each of the production facilities $f \in \mathcal{F} \triangleq \{1, \dots, |\mathcal{F}|\}$, at a cost of p_1 per unit. The demands $\xi_\ell \geq 0$ in each location $\ell \in \mathcal{L} \triangleq \{1, \dots, |\mathcal{L}|\}$ are then observed. The decision maker fulfills these demands by shipping $s_{f\ell} \geq 0$ units from facility $f \in \mathcal{F}$ to location $\ell \in \mathcal{L}$ at a per-unit cost of $c_{f\ell} > 0$. Additionally, after observing demand, the decision maker has the opportunity to produce additional units $y_f \geq 0$ in each facility at a cost of $p_2 > p_1$ per unit. The fulfillment of each unit of demand generates $r > 0$ in revenue. Given the above notation and dynamics, the cost incurred by the decision maker is

$$\begin{aligned}
 c(\boldsymbol{\xi}, \mathbf{x}) = & \sum_{f \in \mathcal{F}} p_1 x_f - \sum_{\ell \in \mathcal{L}} r \xi_\ell + \underset{\mathbf{s} \in \mathbb{R}_+^{\mathcal{L} \times \mathcal{F}}, \mathbf{y} \in \mathbb{R}_+^{\mathcal{F}}}{\text{minimize}} & \sum_{f \in \mathcal{F}} p_2 y_f + \sum_{f \in \mathcal{F}} \sum_{\ell \in \mathcal{L}} c_{f\ell} s_{f\ell} \\
 & \text{subject to} & \sum_{f \in \mathcal{F}} s_{f\ell} \geq \xi_\ell & \forall \ell \in \mathcal{L} \\
 & & \sum_{\ell \in \mathcal{L}} s_{f\ell} \leq x_f + y_f & \forall f \in \mathcal{F}.
 \end{aligned}$$

Experiments. We perform computational experiments using the same parameters and data generation procedure as Bertsimas and Kallus (2020). Specifically, we consider an instance with $|\mathcal{F}| = 4$,

Figure 2 Out-of-sample profit for the shipment planning example.



Note. The profits for SRO and SAA are overlapping.

$|\mathcal{L}| = 12$, $p_1 = 5$, $p_2 = 100$, and $r = 90$. The network topology, transportation costs, and the joint distribution of the side information $\gamma \in \mathbb{R}^3$ and demands $\xi \in \mathbb{R}^{12}$ are the same as Bertsimas and Kallus (2020), with the exception that we generate the side information as i.i.d. samples as opposed to an ARMA process (but with the same marginal distribution).

In our experiments, we compare sample robust optimization with side information, sample average approximation, sample robust optimization, and predictions to prescriptions. For the robust approaches (bottom row of Table 1), we construct the uncertainty sets from Section 3 using the ℓ_1 norm and $\Xi = \mathbb{R}_+^{12}$, solve these problems using the multi-policy approximation with linear decision rules described in Section 6, and consider uncertainty sets with radius $\epsilon \in \{100, 500\}$. For the approaches using side information (right column of Table 1), we used the k_N -nearest neighbors with parameter $k_N = \frac{2N}{5}$. All solutions were evaluated on a test set of size 100 and the results were averaged over 100 independent training sets.

Results. In Figure 2, we present the average out-of-sample profits of the various methods. The results show that the best out-of-sample average profit is attained when using the proposed sample robust optimization with side information. Interestingly, we observe no discernible differences between sample average approximation and sample robust optimization in Figure 2, suggesting the value gained by incorporating side information in this example. Compared to the approach of Bertsimas and Kallus (2020), sample robust optimization with side information achieves a better out-of-sample average performance for each choice of ϵ . Table 5 shows that these differences are

Table 5 Statistical significance for shipment planning problem.

N	ϵ	
	100	500
50	4.6×10^{-13}	5.3×10^{-16}
75	1.3×10^{-14}	6.4×10^{-12}
100	1.2×10^{-13}	1.1×10^{-7}
125	2.6×10^{-15}	1.5×10^{-11}
150	3.4×10^{-12}	1.2×10^{-6}
200	1.4×10^{-12}	1.0×10^{-8}
250	3.4×10^{-10}	1.0×10^{-4}
300	1.8×10^{-6}	5.2×10^{-4}

The p -values from the Wilcoxon signed rank test for comparison with the predictive to prescriptive analytics method (PtP- k NN) and sample robust optimization with side information (SRO- k NN). After adjusting for multiple hypothesis testing, all results are significant at the $\alpha = 0.05$ significance level because all p -values are less than $\frac{0.05}{16} \approx 3.1 \times 10^{-3}$.

statistically significant. This example demonstrates that, in addition to enjoying asymptotic optimality guarantees, sample robust optimization with side information provides meaningful value across various values of N .

8. Conclusion

In this paper, we introduced *sample robust optimization with side information*, a new approach for solving dynamic optimization problems with side information. Through three computational examples, we demonstrated that our method achieves significantly better out-of-sample performance than scenario-based alternatives. We complemented these empirical observations with theoretical analysis, showing our nonparametric method is asymptotically optimal via a new concentration measure result for local learning methods. Finally, we showed our approach inherits the tractability of robust optimization, scaling to problems with many stages via the multi-policy approximation scheme.

Acknowledgements

The authors thank the associate editor and two referees for many helpful suggestions that greatly improved the manuscript.

Endnotes

1. If the random vectors are continuous and $T \geq 2$, it is readily observed that (2) resolves to an optimization problem of the form

$$\underset{\mathbf{x}_1 \in \mathcal{X}_1; \mathbf{x}_2^i \in \mathcal{X}_2, \dots, \mathbf{x}_T^i \in \mathcal{X}_T \forall i}{\text{minimize}} \quad \sum_{i=1}^N w_N^i(\bar{\gamma}) c(\boldsymbol{\xi}_1^i, \dots, \boldsymbol{\xi}_T^i, \mathbf{x}_1, \mathbf{x}_2^i, \dots, \mathbf{x}_T^i).$$

2. To see why this is without loss of generality, consider any other ℓ_p norm where $p \geq 1$. In this case,

$$\|\boldsymbol{\xi} - \boldsymbol{\xi}'\|_p \leq d_\xi^{1/p} \|\boldsymbol{\xi} - \boldsymbol{\xi}'\|_\infty.$$

By the definition of the 1-Wasserstein metric, this implies

$$\mathbf{d}_1^p(\mathbb{P}_{\bar{\gamma}}, \hat{\mathbb{P}}_{\bar{\gamma}}^N) \leq d_\xi^{1/p} \mathbf{d}_1^\infty(\mathbb{P}_{\bar{\gamma}}, \hat{\mathbb{P}}_{\bar{\gamma}}^N),$$

where \mathbf{d}_1^p refers to the 1-Wasserstein metric with the ℓ_p norm. If ϵ_N satisfies Assumption 1, $\epsilon_N/d_\xi^{1/p}$ also satisfies Assumption 1, so the result for all other choices of ℓ_p norms follows from the result with the ℓ_∞ norm.

References

- Anderson E, Philpott A (2019) Improving sample average approximation using distributional robustness
URL http://www.optimization-online.org/DB_FILE/2019/10/7405.pdf.
- Ban GY, Gallien J, Mersereau AJ (2019) Dynamic procurement of new products with covariate information: The residual tree method. *Manufacturing & Service Operations Management* 21(4):798–815.
- Ban GY, Rudin C (2018) The big data newsvendor: practical insights from machine learning. *Operations Research* 67(1):90–108.
- Bertsimas D, Kallus N (2020) From predictive to prescriptive analytics. *Management Science* 66(3):1025–1044.
- Bertsimas D, McCord C (2019) From predictions to prescriptions in multistage optimization problems. *arXiv preprint arXiv:1904.11637* .
- Bertsimas D, Shtern S, Sturt B (2018a) A data-driven approach to multi-stage stochastic linear optimization
URL http://www.optimization-online.org/DB_FILE/2018/11/6907.pdf.
- Bertsimas D, Shtern S, Sturt B (2019) Two-stage sample robust optimization. *arXiv preprint arXiv:1907.07142* .
- Bertsimas D, Sim M, Zhang M (2018b) Adaptive distributionally robust optimization. *Management Science* 65(2):604–618.

