
Data-Driven Stochastic Dual Dynamic Programming: Performance Guarantees and Regularization Schemes

Hyuk Park · Zhuangzhuang Jia · Grani A. Hanasusanto

Received: date / Accepted: date

Abstract We propose a data-driven scheme for multistage stochastic linear programming with Markovian random parameters by extending the stochastic dual dynamic programming (SDDP) algorithm. In our data-driven setting, only a finite number of historical trajectories are available. The proposed SDDP scheme evaluates the cost-to-go functions only at the observed sample points, where the conditional expectations are estimated empirically using kernel regression. The scheme thus avoids the construction of scenario trees, which may incur exponential time complexity during the backward induction step. However, if the training data is sparse, the resulting SDDP algorithm exhibits a high optimistic bias that gives rise to poor out-of-sample performances. To mitigate the small sample effects, we adopt ideas from the distributionally robust optimization (DRO), which replaces the empirical conditional expectation in the cost-to-go function with a worst-case conditional expectation over a polyhedral ambiguity set. We derive the theoretical out-of-sample performance guarantee of the data-driven SDDP scheme and show that the dependence of its sample complexity on the number of time stages is merely polynomial. Finally, we validate the effectiveness of the proposed algorithm and demonstrate its superiority over existing data-driven schemes through extensive numerical experiments.

Keywords stochastic optimization · multistage stochastic linear program · stochastic dual dynamic programming · distributionally robust optimization · Nadaraya-Watson estimator · Markov dependence

1 Introduction

Multistage stochastic linear programming (MSLP) is a framework for sequential decision making under uncertainty with linear objective functions and constraints. At each stage, the decision is made based on realizations of random parameters available up to that stage. This framework has been adopted to model various real-world applications, such as hydrothermal scheduling [17, 45, 46, 60], unit commitment [4, 14, 29, 52], portfolio optimization [11, 12, 16, 27, 42], and manufacturing and capacity planning [2, 3, 23, 62]. Despite such exceptional modeling power, the general consensus is that multistage stochastic programs are very difficult to solve [36, 59].

One popular approach to solve MSLP problems is a sampling-based Benders decomposition method called stochastic dual dynamic programming (SDDP). It was developed by Pereira and Pinto in 1991 as an effort to solve a large-scale hydrothermal scheduling problem [46] and is still considered as the state-of-the-art solution method. This algorithm iteratively approximates the expected (i.e., risk-neutral) cost-to-go functions of the dynamic programming equations by piecewise linear functions called *cuts* during the *forward* and *backward* steps. Since such cuts allow the cost-to-go function to be evaluated at well-chosen candidate solutions, the intractability arising from discretizing the state variables, known as the curse of dimensionality, can be alleviated. Asymptotic convergence of the algorithm is established in [15, 25, 50, 55],

Hyuk Park
Department of Industrial and Enterprise Systems Engineering, University of Illinois Urbana-Champaign, Urbana, IL, 61801.
E-mail: hyukp2@illinois.edu

Zhuangzhuang Jia
Department of Industrial and Enterprise Systems Engineering, University of Illinois Urbana-Champaign, Urbana, IL, 61801.
E-mail: zj12@illinois.edu

Grani A. Hanasusanto
Department of Industrial and Enterprise Systems Engineering, University of Illinois Urbana-Champaign, Urbana, IL, 61801.
E-mail: gah@illinois.edu

while its computational complexity is recently investigated in [37,68]. The standard SDDP framework is extended to risk-averse cases [55,60], and convex [66,67] and nonconvex cost-to-go functions [1,20,69,70].

In the literature on the SDDP framework, it is commonly assumed that the underlying distribution of the uncertain data is known and the data process is stagewise independent. While these assumptions are not entirely realistic, they are crucial for several reasons. First, a known distribution assumption enables to sample scenarios in the forward step and construct a finite representation of the true problem by approximating the underlying stochastic process using a scenario tree [53]. Second, under the stagewise independent data process, we can avoid a scenario tree that grows exponentially with the number of time stages as it collapses to a recombining scenario tree, where nodes between any consecutive stages are fully connected via unconditional transition probabilities. Such a structure allows us to evaluate only one unconditional expected cost-to-go at each stage. In reality, however, there is often a correlation over time in the data process. Under stagewise dependence, the aforementioned approach is no longer valid as we have to take a new set of samples conditional on the whole history of the data process up to the current stage, which makes evaluating the conditional expected cost-to-go functions intractable [?,?,56]. To alleviate such an intractability, several papers instead assume Markov dependence (i.e., interstage dependence) with the purpose of exploiting some degree of correlation in the data process.

In a data-driven setting, where only historical data is available while the true distribution is unknown, there are two approaches to incorporate Markov dependence. One is to reformulate the random parameter as a first-order time series model with additional state variables and assume that the noise term is stagewise independent [32]. However, the time series approach requires adding auxiliary variables to state variables, which aggravates the curse of dimensionality. In addition, from a modeling perspective, this approach is limited in that uncertainty should appear only on the right-hand side of the constraints. Otherwise, the convexity of the cost-to-go is destroyed. A more general approach, known as Markov chain discretization, is to construct a recombining scenario tree and approximate the true conditional distribution using optimal quantization to reflect Markov dependence [10,13,26,38,39,48]. Unlike the time series approach, the Markov chain approach does not require the auxiliary variables and allows uncertainty in any parameters, e.g., objective function coefficients, a recourse matrix, etc, and therefore it can model a broader range of problems, such as portfolio optimization and inventory management. Löhndorf and Shapiro [38] demonstrate that the SDDP based on the Markov chain approach provides tighter lower bounds and improved policies for the hydrothermal scheduling problem compared to the time series approach. Unfortunately, despite the better empirical performance shown in the previous works, there has been no rigorous sample complexity analysis of this method, i.e., it is not clear whether the Markov chain discretization can generalize to out-of-sample data and provide asymptotically consistent policies.

Whichever approach is adopted, such data-driven schemes often suffer from overfitting issues, causing poor out-of-sample performance. The issue is aggravated particularly when the available data is limited. As a remedy for this small sample effect, distributionally robust optimization (DRO) has garnered significant attention recently. DRO relaxes the stringent assumption of known distribution by constructing an ambiguity set of plausible distributions consistent with the available information. Using DRO, an MSLP problem is formulated as a min-max problem at each stage, which yields a policy that performs best under the worst-case distribution that maximizes the expected cost-to-go. Ambiguity sets are commonly categorized into two types: moment-based [18,51,57,71] and discrepancy-based ambiguity sets [5,8,35,40,41,65]. Particularly in the distributionally robust MSLP setting, SDDP frameworks with different discrepancy-based ambiguity sets have been proposed [22,31,49,61]. Similar to this paper, Philpott et al. [49] use the modified χ^2 ambiguity set [7]. They derive a closed-form solution for the inner maximization problem and use the resulting worst-case distribution to generate a cut during the backward step. They assume that the data process is stagewise independent and impose randomness only on the right-hand side. Silva et al. [61] consider a more general distributionally robust MSLP problem where the random recourse matrix has Markov dependence. They use a hidden Markov model (HMM) to capture unobservable states characterizing the distribution of the random parameters (in the recourse matrix) and then use the total variation ambiguity set to account for estimation errors of the transition probabilities among the unobservable states. In Table 1, we compare our proposed method with the existing SDDP algorithms for the distributionally robust MSLP.

This paper focuses on the incorporation of Markov dependence into risk-neutral and risk-averse MSLP problems in a data-driven setting and the sample complexity analysis of the proposed Markov discretization method. From the modeling aspect, our work is similar to [29] where the authors use the Nadaraya-Watson (NW) kernel regression estimator [43,64] for stochastic optimal control problems with endogenous state variables. However, as their stochastic dynamic programming scheme requires discretizing the endogenous state variables, it only works for low-dimensional settings. To deal with large-scale MSLP problems, we propose a data-driven SDDP framework using the NW regression under Markov dependence. Srivastava et al. [63] derive generalization bounds for the NW approximation of a static stochastic program. From a theoretical aspect, therefore, our work extends the out-of-sample guarantee derived in [63] to multistage settings.

Table 1 Comparison of distributionally robust SDDP algorithm variants

Algorithms	Huang et al., 2017 [31]	Philpott et al., 2018 [49]	Duque and Morton, 2020 [22]	Silva et al., 2021 [61]	This Paper
Ambiguity Set	∞ -norm	modified χ^2	Wasserstein	total variation	polyhedral approximation of modified χ^2
Random RHS	✓	✓	✓	✓	✓
Random Technology Matrix	✓	-	-	✓	✓
Random Recourse Matrix	✓	-	-	✓	✓
Random Obj. Func. Coeff.	✓	-	-	✓	✓
Markov Dependence	-	-	-	✓	✓
Stopping Criterion	stochastic gap	maximum iteration	maximum iteration	deterministic gap	deterministic gap
Out-of-Sample Guarantee	-	-	-	-	✓

The main contributions of this paper can be summarized as follows:

1. We propose a data-driven SDDP framework for risk-neutral and risk-averse MSLP problems under the Markov dependence assumption. To incorporate Markov dependence, the true conditional probability is estimated using the Nadaraya-Watson kernel regression. Unlike other Markov chain discretization approaches that do not have convergence guarantees, the proposed discretization provides convergence to the true optimal policy as the number of data tends to infinity.
2. Leveraging the generalization bounds for a static stochastic program established in [63], we derive for the first time a theoretical out-of-sample guarantee for the data-driven risk-neutral MSLP problem under Markov dependence. Our result indicates that the out-of-sample suboptimality bound is at most $\tilde{O}_p(T^{\frac{3}{2}}/N^{\frac{2}{p+4}})$,² where T is the number of time stages. This mild (polynomial) dependence on T suggests that our scheme is applicable to solve MSLP problems with a large number of stages. The result is surprising since existing analysis [59] suggests that the theoretical suboptimality bound drastically worsens as the number of stages increases.
3. Our theoretical guarantee suggests the use of a variance-based regularization scheme for improving out-of-sample performance. Unfortunately, the scheme is intractable due to nonconvexity. As a tractable alternative, we develop a conservative MSLP problem using a DRO formulation with a polyhedral outer approximation of the modified χ^2 ambiguity set. We further prove that the DRO formulation asymptotically converges to the variance regularization scheme.
4. Numerical experiments in the context of portfolio optimization and wind energy commitment problems demonstrate that our data-driven schemes are superior to the stagewise independent scheme and the benchmark schemes proposed in the literature in terms of out-of-sample performance.

Notation and terminology. We use bold letters for vectors and regular fonts for scalars. We define \mathbf{e} as the vector of all ones—its dimension will be clear from the context. The tilde symbol is used to denote random variables (e.g., $\tilde{\xi}_t$) to differentiate from their realizations (e.g., ξ_t). For any $t \in \mathbb{N}$, we define $[t]$ as the index set $\{1, \dots, t\}$. A sequence of realizations of the random data process up to stage t is denoted as $\xi_{[t]} = (\xi_1, \dots, \xi_t)$. For a random variable \tilde{Z} , $\mathbb{E}[\tilde{Z}]$ and $\mathbb{V}[\tilde{Z}]$ denote its expectation and variance, respectively. The indicator function of a subset \mathcal{X} is defined through $\mathbb{1}_{\mathcal{X}}(\mathbf{x}) = 0$ if $\mathbf{x} \in \mathcal{X}$ and $\mathbb{1}_{\mathcal{X}}(\mathbf{x}) = +\infty$ otherwise. The Dirac distribution, which assigns unit mass on ξ_t , is denoted by δ_{ξ_t} . In asymptotic analysis, we use the standard O notations for the convergence of sets of ordinary numbers and O_p for the convergence of sets of random variables. In addition, \tilde{O}_p is used to suppress multiplicative terms with logarithmic dependence on n in the sense of convergence in probability.

2 Problem Statement

Consider the following risk-neutral MSLP problem under *Markov dependence*:

$$\min_{\mathbf{x}_1 \in \mathcal{X}_1(\mathbf{x}_0, \xi_1)} \mathbf{c}_1^\top \mathbf{x}_1 + \mathbb{E} \left[\min_{\mathbf{x}_2 \in \mathcal{X}_2(\mathbf{x}_1, \xi_2)} \tilde{\mathbf{c}}_2^\top \mathbf{x}_2 + \mathbb{E} \left[\dots + \mathbb{E} \left[\min_{\mathbf{x}_T \in \mathcal{X}_T(\mathbf{x}_{T-1}, \tilde{\xi}_T)} \tilde{\mathbf{c}}_T^\top \mathbf{x}_T \mid \tilde{\xi}_{T-1} \right] \dots \mid \tilde{\xi}_2 \right] \mid \xi_1 \right], \quad (2.1)$$

Here, the problem parameters are summarized by the random vectors $\tilde{\xi}_t = (\tilde{\mathbf{c}}_t, \tilde{\mathbf{b}}_t, \tilde{\mathbf{A}}_t, \tilde{\mathbf{B}}_t)$ for every $t \in [T] \setminus \{1\}$ governed by an (unknown) continuous joint distribution while the first stage parameters $\xi_1 = (\mathbf{c}_1, \mathbf{b}_1, \mathbf{A}_1, \mathbf{B}_1)$ are deterministic with the initial state \mathbf{x}_0 is given as input. Here, the decision space is defined as the polytope $\mathcal{X}_t(\mathbf{x}_{t-1}, \tilde{\xi}_t) = \{\mathbf{x}_t \in \mathbb{R}_+^{d_t} : \tilde{\mathbf{A}}_t \mathbf{x}_t + \tilde{\mathbf{B}}_t \mathbf{x}_{t-1} = \tilde{\mathbf{b}}_t\}$ and the risk measures are given by conditional expectations $\mathbb{E}[\cdot \mid \xi_t]$. In this setting, the decision maker takes sequential decisions in a

² N denotes the number of samples and p is the dimension of random parameters in an MSLP problem.

nonanticipative manner as the realization of the random parameter vector $\tilde{\boldsymbol{\xi}}_t = (\tilde{\mathbf{c}}_t, \tilde{\mathbf{b}}_t, \tilde{\mathbf{A}}_t, \tilde{\mathbf{B}}_t)$ is revealed at each time period. More precisely, when the decision maker takes the first-stage decision \mathbf{x}_1 from the deterministic feasible region $\mathcal{X}_1(\mathbf{x}_0, \boldsymbol{\xi}_1)$, the decision incurs an immediate cost $\mathbf{c}_1^\top \mathbf{x}_1$. Next, at the beginning of the second stage, the realization $\boldsymbol{\xi}_2$ of the random parameter vector $\tilde{\boldsymbol{\xi}}_2$ is revealed and the second-stage decision \mathbf{x}_2 is chosen from the feasible region $\mathcal{X}_2(\mathbf{x}_1, \boldsymbol{\xi}_2)$ with a cost $\mathbf{c}_2^\top \mathbf{x}_2$. The process continues until the terminal stage T . The goal of the problem is to optimize a sequence of functions $\{\mathbf{x}_t(\boldsymbol{\xi}_{[t]})\}_{t=1}^T$ called policies that minimize the expected total cost over the T stages. Here, each $\mathbf{x}_t(\boldsymbol{\xi}_{[t]})$ is a function of the realizations $\boldsymbol{\xi}_{[t]} = (\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_t)$ of random parameter vectors up to stage t . In our setting, we refer to (2.1) as the *true* problem where the random process $\tilde{\boldsymbol{\xi}}_{[T]}$ (i.e., the evolution of the random parameter vector) is governed by a continuous state Markov process. Using the Bellman's optimality principle [6], the true problem (2.1) at stage $t \in [T]$ can equivalently be expressed as the dynamic program

$$\begin{aligned} Q_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) &= \min \mathbf{c}_t^\top \mathbf{x}_t + Q_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_t) \\ \text{s.t. } \mathbf{x}_t &\in \mathbb{R}_+^d, \\ \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{x}_{t-1} &= \mathbf{b}_t, \end{aligned} \quad (2.2)$$

where $Q_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_t) = \mathbb{E}[Q_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t]$ is the expected cost-to-go function. Here, for $t \in [T]$, the cost-to-go function $Q_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$ in (2.2) represents the minimum future expected cost accrued from stage t to the terminal stage T given that the previous decision \mathbf{x}_{t-1} is taken and the current stage parameter vector $\boldsymbol{\xi}_t$ is realized. For simplicity, we assume the terminal cost $Q_{T+1}(\cdot) \equiv 0$. In the remainder of this paper, we focus on the dynamic programming equation since the solution scheme proposed in Section 3 is based on a cutting plane approximation of the expected cost-to-go functions in (2.2).

2.1 Assumptions

In the previous works studying MSLP, it is commonly assumed that i) the probability distribution of random parameters at each stage is known and ii) the data process is stagewise independent. In this paper, we depart from those assumptions with the purpose of developing a practically meaningful data-driven scheme for the MSLP problem (2.2). We make the following assumptions throughout the main paper:

- (A1) Markovian Random Parameters.** The first stage parameters $\boldsymbol{\xi}_1 = (\mathbf{c}_1, \mathbf{b}_1, \mathbf{A}_1, \mathbf{B}_1)$ are deterministic, whereas the subsequent stage parameters $\tilde{\boldsymbol{\xi}}_t = (\tilde{\mathbf{c}}_t, \tilde{\mathbf{b}}_t, \tilde{\mathbf{A}}_t, \tilde{\mathbf{B}}_t)$ for $t \in [T] \setminus \{1\}$ are stochastic with support $\Xi_t \subset \mathbb{R}^p$. Furthermore, the stochastic process $\boldsymbol{\xi}_{[T]} = (\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_T)$ satisfies the Markov property. That is, the conditional distribution of $\boldsymbol{\xi}_{t+1}$ given the data process $\boldsymbol{\xi}_{[t]} = (\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_t)$ up to time t does not depend on $\boldsymbol{\xi}_{[t-1]} = (\boldsymbol{\xi}_1, \dots, \boldsymbol{\xi}_{t-1})$.
- (A2) Unknown Distribution.** The true joint distribution of $\tilde{\boldsymbol{\xi}}_{[T]} = (\tilde{\boldsymbol{\xi}}_1, \tilde{\boldsymbol{\xi}}_2, \dots, \tilde{\boldsymbol{\xi}}_T)$ is unknown and only N independent and identically distributed (i.i.d.) historical trajectories $\boldsymbol{\xi}_{[T]}^i = (\boldsymbol{\xi}_1^i, \boldsymbol{\xi}_2^i, \dots, \boldsymbol{\xi}_T^i)$, $i \in [N]$, from the stochastic process are available. In this paper, we denote the empirical uncertainty sets by $\hat{\Xi}_1 = \{\boldsymbol{\xi}_1^i\}$ and $\hat{\Xi}_t = \{\boldsymbol{\xi}_t^i : i \in [N]\}$ for every $t \in [T] \setminus \{1\}$.
- (A3) Relatively Complete Recourse and Compactness.** For each stage $t \in [T]$, the feasible region $\mathcal{X}_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$ is nonempty and compact for any feasible \mathbf{x}_{t-1} and any $\boldsymbol{\xi}_t \in \Xi_t$.

2.2 Discretized Problem

It is impossible to solve the true problem (2.2) exactly under assumptions (A1) and (A2) since the expected cost-to-go functions $Q_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_t)$ for every $t \in [T-1]$ cannot be evaluated without full knowledge of the underlying (conditional) distribution. Even if we knew the true distribution, evaluating the multivariate (conditional) expectation and optimizing over the continuum of the state variables \mathbf{x}_t make the true problem intractable to solve. To address these fundamental challenges, in this paper, we propose to estimate the true conditional probability using the historical data via the Nadaraya-Watson (NW) kernel regression [43, 64] defined as follows:

$$\hat{w}_t^i(\boldsymbol{\xi}_t) = \frac{\mathcal{K}\left(\frac{\boldsymbol{\xi}_t - \boldsymbol{\xi}_t^i}{h}\right)}{\sum_{j \in [N]} \mathcal{K}\left(\frac{\boldsymbol{\xi}_t - \boldsymbol{\xi}_t^j}{h}\right)} \quad \forall i \in [N]. \quad (2.3)$$

Here, \mathcal{K} is a kernel function of choice and $h > 0$ is a smoothing parameter called bandwidth. In this paper, we use the exponential kernel function of the form $\mathcal{K}(\boldsymbol{\rho}) = \exp(-\|\boldsymbol{\rho}\|_2)$. The bandwidth parameter h controls the smoothness of the estimator (2.3). A too small h leads to undersmoothing, meaning that most weights are assigned to points close to $\boldsymbol{\xi}_t$. On the other hand, an extremely large h reduces the weights $\hat{w}_t^i(\boldsymbol{\xi}_t)$ to

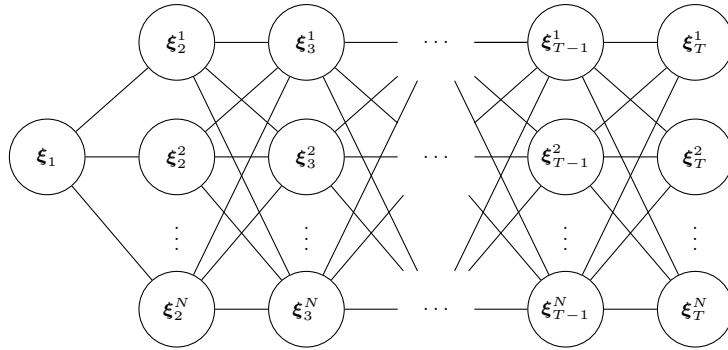


Fig. 1 Recombining scenario tree using N historical trajectories.

$1/N$, $\forall i \in [N]$. Györfi et al. [28] establish that the bandwidth scaling $h = O(1/N^{\frac{1}{p+4}})$ minimizes the expected estimation errors, which we adopt throughout the paper.

In our data-driven setting, the sample space of the true random process is discretized by i.i.d. historical trajectories $\{\xi_t^i\}_{t=1}^T$ for every $i \in [N]$, giving rise to a recombining scenario tree depicted in Figure 1. Here, the true conditional distribution of ξ_{t+1}^j given ξ_t^i is approximated by the NW estimator (2.3). We refer to the MSLP problem represented by Figure 1 as the *discretized problem*. Similar to the true problem (2.2), using the dynamic programming equation, we present the discretized problem for stage $t \in [T]$ and each $\xi_t \in \hat{\Xi}_t$ as

$$\begin{aligned} \hat{Q}_t(\mathbf{x}_{t-1}, \xi_t) = \min & \mathbf{c}_t^\top \mathbf{x}_t + \hat{Q}_{t+1}(\mathbf{x}_t, \xi_t) \\ \text{s.t. } & \mathbf{x}_t \in \mathbb{R}_+^{d_t}, \\ & \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{x}_{t-1} = \mathbf{b}_t, \end{aligned} \quad (2.4)$$

and $\hat{Q}_{T+1}(\cdot) \equiv 0$. Here, the approximate expected cost-to-go function $\hat{Q}_{t+1}(\mathbf{x}_t, \xi_t)$ is defined as

$$\hat{Q}_{t+1}(\mathbf{x}_t, \xi_t) = \hat{\mathbb{E}} \left[\hat{Q}_{t+1}(\mathbf{x}_t, \tilde{\xi}_{t+1}) \mid \xi_t \right] = \sum_{i \in [N]} \hat{w}_t^i(\xi_t) \cdot \hat{Q}_{t+1}(\mathbf{x}_t, \xi_{t+1}^i). \quad (2.5)$$

Note that the approximate conditional expectation $\hat{Q}_{t+1}(\mathbf{x}_t, \xi_t)$ is taken with respect to N data points, using the NW estimator in (2.3) and this weighted sum replaces the true expected cost-to-go function $Q_{t+1}(\mathbf{x}_t, \xi_t)$ in (2.2). This implies that the results from solving the discretized problem are not necessarily valid for the true problem. In addition, even with this drastic simplification, the discretized problem remains hard since the approximate cost-to-go functions still need to be evaluated for every $\mathbf{x}_t \in \mathcal{X}_t(\mathbf{x}_{t-1}, \xi_t)$, which makes the discretized problem intractable. In Section 3, we discuss the stochastic dual dynamic programming algorithm that can deal with this kind of intractability.

2.3 Out-of-Sample Performance Guarantee

Data-driven solutions obtained from solving the discretized problem (2.4) can be suboptimal for the true problem due to the estimation errors from the NW estimator. In this paper, we are interested in understanding how well the data-driven scheme approximates the true problem and provides a reasonably good solution with a high probability.

Srivastava et al. [63] investigate the out-of-sample guarantee for the static stochastic problem with side information. In this section, we extend their result to the multi-stage setting. We state the following mild regularity conditions required for the generalization bound.

- (A4) **Compact Uncertainty Set.** For each stage $t \in [T]$, the random parameter vector $\tilde{\xi}_t$ is supported on a compact set Ξ_t .
- (A5) **Differentiability.** For each stage $t \in [T]$, the joint density function $f(\xi_t, \xi_{t-1})$ is twice differentiable with continuous and bounded partial derivatives and the marginal density $f(\xi_t \mid \xi_{t-1})$ is non-zero for every $\xi_{t-1} \in \hat{\Xi}_{t-1}$.
- (A6) **Bandwidth.** The bandwidth parameter h for the kernel function \mathcal{K} is scaled such that $\lim_{N \rightarrow \infty} h_N = 0$ and $\lim_{N \rightarrow \infty} N h_N^p = \infty$.

The assumptions (A4)-(A6) are common in kernel regression estimation. The condition about the marginal density in the assumption (A5) ensures that the conditional probability is always positive throughout the time horizon. The condition about the bandwidth h in the assumption (A6) guarantees that the NW

estimates asymptotically converge to the true conditional distribution. Moreover, in view of the assumption (A4), we can establish the following Lipschitz continuity condition for the cost-to-go functions.

Lemma 1 (Lipschitz Cost-to-go Function) *For each stage $t \in [T]$, the true cost-to-go function $Q_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$ and the approximate cost-to-go function $\hat{Q}_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$ are L_t -Lipschitz continuous in \mathbf{x}_{t-1} . That is, there exists a constant $L_t > 0$ such that*

$$\begin{aligned} |Q_t(\mathbf{x}, \boldsymbol{\xi}_t) - Q_t(\mathbf{x}', \boldsymbol{\xi}_t)| &\leq L_t \|\mathbf{x} - \mathbf{x}'\| & \forall \mathbf{x}, \mathbf{x}' \in \mathcal{X}_{t-1}(\mathbf{x}_{t-2}, \boldsymbol{\xi}_{t-1}) \quad \forall \boldsymbol{\xi}_t \in \Xi_t, \\ |\hat{Q}_t(\mathbf{x}, \boldsymbol{\xi}_t) - \hat{Q}_t(\mathbf{x}', \boldsymbol{\xi}_t)| &\leq L_t \|\mathbf{x} - \mathbf{x}'\| & \forall \mathbf{x}, \mathbf{x}' \in \mathcal{X}_{t-1}(\mathbf{x}_{t-2}, \boldsymbol{\xi}_{t-1}) \quad \forall \boldsymbol{\xi}_t \in \Xi_t. \end{aligned}$$

The proof of Lemma 1 can be found in Appendix A.1. Based on the result, we obtain the following corollary on the Lipschitz continuity of the expected cost-to-go functions.

Lemma 2 *At stage $t \in [T] \setminus \{1\}$, for any $\boldsymbol{\xi}_{t-1}$ we have*

$$\begin{aligned} |\mathbb{E}[Q_t(\mathbf{x}, \tilde{\boldsymbol{\xi}}_t) \mid \boldsymbol{\xi}_{t-1}] - \mathbb{E}[Q_t(\mathbf{x}', \tilde{\boldsymbol{\xi}}_t) \mid \boldsymbol{\xi}_{t-1}]| &\leq L_t \|\mathbf{x} - \mathbf{x}'\| & \forall \mathbf{x}, \mathbf{x}' \in \mathcal{X}_{t-1}(\mathbf{x}_{t-2}, \boldsymbol{\xi}_{t-1}), \\ |\hat{\mathbb{E}}[Q_t(\mathbf{x}, \tilde{\boldsymbol{\xi}}_t) \mid \boldsymbol{\xi}_{t-1}] - \hat{\mathbb{E}}[Q_t(\mathbf{x}', \tilde{\boldsymbol{\xi}}_t) \mid \boldsymbol{\xi}_{t-1}]| &\leq L_t \|\mathbf{x} - \mathbf{x}'\| & \forall \mathbf{x}, \mathbf{x}' \in \mathcal{X}_{t-1}(\mathbf{x}_{t-2}, \boldsymbol{\xi}_{t-1}), \\ |\mathbb{E}[\hat{Q}_t(\mathbf{x}, \tilde{\boldsymbol{\xi}}_t) \mid \boldsymbol{\xi}_{t-1}] - \mathbb{E}[\hat{Q}_t(\mathbf{x}', \tilde{\boldsymbol{\xi}}_t) \mid \boldsymbol{\xi}_{t-1}]| &\leq L_t \|\mathbf{x} - \mathbf{x}'\| & \forall \mathbf{x}, \mathbf{x}' \in \mathcal{X}_{t-1}(\mathbf{x}_{t-2}, \boldsymbol{\xi}_{t-1}), \\ |\hat{\mathbb{E}}[\hat{Q}_t(\mathbf{x}, \tilde{\boldsymbol{\xi}}_t) \mid \boldsymbol{\xi}_{t-1}] - \hat{\mathbb{E}}[\hat{Q}_t(\mathbf{x}', \tilde{\boldsymbol{\xi}}_t) \mid \boldsymbol{\xi}_{t-1}]| &\leq L_t \|\mathbf{x} - \mathbf{x}'\| & \forall \mathbf{x}, \mathbf{x}' \in \mathcal{X}_{t-1}(\mathbf{x}_{t-2}, \boldsymbol{\xi}_{t-1}). \end{aligned}$$

Lemma 2 can be verified by directly applying Lemma 1. Based on the assumption (A4), we also have the following lemma that plays an important role in deriving our generalization bound in Theorem 3.