- Bertsimas D, Van Parys B (2017) Bootstrap robust prescriptive analytics. *arXiv preprint arXiv:1711.09974* .
- Biau G, Devroye L (2015) *Lectures on the Nearest Neighbor Method* (Springer).
- Blanchet J, Murthy K (2019) Quantifying distributional model risk via optimal transport. *Mathematics of Operations Research* 44(2):565–600.
- Breiman L (2001) Random forests. *Machine Learning* 45(1):5–32.
- Breiman L, Friedman JH, Olshen RA, Stone CJ (1984) *Classification and Regression Trees* (Chapman & Hall/CRC).
- Chen X, Zhang Y (2009) Uncertain linear programs: extended affinely adjustable robust counterparts. *Operations Research* 57(6):1469–1482.
- Chen Z, Sim M, Xiong P (2020) Robust stochastic optimization made easy with rsome. *Forthcoming in Management Science* .
- Clement P, Desch W (2008) An elementary proof of the triangle inequality for the wasserstein metric. *Proceedings of the American Mathematical Society* 136(1):333–339.
- Delage E, Ye Y (2010) Distributionally robust optimization under moment uncertainty with application to data-driven problems. *Operations Research* 58(3):595–612.
- Elmachtoub AN, Grigas P (2017) Smart “predict, then optimize”. *arXiv preprint arXiv:1710.08005* .
- Fournier N, Guillin A (2015) On the rate of convergence in wasserstein distance of the empirical measure. *Probability Theory and Related Fields* 162(3-4):707–738.
- Friedman J, Hastie T, Tibshirani R (2001) *The Elements of Statistical Learning*, volume 1 (Springer).
- Gao R, Kleywegt AJ (2016) Distributionally robust stochastic optimization with wasserstein distance. *arXiv preprint arXiv:1604.02199* .
- Georghiou A, Wiesemann W, Kuhn D (2015) Generalized decision rule approximations for stochastic programming via liftings. *Mathematical Programming* 152(1-2):301–338.
- Gotoh Jy, Kim MJ, Lim AE (2018) Robust empirical optimization is almost the same as mean–variance optimization. *Operations Research Letters* 46(4):448–452.
- Hanasusanto GA, Kuhn D (2018) Conic programming reformulations of two-stage distributionally robust linear programs over wasserstein balls. *Operations Research* 66(3):849–869.
- Hanasusanto GA, Kuhn D, Wiesemann W (2016) K-adaptability in two-stage distributionally robust binary programming. *Operations Research Letters* 44(1):6–11.
- Hannah L, Powell W, Blei DM (2010) Nonparametric density estimation for stochastic optimization with an observable state variable. *Advances in Neural Information Processing Systems*, 820–828.

- Ho CP, Hanasusanto G (2019) On data-driven prescriptive analytics with side information: a regularized nadaraya-watson approach URL http://www.optimization-online.org/DB_FILE/2019/01/7043.pdf.
- Kantorovich L, Rubinstein G (1958) On a space of totally additive functions. *Vestn Lening. Univ* 13:52–59.
- Kullback S, Leibler RA (1951) On information and sufficiency. *The Annals of Mathematical Statistics* 22(1):79–86.
- Mohajerin Esfahani P, Kuhn D (2018) Data-driven distributionally robust optimization using the wasserstein metric: Performance guarantees and tractable reformulations. *Mathematical Programming* 171(1):115–166.
- Natarajan K, Teo CP, Zheng Z (2011) Mixed 0-1 linear programs under objective uncertainty: A completely positive representation. *Operations Research* 59(3):713–728.
- Rockafellar RT, Uryasev S (2000) Optimization of conditional value-at-risk. *Journal of Risk* 2:21–42.
- Van Parys BP, Esfahani PM, Kuhn D (2017) From data to decisions: distributionally robust optimization is optimal. *arXiv preprint arXiv:1704.04118* .
- Vershynin R (2018) *High-Dimensional Probability: An Introduction with Applications in Data Science*, volume 47 (Cambridge University Press).
- Walk H (2010) Strong laws of large numbers and nonparametric estimation. *Recent Developments in Applied Probability and Statistics*, 183–214 (Springer).
- Xu H, Caramanis C, Mannor S (2012) A distributional interpretation of robust optimization. *Mathematics of Operations Research* 37(1):95–110.

Electronic Companion

EC.1. Properties of Weight Functions

In this section, we show that the k -nearest neighbor and kernel regression weight functions satisfy several guarantees. These results are used in the proof of Theorem 2, found in Section 4.3. The main result of this section is the following. For convenience, the equations below are numbered the same as in the proof of Theorem 2.

THEOREM EC.1. *If Assumptions 1 and 4 hold, then*

$$\{w_N^i(\bar{\gamma})\} \text{ are not functions of } \xi^1, \dots, \xi^N; \quad (5)$$

$$\sum_{i=1}^N w_N^i(\bar{\gamma}) = 1 \text{ and } w_N^1(\bar{\gamma}), \dots, w_N^N(\bar{\gamma}) \geq 0, \quad \forall N \in \mathbb{N}. \quad (6)$$

Moreover, there exists constants $k_2 > 0$ and $\eta > p(2 + d_\xi)$ such that

$$\lim_{N \rightarrow \infty} \frac{1}{\epsilon_N} \sum_{i=1}^N w_N^i(\bar{\gamma}) \|\gamma^i - \bar{\gamma}\| = 0, \quad \mathbb{P}^\infty\text{-almost surely}; \quad (8)$$

$$\mathbb{E}_{\mathbb{P}^N} \left[\exp \left(\frac{-\theta}{\sum_{i=1}^N w_N^i(\bar{\gamma})^2} \right) \right] \leq \exp(-k_2 \theta N^\eta), \quad \forall \theta \in (0, 1), N \in \mathbb{N}. \quad (9)$$

Proof. We observe that (5) and (6) follow directly from the definitions of the weight functions. The proofs of (8) and (9) are split into two parts, one for the k -nearest neighbor weights and one for kernel regression weights.

k-Nearest Neighbors: For the proof of (8), we note

$$\sum_{i=1}^N w_N^i(\bar{\gamma}) \|\gamma^i - \bar{\gamma}\| \leq \|\gamma^{(k_N)}(\bar{\gamma}) - \bar{\gamma}\|,$$

where $\gamma^{(k_N)}(\bar{\gamma})$ denotes the k_N th nearest neighbor of $\bar{\gamma}$ out of $\gamma^1, \dots, \gamma^N$. Therefore, for any $\lambda > 0$,

$$\begin{aligned} \mathbb{P}^N \left(\sum_{i=1}^N w_N^i(\bar{\gamma}) \|\gamma^i - \bar{\gamma}\| > \lambda \epsilon_N \right) &\leq \mathbb{P}^N (\|\gamma^{(k_N)}(\bar{\gamma}) - \bar{\gamma}\| > \lambda \epsilon_N) \\ &\leq \mathbb{P}^N (|\{i : \|\gamma^i - \bar{\gamma}\| \leq \lambda \epsilon_N\}| \leq k_N - 1). \end{aligned}$$

By Assumption 4, this probability is upper bounded by $\mathbb{P}(\beta \leq k_N - 1)$, where $\beta \sim \text{Binom}(N, g(\lambda \epsilon_N)^{d_\gamma})$. By Hoeffding's inequality,

$$\mathbb{P}^N \left(\sum_{i=1}^N w_N^i(\bar{\gamma}) \|\gamma^i - \bar{\gamma}\| > \lambda \epsilon_N \right) \leq \exp \left(\frac{-2(Ng(\lambda k_1/N^p)^{d_\gamma} - k_N + 1)^2}{N} \right),$$

for $k_N \leq Ng(\lambda k_1/N^p)^{d_\gamma} + 1$. We note that this condition on k_N is satisfied for N sufficiently large because $\delta + pd_\gamma < 1$ by Assumption 1. Because the right hand side in the above inequality has a finite sum over N , (8) follows by the Borel Cantelli lemma.