Lemma 3 (Bounded Conditional Variance) *For each stage $t \in [T]$, there exists a constant $\sigma_{t+1}^2 > 0$ such that for any $\mathbf{x}_t \in \mathcal{X}_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$ and all feasible \mathbf{x}_{t-1} , and all $\boldsymbol{\xi}_t$, it holds that $\mathbb{V}[Q_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t] \leq \sigma_{t+1}^2$.*

Under these assumptions and lemmas, we can use Theorem 1 in [63] to obtain the following generalization bound on the error of the NW estimate for any fixed state variables.

Theorem 1 (Generalization Bound for Fixed State Variables \mathbf{x}_t) *At stage $t \in [T]$, for any fixed $\mathbf{x}_t \in \mathcal{X}_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$ and any $\boldsymbol{\xi}_t \in \Xi_t$, we have*

$$\left| \mathbb{E}[Q_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t] - \hat{\mathbb{E}}[Q_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t] \right| \leq \sqrt{\frac{\mathbb{V}[Q_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t]}{O(N^{\frac{4}{p+4}})(1+o(1))g(\boldsymbol{\xi}_t)} \log\left(\frac{1}{\delta_{t+1}}\right)}, \quad (2.6)$$

with probability at least $1 - \delta_{t+1}$. Here,

$$g(\boldsymbol{\xi}_t) = \frac{f(\boldsymbol{\xi}_{t+1} \mid \boldsymbol{\xi}_t)}{2 \int_{\mathbb{R}^p} \mathcal{K}^2(\boldsymbol{\rho}) d\boldsymbol{\rho}}$$

is the scaled marginal density of $\tilde{\boldsymbol{\xi}}_t$.

Theorem 1 provides some insights about the error bound of the NW estimator: a small conditional variance $\mathbb{V}[Q_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t]$ or a large scaled density $g(\boldsymbol{\xi}_t)$ can provide us a better bound. Now we extend the result to obtain a uniform generalization bound for every state variable in a continuous and bounded set under the assumptions (A4)-(A6).

Theorem 2 (Generalization Bound for a Continuous and Bounded Feasible Region) *At stage $t \in [T]$, for any fixed tolerance level $\eta > 0$ and $\boldsymbol{\xi}_t \in \Xi_t$, we have*

$$\begin{aligned} &\left| \mathbb{E}[Q_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t] - \hat{\mathbb{E}}[Q_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t] \right| \\ &\leq \sqrt{\sigma_{t+1}^2 \frac{\log\left(\frac{O(1)(D_t/\eta)^{d_t}}{\delta_{t+1}}\right)}{O(N^{\frac{4}{p+4}})(1+o(1))g(\boldsymbol{\xi}_t)}} + 2L_{t+1}\eta \quad \forall \mathbf{x}_t \in \mathcal{X}_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t), \quad (2.7) \end{aligned}$$

with probability at least $1 - \delta_{t+1}$. Here, L_{t+1} is a Lipschitz constant for $Q_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1})$ defined in Lemma 1 and D_t is a positive constant satisfying

$$\sup_{\mathbf{x}_t, \mathbf{x}'_t \in \mathcal{X}_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)} \|\mathbf{x}_t - \mathbf{x}'_t\| \leq D_t \quad \forall \mathbf{x}_{t-1} \quad \forall \boldsymbol{\xi}_t.$$

We defer the proof of Theorem 2 to Appendix A.2. We note that Theorems 1 and 2 only provide bounds on errors purely due to the NW estimator. In this paper, we want to determine the errors between the true expected cost-to-go and the approximate expected cost-to-go that result from solving the discretized problem (2.4). Let \mathbf{x}_1^* and $\hat{\mathbf{x}}_1$ be optimal solutions to the true problem and the discretized problem, respectively. By recursively applying Theorem 2 from the terminal stage to the first stage, we derive the desired suboptimality bound in the following theorem.

Theorem 3 (Out-of-Sample Performance Guarantee) *For a fixed tolerance level $\eta > 0$, with probability at least $1 - \sum_{t=2}^T \delta_t$, we have*

$$\begin{aligned} & \left(\mathbf{c}_1^\top \hat{\mathbf{x}}_1 + \mathbb{E} [Q_2(\hat{\mathbf{x}}_1, \tilde{\xi}_2) \mid \xi_1] \right) - \left(\mathbf{c}_1^\top \mathbf{x}_1^* + \mathbb{E} [Q_2(\mathbf{x}_1^*, \tilde{\xi}_2) \mid \xi_1] \right) \\ & \leq 2 \sum_{t=2}^T \left(\sqrt{\sigma_t^2 \frac{\log \left(\frac{O(1)N^{t-2} \prod_{s=1}^{t-1} (D_s/\eta)^{d_s}}{\delta_t} \right)}{O(N^{\frac{4}{p+4}})(1+o(1))g_{t-1}}} + 2L_t\eta \right), \end{aligned} \quad (2.8)$$

where L_t is a Lipschitz constant for $Q_t(\mathbf{x}_{t-1}, \xi_t)$ defined in Lemma 1, D_t is the finite diameter defined in Theorem 2, and $g_{t-1} = \min_{i \in [N]} g(\xi_{t-1}^i)$.

Proof We proceed by backward induction in the stages.

Stage $t = T$: We have $Q_T(\cdot, \xi_T) = \hat{Q}_T(\cdot, \xi_T)$ for every ξ_T due to the fact that $Q_{T+1}(\cdot) \equiv 0$. From Theorem 2, for any fixed \mathbf{x}_{T-2} and ξ_{T-1} we have

$$\begin{aligned} \left| \mathbb{E} [Q_T(\mathbf{x}_{T-1}, \tilde{\xi}_T) \mid \xi_{T-1}] - \hat{\mathbb{E}} [\hat{Q}_T(\mathbf{x}_{T-1}, \tilde{\xi}_T) \mid \xi_{T-1}] \right| & \leq \sqrt{\sigma_T^2 \frac{\log \left(\frac{O(1)(D_{T-1}/\eta)^{d_{T-1}}}{\delta_T} \right)}{O(N^{\frac{4}{p+4}})(1+o(1))g(\xi_{T-1})}} + 2L_T\eta \\ & \quad \forall \mathbf{x}_{T-1} \in \mathcal{X}_{T-1}(\mathbf{x}_{T-2}, \xi_{T-1}) \end{aligned}$$

with probability at least $1 - \delta_T$. Taking minimization over \mathcal{X}_{T-1} , for any fixed \mathbf{x}_{T-2} and ξ_{T-1} we have

$$\begin{aligned} & \left| \min_{\mathbf{x}_{T-1} \in \mathcal{X}_{T-1}(\mathbf{x}_{T-2}, \xi_{T-1})} \mathbf{c}_{T-1}^\top \mathbf{x}_{T-1} + \mathbb{E} [Q_T(\mathbf{x}_{T-1}, \tilde{\xi}_T) \mid \xi_{T-1}] \right. \\ & \quad \left. - \min_{\mathbf{x}_{T-1} \in \mathcal{X}_{T-1}(\mathbf{x}_{T-2}, \xi_{T-1})} \mathbf{c}_{T-1}^\top \mathbf{x}_{T-1} + \hat{\mathbb{E}} [\hat{Q}_T(\mathbf{x}_{T-1}, \tilde{\xi}_T) \mid \xi_{T-1}] \right| \\ & = \left| Q_{T-1}(\mathbf{x}_{T-2}, \xi_{T-1}) - \hat{Q}_{T-1}(\mathbf{x}_{T-2}, \xi_{T-1}) \right| \leq \sqrt{\sigma_T^2 \frac{\log \left(\frac{O(1)(D_{T-1}/\eta)^{d_{T-1}}}{\delta_T} \right)}{O(N^{\frac{4}{p+4}})(1+o(1))g(\xi_{T-1})}} + 2L_T\eta \end{aligned}$$

with probability at least $1 - \delta_T$. Applying union bound over the N sample points $\xi_{T-1}^i \in \hat{\Xi}_{T-1}$, $i \in [N]$, for any fixed \mathbf{x}_{T-2} , we have

$$\begin{aligned} \left| Q_{T-1}(\mathbf{x}_{T-2}, \xi_{T-1}^i) - \hat{Q}_{T-1}(\mathbf{x}_{T-2}, \xi_{T-1}^i) \right| & \leq \sqrt{\sigma_T^2 \frac{\log \left(\frac{O(1)N(D_{T-1}/\eta)^{d_{T-1}}}{\delta_T} \right)}{O(N^{\frac{4}{p+4}})(1+o(1))g_{T-1}}} + 2L_T\eta \\ & \quad \forall \xi_{T-1}^i \in \hat{\Xi}_{T-1}, i \in [N] \end{aligned} \quad (2.9)$$

with probability at least $1 - \delta_T$ and $g_{T-1} = \min_{i \in [N]} g(\xi_{T-1}^i)$.

Stage $t = T - 1$: From Theorem 1, for any fixed \mathbf{x}_{T-2} and ξ_{T-2} , we have

$$\left| \mathbb{E} [Q_{T-1}(\mathbf{x}_{T-2}, \tilde{\xi}_{T-1}) \mid \xi_{T-2}] - \hat{\mathbb{E}} [Q_{T-1}(\mathbf{x}_{T-2}, \tilde{\xi}_{T-1}) \mid \xi_{T-2}] \right| \leq \sqrt{\sigma_{T-1}^2 \frac{\log \left(\frac{1}{\delta_{T-1}} \right)}{O(N^{\frac{4}{p+4}})(1+o(1))g(\xi_{T-2})}} \quad (2.10)$$

with probability at least $1 - \delta_{T-1}$. Note that $\hat{\mathbb{E}} [Q_{T-1}(\mathbf{x}_{T-2}, \tilde{\xi}_{T-1}) \mid \xi_{T-2}]$ in (2.10) is not the approximate expected cost-to-go defined in (2.5). Therefore, we use (2.9) to replace $Q_{T-1}(\mathbf{x}_{T-2}, \xi_{T-1})$ with

$\hat{Q}_{T-1}(\mathbf{x}_{T-2}, \boldsymbol{\xi}_{T-1})$. Applying union bound to (2.9) and (2.10), for any fixed \mathbf{x}_{T-2} and $\boldsymbol{\xi}_{T-2}$, we have

$$\begin{aligned} & \left| \mathbb{E} [Q_{T-1}(\mathbf{x}_{T-2}, \tilde{\boldsymbol{\xi}}_{T-1}) \mid \boldsymbol{\xi}_{T-2}] - \hat{\mathbb{E}} \left[\hat{Q}_{T-1}(\mathbf{x}_{T-2}, \tilde{\boldsymbol{\xi}}_{T-1}) \mid \boldsymbol{\xi}_{T-2} \right] \right| \\ & \leq \sqrt{\sigma_T^2 \frac{\log \left(\frac{O(1)N(D_{T-1}/\eta)^{d_{T-1}}}{\delta_T} \right)}{O(N^{\frac{4}{p+4}})(1+o(1))g_{T-1}}} + \sqrt{\sigma_{T-1}^2 \frac{\log \left(\frac{1}{\delta_{T-1}} \right)}{O(N^{\frac{4}{p+4}})(1+o(1))g(\boldsymbol{\xi}_{T-2})}} + 2L_T\eta \end{aligned} \quad (2.11)$$

with probability at least $1 - \delta_T - \delta_{T-1}$. Since the cost-to-go function $Q_{T-1}(\mathbf{x}_{T-2}, \boldsymbol{\xi}_{T-1})$ is L_{T-1} -Lipschitz continuous in \mathbf{x}_{T-2} , Lemma 2 implies that for any $\mathbf{x}_{T-2} \in \mathcal{X}_{T-2}(\mathbf{x}_{T-3}, \boldsymbol{\xi}_{T-2})$, there exists $\mathbf{x}'_{T-2} \in \mathcal{X}_{T-2}^\eta(\mathbf{x}_{T-3}, \boldsymbol{\xi}_{T-2})$, $\|\mathbf{x}_{T-2} - \mathbf{x}'_{T-2}\| \leq \eta$, such that:

$$\left| \mathbb{E} [Q_{T-1}(\mathbf{x}_{T-2}, \tilde{\boldsymbol{\xi}}_{T-1}) \mid \boldsymbol{\xi}_{T-2}] - \mathbb{E} [Q_{T-1}(\mathbf{x}'_{T-2}, \tilde{\boldsymbol{\xi}}_{T-1}) \mid \boldsymbol{\xi}_{T-2}] \right| \leq L_{T-1}\eta \quad (2.12)$$