For the proof of (9), it follows from Assumption 1 that

$$\sum_{i=1}^N w_N^i(\bar{\gamma})^2 \leq k_3^{-2} N^{1-2\delta}$$

deterministically (for all sufficiently large N such that $\lceil k_3 N^\delta \rceil \leq N - 1$) and $2\delta - 1 > p(d_\xi + 2)$. Thus, (9) follows with $\eta = 2\delta - 1$.

Kernel regression: Assumption 1 stipulates that the kernel function $K(\cdot)$ is Gaussian, triangular, or Epanechnikov, which are defined in Section 3. It is easy to verify that these kernel functions satisfy the following:

1. K is nonnegative, finite valued, and monotonically decreasing (for nonnegative inputs).
2. $u^\alpha K(u) \rightarrow 0$ as $u \rightarrow \infty$ for any $\alpha \in \mathbb{R}$.
3. $\exists u^* > 0$ such that $K(u^*) > 0$.

For the proof of (8), define $q > 0$ such that $p < q < \delta$. Letting D be the diameter of Γ and $g_N(\bar{\gamma}) = \sum_{i=1}^N K(\|\gamma^i - \bar{\gamma}\|/h_N)$, we have

$$\begin{aligned} & \sum_{i=1}^N w_N^i(\bar{\gamma}) \|\gamma^i - \bar{\gamma}\| \\ &= \sum_{i=1}^N w_N^i(\bar{\gamma}) \mathbf{1}\{\|\gamma^i - \bar{\gamma}\| \leq N^{-q}\} \|\gamma^i - \bar{\gamma}\| + \frac{1}{g_N(\bar{\gamma})} \sum_{i=1}^N K\left(\frac{\|\gamma^i - \bar{\gamma}\|}{h_N}\right) \mathbf{1}\{\|\gamma^i - \bar{\gamma}\| > N^{-q}\} \|\gamma^i - \bar{\gamma}\| \\ &\leq N^{-q} + \frac{NDK(N^{-q}/h_N)}{g_N(\bar{\gamma})}, \end{aligned}$$

where the inequality follows from the monotonicity of K . By construction, $N^{-q}/\epsilon_N \rightarrow 0$, so we just need to handle the second term. We note, for any $\lambda > 0$,

$$\mathbb{P}^N \left(\frac{NDK(N^{-q}/h_N)}{g_N(\bar{\gamma})} > \lambda \epsilon_N \right) \leq \mathbb{P}^N \left(\sum_{i=1}^N Z_i^N K(u^*) < \frac{NDK(N^{-q}/h_N)}{\lambda \epsilon_N} \right),$$

where $Z_i^N = \mathbf{1}\{\|\gamma^i - \bar{\gamma}\| \leq u^* h_N\}$. To achieve this inequality, we lower bounded each term in $g_N(\bar{\gamma})$ by $K(u^*)$ or 0, because of the monotonicity of K . By Hoeffding's inequality,

$$\mathbb{P}^N \left(\sum_{i=1}^N Z_i^N K(u^*) < \frac{NDK(N^{-q}/h_N)}{\lambda \epsilon_N} \right) \leq \exp \left(- \frac{2 \left(N \mathbb{E} Z_i^N - \frac{ND}{\lambda \epsilon_N K(u^*)} K(N^{-q}/h_N) \right)_+^2}{N} \right)$$

$$\begin{aligned} &\leq \exp\left(-\frac{2\left(Ng(u^*h_N)^{d_\gamma} - \frac{ND}{\lambda\epsilon_N K(u^*)}K(N^{-q}/h_N)\right)_+^2}{N}\right) \\ &= \exp\left(-\left(k_5N^{1/2-\delta d_\gamma} - k_6N^{1/2+p}K(k_4N^{-q+\delta})\right)_+^2\right), \end{aligned}$$

for some constants $k_5, k_6 > 0$ that do not depend on N . We used Assumption 4 for the second inequality. Because $\delta > q$, the second kernel property implies $N^{1/2+p}K(k_4N^{-q+\delta})$ goes to 0 as N goes to infinity, so that term is irrelevant. Because $1/2 - \delta d_\gamma > 0$ by Assumption 1, the right hand side of the inequality has a finite sum over N , and thus (8) follows from the Borel Cantelli lemma.

For the proof of (9), define

$$v^N = \begin{pmatrix} K(\|\gamma^1 - \bar{\gamma}\|/h_N) \\ \vdots \\ K(\|\gamma^N - \bar{\gamma}\|/h_N) \end{pmatrix}.$$

We note that

$$\sum_{i=1}^N w_N^i(\bar{\gamma})^2 = \frac{\|v^N\|_2^2}{\|v^N\|_1^2} \leq \frac{\|v^N\|_\infty}{\|v^N\|_1} \leq \frac{K(0)}{K(u^*) \sum_{i=1}^N Z_i^N},$$

where Z_i^N is defined above. The first inequality follows from Holder's inequality, and the second inequality follows from the monotonicity of K . Next, we define \bar{Z}_i^N to be a Bernoulli random variable with parameter $g(u^*h_N)^{d_\gamma}$ for each i . For any $\theta \in (0, 1)$,

$$\begin{aligned} \mathbb{E}_{\mathbb{P}^N} \left[\exp\left(\frac{-\theta}{\sum_{i=1}^N w_N^i(\bar{\gamma})^2}\right) \right] &\leq \mathbb{E}_{\mathbb{P}^N} \left[\exp\left(\frac{-\theta K(u^*) \sum_{i=1}^N \bar{Z}_i^N}{K(0)}\right) \right] \\ &= \left(1 - g(u^*h_N)^{d_\gamma} + g(u^*h_N)^{d_\gamma} \exp(-\theta K(u^*)/K(0))\right)^N \\ &\leq \exp(-Ng(u^*h_N)^{d_\gamma} (1 - \exp(-\theta K(u^*)/K(0)))) \\ &\leq \exp\left(-Ng(u^*h_N)^{d_\gamma} \frac{\theta K(u^*)}{2K(0)}\right) \\ &= \exp\left(-\frac{\theta K(u^*)g(k_4u^*)^{d_\gamma} N^{1-\delta d_\gamma}}{2K(0)}\right). \end{aligned}$$

The first inequality follows because $g(u^*h_N)^{d_\gamma}$ is an upper bound on $\mathbb{P}(\|\gamma^i - \bar{\gamma}\| \leq u^*h_N)$ by Assumption 4. The first equality follows from the definition of the moment generating function for a binomial random variable. The next line follows from the inequality $e^x \geq 1 + x$ and the following from the inequality $1 - e^{-x} \geq x/2$ for $0 \leq x \leq 1$. Because $1 - \delta d_\gamma > p(2 + d_\xi)$, this completes the proof of (9) with $\eta = 1 - \delta d_\gamma$ and $k_2 = K(u^*)g(k_4u^*)^{d_\gamma}/2K(0)$. \square

EC.2. Proof of Theorem 1

In this section, we present our proof of Theorem 1. We make use of the following result from Bertsimas et al. (2018a) (their Lemma EC.2), which bounds the difference in worst case objective

values between distributionally robust optimization with the type-1 Wasserstein ambiguity set and sample robust optimization³ problems. We note that [Bertsimas et al. \(2018a\)](#) proved the following result for the case that \mathbb{Q}' is the unweighted empirical measure, but their proof carries through for the case here in which \mathbb{Q}' is a weighted empirical measure.

LEMMA EC.1. *Let $\mathcal{Z} \subseteq \mathbb{R}^d$, $f: \mathcal{Z} \rightarrow \mathbb{R}$ be measurable, and $\zeta^1, \dots, \zeta^N \in \mathcal{Z}$. Suppose that*

$$\mathbb{Q}' = \sum_{i=1}^N w^i \delta_{\zeta^i}$$

for given weights $w^1, \dots, w^N \geq 0$ that sum to one. If $\theta_2 \geq 2\theta_1 \geq 0$, then

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

We now restate and prove the main result, which combines the new measure concentration result from this paper with similar proof techniques as [Bertsimas et al. \(2018a\)](#) and [Mohajerin Esfahani and Kuhn \(2018\)](#).

THEOREM 1. *Suppose the weight function and uncertainty sets satisfy Assumption 1, the joint probability distribution of (γ, ξ) satisfies Assumptions 2-4 from Section 4.3, and the cost function satisfies Assumptions 5-6 from Section 4.4. Then, for every $\bar{\gamma} \in \Gamma$,*

$$\lim_{N \rightarrow \infty} \hat{v}^N(\bar{\gamma}) = v^*(\bar{\gamma}), \quad \mathbb{P}^\infty\text{-almost surely.}$$

Proof. We break the limit into upper and lower parts. The proof of the lower part follows from an argument similar to that used by [Bertsimas et al. \(2018a\)](#). The proof of the upper part follows from the argument used by [Mohajerin Esfahani and Kuhn \(2018\)](#).

Lower bound. We first show that

$$\liminf_{N \rightarrow \infty} \hat{v}^N(\bar{\gamma}) \geq v^*(\bar{\gamma}), \quad \mathbb{P}^\infty\text{-almost surely.} \quad (\text{EC.1})$$

Indeed, it follows from Assumptions 1-2 and the union bound that there exists $N_0 \in \mathbb{N}$ such that

$$\mathbb{P}^N \left(\sup_{\zeta \in \cup_{i=1}^N \mathcal{U}_N^i} \|\zeta\| > \log N \right) < \exp(-(\log N)^{1.99}), \quad \forall N \geq N_0.$$

Therefore, the Borel-Cantelli lemma implies that there exists $N_1 \in \mathbb{N}$, \mathbb{P}^∞ -almost surely, such that

$$\cup_{i=1}^N \mathcal{U}_N^i \subseteq D_N \triangleq \{\zeta : \|\zeta\| \leq \log N\}, \quad \forall N \geq N_1. \quad (\text{EC.2})$$

Consider any $r > 0$ such that $\epsilon_N N^{-r}$ satisfies Assumption 1, and let Π^N denote the set of decision rules which satisfy the conditions of Assumption 5. Then, the following holds for all $N \geq N_1 \triangleq \max\{N_0, 2^{\frac{1}{r}}\}$ and $\boldsymbol{\pi} \in \Pi^N$:

$$\begin{aligned}
& \sup_{\mathbb{Q} \in \mathcal{P}(D_N \cap \Xi): d_1(\mathbb{Q}, \hat{\mathbb{P}}_{\bar{\gamma}}^N) \leq \frac{\epsilon_N}{N^r}} \mathbb{E}_{\boldsymbol{\xi} \sim \mathbb{Q}}[c^\pi(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T)] \\
& \leq \sum_{i=1}^N w_i^N(\bar{\gamma}) \sup_{\boldsymbol{\zeta} \in D_N \cap \Xi: \|\boldsymbol{\zeta} - \boldsymbol{\xi}^i\| \leq \epsilon_N} c^\pi(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_T) + \frac{4}{N^r} \sup_{\boldsymbol{\zeta} \in D_N \cap \Xi} |c^\pi(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_T)| \\
& = \sum_{i=1}^N w_i^N(\bar{\gamma}) \sup_{\boldsymbol{\zeta} \in \mathcal{U}_N^i} c^\pi(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_T) + \frac{4}{N^r} \sup_{\boldsymbol{\zeta} \in D_N \cap \Xi} |c^\pi(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_T)| \\
& \leq \sum_{i=1}^N w_i^N(\bar{\gamma}) \sup_{\boldsymbol{\zeta} \in \mathcal{U}_N^i} c^\pi(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_T) + \frac{4}{N^r} M \left(1 + \max \left\{ \|\boldsymbol{\zeta}\|, \sup_{\boldsymbol{\zeta}' \in \cup_{i=1}^N \mathcal{U}_N^i} \|\boldsymbol{\zeta}'\| \right\} \right) \\
& \leq \sum_{i=1}^N w_i^N(\bar{\gamma}) \sup_{\boldsymbol{\zeta} \in \mathcal{U}_N^i} c^\pi(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_T) + \frac{4M}{N^r} (1 + \log N). \tag{EC.3}
\end{aligned}$$

Indeed, the first inequality follows from Lemma EC.1 since $N \geq 2^{\frac{1}{r}}$, the equality follows from $N \geq N_1$, the second inequality holds because $\boldsymbol{\pi} \in \Pi^N$, and the third and final inequality follows from the definition of D_N and $N \geq N_1$. We observe that the second term in (EC.3) converges to zero as $N \rightarrow \infty$.

We now observe that

$$\begin{aligned}
\mathbb{E}[c^\pi(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T) \mid \gamma = \bar{\gamma}] & \triangleq \mathbb{E}_{\boldsymbol{\xi} \sim \mathbb{P}_{\bar{\gamma}}} [c^\pi(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T)] \\
& = \mathbb{E}_{\boldsymbol{\xi} \sim \mathbb{P}_{\bar{\gamma}}} [c^\pi(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T) \mathbf{1}\{\boldsymbol{\xi} \notin D_N\}] + \mathbb{E}_{\boldsymbol{\xi} \sim \mathbb{P}_{\bar{\gamma}}} [c^\pi(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T) \mathbf{1}\{\boldsymbol{\xi} \in D_N\}].
\end{aligned}$$

We handle the first term with the Cauchy-Schwartz inequality,

$$\mathbb{E}_{\boldsymbol{\xi} \sim \mathbb{P}_{\bar{\gamma}}} [c^\pi(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T) \mathbf{1}\{\boldsymbol{\xi} \notin D_N\}] \leq \sqrt{\mathbb{E}_{\boldsymbol{\xi} \sim \mathbb{P}_{\bar{\gamma}}} [c^\pi(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T)^2] \mathbb{P}_{\bar{\gamma}}(\boldsymbol{\xi} \notin D_N)}.$$

By Assumption 2, the above bound is finite and converges to zero as $N \rightarrow \infty$ uniformly over $\boldsymbol{\pi} \in \Pi^N$.