Furthermore, applying union bound to (2.11) over $\mathbf{x}'_{T-2} \in \mathcal{X}_{T-2}^\eta(\mathbf{x}_{T-3}, \boldsymbol{\xi}_{T-2})$, we get

$$\begin{aligned} & \left| \mathbb{E} [Q_{T-1}(\mathbf{x}'_{T-2}, \tilde{\boldsymbol{\xi}}_{T-1}) \mid \boldsymbol{\xi}_{T-2}] - \hat{\mathbb{E}} \left[\hat{Q}_{T-1}(\mathbf{x}'_{T-2}, \tilde{\boldsymbol{\xi}}_{T-1}) \mid \boldsymbol{\xi}_{T-2} \right] \right| \\ & \leq \sqrt{\sigma_T^2 \frac{\log \left(\frac{O(1)N(D_{T-1}/\eta)^{d_{T-1}}(D_{T-2}/\eta)^{d_{T-2}}}{\delta_T} \right)}{O(N^{\frac{4}{p+4}})(1+o(1))g_{T-1}}} + \sqrt{\sigma_{T-1}^2 \frac{\log \left(\frac{O(1)(D_{T-2}/\eta)^{d_{T-2}}}{\delta_{T-1}} \right)}{O(N^{\frac{4}{p+4}})(1+o(1))g(\boldsymbol{\xi}_{T-2})}} + 2L_T\eta \\ & \quad \forall \mathbf{x}'_{T-2} \in \mathcal{X}_{T-2}^\eta(\mathbf{x}_{T-3}, \boldsymbol{\xi}_{T-2}) \end{aligned} \quad (2.13)$$

with probability at least $1 - \delta_T - \delta_{T-1}$. Using the Lipschitz continuity of $Q_{T-1}(\mathbf{x}_{T-2}, \boldsymbol{\xi}_{T-1})$ and $\hat{Q}_{T-1}(\mathbf{x}_{T-2}, \boldsymbol{\xi}_{T-1})$, we obtain

$$\begin{aligned} & \left| \mathbb{E} [Q_{T-1}(\mathbf{x}_{T-2}, \tilde{\boldsymbol{\xi}}_{T-1}) \mid \boldsymbol{\xi}_{T-2}] - \hat{\mathbb{E}} \left[\hat{Q}_{T-1}(\mathbf{x}_{T-2}, \tilde{\boldsymbol{\xi}}_{T-1}) \mid \boldsymbol{\xi}_{T-2} \right] \right| \\ & \leq \sqrt{\sigma_T^2 \frac{\log \left(\frac{O(1)N(D_{T-1}/\eta)^{d_{T-1}}(D_{T-2}/\eta)^{d_{T-2}}}{\delta_T} \right)}{O(N^{\frac{4}{p+4}})(1+o(1))g_{T-1}}} + \sqrt{\sigma_{T-1}^2 \frac{\log \left(\frac{O(1)(D_{T-2}/\eta)^{d_{T-2}}}{\delta_{T-1}} \right)}{O(N^{\frac{4}{p+4}})(1+o(1))g(\boldsymbol{\xi}_{T-2})}} + 2L_T\eta + 2L_{T-1}\eta \\ & \quad \forall \mathbf{x}_{T-2} \in \mathcal{X}_{T-2}(\mathbf{x}_{T-3}, \boldsymbol{\xi}_{T-2}) \end{aligned}$$

with probability at least $1 - \delta_T - \delta_{T-1}$. Then minimizing over \mathcal{X}_{T-2} , for any fixed $\boldsymbol{\xi}_{T-2}$ and \mathbf{x}_{T-3} , we have

$$\begin{aligned} & \left| \min_{\mathbf{x}_{T-2} \in \mathcal{X}_{T-2}(\mathbf{x}_{T-3}, \boldsymbol{\xi}_{T-2})} \mathbf{c}_{T-2}^\top \mathbf{x}_{T-2} + \mathbb{E} [Q_{T-1}(\mathbf{x}_{T-2}, \tilde{\boldsymbol{\xi}}_{T-1}) \mid \boldsymbol{\xi}_{T-2}] \right. \\ & \quad \left. - \min_{\mathbf{x}_{T-2} \in \mathcal{X}_{T-2}(\mathbf{x}_{T-3}, \boldsymbol{\xi}_{T-2})} \mathbf{c}_{T-2}^\top \mathbf{x}_{T-2} + \hat{\mathbb{E}} \left[\hat{Q}_{T-1}(\mathbf{x}_{T-2}, \tilde{\boldsymbol{\xi}}_{T-1}) \mid \boldsymbol{\xi}_{T-2} \right] \right| \\ & = \left| Q_{T-2}(\mathbf{x}_{T-3}, \boldsymbol{\xi}_{T-2}) - \hat{Q}_{T-2}(\mathbf{x}_{T-3}, \boldsymbol{\xi}_{T-2}) \right| \\ & \leq \sqrt{\sigma_T^2 \frac{\log \left(\frac{O(1)N(D_{T-1}/\eta)^{d_{T-1}}(D_{T-2}/\eta)^{d_{T-2}}}{\delta_T} \right)}{O(N^{\frac{4}{p+4}})(1+o(1))g_{T-1}}} + \sqrt{\sigma_{T-1}^2 \frac{\log \left(\frac{O(1)(D_{T-2}/\eta)^{d_{T-2}}}{\delta_{T-1}} \right)}{O(N^{\frac{4}{p+4}})(1+o(1))g(\boldsymbol{\xi}_{T-2})}} + 2L_T\eta + 2L_{T-1}\eta \end{aligned}$$

with probability at least $1 - \delta_T - \delta_{T-1}$. Applying union bound over the N sample points $\boldsymbol{\xi}_{T-2}^i \in \hat{\Xi}_{T-2}$, $i \in [N]$, for any fixed \mathbf{x}_{T-3} , we have

$$\begin{aligned} & \left| Q_{T-2}(\mathbf{x}_{T-3}, \boldsymbol{\xi}_{T-2}^i) - \hat{Q}_{T-2}(\mathbf{x}_{T-3}, \boldsymbol{\xi}_{T-2}^i) \right| \\ & \leq \sqrt{\sigma_T^2 \frac{\log \left(\frac{O(1)N^2(D_{T-1}/\eta)^{d_{T-1}}(D_{T-2}/\eta)^{d_{T-2}}}{\delta_T} \right)}{O(N^{\frac{4}{p+4}})(1+o(1))g_{T-1}}} + \sqrt{\sigma_{T-1}^2 \frac{\log \left(\frac{O(1)N(D_{T-2}/\eta)^{d_{T-2}}}{\delta_{T-1}} \right)}{O(N^{\frac{4}{p+4}})(1+o(1))g_{T-2}}} + 2L_T\eta + 2L_{T-1}\eta \\ & \quad \forall \boldsymbol{\xi}_{T-2}^i \in \hat{\Xi}_{T-2}, i \in [N] \end{aligned}$$

with probability at least $1 - \delta_T - \delta_{T-1}$ and $g_{T-2} = \min_{i \in [N]} g(\xi_{T-2}^i)$.
 For stage $t = 2$, by backward induction, we obtain

$$\left| \mathbb{E} [Q_2(\mathbf{x}_1, \tilde{\xi}_2) \mid \xi_1] - \hat{\mathbb{E}} [\hat{Q}_2(\mathbf{x}_1, \tilde{\xi}_2) \mid \xi_1] \right| \leq \underbrace{\sum_{t=2}^T \left(\sqrt{\sigma_t^2 \frac{\log \left(\frac{O(1)N^{t-2} \prod_{s=1}^{t-1} (D_s/\eta)^{d_s}}{\delta_t} \right)}{O(N^{\frac{4}{p+4}})(1+o(1))g_{t-1}} + 2L_t\eta} \right)}_{\equiv \epsilon(\delta_2, \delta_3, \dots, \delta_T)} \quad (2.14)$$

$\forall \mathbf{x}_1 \in \mathcal{X}_1(\mathbf{x}_0, \xi_1)$

with probability at least $1 - \sum_{t=2}^T \delta_t$. To ease the notation, let the error term in (2.14) be a function of $\delta_2, \delta_3, \dots, \delta_T$, namely $\epsilon(\delta_2, \delta_3, \dots, \delta_T)$. Then, for $\mathbf{x}_1 = \hat{\mathbf{x}}_1$, we have

$$\mathbf{c}_1^\top \hat{\mathbf{x}}_1 + \mathbb{E} [Q_2(\hat{\mathbf{x}}_1, \tilde{\xi}_2) \mid \xi_1] \leq \left(\mathbf{c}_1^\top \hat{\mathbf{x}}_1 + \hat{\mathbb{E}} [\hat{Q}_2(\hat{\mathbf{x}}_1, \tilde{\xi}_2) \mid \xi_1] \right) + \epsilon(\delta_2, \delta_3, \dots, \delta_T) \quad (2.15a)$$

$$\leq \left(\mathbf{c}_1^\top \mathbf{x}_1^* + \hat{\mathbb{E}} [\hat{Q}_2(\mathbf{x}_1^*, \tilde{\xi}_2) \mid \xi_1] \right) + \epsilon(\delta_2, \delta_3, \dots, \delta_T), \quad (2.15b)$$

where the second inequality (2.15b) holds since \mathbf{x}_1^* is suboptimal to the discretized problem. Similarly, for $\mathbf{x}_1 = \mathbf{x}_1^*$, we have

$$-\left(\mathbf{c}_1^\top \mathbf{x}_1^* + \mathbb{E} [Q_2(\mathbf{x}_1^*, \tilde{\xi}_2) \mid \xi_1] \right) + \left(\mathbf{c}_1^\top \mathbf{x}_1^* + \hat{\mathbb{E}} [\hat{Q}_2(\mathbf{x}_1^*, \tilde{\xi}_2) \mid \xi_1] \right) \leq \epsilon(\delta_2, \delta_3, \dots, \delta_T) \quad (2.16)$$

Applying union bound to (2.15b) and (2.16), we establish the desired result:

$$\left(\mathbf{c}_1^\top \hat{\mathbf{x}}_1 + \mathbb{E} [Q_2(\hat{\mathbf{x}}_1, \tilde{\xi}_2) \mid \xi_1] \right) - \left(\mathbf{c}_1^\top \mathbf{x}_1^* + \mathbb{E} [Q_2(\mathbf{x}_1^*, \tilde{\xi}_2) \mid \xi_1] \right) \leq 2\epsilon(\delta_2, \delta_3, \dots, \delta_T)$$

with probability at least $1 - \sum_{t=2}^T \delta_t$. Thus, the claim follows. \blacksquare

To gain insight from Theorem 3, we derive the following simple upper bound.

Corollary 1 Define $L_{\max} = \max_{t \in [T-1]} L_{t+1}$, $D_{\max} = \max_{t \in [T-1]} (D_{t+1}/\eta)^{d_{t+1}}$, $\delta_{\min} = \min_{t \in [T-1]} \delta_{t+1}$, $\sigma_{\max} = \max_{t \in [T-1]} \sigma_{t+1}^2$, and $g_{\min} = \min_{t \in [T-1]} g_t$. Then, we have

$$\begin{aligned} & \left(\mathbf{c}_1^\top \hat{\mathbf{x}}_1 + \mathbb{E} [Q_2(\hat{\mathbf{x}}_1, \tilde{\xi}_2) \mid \xi_1] \right) - \left(\mathbf{c}_1^\top \mathbf{x}_1^* + \mathbb{E} [Q_2(\mathbf{x}_1^*, \tilde{\xi}_2) \mid \xi_1] \right) \\ & \leq \frac{2T}{O(N^{\frac{2}{p+4}})} \sqrt{\frac{\sigma_{\max}^2}{(1+o(1))g_{\min}} \log \left(\frac{O(1)N^{T-2} D_{\max}^{T-1}}{\delta_{\min}} \right)} + 4L_{\max}T\eta \quad (2.17) \end{aligned}$$

with probability at least $1 - \sum_{t=2}^T \delta_t$.

Proof Since $L_{\max} = \max_{t \in [T-1]} L_{t+1}$, $D_{\max} = \max_{t \in [T-1]} (D_{t+1}/\eta)^{d_{t+1}}$, $\delta_{\min} = \min_{t \in [T-1]} \delta_{t+1}$, $\sigma_{\max}^2 = \max_{t \in [T-1]} \sigma_t^2$, and $g_{\min} = \min_{t \in [T-1]} g_t$, Theorem 3 implies

$$\begin{aligned} & \left(\mathbf{c}_1^\top \hat{\mathbf{x}}_1 + \mathbb{E} [Q_2(\hat{\mathbf{x}}_1, \tilde{\xi}_2) \mid \xi_1] \right) - \left(\mathbf{c}_1^\top \mathbf{x}_1^* + \mathbb{E} [Q_2(\mathbf{x}_1^*, \tilde{\xi}_2) \mid \xi_1] \right) \\ & \leq 2 \sum_{t=2}^T \left(\sqrt{\sigma_t^2 \frac{\log \left(\frac{O(1)N^{t-2} \prod_{s=1}^{t-1} (D_s/\eta)^{d_s}}{\delta_t} \right)}{O(N^{\frac{4}{p+4}})(1+o(1))g_{t-1}} + 2L_t\eta} \right) \leq \frac{2T}{O(N^{\frac{2}{p+4}})} \sqrt{\frac{\sigma_{\max}^2}{(1+o(1))g_{\min}} \log \left(\frac{O(1)N^{T-2} D_{\max}^{T-1}}{\delta_{\min}} \right)} + 4L_{\max}T\eta \end{aligned}$$

with probability at least $1 - \sum_{t=2}^T \delta_t$. \blacksquare

Corollary 1 asserts that the generalization bound in Theorem 3 is at most $\tilde{O}_p(T^{\frac{3}{2}}/N^{\frac{2}{p+4}})$. The result is meaningful from both theoretical and practical perspectives since i) it provides the first statistical analysis of the Markov discretization method and ii) the polynomial dependence on T suggests that our scheme is applicable to problems that require a large number of time stages. However, the bound also depends on the dimension p of the random parameter vector $\tilde{\xi}_t$, implying large out-of-sample errors when p is large. The estimation errors from the high dimensional regression are a relatively common problem. To mitigate the issue, one can utilize a dimensionality reduction algorithm. Using the procedure, one can improve the decay rate of the errors to $\tilde{O}_p(T^{\frac{3}{2}}/N^{\frac{2}{p'+4}})$ when the effective dimensionality p' of the random parameters is much smaller than the dimensionality p of the ambient space.

Remark 1 Shapiro [54] investigates MSLP problems with stagewise independence and known uncertainty distribution. He shows that using the popular sample average approximation (SAA) method, the errors diminish slowly at the rate of $\tilde{O}_p(T^{\frac{1}{2}}/N^{\frac{1}{2T}})$. Thus, the sample complexity of solving the true problem grows exponentially in the number of stages, which makes SAA practically inapplicable for MSLP problems with a large T . ■

3 Data-Driven SDDP

Our solution method to solve the discretized problem in Section 2.2 is based on stochastic dual dynamic programming (SDDP) algorithm introduced by Pereira and Pinto [46], which is a sampling-based variant of the nested Benders decomposition method by Birge [9]. SDDP is an iterative cutting plane approximation algorithm [33] where each iteration consists of two steps called the *forward* and *backward* steps. In the forward step, M scenarios—sample paths in the scenario tree in Figure 1 from the first stage to the terminal stage—are sampled, and a sequence of candidate solutions corresponding to each scenario is obtained under the current approximation of the expected cost-to-go function $\hat{Q}_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_t)$ in (2.4). In the subsequent backward step, the lower and upper bounds of the expected cost-to-go functions are updated starting from the terminal stage to the first stage along the candidate solutions obtained in the forward step. As the lower and upper approximations are evaluated only at the candidate solutions, SDDP can alleviate the intractability arising from discretizing state variables.