We handle the second term by the new concentration measure from this paper. Specifically, it follows from Theorem 2 that there exists an $N_2 \geq N_1$, \mathbb{P}^∞ -almost surely, such that

$$d_1(\mathbb{P}_{\bar{\gamma}}, \hat{\mathbb{P}}_{\bar{\gamma}}^N) \leq \frac{\epsilon_N}{N^r} \quad \forall N \geq N_2.$$

Therefore, for all $N \geq N_2$ and decision rules $\boldsymbol{\pi} \in \Pi^N$:

$$\begin{aligned}
& \mathbb{E}_{\boldsymbol{\xi} \sim \mathbb{P}_{\bar{\gamma}}} [c^\pi(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T) \mathbf{1}\{\boldsymbol{\xi} \in D_N\}] \\
& = \mathbb{E}_{\boldsymbol{\xi} \sim \mathbb{P}_{\bar{\gamma}}} \left[\left(c^\pi(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T) - \inf_{\boldsymbol{\zeta} \in D_N \cap \Xi} c^\pi(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_T) \right) \mathbf{1}\{\boldsymbol{\xi} \in D_N\} \right] + \underbrace{\mathbb{P}_{\bar{\gamma}}(\boldsymbol{\xi} \in D_N) \inf_{\boldsymbol{\zeta} \in D_N \cap \Xi} c^\pi(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_T)}_{\alpha_N}
\end{aligned}$$

$$\begin{aligned}
&\leq \sup_{\mathbb{Q} \in \mathcal{P}(\Xi): d_1(\mathbb{Q}, \hat{\mathbb{P}}_{\bar{\gamma}}^N) \leq \frac{\epsilon_N}{N^*}} \mathbb{E}_{\xi \sim \mathbb{Q}} \left[\left(c^\pi(\xi_1, \dots, \xi_T) - \inf_{\zeta \in D_N \cap \Xi} c^\pi(\zeta_1, \dots, \zeta_T) \right) \mathbb{1}\{\xi \in D_N\} \right] + \alpha_N \\
&= \sup_{\mathbb{Q} \in \mathcal{P}(\Xi \cap D_N): d_1(\mathbb{Q}, \hat{\mathbb{P}}_{\bar{\gamma}}^N) \leq \frac{\epsilon_N}{N^*}} \mathbb{E}_{\xi \sim \mathbb{Q}} \left[c^\pi(\xi_1, \dots, \xi_T) - \inf_{\zeta \in D_N \cap \Xi} c^\pi(\zeta_1, \dots, \zeta_T) \right] + \alpha_N \\
&= \sup_{\mathbb{Q} \in \mathcal{P}(\Xi \cap D_N): d_1(\mathbb{Q}, \hat{\mathbb{P}}_{\bar{\gamma}}^N) \leq \frac{\epsilon_N}{N^*}} \mathbb{E}_{\xi \sim \mathbb{Q}} [c^\pi(\xi_1, \dots, \xi_T)] - \mathbb{P}_{\bar{\gamma}}(\xi \notin D_N) \inf_{\zeta \in D_N \cap \Xi} c^\pi(\zeta_1, \dots, \zeta_T),
\end{aligned}$$

where the inequality follows from $N \geq N_2$. It follows from (EC.2) that the second term in the final equality converges to zero as $N \rightarrow \infty$ uniformly over $\pi \in \Pi^N$.

Combining the above, we conclude that

$$\begin{aligned}
\liminf_{N \rightarrow \infty} \hat{v}^N(\bar{\gamma}) &= \liminf_{N \rightarrow \infty} \inf_{\pi \in \Pi} \sum_{i=1}^N w_i^N(\bar{\gamma}) \sup_{\zeta \in \mathcal{U}_N^i} c^\pi(\zeta_1, \dots, \zeta_T) \\
&= \liminf_{N \rightarrow \infty} \inf_{\pi \in \Pi^N} \sum_{i=1}^N w_i^N(\bar{\gamma}) \sup_{\zeta \in \mathcal{U}_N^i} c^\pi(\zeta_1, \dots, \zeta_T) \tag{EC.4}
\end{aligned}$$

$$\begin{aligned}
&\geq \liminf_{N \rightarrow \infty} \inf_{\pi \in \Pi^N} \mathbb{E}[c^\pi(\xi_1, \dots, \xi_T) \mid \gamma = \bar{\gamma}], \quad \mathbb{P}^\infty\text{-almost surely} \\
&\geq \inf_{\pi \in \Pi} \mathbb{E}[c^\pi(\xi_1, \dots, \xi_T) \mid \gamma = \bar{\gamma}] \tag{EC.5} \\
&= v^*(\bar{\gamma}),
\end{aligned}$$

where (EC.4) follows from Assumption 5 and (EC.5) follows because $\Pi^N \subseteq \Pi$ for all $N \in \mathbb{N}$. This completes the proof of (EC.1).

Upper bound. We now prove that

$$\limsup_{N \rightarrow \infty} \hat{v}^N(\bar{\gamma}) \leq v^*(\bar{\gamma}), \quad \mathbb{P}^\infty\text{-almost surely.} \tag{EC.6}$$

Indeed, for any arbitrary $\delta > 0$, let $\pi_\delta \in \Pi$ be a δ -optimal solution for (1). Moreover, without any loss of generality, we assume that the decision rule is chosen to satisfy the conditions of Assumption 6. Then it follows from [Mohajerin Esfahani and Kuhn \(2018, Lemma A.1\)](#) that there exists a non-increasing sequence of functions $f^j(\zeta_1, \dots, \zeta_T)$, $j \in \mathbb{N}$, such that

$$\lim_{j \rightarrow \infty} f^j(\zeta_1, \dots, \zeta_T) = c^{\pi_\delta}(\zeta_1, \dots, \zeta_T), \quad \forall \zeta \in \Xi$$

and f^j is L_j -Lipschitz continuous. Furthermore, for each $N \in \mathbb{N}$, choose any probability distribution $\hat{\mathbb{Q}}^N \in \mathcal{P}(\Xi)$ such that $d_1(\hat{\mathbb{Q}}^N, \hat{\mathbb{P}}_{\bar{\gamma}}^N) \leq \epsilon_N$ and

$$\sup_{\mathbb{Q} \in \mathcal{P}(\Xi): d_1(\mathbb{Q}, \hat{\mathbb{P}}_{\bar{\gamma}}^N) \leq \epsilon_N} \mathbb{E}_{\xi \sim \mathbb{Q}} [c^{\pi_\delta}(\xi_1, \dots, \xi_T)] \leq \mathbb{E}_{\xi \sim \hat{\mathbb{Q}}^N} [c^{\pi_\delta}(\xi_1, \dots, \xi_T)] + \delta.$$

For any $j \in \mathbb{N}$,

$$\begin{aligned}
\limsup_{N \rightarrow \infty} \hat{v}^N(\bar{\gamma}) &\leq \limsup_{N \rightarrow \infty} \sum_{i=1}^N w_i^N(\bar{\gamma}) \sup_{\boldsymbol{\zeta} \in \mathcal{U}_N^i} c^{\pi^\delta}(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_T) \\
&= \limsup_{N \rightarrow \infty} \sup_{\mathbb{Q} \in \mathcal{P}(\Xi): d_\infty(\mathbb{Q}, \hat{\mathbb{P}}_{\bar{\gamma}}^N) \leq \epsilon_N} \mathbb{E}_{\boldsymbol{\xi} \sim \mathbb{Q}}[c^{\pi^\delta}(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T)] \\
&\leq \limsup_{N \rightarrow \infty} \sup_{\mathbb{Q} \in \mathcal{P}(\Xi): d_1(\mathbb{Q}, \hat{\mathbb{P}}_{\bar{\gamma}}^N) \leq \epsilon_N} \mathbb{E}_{\boldsymbol{\xi} \sim \mathbb{Q}}[c^{\pi^\delta}(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T)] \\
&\leq \limsup_{N \rightarrow \infty} \mathbb{E}_{\boldsymbol{\xi} \sim \hat{\mathbb{Q}}^N}[c^{\pi^\delta}(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T)] + \delta \\
&\leq \limsup_{N \rightarrow \infty} \mathbb{E}_{\boldsymbol{\xi} \sim \hat{\mathbb{Q}}^N}[f^j(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T)] + \delta \\
&\leq \limsup_{N \rightarrow \infty} \mathbb{E}_{\boldsymbol{\xi} \sim \mathbb{P}_{\bar{\gamma}}}[f^j(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T)] + L_j d_1(\mathbb{P}_{\bar{\gamma}}, \hat{\mathbb{Q}}^N) + \delta \\
&\leq \limsup_{N \rightarrow \infty} \mathbb{E}_{\boldsymbol{\xi} \sim \mathbb{P}_{\bar{\gamma}}}[f^j(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T)] + L_j(d_1(\mathbb{P}_{\bar{\gamma}}, \hat{\mathbb{P}}_{\bar{\gamma}}^N) + d_1(\hat{\mathbb{Q}}^N, \hat{\mathbb{P}}_{\bar{\gamma}}^N)) + \delta \\
&\leq \limsup_{N \rightarrow \infty} \mathbb{E}_{\boldsymbol{\xi} \sim \mathbb{P}_{\bar{\gamma}}}[f^j(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T)] + L_j(d_1(\mathbb{P}_{\bar{\gamma}}, \hat{\mathbb{P}}_{\bar{\gamma}}^N) + \epsilon_N) + \delta \\
&= \mathbb{E}_{\mathbb{P}_{\bar{\gamma}}}[f^j(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T)] + \delta, \quad \mathbb{P}^\infty\text{-almost surely,}
\end{aligned}$$

where we have used the relationship between sample robust optimization and distributionally robust optimization with the type- ∞ Wasserstein ambiguity set for the first equality (Bertsimas et al. 2018a, Section 6), the fact $d_1(\mathbb{P}, \mathbb{Q}) \leq d_\infty(\mathbb{P}, \mathbb{Q})$ for the second inequality, the dual form of the 1-Wasserstein metric for the fifth inequality (because f^j is L_j -Lipschitz), and Theorem 2 for the equality. Taking the limit as $j \rightarrow \infty$, and applying the monotone convergence theorem (which is allowed because $\mathbb{E}_{\boldsymbol{\xi} \sim \mathbb{P}_{\bar{\gamma}}}|f^1(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T)| \leq L_1 \mathbb{E}_{\boldsymbol{\xi} \sim \mathbb{P}_{\bar{\gamma}}}\|\boldsymbol{\xi}\| + |f^1(0)| < \infty$ by Assumption 4), gives

$$\limsup_{N \rightarrow \infty} \hat{v}^N(\bar{\gamma}) \leq \mathbb{E}_{\boldsymbol{\xi} \sim \mathbb{P}_{\bar{\gamma}}}[c^{\pi^\delta}(\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T)] + \delta \leq v^*(\bar{\gamma}) + 2\delta, \quad \mathbb{P}^\infty\text{-almost surely.}$$

Since $\delta > 0$ was chosen arbitrarily, the proof of (EC.6) is complete. \square

EC.3. Proof of Theorem 4

In this section, we present our proof of Theorem 4 from Section 6. We restate the theorem here for convenience.

THEOREM 4. For cost functions of the form (15), $\tilde{v}^N(\bar{\gamma}) = \hat{v}^N(\bar{\gamma})$.

Proof. We first show that $\tilde{v}^N(\bar{\gamma}) \geq \hat{v}^N(\bar{\gamma})$. Indeed, consider any primary decision rule $\bar{\pi}$ and auxiliary decision rules $\bar{\mathbf{y}}_1^i, \dots, \bar{\mathbf{y}}_T^i$ for each $i \in \{1, \dots, N\}$ which are optimal for (16).⁴ Then, it follows from feasibility to (16) that

$$\mathbf{h}_t^\top \bar{\mathbf{y}}_t^i(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_t) \geq \min_{\mathbf{y}_t \in \mathbb{R}^{d_t^i}} \left\{ \mathbf{h}_t^\top \mathbf{y}_t : \sum_{s=1}^t \mathbf{A}_{t,s} \bar{\pi}_s(\boldsymbol{\zeta}_1, \dots, \boldsymbol{\zeta}_{s-1}) + \sum_{s=1}^t \mathbf{B}_{t,s} \boldsymbol{\zeta}_s + \mathbf{C}_t \mathbf{y}_t \leq \mathbf{d}_t \right\}$$

for each $i \in \{1, \dots, N\}$, $\zeta \in \mathcal{U}_N^i$, and $t \in \{1, \dots, T\}$. Thus,

$$\begin{aligned} \hat{v}^N(\bar{\gamma}) &= \min_{\bar{\pi} \in \Pi} \sum_{i=1}^N w_N^i(\bar{\gamma}) c^{\bar{\pi}}(\zeta_1, \dots, \zeta_T) \\ &\leq \sum_{i=1}^N w_N^i(\bar{\gamma}) c^{\bar{\pi}}(\zeta_1, \dots, \zeta_T) \\ &\leq \sum_{i=1}^N w_N^i(\bar{\gamma}) \sup_{\zeta \in \mathcal{U}_N^i} \sum_{t=1}^T (\mathbf{f}_t^\top \bar{\pi}_t(\zeta_1, \dots, \zeta_{t-1}) + \mathbf{g}_t^\top \zeta_t + \mathbf{h}_t^\top \bar{\mathbf{y}}_t^i(\zeta_1, \dots, \zeta_t)) = \tilde{v}^N(\bar{\gamma}). \end{aligned}$$

The other side of the inequality follows from similar reasoning. Indeed, let $\bar{\pi}$ be an optimal solution to (4). For each $i \in \{1, \dots, N\}$ and $t \in \{1, \dots, T\}$, define $\bar{\mathbf{y}}_t^i \in \mathcal{R}_t$ as any decision rule that satisfies

$$\bar{\mathbf{y}}_t^i(\zeta_1, \dots, \zeta_t) \in \arg \min_{\mathbf{y}_t \in \mathbb{R}^{d_y^t}} \left\{ \mathbf{h}_t^\top \mathbf{y}_t : \sum_{s=1}^t \mathbf{A}_{t,s} \bar{\pi}_s(\zeta_1, \dots, \zeta_{s-1}) + \sum_{s=1}^t \mathbf{B}_{t,s} \zeta_s + \mathbf{C}_t \mathbf{y}_t \leq \mathbf{d}_t \right\}$$

for every $\zeta \in \mathcal{U}_N^i$. Then,

$$\begin{aligned} \tilde{v}^N(\bar{\gamma}) &\leq \sum_{i=1}^N w_N^i(\bar{\gamma}) \sup_{\zeta \in \mathcal{U}_N^i} \sum_{t=1}^T (\mathbf{f}_t^\top \bar{\pi}_t(\zeta_1, \dots, \zeta_{t-1}) + \mathbf{g}_t^\top \zeta_t + \mathbf{h}_t^\top \bar{\mathbf{y}}_t^i(\zeta_1, \dots, \zeta_t)) \\ &= \sum_{i=1}^N w_N^i(\bar{\gamma}) \sup_{\zeta \in \mathcal{U}_N^i} c^{\bar{\pi}}(\zeta_1, \dots, \zeta_T) = \hat{v}^N(\bar{\gamma}). \end{aligned}$$

Combining the above inequalities, the proof is complete. \square

EC.4. Tractable Reformulation of the Multi-Policy Approximation

For completeness, we now show how to reformulate the multi-policy approximation scheme with linear decision rules from Section 6 into a deterministic optimization problem using standard techniques from robust optimization.