3.1 Lower Bound on $\hat{Q}_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_t)$

At the beginning of the k -th iteration of the SDDP algorithm, we consider the following lower bound problem for every $t \in [T]$ and $\boldsymbol{\xi}_t \in \hat{\Xi}_t$:

$$\begin{aligned} \underline{Q}_t^{k-1}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) = \min & \mathbf{c}_t^\top \mathbf{x}_t + \underline{Q}_{t+1}^{k-1}(\mathbf{x}_t, \boldsymbol{\xi}_t) \\ \text{s.t. } & \mathbf{x}_t \in \mathbb{R}_+^{d_t}, \quad \mathbf{z}_t \in \mathbb{R}_+^{d_{t-1}}, \\ & \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{z}_t = \mathbf{b}_t, \\ & \mathbf{z}_t = \mathbf{x}_{t-1}, \end{aligned} \quad (3.1)$$

Here, $\underline{Q}_{t+1}^{k-1}(\mathbf{x}_t, \boldsymbol{\xi}_t) = \sum_{i \in [N]} \hat{w}_t^i(\boldsymbol{\xi}_t) \cdot \underline{Q}_{t+1}^{k-1}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i)$. Here, the expected cost-to-go function $\hat{Q}_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_t)$ in (2.4) is replaced by the lower bound approximation $\underline{Q}_{t+1}^{k-1}(\mathbf{x}_t, \boldsymbol{\xi}_t)$, which is the maximum of a collection of linear functions, known as *cuts*, derived at each iteration. Note that we add the equality constraints $\mathbf{z}_t = \mathbf{x}_{t-1}$ where the auxiliary variables \mathbf{z}_t acting as a copy of the previous state variables \mathbf{x}_{t-1} . This has the drawback of adding one extra variable and constraint for each state variable but simplifies the implementation of the algorithm by avoiding the matrix multiplication for cut gradients [21, 25].

In the forward step, a scenario² is sampled using the approximate conditional distributions from the NW regression. More precisely, conditional on $\boldsymbol{\xi}_1$, we sample $\hat{\boldsymbol{\xi}}_2$ from the empirical uncertainty set $\hat{\Xi}_2$ according to the probabilities $\{\hat{w}_2^i(\boldsymbol{\xi}_1)\}_{i=1}^N$. Given $\hat{\boldsymbol{\xi}}_2$, we then sample $\hat{\boldsymbol{\xi}}_3$ from $\hat{\Xi}_3$ according to the probabilities $\{\hat{w}_3^i(\hat{\boldsymbol{\xi}}_2)\}_{i=1}^N$. This process is repeated until the terminal stage to generate a scenario. For $t = 1, \dots, T-1$, the current lower bound problem (3.1) is solved at the previous stage solution $\bar{\mathbf{x}}_{t-1}^k$ and the scenario $\boldsymbol{\xi}_t$. The main output in the forward step is a sequence of optimal solutions $\bar{\mathbf{x}}_{[T-1]}^k$, known as *candidate solutions*. In the following backward step, for $t = T, \dots, 2$, the lower bound problem (3.1) is solved at candidate solution $\bar{\mathbf{x}}_{t-1}^k$ and every sample $\boldsymbol{\xi}_t \in \hat{\Xi}_t$. Then, the objective values V_t^i and dual solutions $\boldsymbol{\pi}_t^i$, $i \in [N]$, corresponding to $\mathbf{z}_t = \mathbf{x}_{t-1}$ are used to generate a new cut for updating the k -th approximation $\underline{Q}_t^k(\mathbf{x}_{t-1}, \boldsymbol{\xi}_{t-1})$ for every $\boldsymbol{\xi}_{t-1} \in \hat{\Xi}_{t-1}$, as follows:

$$\underline{Q}_t^k(\mathbf{x}_{t-1}, \boldsymbol{\xi}_{t-1}) \leftarrow \max \left\{ \underline{Q}_t^{k-1}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_{t-1}), \left(\sum_{i \in [N]} \hat{w}_{t-1}^i(\boldsymbol{\xi}_{t-1}) \cdot \boldsymbol{\pi}_t^i \right)^\top (\mathbf{x}_{t-1} - \bar{\mathbf{x}}_{t-1}^k) + \sum_{i \in [N]} \hat{w}_{t-1}^i(\boldsymbol{\xi}_{t-1}) \cdot V_t^i \right\}.$$

Due to the convexity of $\hat{Q}_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_t)$ in \mathbf{x}_t , which we will discuss in detail in Section 4.3, $\underline{Q}_t^k(\cdot)$ is a valid lower approximation for $\hat{Q}_{t+1}(\cdot)$ for every $t \in [T-1] \setminus \{1\}$ in any iteration k . With these updated approximations, we solve $\min_{\mathbf{x}_1 \in \mathcal{X}_1(\mathbf{x}_0, \boldsymbol{\xi}_1)} \mathbf{c}_1^\top \mathbf{x}_1 + \underline{Q}_1^k(\mathbf{x}_0, \boldsymbol{\xi}_1)$ to obtain a deterministic lower bound on the value of optimal policies of the discretized problem.

²More than one scenario can be sampled in each iteration. In this case, multiple sequences of candidate solutions can be used to generate cuts in the backward step.

3.2 Upper Bound on $\hat{Q}_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_t)$

In the standard SDDP, the algorithm terminates if a gap between the *deterministic* lower bound and the *stochastic* upper bound becomes within a predetermined tolerance level; see [55] for more detail. Although a small gap may indicate that the current lower bound problem provides a good approximation of the true problem, it may not be the case if the randomness in the upper bound construction causes an early termination of the algorithm. Moreover, it fails to provide a deterministic quality of the current approximation, which would be of interest to the decision maker. We avoid such a stochastic upper bound scheme; instead, we utilize the deterministic upper bound proposed in [24,47,61].

At the k -th iteration, we consider the following upper bound problem for every $t \in [T]$ and $\boldsymbol{\xi}_t \in \hat{\Xi}_t$:

$$\begin{aligned} \bar{Q}_t^{k-1}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) = \min \quad & \mathbf{c}_t^\top \mathbf{x}_t + \bar{Q}_{t+1}^{k-1}(\mathbf{x}_t, \boldsymbol{\xi}_t) + M_t \|\mathbf{y}_t\|_1 \\ \text{s.t.} \quad & \mathbf{x}_t \in \mathbb{R}_+^{d_t}, \quad \boldsymbol{\theta} \in \mathbb{R}_+^{k-1}, \quad \mathbf{y}_t \in \mathbb{R}^{d_t}, \\ & \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{x}_{t-1} = \mathbf{b}_t, \\ & \sum_{j \in [k-1]} \theta_j \bar{\mathbf{x}}_t^j + \mathbf{y}_t = \mathbf{x}_t, \\ & \mathbf{e}^\top \boldsymbol{\theta} = 1. \end{aligned} \tag{3.2}$$

Here, $\bar{Q}_{t+1}^{k-1}(\mathbf{x}_t, \boldsymbol{\xi}_t) = \sum_{i \in [N]} \hat{w}_t^i(\boldsymbol{\xi}_t) \cdot \bar{Q}_{t+1}^{k-1}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i)$. Similar to the lower bound, the upper bound problem is obtained by replacing $\hat{Q}_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_t)$ in (2.4) with the upper bound approximation $\bar{Q}_{t+1}^{k-1}(\mathbf{x}_t, \boldsymbol{\xi}_t)$. Here, $\bar{Q}_{t+1}^{k-1}(\mathbf{x}_t, \boldsymbol{\xi}_t)$ is a lower convex envelope defined on a convex combination of the previously obtained $k-1$ candidate solutions, namely, $\sum_{j \in [k-1]} \theta_j \bar{\mathbf{x}}_t^j$ for every $\boldsymbol{\theta} \in \mathbb{R}_+^{k-1}$ such that $\mathbf{e}^\top \boldsymbol{\theta} = 1$. Note that the auxiliary variables \mathbf{y}_t are introduced in (3.2) to ensure the feasibility of the upper bound problem because the relatively complete recourse in the assumption (A3) alone does not guarantee a feasible solution $\mathbf{x}_t \in \mathcal{X}_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$ that also belongs to the convex combination [61]. That is, \mathbf{y}_t becomes nonzero if there does not exist $\boldsymbol{\theta}$ that satisfies the constraints $\sum_{j \in [k-1]} \theta_j \bar{\mathbf{x}}_t^j = \mathbf{x}_t$ and $\mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{x}_{t-1} = \mathbf{b}_t$. Hence, a penalty proportional to a sufficiently large scalar M_t is imposed whenever \mathbf{y}_t is nonzero. This, in turn, expands the state space of $\bar{Q}_{t+1}^{k-1}(\cdot, \boldsymbol{\xi}_t)$ by adding a new candidate solution $\bar{\mathbf{x}}_t^k$ to the current convex combination.

In the backward step, for $t = T, \dots, 2$, the current upper bound problem (3.2) is solved at every $\boldsymbol{\xi}_t \in \hat{\Xi}_t$ and candidate solution $\bar{\mathbf{x}}_{t-1}^k$, and then the objective values \bar{V}_t^i for every $i \in [N]$ and $\bar{\mathbf{x}}_{t-1}^k$ are used to update the k -th approximation $\bar{Q}_t^k(\mathbf{x}_{t-1}, \boldsymbol{\xi}_{t-1})$ for every $\boldsymbol{\xi}_{t-1} \in \hat{\Xi}_{t-1}$, as follows:

$$\bar{Q}_t^k(\mathbf{x}_{t-1}, \boldsymbol{\xi}_{t-1}) \leftarrow \text{env} \left(\min \left\{ \bar{Q}_t^{k-1}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_{t-1}), \sum_{i \in [N]} \hat{w}_{t-1}^i(\boldsymbol{\xi}_{t-1}) \cdot \bar{V}_t^i + \mathbb{1}_{\{\bar{\mathbf{x}}_{t-1}^k\}}(\mathbf{x}_{t-1}) \right\} \right).$$

Here, ‘env’ represents the lower convex envelope where $\mathbb{1}_{\{\bar{\mathbf{x}}_{t-1}^k\}}(\cdot)$ is an indicator function for the singleton set $\{\bar{\mathbf{x}}_{t-1}^k\}$ defined as

$$\mathbb{1}_{\{\bar{\mathbf{x}}_{t-1}^k\}}(\mathbf{x}_{t-1}) = \begin{cases} 0 & \text{if } \mathbf{x}_{t-1} = \bar{\mathbf{x}}_{t-1}^k \\ +\infty & \text{otherwise.} \end{cases}$$

The convexity of $\hat{Q}_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_t)$ in \mathbf{x}_t guarantees that this convex envelope is a valid upper bound in any iteration k as discussed in Section 4.3. As the approximations $\bar{Q}_t^k(\cdot)$ for every $t \in [T] \setminus \{1\}$ are refined along the new candidate solutions in the backward step, the objective value of $\bar{Q}_1^k(\mathbf{x}_0, \boldsymbol{\xi}_1)$ provides a tighter deterministic upper bound on the value of optimal policies of the discretized problem. Algorithm 1 shows our data-driven SDDP (DD-SDDP) scheme where we adopt the deterministic upper bound in lieu of the standard stochastic one.

3.3 Convergence

We defer the convergence proof to Section 4.3 since the data-driven distributionally robust SDDP (DDR-SDDP) scheme introduced in the next section generalizes DD-SDDP. We emphasize that the convergence indicates the lower and upper bound problems converge to the discretized problem (2.4), not to the true problem (2.2). This implies that our DD-SDDP scheme may suffer from poor out-of-sample performance discussed in Section 2.3 if the discretized problem does not provide a good approximation of the true problem.

Algorithm 1: Data-Driven SDDP (DD-SDDP)

Input: Given observation histories $\hat{\Xi}_t = \{\xi_t^i\}_{i=1}^N \forall t \in [T]$ where $\xi_t^i = (\alpha_t^i, \mathbf{A}_t^i, \mathbf{B}_t^i, \mathbf{b}_t^i)$, tolerance level ε

Initialization: $\underline{Q}_{t+1}^0(\cdot, \xi_t^i) = +\infty$, $\overline{Q}_{t+1}^0(\cdot, \xi_t^i) = -\infty \forall t \in [T-1] \forall i \in [N]$
 $UB = +\infty$, $LB = -\infty$, Iteration $k = 1$