We begin by transforming (16) with linear decision rules into a more compact representation. First, we combine the primary linear decision rules across stages as

$$\mathbf{x}_0 = \begin{bmatrix} \mathbf{x}_{1,0} \\ \vdots \\ \mathbf{x}_{T,0} \end{bmatrix} \in \mathbb{R}^{d_x}, \quad \mathbf{X} = \begin{bmatrix} \mathbf{0} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{X}_{2,1} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{X}_{3,1} & \mathbf{X}_{3,2} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ \mathbf{X}_{T-2,1} & \mathbf{X}_{T-2,2} & \mathbf{X}_{T-2,3} & \cdots & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{X}_{T-1,1} & \mathbf{X}_{T-1,2} & \mathbf{X}_{T-1,3} & \cdots & \mathbf{X}_{T-1,T-2} & \mathbf{0} & \mathbf{0} \\ \mathbf{X}_{T,1} & \mathbf{X}_{T,2} & \mathbf{X}_{T,3} & \cdots & \mathbf{X}_{T,T-2} & \mathbf{X}_{T,T-1} & \mathbf{0} \end{bmatrix} \in \mathbb{R}^{d_x \times d_\xi}.$$

We note that the zero entries in the above matrix are necessary to ensure that the linear decision rules are non-anticipative. Similarly, for each $i \in \{1, \dots, N\}$, we represent the auxiliary linear decision rules as

$$\mathbf{y}_0^i = \begin{bmatrix} \mathbf{y}_{1,0}^i \\ \vdots \\ \mathbf{y}_{T,0}^i \end{bmatrix} \in \mathbb{R}^{d_y}, \quad \mathbf{Y}^i = \begin{bmatrix} \mathbf{Y}_{1,1}^i & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ \mathbf{Y}_{2,1}^i & \mathbf{Y}_{2,2}^i & \cdots & \mathbf{0} & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{Y}_{T-1,1}^i & \mathbf{Y}_{T-1,2}^i & \cdots & \mathbf{Y}_{T-1,T-1}^i & \mathbf{0} \\ \mathbf{Y}_{T,1}^i & \mathbf{Y}_{T,2}^i & \cdots & \mathbf{Y}_{t,t-1}^i & \mathbf{Y}_{T,T}^i \end{bmatrix} \in \mathbb{R}^{d_y \times d_\xi}.$$

We now combine the problem parameters. Let $\mathbf{d} = (\mathbf{d}_1, \dots, \mathbf{d}_T) \in \mathbb{R}^m$ and

$$\begin{aligned} \mathbf{f} &= \begin{bmatrix} \mathbf{f}_1 \\ \vdots \\ \mathbf{f}_T \end{bmatrix} \in \mathbb{R}^{d_x}, & \mathbf{A} &= \begin{bmatrix} \mathbf{A}_{1,1} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ \mathbf{A}_{2,1} & \mathbf{A}_{2,2} & \cdots & \mathbf{0} & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{A}_{T-1,1} & \mathbf{A}_{T-1,2} & \cdots & \mathbf{A}_{T-1,T-1} & \mathbf{0} \\ \mathbf{A}_{T,1} & \mathbf{A}_{T,2} & \cdots & \mathbf{A}_{t,t-1} & \mathbf{A}_{T,T} \end{bmatrix} \in \mathbb{R}^{m \times d_x}, \\ \mathbf{g} &= \begin{bmatrix} \mathbf{g}_1 \\ \vdots \\ \mathbf{g}_T \end{bmatrix} \in \mathbb{R}^{d_\xi}, & \mathbf{B} &= \begin{bmatrix} \mathbf{B}_{1,1} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ \mathbf{B}_{2,1} & \mathbf{B}_{2,2} & \cdots & \mathbf{0} & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{B}_{T-1,1} & \mathbf{B}_{T-1,2} & \cdots & \mathbf{B}_{T-1,T-1} & \mathbf{0} \\ \mathbf{B}_{T,1} & \mathbf{B}_{T,2} & \cdots & \mathbf{B}_{t,t-1} & \mathbf{B}_{T,T} \end{bmatrix} \in \mathbb{R}^{m \times d_x}, \\ \mathbf{h} &= \begin{bmatrix} \mathbf{h}_1 \\ \vdots \\ \mathbf{h}_T \end{bmatrix} \in \mathbb{R}^{d_y}, & \mathbf{C} &= \begin{bmatrix} \mathbf{C}_{1,1} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_{2,2} & \cdots & \mathbf{0} & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{C}_{T-1,T-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} & \mathbf{C}_{T,T} \end{bmatrix} \in \mathbb{R}^{m \times d_x}. \end{aligned}$$

Therefore, using the above compact notation, we can rewrite the multi-policy approximation with linear decision rules as

$$\begin{aligned} & \underset{\substack{\mathbf{x}_0 \in \mathbb{R}^{d_x}, \mathbf{X} \in \mathbb{R}^{d_x \times d_\xi} \\ \mathbf{y}_0^i \in \mathbb{R}^{d_y}, \mathbf{Y}^i \in \mathbb{R}^{d_y \times d_\xi}}}{\text{minimize}} \sum_{i=1}^N w_N^i(\bar{\gamma}) \sup_{\zeta \in \mathcal{U}_N^i} \{ \mathbf{f}^\top (\mathbf{x}_0 + \mathbf{X}\zeta) + \mathbf{g}^\top \zeta + \mathbf{h}^\top (\mathbf{y}_0^i + \mathbf{Y}^i \zeta) \} \\ & \text{subject to} \quad \mathbf{A}(\mathbf{x}_0 + \mathbf{X}\zeta) + \mathbf{B}\zeta + \mathbf{C}(\mathbf{y}_0^i + \mathbf{Y}^i \zeta) \leq \mathbf{d} \tag{EC.7} \\ & \quad \mathbf{x}_0 + \mathbf{X}\zeta \in \mathcal{X} \\ & \quad \forall \zeta \in \mathcal{U}_N^i, i \in \{1, \dots, N\}, \end{aligned}$$

where $\mathcal{X} \triangleq \mathcal{X}_1 \times \cdots \times \mathcal{X}_T$ and the matrices \mathbf{X} and \mathbf{Y} are non-anticipative. Note that the linear decision rules in the above optimization problem are represented using $O(d_\xi \max\{d_x, Nd_y\})$ decision variables, where $d_x \triangleq d_x^1 + \cdots + d_x^T$ and $d_y \triangleq d_y^1 + \cdots + d_y^T$. Thus, the complexity of representing the primary and auxiliary linear decision rules scales efficiently both in the size of the dataset and the number of stages. For simplicity, we present the reformulation for the case in which there are no constraints on the decision variables and nonnegativity constraints on the random variables.

THEOREM EC.2. Suppose $\Xi = \mathbb{R}_+^{d_\xi}$ and $\mathcal{X} = \mathbb{R}^{d_x}$. Then, (EC.7) is equivalent to

$$\begin{aligned} & \underset{\substack{\mathbf{x}_0 \in \mathbb{R}^{d_x}, \mathbf{X} \in \mathbb{R}^{d_x \times d_\xi} \\ \mathbf{y}_0^i \in \mathbb{R}^{d_y}, \mathbf{Y}^i \in \mathbb{R}^{d_y \times d_\xi} \\ \Lambda^i \in \mathbb{R}_+^{m \times d_\xi}, \mathbf{s}^i \in \mathbb{R}_+^{d_\xi}}}{\text{minimize}} & \sum_{i=1}^N w_N^i(\bar{\gamma}) \left(\mathbf{f}^\top (\mathbf{x}_0 + \mathbf{X}\boldsymbol{\xi}^i) + \mathbf{g}^\top \boldsymbol{\xi}^i + \mathbf{h}^\top (\mathbf{y}_0^i + \mathbf{Y}^i \boldsymbol{\xi}^i) + (\mathbf{s}^i)^\top \boldsymbol{\xi}^i + \epsilon_N \|\mathbf{X}^\top \mathbf{f} + \mathbf{g} + (\mathbf{Y}^i)^\top \mathbf{h} + \mathbf{s}^i\|_* \right) \\ \text{subject to} & \quad \mathbf{A} (\mathbf{x}_0 + \mathbf{X}\boldsymbol{\xi}^i) + \mathbf{B}\boldsymbol{\xi}^i + \mathbf{C} (\mathbf{y}_0^i + \mathbf{Y}^i \boldsymbol{\xi}^i) + \Lambda^i \boldsymbol{\xi}^i + \epsilon_N \|\mathbf{A}\mathbf{X} + \mathbf{B} + \mathbf{C}\mathbf{Y}^i + \Lambda^i\|_* \leq \mathbf{d} \\ & \quad \forall i \in \{1, \dots, N\}. \end{aligned}$$

where $\|\mathbf{Z}\|_* \triangleq (\|\mathbf{z}_1\|_*, \dots, \|\mathbf{z}_r\|_*) \in \mathbb{R}^r$ for any matrix $\mathbf{Z} \in \mathbb{R}^{r \times n}$.

Proof. For any $\mathbf{c} \in \mathbb{R}^{d_\xi}$ and $\boldsymbol{\xi} \in \Xi$, it follows directly from strong duality for conic optimization that

$$\max_{\boldsymbol{\zeta} \geq \mathbf{0}} \{ \mathbf{c}^\top \boldsymbol{\zeta} : \|\boldsymbol{\zeta} - \boldsymbol{\xi}\| \leq \epsilon \} = \min_{\boldsymbol{\lambda} \geq \mathbf{0}} \{ (\mathbf{c} + \boldsymbol{\lambda})^\top \boldsymbol{\xi} + \epsilon \|\mathbf{c} + \boldsymbol{\lambda}\|_* \}.$$

We use this result to reformulate the objective and constraints of (EC.7). First, let the j -th rows of $\mathbf{A}, \mathbf{B}, \mathbf{C}$ and the j -th element of \mathbf{d} be denoted by $\mathbf{a}_j \in \mathbb{R}^{d_x}$, $\mathbf{b}_j \in \mathbb{R}^{d_\xi}$, $\mathbf{c}_j \in \mathbb{R}^{d_y}$, and $d_j \in \mathbb{R}$. Then, each robust constraint has the form

$$\mathbf{a}_j^\top (\mathbf{x}_0 + \mathbf{X}\boldsymbol{\zeta}) + \mathbf{b}_j^\top \boldsymbol{\zeta} + \mathbf{c}_j^\top (\mathbf{y}_0^i + \mathbf{Y}^i \boldsymbol{\zeta}) \leq d_j \quad \forall \boldsymbol{\zeta} \in \mathcal{U}_N^i.$$

Rearranging terms,

$$(\mathbf{a}_j^\top \mathbf{X} + \mathbf{b}_j^\top + \mathbf{c}_j^\top \mathbf{Y}^i) \boldsymbol{\zeta} \leq d_j - \mathbf{a}_j^\top \mathbf{x}_0 - \mathbf{c}_j^\top \mathbf{y}_0^i \quad \forall \boldsymbol{\zeta} \in \mathcal{U}_N^i,$$

which applying duality becomes

$$\exists \boldsymbol{\lambda}_j^i \geq \mathbf{0} : (\mathbf{X}^\top \mathbf{a}_j + \mathbf{b}_j + (\mathbf{Y}^i)^\top \mathbf{c}_j + \boldsymbol{\lambda}_j^i)^\top \boldsymbol{\xi}^i + \epsilon_N \|\mathbf{X}^\top \mathbf{a}_j + \mathbf{b}_j + (\mathbf{Y}^i)^\top \mathbf{c}_j + \boldsymbol{\lambda}_j^i\|_* \leq d_j - \mathbf{a}_j^\top \mathbf{x}_0 - \mathbf{c}_j^\top \mathbf{y}_0^i.$$

Rearranging terms, the robust constraints for each $i \in \{1, \dots, N\}$ are satisfied if and only if

$$\exists \Lambda^i \geq \mathbf{0} : \mathbf{A} (\mathbf{x}_0 + \mathbf{X}\boldsymbol{\xi}^i) + \mathbf{B}\boldsymbol{\xi}^i + \mathbf{C} (\mathbf{y}_0^i + \mathbf{Y}^i \boldsymbol{\xi}^i) + \Lambda^i \boldsymbol{\xi}^i + \epsilon_N \|\mathbf{A}\mathbf{X} + \mathbf{B} + \mathbf{C}\mathbf{Y}^i + \Lambda^i\|_* \leq \mathbf{d},$$

where the dual norm for a matrix is applied separately for each row. Similarly, the objective function takes the form

$$\begin{aligned} & \sum_{i=1}^N w_N^i(\bar{\gamma}) \sup_{\boldsymbol{\zeta} \in \mathcal{U}_N^i} \{ \mathbf{f}^\top (\mathbf{x}_0 + \mathbf{X}\boldsymbol{\zeta}) + \mathbf{g}^\top \boldsymbol{\zeta} + \mathbf{h}^\top (\mathbf{y}_0^i + \mathbf{Y}^i \boldsymbol{\zeta}) \} \\ &= \sum_{i=1}^N w_N^i(\bar{\gamma}) \left(\mathbf{f}^\top \mathbf{x}_0 + \mathbf{h}^\top \mathbf{y}_0^i + \sup_{\boldsymbol{\zeta} \in \mathcal{U}_N^i} (\mathbf{f}^\top \mathbf{X} + \mathbf{g}^\top + \mathbf{h}^\top \mathbf{Y}^i) \boldsymbol{\zeta} \right) \\ &= \sum_{i=1}^N w_N^i(\bar{\gamma}) \left(\mathbf{f}^\top \mathbf{x}_0 + \mathbf{h}^\top \mathbf{y}_0^i + \inf_{\mathbf{s}^i \geq \mathbf{0}} \{ (\mathbf{X}^\top \mathbf{f} + \mathbf{g} + (\mathbf{Y}^i)^\top \mathbf{h} + \mathbf{s}^i)^\top \boldsymbol{\xi}^i + \epsilon_N \|\mathbf{X}^\top \mathbf{f} + \mathbf{g} + (\mathbf{Y}^i)^\top \mathbf{h} + \mathbf{s}^i\|_* \} \right) \\ &= \sum_{i=1}^N w_N^i(\bar{\gamma}) \left(\mathbf{f}^\top (\mathbf{x}_0 + \mathbf{X}\boldsymbol{\xi}^i) + \mathbf{g}^\top \boldsymbol{\xi}^i + \mathbf{h}^\top (\mathbf{y}_0^i + \mathbf{Y}^i \boldsymbol{\xi}^i) + \inf_{\mathbf{s}^i \geq \mathbf{0}} \{ (\mathbf{s}^i)^\top \boldsymbol{\xi}^i + \epsilon_N \|\mathbf{X}^\top \mathbf{f} + \mathbf{g} + (\mathbf{Y}^i)^\top \mathbf{h} + \mathbf{s}^i\|_* \} \right). \end{aligned}$$

Combining the reformulations above, we obtain the desired reformulation. \square