Output: Optimization problems with added cuts

- 1 **while** $|UB - LB| > \varepsilon \cdot \min\{|UB|, |LB|\}$ **do**
- 2 **(forward step)**
- 3 **for** $t = 1, \dots, T-1$ **do**
- 4 **if** $t = 1$ **then**
- 5 Solve $\underline{Q}_1^{k-1}(\mathbf{x}_0, \xi_1)$ in (3.1) and obtain $\bar{\mathbf{x}}_1^k$
- 6 **else**
- 7 Generate a sample $\hat{\xi}_t$ from $\hat{\Xi}_t$ with probability $\{\hat{w}_{t-1}^i(\hat{\xi}_{t-1})\}_{i=1}^N$
- 8 Solve $\underline{Q}_t^{k-1}(\bar{\mathbf{x}}_{t-1}^k, \hat{\xi}_t)$ in (3.1) and obtain $\bar{\mathbf{x}}_t^k$.
- 9 **(backward step)**
- 10 **for** $t = T, \dots, 2$ **do**
- 11 **for** $i = 1, \dots, N$ **do**
- 12 **(LB)** Solve $\underline{Q}_t^k(\bar{\mathbf{x}}_{t-1}^k, \xi_t^i)$ in (3.1). Then obtain the optimal value \underline{V}_t^i and the dual solution π_t^i corresponding to $\mathbf{z}_t = \bar{\mathbf{x}}_{t-1}^k$
- 13 **(UB)** Solve $\overline{Q}_t^k(\bar{\mathbf{x}}_{t-1}^k, \xi_t^i)$ in (3.2) and obtain the optimal value \overline{V}_t^i
- 14 **if** $t > 2$ **then**
- 15 **for** $j = 1, \dots, N$ **do**
- 16 **(Update Lower Bound)**
- 17 Compute $\underline{\mathcal{Y}}_{tj}^k = \sum_{i \in [N]} \hat{w}_{t-1}^i(\xi_{t-1}^j) \underline{V}_t^i$ and $\underline{\mathcal{G}}_{tj}^k = \sum_{i \in [N]} \hat{w}_{t-1}^i(\xi_{t-1}^j) \pi_t^i$
- 18 $\underline{Q}_t^k(\mathbf{x}_{t-1}, \xi_{t-1}^j) \leftarrow \max\{\underline{Q}_t^{k-1}(\mathbf{x}_{t-1}, \xi_{t-1}^j), \underline{\mathcal{G}}_{tj}^{k \top}(\mathbf{x}_{t-1} - \bar{\mathbf{x}}_{t-1}^k) + \underline{\mathcal{Y}}_{tj}^k\}$
- 19 **(Update Upper Bound)**
- 20 Compute $\overline{\mathcal{Y}}_{tj}^k = \sum_{i \in [N]} \hat{w}_{t-1}^i(\xi_{t-1}^j) \overline{V}_t^i$
- 21 $\overline{Q}_t^k(\mathbf{x}_{t-1}, \xi_{t-1}^j) \leftarrow \text{env}(\min\{\overline{Q}_t^{k-1}(\mathbf{x}_{t-1}, \xi_{t-1}^j), \overline{\mathcal{Y}}_{tj}^k + \mathbb{1}_{\{\bar{\mathbf{x}}_{t-1}^k\}}(\mathbf{x}_{t-1})\})$
- 22 **else**
- 23 **(Update Lower Bound)**
- 24 Compute $\underline{\mathcal{Y}}_2^k = \sum_{i \in [N]} \underline{V}_2^i / N$ and $\underline{\mathcal{G}}_2^k = \sum_{i \in [N]} \pi_2^i / N$
- 25 $\underline{Q}_2^k(\mathbf{x}_1, \xi_1) \leftarrow \max\{\underline{Q}_2^{k-1}(\mathbf{x}_1, \xi_1), \underline{\mathcal{G}}_2^{k \top}(\mathbf{x}_1 - \bar{\mathbf{x}}_1^k) + \underline{\mathcal{Y}}_2^k\}$
- 26 **(Update Upper Bound)**
- 27 Compute $\overline{\mathcal{Y}}_2^k = \sum_{i \in [N]} \overline{V}_2^i / N$
- 28 $\overline{Q}_2^k(\mathbf{x}_1, \xi_1) \leftarrow \text{env}(\min\{\overline{Q}_2^{k-1}(\mathbf{x}_1, \xi_1), \mathbb{1}_{\{\bar{\mathbf{x}}_1^k\}}(\mathbf{x}_1) + \overline{\mathcal{Y}}_2^k\})$
- 29 Solve $\underline{Q}_1^k(\mathbf{x}_0, \xi_1)$ in (3.1) and obtain the optimal value LB
- 30 Solve $\overline{Q}_1^k(\mathbf{x}_0, \xi_1)$ in (3.2) and obtain the optimal value UB
- 31 $k = k + 1$

4 Regularization Schemes

Based on Theorem 2, the out-of-sample performance may be poor if the constant upper bound σ_{t+1}^2 on the conditional variance in (2.7) is large. In principle, σ_{t+1}^2 can be replaced with the true conditional variance $\mathbb{V}[Q_{t+1}(\mathbf{x}_t, \tilde{\xi}_{t+1}) \mid \xi_t]$ to obtain a tighter generalization bound. We make use of a regularization scheme involving the conditional variance term that provides a better out-of-sample performance guarantee. The regularization scheme, however, is intractable due to nonconvexity. In this section, we propose the data-driven distributionally robust SDDP (DDR-SDDP) scheme as a tractable approximation.

4.1 Variance Regularization Scheme

In Theorem 2, the constant upper bound σ_{t+1}^2 on the true conditional variance $\mathbb{V}[Q_{t+1}(\mathbf{x}_t, \tilde{\xi}_{t+1}) \mid \xi_t]$ is used as a proxy. Since $\sigma_{t+1}^2 \geq \mathbb{V}[Q_{t+1}(\mathbf{x}_t, \tilde{\xi}_{t+1}) \mid \xi_t]$ for every feasible \mathbf{x}_t and $\xi_t \in \Xi_t$, replacing σ_{t+1}^2 with

$\mathbb{V}[Q_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t]$ gives rise to a tighter generalization bound

$$\begin{aligned} & \left| \mathbb{E} [Q_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t] - \hat{\mathbb{E}} [Q_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t] \right| \\ & \leq \sqrt{\frac{\mathbb{V} [Q_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t]}{g(\boldsymbol{\xi}_t)(1+o(1))N^{\frac{4}{p+4}}} \log \left(\frac{O(1)(D_t/\eta)^{d_t}}{\delta_{t+1}} \right)} + L_{t+1}\eta \left(1 + \sqrt{\frac{\log \left(\frac{O(1)(D_t/\eta)^{d_t}}{\delta_{t+1}} \right)}{g(\boldsymbol{\xi}_t)(1+o(1))N^{\frac{4}{p+4}}}} \right) \end{aligned} \quad (4.1)$$

with probability at least $1 - \delta_{t+1}$. Note that the true variance term provides a more faithful characterization of the generalization errors compared to the bound in Theorem 2. Thus, it implies that a regularization scheme involving the true conditional variance yields good out-of-sample performance, in view of the bound. Nevertheless, the true conditional variance is unknown under the assumption (A2), so we utilize the empirical conditional variance given by

$$\begin{aligned} \hat{\mathbb{V}} \left[\hat{Q}_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t \right] & \equiv \hat{\mathbb{E}} \left[\left(\hat{Q}_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) - \hat{\mathbb{E}} \left[\hat{Q}_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t \right] \right)^2 \mid \boldsymbol{\xi}_t \right] \\ & = \hat{\mathbb{E}} \left[\hat{Q}_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1})^2 \mid \boldsymbol{\xi}_t \right] - \hat{\mathbb{E}} \left[\hat{Q}_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t \right]^2. \end{aligned} \quad (4.2)$$

Using this approximation, we obtain the variance regularized formulation of the data-driven dynamic program (2.4), as follows:

$$\begin{aligned} \hat{Q}_t^{\mathcal{V}\mathcal{R}}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) & = \min \mathbf{c}_t^\top \mathbf{x}_t + \hat{Q}_{t+1}^{\mathcal{V}\mathcal{R}}(\mathbf{x}_t, \boldsymbol{\xi}_t) + \lambda \sqrt{\hat{\mathbb{V}} \left[\hat{Q}_{t+1}^{\mathcal{V}\mathcal{R}}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t \right]} \\ \text{s.t. } \mathbf{x}_t & \in \mathbb{R}_+^{d_t}, \\ \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{x}_{t-1} & = \mathbf{b}_t, \end{aligned} \quad (4.3)$$

where $\hat{Q}_{t+1}^{\mathcal{V}\mathcal{R}}(\mathbf{x}_t, \boldsymbol{\xi}_t) = \hat{\mathbb{E}}[\hat{Q}_{t+1}^{\mathcal{V}\mathcal{R}}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t] = \sum_{i \in [N]} \hat{w}_t^i(\boldsymbol{\xi}_t) \cdot \hat{Q}_{t+1}^{\mathcal{V}\mathcal{R}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i)$ and $\lambda \geq 0$ is a tuning parameter that controls the degree of regularization.

4.2 Data-Driven Distributionally Robust SDDP

Unfortunately, solving the regularization problem (4.3) is intractable since the empirical conditional variance term $\hat{\mathbb{V}}[\hat{Q}_{t+1}^{\mathcal{V}\mathcal{R}}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t]$ in the objective function is nonconvex in \mathbf{x}_t . Instead, we apply the distributionally robust optimization (DRO) methodology to the discretized problem (2.2) as a tractable alternative to the variance regularization scheme. Specifically, we propose the following distributionally robust dynamic program for $t \in [T]$ and $\boldsymbol{\xi}_t \in \hat{\Xi}_t$,

$$\begin{aligned} \hat{Q}_t^{\mathcal{D}\mathcal{R}\mathcal{O}}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) & = \min \mathbf{c}_t^\top \mathbf{x}_t + \max_{\mathbb{P}_t \in \mathcal{P}_t^\lambda(\hat{\mathbb{P}}_t)} \mathbb{E}_{\mathbb{P}_t} \left[\hat{Q}_{t+1}^{\mathcal{D}\mathcal{R}\mathcal{O}}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t \right] \\ \text{s.t. } \mathbf{x}_t & \in \mathbb{R}_+^{d_t}, \\ \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{x}_{t-1} & = \mathbf{b}_t, \end{aligned} \quad (4.4)$$

where $\hat{Q}_{T+1}^{\mathcal{D}\mathcal{R}\mathcal{O}}(\cdot) \equiv 0$. Here, the ambiguity set $\mathcal{P}_t^\lambda(\hat{\mathbb{P}}_t)$ with radius parameter $\lambda \geq 0$ is defined as

$$\mathcal{P}_t^\lambda(\hat{\mathbb{P}}_t) = \left\{ \mathbb{P}_t = \sum_{i \in [N]} w_t^i \delta_{\boldsymbol{\xi}_t^i} : \sum_{i \in [N]} \left| \frac{w_t^i - \hat{w}_t^i(\boldsymbol{\xi}_t)}{\sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)}} \right| \leq \sqrt{N}\lambda, \max_{i \in [N]} \left| \frac{w_t^i - \hat{w}_t^i(\boldsymbol{\xi}_t)}{\sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)}} \right| \leq \lambda, \mathbf{e}^\top \mathbf{w}_t = 1, \mathbf{w}_t \in \mathbb{R}_+^N \right\}, \quad (4.5)$$

where each candidate distribution $\mathbb{P}_t \in \mathcal{P}_t^\lambda(\hat{\mathbb{P}}_t)$ is supported on the N observed data points with probabilities close to the nominal weights $\{\hat{w}_t^i(\boldsymbol{\xi}_t)\}_{i=1}^N$ up to a certain distance λ with respect to the measure described in the ambiguity set (4.5). The distance measure of our choice is a polyhedral outer approximation to the *modified* χ^2 ambiguity set; see [5, 7]. More precisely, we construct the ambiguity set using 1-norm and ∞ -norm to conservatively approximate the 2-norm in the description of the modified χ^2 distance. While utilizing the polyhedral approximation is mainly motivated by the computational benefits, it also provides a connection between the variance regularization formulation (4.3) and the DRO formulation (4.4) as stated in the following proposition.

Proposition 1 *For any \mathbf{x}_t and $\boldsymbol{\xi}_t$, we have*

$$\hat{\mathbb{E}} \left[\hat{Q}_{t+1}^{\mathcal{V}\mathcal{R}}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t \right] + \lambda \sqrt{\hat{\mathbb{V}} \left[\hat{Q}_{t+1}^{\mathcal{V}\mathcal{R}}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t \right]} \leq \max_{\mathbb{P}_t \in \mathcal{P}_t^\lambda(\hat{\mathbb{P}}_t)} \mathbb{E}_{\mathbb{P}_t} \left[\hat{Q}_{t+1}^{\mathcal{D}\mathcal{R}\mathcal{O}}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t \right] + O(\lambda^2).$$

The proof of Proposition 1 can be found in Appendix B.1. Proposition 1 implies that the DRO model provides a more accurate approximation of the variance regularization scheme as we observe more sample trajectories and decrease λ accordingly. Specifically, the suboptimality bound for the variance regularization scheme in [63] suggests to scale the regularization parameter according to $\lambda = O(1/N^{\frac{2}{p+4}})$, which converges to 0 as $N \rightarrow \infty$. In the following proposition, we present a tractable single-level reformulation for (4.4).

Proposition 2 (A Tractable Reformulation) *At stage $t \in [T]$, the DRO problem (4.4) can be reformulated as the linear program*

$$\begin{aligned}
\hat{Q}_t^{\text{DRO}}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) = \min \quad & \mathbf{c}_t^\top \mathbf{x}_t + \gamma + \lambda \left(\beta + \sum_{i \in [N]} \psi_i \right) + \sum_{i \in [N]} \sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)} (\mu_i - \zeta_i) \\
\text{s.t.} \quad & \mathbf{x}_t \in \mathbb{R}_+^{d_t}, \quad \gamma \in \mathbb{R}, \quad \beta \in \mathbb{R}_+, \quad \boldsymbol{\mu}, \boldsymbol{\zeta}, \boldsymbol{\psi} \in \mathbb{R}_+^N, \\
& \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{x}_{t-1} = \mathbf{b}_t, \\
& \hat{Q}_{t+1}^{\text{DRO}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i) \leq \gamma + \frac{\mu_i - \zeta_i}{\sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)}} \quad \forall i \in [N], \\
& \mu_i + \zeta_i = \psi_i + \frac{\beta}{\sqrt{N}} \quad \forall i \in [N].
\end{aligned} \tag{4.6}$$

The proof of Proposition 2 is deferred to Appendix B.2.

As a solution scheme for the DRO problem (4.6), we propose our data-driven distributionally robust SDDP (DDR-SDDP) scheme in Algorithm 2. Similar to DD-SDDP in Section 3, we construct the lower and upper bound problems, as follows:

$$\begin{aligned}
\underline{Q}_{t,k-1}^{\text{DRO}}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) = \min \quad & \mathbf{c}_t^\top \mathbf{x}_t + \gamma + \lambda \left(\beta + \sum_{i \in [N]} \psi_i \right) + \sum_{i \in [N]} \sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)} (\mu_i - \zeta_i) \\
\text{s.t.} \quad & \mathbf{x}_t \in \mathbb{R}_+^{d_t}, \quad \mathbf{z}_t \in \mathbb{R}_+^{d_{t-1}}, \quad \gamma \in \mathbb{R}, \quad \beta \in \mathbb{R}_+, \quad \boldsymbol{\mu}, \boldsymbol{\zeta}, \boldsymbol{\psi} \in \mathbb{R}_+^N \\
& \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{z}_t = \mathbf{b}_t, \\
& \mathbf{z}_t = \mathbf{x}_{t-1}, \\
& \underline{Q}_{t+1,k-1}^{\text{DRO}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i) \leq \gamma + \frac{\mu_i - \zeta_i}{\sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)}} \quad \forall i \in [N], \\
& \mu_i + \zeta_i = \psi_i + \frac{\beta}{\sqrt{N}} \quad \forall i \in [N].
\end{aligned} \tag{4.7}$$

$$\begin{aligned}
\overline{Q}_{t,k-1}^{\text{DRO}}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) = \min \quad & \mathbf{c}_t^\top \mathbf{x}_t + \gamma + \lambda \left(\beta + \sum_{i \in [N]} \psi_i \right) + \sum_{i \in [N]} \sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)} (\mu_i - \zeta_i) \\
\text{s.t.} \quad & \mathbf{x}_t, \mathbf{y}_t^i \in \mathbb{R}_+^{d_t}, \quad \gamma \in \mathbb{R}, \quad \beta \in \mathbb{R}_+, \quad \boldsymbol{\mu}, \boldsymbol{\zeta}, \boldsymbol{\psi} \in \mathbb{R}_+^N, \quad \boldsymbol{\theta}^i \in \mathbb{R}_+^{k-1} \quad \forall i \in [N], \\
& \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{x}_{t-1} = \mathbf{b}_t, \\
& \overline{Q}_{t+1}^{\text{DRO}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i) \leq \gamma + \frac{\mu_i - \zeta_i}{\sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)}} \quad \forall i \in [N], \\
& \sum_{j \in [k]} \theta_j^i \overline{Q}_{t+1,j}^{\text{DRO}}(\bar{\mathbf{x}}_t^j, \boldsymbol{\xi}_{t+1}^i) + M_{t+1} \|\mathbf{y}_t^i\|_1 = \overline{Q}_{t+1}^{\text{DRO}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i) \quad \forall i \in [N], \\
& \sum_{j \in [k]} \theta_j^i \bar{\mathbf{x}}_t^j + \mathbf{y}_t^i = \mathbf{x}_t \quad \forall i \in [N], \\
& \mu_i + \zeta_i = \psi_i + \frac{\beta}{\sqrt{N}} \quad \forall i \in [N], \\
& \mathbf{e}^\top \boldsymbol{\theta}^i = 1 \quad \forall i \in [N].
\end{aligned} \tag{4.8}$$

In DDR-SDDP, we use a multicut version of SDDP, i.e., for each $\boldsymbol{\xi}_t \in \hat{\Xi}_t$, the cost-to-go function $\hat{Q}_t^{\text{DRO}}(\cdot, \boldsymbol{\xi}_t)$ is approximated by $\underline{Q}_{t,k}^{\text{DRO}}(\cdot, \boldsymbol{\xi}_t)$ and $\overline{Q}_{t,k}^{\text{DRO}}(\cdot, \boldsymbol{\xi}_t)$. Each iteration of DDR-SDDP is tractable since the upper and lower bound problems are linear programs, which are efficiently solvable using off-the-shelf solvers.

Algorithm 2: Data-Driven Distributionally Robust SDDP (DDR-SDDP)

Input: Given observation histories $\hat{\Xi}_t = \{\xi_t^i\}_{i=1}^N \forall t \in [T]$ where $\xi_t^i = (\alpha_t^i, \mathbf{A}_t^i, \mathbf{B}_t^i, \mathbf{b}_t^i)$, tolerance level ε

Initialization: $\bar{Q}_{t+1,0}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \xi_{t+1}^i) = +\infty$, $Q_{t+1,0}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \xi_{t+1}^i) = -\infty \forall t \in [T-1] \forall i \in [N]$
 $UB = +\infty$, $LB = -\infty$, Iteration $k = 1$

Output: Optimization problems with added cuts

- 1 **while** $|UB - LB| > \varepsilon \cdot \min\{|UB|, |LB|\}$ **do**
- (forward step)
- for** $t = 1, \dots, T-1$ **do**
- if** $t = 1$ **then**
- Solve $Q_{1,k-1}^{\text{DR}\mathcal{O}}(\mathbf{x}_0, \xi_1)$ in (4.7) and obtain $\bar{\mathbf{x}}_1^k$
- else**
- Generate a sample $\hat{\xi}_t$ from $\hat{\Xi}_t$ with probability $\{\hat{w}_{t-1}^i(\hat{\xi}_{t-1})\}_{i=1}^N$
- Solve $Q_{t,k-1}^{\text{DR}\mathcal{O}}(\bar{\mathbf{x}}_{t-1}^k, \hat{\xi}_t)$ in (4.7) and obtain $\bar{\mathbf{x}}_t^k$.
- (backward step)
- for** $t = T, \dots, 2$ **do**
- for** $i = 1, \dots, N$ **do**
- (LB) Solve $Q_{t,k}^{\text{DR}\mathcal{O}}(\bar{\mathbf{x}}_{t-1}^k, \xi_t^i)$ in (4.7). Then obtain the optimal value \underline{V}_t^i and the dual solution π_t^i corresponding to $\mathbf{z}_t = \bar{\mathbf{x}}_{t-1}^k$
- (UB) Solve $\bar{Q}_{t,k}^{\text{DR}\mathcal{O}}(\bar{\mathbf{x}}_{t-1}^k, \xi_t^i)$ in (4.8) and obtain the optimal value \bar{V}_t^i
- for** $i = 1, \dots, N$ **do**
- (Update Lower Bound)
- $Q_{t,k}^{\text{DR}\mathcal{O}}(\mathbf{x}_{t-1}, \xi_t^i) \leftarrow \max\{Q_{t,k-1}^{\text{DR}\mathcal{O}}(\mathbf{x}_{t-1}, \xi_t^i), \pi_t^{i\top}(\mathbf{x}_{t-1} - \bar{\mathbf{x}}_{t-1}^k) + \underline{V}_t^i\}$
- (Update Upper Bound)
- $\bar{Q}_{t,k}^{\text{DR}\mathcal{O}}(\mathbf{x}_{t-1}, \xi_t^i) \leftarrow \text{env}(\min\{\bar{Q}_{t,k-1}^{\text{DR}\mathcal{O}}(\mathbf{x}_{t-1}, \xi_t^i), \bar{V}_t^i + \mathbb{1}_{\{\bar{\mathbf{x}}_{t-1}^k\}}(\mathbf{x}_{t-1})\})$
- 15 Solve $Q_{1,k}^{\text{DR}\mathcal{O}}(\mathbf{x}_0, \xi_1)$ in (4.7) and obtain the optimal value LB
- 16 Solve $\bar{Q}_{1,k}^{\text{DR}\mathcal{O}}(\mathbf{x}_0, \xi_1)$ in (4.8) and obtain the optimal value UB
- 17 $k = k + 1$

4.3 Convergence

To establish the convergence of the DDR-SDDP scheme presented in Algorithm 2, we make the following assumption.

(B1) Basic Feasible Solution. Candidate solution $\bar{\mathbf{x}}_t^k$ obtained during the forward step and dual solutions $\pi_{t,k}^i$ for every $i \in [N]$ obtained during the backward step are basic feasible solutions for any iteration k .

(B1) is a mild assumption since it can be easily satisfied by using the simplex method to solve the linear programs during the forward and backward steps.

As stated earlier, the DRO problem (4.4) reduces to the discretized problem (2.4) if λ is set to 0. That is, the expected cost-to-go function $\sum_{i \in [N]} \hat{w}_t^i(\xi_t) \cdot \hat{Q}_{t+1}(\mathbf{x}_t, \xi_{t+1}^i)$ in (2.4) replaces the inner maximization problem in (4.4) if the ambiguity set $\mathcal{P}_t^\lambda(\hat{\mathbb{P}}_t)$ contains only the nominal distribution $\hat{\mathbb{P}}_t$. Thus, the convergence proof presented in this section holds for DD-SDDP in Algorithm 1 as well. Before we present the convergence proof, we show that the cost-to-go function $\hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \xi_{t+1})$ is piecewise linear convex in the state variables \mathbf{x}_t .

Lemma 4 *At each stage $t \in [T-1]$, for any $\xi_t \in \hat{\Xi}_t$, $\xi_{t+1} \in \hat{\Xi}_{t+1}$, and feasible \mathbf{x}_{t-1} , the cost-to-go function $\hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \xi_{t+1})$ is piecewise linear convex in $\mathbf{x}_t \in \mathcal{X}_t(\mathbf{x}_{t-1}, \xi_t)$ with a finite number of pieces.*

The proof of Lemma 4 can be found in Appendix B.3. Lemma 4 implies that the cost-to-go function $\hat{Q}_t^{\text{DR}\mathcal{O}}(\mathbf{x}_{t-1}, \xi_t)$ for every $\xi_t \in \hat{\Xi}_t$ can be restored by a finite number of supporting hyperplanes, although the number of such hyperplanes can be prohibitively large. Furthermore, in the following lemma, we establish that the lower and upper bound problems provide valid lower and upper bounds for the discretized problem, respectively.

Lemma 5 *For the upper bound problem in (4.8), suppose $M_{t+1} \geq L_{t+1}^{\text{DR}\mathcal{O}} \forall t \in [T-1]$, where $L_{t+1}^{\text{DR}\mathcal{O}}$ is a Lipschitz constant for $\hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \xi_{t+1})$ under the 1-norm. Then, for any iteration k , the upper and lower bound problems provide valid bounds on the optimal value of the discretized problem, i.e.,*

$$\underline{Q}_{t+1,k}^{\text{DR}\mathcal{O}}(\cdot, \xi_{t+1}) \leq \hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\cdot, \xi_{t+1}) \leq \bar{Q}_{t+1,k}^{\text{DR}\mathcal{O}}(\cdot, \xi_{t+1}) \quad \forall t \in [T-1] \forall \xi_{t+1} \in \hat{\Xi}_{t+1}.$$

The proof of Lemma 5 is deferred to Appendix B.4.

An almost sure finite convergence of the standard SDDP algorithm under stagewise independence has been established in the existing works [50, 55]. These papers, however, only consider the lower bound convergence. It is also worth mentioning that they have different settings on uncertain parameters. Philpott and Guan [50] consider randomness only on the right-hand side of the constraints of the linear program, whereas Shapiro [55] considers a more general MSLP problem where uncertainty can enter any parameters, e.g., objective function coefficients, a recourse matrices, etc. The following theorem presents the convergence of the data-driven distributionally robust SDDP algorithm.

Theorem 4 *Consider the DDR-SDDP scheme described in Algorithm 2 without the termination condition in line 1. Then, the algorithm converges to an optimal solution of the following DRO problem,*

$$\min_{\mathbf{x}_1 \in \mathcal{X}_1(\mathbf{x}_0, \boldsymbol{\xi}_1)} \mathbf{c}_1^\top \mathbf{x}_1 + \max_{\mathbb{P}_1 \in \mathcal{P}_1^\lambda(\hat{\mathbb{P}}_1)} \mathbb{E}_{\mathbb{P}_1} \left[\hat{Q}_2^{\text{DRO}}(\mathbf{x}_1, \tilde{\boldsymbol{\xi}}_2) \mid \boldsymbol{\xi}_1 \right]$$

in a finite number of iterations with probability one.

We defer convergence analysis to Appendix B.5, where we propose our upper bound convergence proof while the lower bound convergence relies on [55].

5 Numerical Experiments

In this section, we conduct numerical experiments to assess the performance of our proposed schemes. All optimization problems are solved using Python 3.7 with GUROBIPY 9.5.2 on a 6-core, 2.3GHz Intel Core i7 CPU laptop with 16GB RAM.

5.1 Portfolio Optimization

We consider the classical multistage portfolio optimization problem where an investor aims to maximize his/her utility at the terminal stage by re-balancing the portfolio at each time stage. The portfolio can be selected from K risky assets and one risk-free asset with a fixed return rate r_f . Initially, the investor has \$1 available in the risk-free asset. At each stage $t \in [T-1]$, the investor can either hold his/her position, buy more, or sell off part (or all) of asset $i \in [K]$ before observing the returns $\boldsymbol{\xi}_{t+1} \in \mathbb{R}^K$ of risky assets at stage $t+1$. We denote by $\mathbf{u}_t^+ \in \mathbb{R}_+^K$ the amount of risky assets bought and by $\mathbf{u}_t^- \in \mathbb{R}_+^K$ the amount of risky assets sold at stage t . At the end of stage t , the value of asset i , $s_{t,i}$, equals the previous value $s_{t-1,i}$, plus the realized return $\xi_{t,i}s_{t-1,i}$ during the period, plus the newly bought amount $u_{t,i}^+$, minus the newly sold amount $u_{t,i}^-$. We use f_b and f_s to denote per-unit transaction costs for buying and selling a unit asset, respectively.

The problem can be solved via the (true) dynamic program for $t \in [T-1]$:

$$\begin{aligned} Q_t(\mathbf{s}_{t-1}, \boldsymbol{\xi}_t) = & \max \mathbb{E} [Q_{t+1}(\mathbf{s}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t] \\ \text{s.t. } & \mathbf{s}_t \in \mathbb{R}_+^{K+1}, \quad \mathbf{u}_t^+, \mathbf{u}_t^- \in \mathbb{R}_+^K, \\ & s_{t,i} = \xi_{t,i}s_{t-1,i} + u_{t,i}^+ - u_{t,i}^- \quad \forall i \in [K], \\ & s_{t,K+1} = r_f s_{t-1,K+1} - (1 + f_b)\mathbf{e}^\top \mathbf{u}_t^+ + (1 - f_s)\mathbf{e}^\top \mathbf{u}_t^-. \end{aligned} \quad (5.1)$$

Here, $s_{0,i} = 0$ for every $i \in [K]$, $s_{0,K+1} = 1$, and the cost-to-go function at the terminal stage T is

$$\begin{aligned} Q_T(\mathbf{s}_{T-1}, \boldsymbol{\xi}_T) = & \max U(\boldsymbol{\xi}_T^\top \mathbf{x}_T) \\ \text{s.t. } & \mathbf{s}_T \in \mathbb{R}_+^{K+1}, \\ & s_{T,i} = \xi_{T,i}s_{T-1,i} \quad \forall i \in [K], \\ & s_{T,K+1} = r_f s_{T-1,K+1}. \end{aligned} \quad (5.2)$$

Here, the function $U : \mathbb{R}_{++} \rightarrow \mathbb{R}_+$ is a linear approximation to the log utility function. Note that no immediate cost is imposed at stage $t \in [T-1]$ since our goal is to maximize the utility of cumulative wealth at the terminal stage T .

We derive a single-level DRO reformulation for the true problem based on (4.6) and compare our data-driven SDDP scheme (DD-SDDP) in Algorithm 1 and the distributionally robust version (DDR-SDDP) in Algorithm 2 against the stagewise independent scheme (Independent), the equally weighted portfolio (Equal), and the hidden Markov model based SDDP scheme (HMM-SDDP) proposed by Silva et al. [61]. The stagewise independent scheme is an SDDP scheme that ignores correlation over time in the random

return process. Thus, we replace the NW estimates for the conditional distribution with discrete uniform distribution. The equally weighted portfolio is to allocate an equal amount of the current wealth in all assets at each stage, hence known as the $1/n$ portfolio strategy. DeMiguel et al. [19] show that this seemingly naïve strategy is optimal under certain conditions. The hidden Markov model based SDDP scheme assumes three unobservable states governing random returns ξ_t , namely *bull*, *regular*, and *bear* states and apply DRO to account for estimation errors of the transition probabilities among those states. Similar to our schemes, the resulting MSLP problem is then solved by SDDP.

We test the schemes with the historical weekly returns of the following data sets from December 2003 to January 2023: the 10 Industry Portfolios and 12 Industry Portfolios from the Fama-French online data library³, which include US stock portfolios categorized by industries; and the iShares Exchange-Trades Funds (iShares) data set downloaded from *yfinance*⁹, which includes the following nine funds: EWG (Germany), EWH (Hong Kong), EWI (Italy), EWK (Belgium), EWL (Switzerland), EWN (Netherlands), EWP (Spain), EWQ (France), and EWU (United Kingdom).

We obtain 120 historical trajectories of the weekly closing prices from each data set by setting the time horizon to $T = 8$ weeks. The first 50 trajectories are used as training data for our models, while the remaining 70 trajectories are used for the out-of-sample tests. In our experiment, we utilize a five-fold cross-validation procedure to determine the radius of the ambiguity set λ for DDR-SDDP and HMM-SDDP. In each trial, we split the training data into five equal-sized subsets where four of the five subsets are put together to train the model. The resulting policy is then tested on the remaining set for λ in $[10^{-3}, 10^1]/N^{2/(p+4)}$ on a logarithm search grid with 10 equidistant points. This process is repeated five times for different partitions of the data to choose λ that performs best overall.

First, we show the convergence of our data-driven schemes. Figure 2 depicts the evolution of the upper and lower bounds using the 10 Industry Portfolios training data set. We observe that both DD-SDDP and DDR-SDDP are able to close the optimality gap to less than 3% after 800 iterations, and the gap for DDR-SDDP is about three times smaller than one for DD-SDDP.

Table 2 reports the out-of-sample performance for five different schemes. We fixed the number of iterations $k = 200$ for every SDDP scheme. It took about 30 minutes to solve an MSLP problem instance in this setup. Along with other statistics, we include the Sharpe ratio to measure the performance of a portfolio relative to a risk-free asset, after adjusting for its risk. Here, the Sharpe ratio is defined as

$$\text{Sharpe ratio} = \frac{\text{mean}\{R_i\}_{i=1}^{70} - R_f}{\text{std}\{R_i\}_{i=1}^{70}},$$

where R_i is the return of the portfolio along test trajectory i over 8 weeks and R_f is the return of the risk-free asset over 8 weeks. The results indicate that the DDR-SDDP scheme performs favorably relative to its competitors: it achieves the largest utility, the largest mean return, and the largest Sharpe ratio over all data sets. In addition, compared to DD-SDDP, DDR-SDDP provides a less risky policy in terms of standard deviation for every data set, illustrating the connection with the variance regularization scheme discussed in Section 4. Meanwhile, we observe that DD-SDDP also performs significantly better than the stagewise independent scheme, showing the benefit of incorporating stagewise dependence into the model.

Table 2 Out-of-sample statistics of different schemes

Data set	Model	Utility	Mean return	Std. dev.	Sharpe ratio
10 Industry Fama-French	DDR-SDDP	0.03190	0.03411	0.06348	0.47656
	DD-SDDP	0.02654	0.02904	0.06957	0.36176
	Independent	0.01582	0.01762	0.06228	0.22088
	HMM-SDDP	0.02428	0.02540	0.04697	0.45878
	Equal	-0.00215	-0.00197	0.02743	-0.03518
12 Industry Fama-French	DDR-SDDP	0.02486	0.02669	0.05974	0.38238
	DD-SDDP	0.02348	0.02591	0.06899	0.31962
	Independent	0.02014	0.02203	0.06316	0.28780
	HMM-SDDP	0.01968	0.02086	0.04970	0.34211
	Equal	-0.00152	-0.00136	0.02715	-0.03166
iShares	DDR-SDDP	0.01380	0.01425	0.03605	0.28860
	DD-SDDP	0.00734	0.00945	0.06855	0.08160
	Independent	0.00343	0.00534	0.06511	0.02283
	HMM-SDDP	0.00568	0.00609	0.03433	0.06526
	Equal	-0.05000	-0.04865	0.02729	-0.31754

³https://mba.tuck.dartmouth.edu/pages/faculty/ken.french/data_library.html

⁹<https://pypi.org/project/yfinance/>

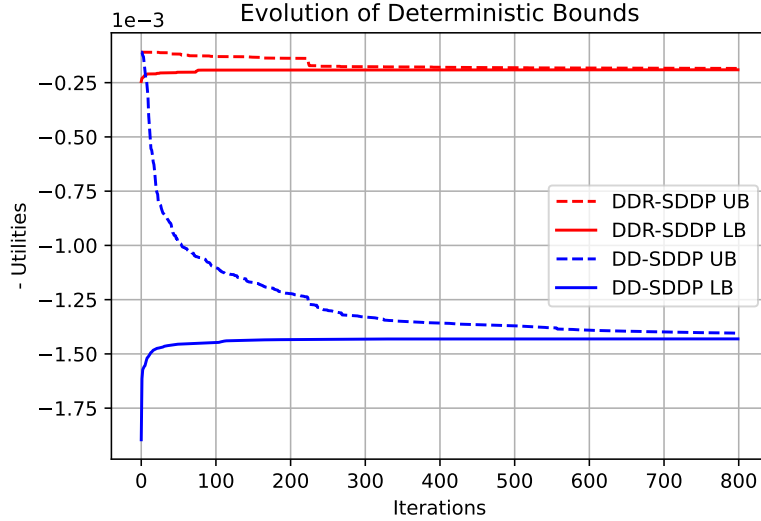


Fig. 2 Evolution of both upper and lower bounds for DDR-SDDP and DD-SDDP over 800 iterations using the 10 Industry Portfolios data set.

5.2 Day-Ahead Wind Energy Commitment

In this experiment, we explore the wind energy commitment problem on the day-ahead market considered in [30] and [34]. At the beginning of day t , the day-ahead electricity prices $\mathbf{p}_t \in \mathbb{R}_+^{24}$ for the next 24 hours are known to a wind energy producer. The producer determines how much wind energy to commit for each of the next 24 hours without knowing the actual hourly amounts of wind energy $\mathbf{w}_{t+1} \in \mathbb{R}_+^{24}$. That is, the energy commitment levels $\mathbf{u}_t \in \mathbb{R}_+^{24}$ are determined after the observation of the day-ahead prices \mathbf{p}_t but before the actual amounts of wind energy \mathbf{w}_{t+1} are realized. On day $t+1$, the actual production can be used to satisfy the scheduled commitment \mathbf{u}_t or charge three storage devices indexed by $l \in \{1, 2, 3\}$. It can also be dumped if all three storage devices are fully charged. The hourly commitment $u_{t,h}$ can be satisfied directly from the newly generated wind energy $w_{t+1,h}$ or by discharging the storage devices. Otherwise, there is a penalty of twice the respective day-ahead price $p_{t,h}$ for the unsatisfied commitment. Each of the devices has a different capacity \bar{s}^l , hourly leakage λ^l , charging efficiency λ_c^l and discharging efficiency λ_d^l . We denote by $s_t^l \in \mathbb{R}_+^{24}$ the hourly levels of storage l over the next 24 hours. The state variables $\mathbf{s}_t = (s_{t,24}^1, s_{t,24}^2, s_{t,24}^3) \in \mathbb{R}_+^3$ represent the storage levels at the end of day t while the random parameters include the day-ahead prices \mathbf{p}_t and the wind energy production levels \mathbf{w}_t during the day, i.e., $\boldsymbol{\xi}_t = (\mathbf{p}_t, \mathbf{w}_t)$. The wind producer's objective is to maximize the expected profit over the time horizon of $T = 7$ days. The optimal bidding and storage strategy for the wind energy producer is obtained by solving the (true) dynamic program for $t \in [T]$

$$\begin{aligned}
Q_t(\mathbf{s}_t, \boldsymbol{\xi}_t) &= \max \mathbf{p}_t^\top \mathbf{u}_t - 2\mathbf{p}_t^\top \mathbb{E}[e_{t+1}^u | \boldsymbol{\xi}_t] + \mathbb{E}[Q_{t+1}(\mathbf{s}_{t+1}, \tilde{\boldsymbol{\xi}}_{t+1}) | \boldsymbol{\xi}_t] \\
&\text{s.t. } \mathbf{u}_t, e_{t+1}^{\{s,u,d\}} \in \mathbb{R}_+^{24}, \quad e_{t+1}^{\{+,-\},l}, s_{t+1}^l \in \mathbb{R}_+^{24} \quad \forall l \in [3], \\
&\quad w_{t+1,h} = e_{t+1,h}^s + e_{t+1,h}^{+,1} + e_{t+1,h}^{+,2} + e_{t+1,h}^{+,3} + e_{t+1,h}^d \quad \forall h \in [24], \\
&\quad u_{t,h} = e_{t+1,h}^s + e_{t+1,h}^{-,1} + e_{t+1,h}^{-,2} + e_{t+1,h}^{-,3} + e_{t+1,h}^u \quad \forall h \in [24], \\
&\quad s_{t+1,h}^l = \lambda^l s_{t+1,h-1}^l + \lambda_c^l e_{t+1,h}^{+,l} - \frac{1}{\lambda_d^l} e_{t+1,h}^{-,l} \quad \forall h \in [24] \quad \forall l \in [3], \\
&\quad s_{t+1,h}^l \leq \bar{s}^l \quad \forall h \in [24] \quad \forall l \in [3],
\end{aligned} \tag{5.3}$$

where $Q_{T+1}(\cdot) \equiv 0$. Here, e_{t+1}^s and e_{t+1}^u represent the amounts of satisfied and unsatisfied energy commitments, respectively, while e_{t+1}^d represents the amounts of dumped wind energy. In addition, $e_{t+1}^{+,l}$ represents the amounts of wind energy used to charge storage l and $e_{t+1}^{-,l}$ the amounts of energy discharged from storage l to meet the commitments.

We obtain the hourly day-ahead prices in the PJM market and the hourly wind energy from 2002 to 2011 at the following locations: **Ohio** (41.8125N, 81.5625W) and **North Carolina** (33.9375N, 77.9375W). By setting $T = 7$ days, we obtain 520 historical trajectories for each location. As the wind energy and the day-ahead prices show clear seasonality patterns, we divide the trajectories into four parts according to different seasons and evaluate the out-of-sample performance for each of them separately. We perform principal component analysis on the high dimensional data $\boldsymbol{\xi}_t = (\mathbf{p}_t, \mathbf{w}_t) \in \mathbb{R}_+^{48}$ to obtain a 6-dimensional

subspace that explains more than 90% of the variability of the historical observations, which mitigates large out-of-sample errors as discussed in Section 2.3.

Table 3 presents the out-of-sample performance for DDR-SDDP, DD-SDDP, and the stagewise independent scheme (Independent). Similar to the previous example, our data-driven schemes outperform the stagewise independent scheme in all criteria. We observe that DDR-SDDP wins in all the categories. Particularly, it is more robust in terms of the 10th percentile compared to other schemes while the stagewise independent scheme has a significant risk of incurring a loss (negative profits).

Table 3 Out-of-sample statistics of profit (in \$100,000)

Data Set	Model	Mean	Variance	10th pct.
Ohio	DDR-SDDP	7.12	18.61	2.72
	DD-SDDP	6.71	22.57	2.07
	Independent	5.61	21.58	0.78
North Carolina	DDR-SDDP	8.62	42.82	1.47
	DD-SDDP	8.27	53.47	1.07
	Independent	7.55	57.12	-0.30

6 Conclusion

In this paper, we introduced a data-driven SDDP (DD-SDDP) scheme for solving a risk-neutral multistage stochastic linear programming (MSLP) problem under an unknown underlying Markov process for random parameters. We utilized the NW kernel regression to estimate the true conditional distribution and established for the first time the theoretical out-of-sample guarantee for the data-driven MSLP problem. This theoretical result inspired us to develop a data-driven distributionally robust SDDP (DDR-SDDP) scheme as a tractable regularization scheme. The numerical experiments demonstrated that our robust scheme outperformed all other benchmarks in real-world applications. In the future, it would be interesting to extend our scheme to solve general convex and nonconvex problems as well as problems with discrete decisions.

Funding This research was supported by the National Science Foundation grants no. 1752125 and 2153606.

Conflict of interest The authors declare that they have no conflict of interest.

References

- Ahmed, S., Cabral, F.G., Freitas Paulo da Costa, B.: Stochastic lipschitz dynamic programming. *Mathematical Programming* **191**(2), 755–793 (2022)
- Ahmed, S., King, A.J., Parija, G.: A multi-stage stochastic integer programming approach for capacity expansion under uncertainty. *Journal of Global Optimization* **26**(1), 3–24 (2003)
- Ahmed, S., Sahinidis, N.V.: An approximation scheme for stochastic integer programs arising in capacity expansion. *Operations Research* **51**(3), 461–471 (2003)
- Barth, R., Brand, H., Meibom, P., Weber, C.: A stochastic unit-commitment model for the evaluation of the impacts of integration of large amounts of intermittent wind power. In: *International Conference on Probabilistic Methods Applied to Power Systems*, pp. 1–8. IEEE (2006)
- Bayraksan, G., Love, D.K.: Data-driven stochastic programming using phi-divergences. In: *The Operations Research Revolution*, pp. 1–19. INFORMS (2015)
- Bellman, R.: *Dynamic Programming*. Princeton University Press (1957)
- Ben-Tal, A., Den Hertog, D., De Waegenaere, A., Melenberg, B., Rennen, G.: Robust solutions of optimization problems affected by uncertain probabilities. *Management Science* **59**(2), 341–357 (2013)
- Bertsimas, D., Gupta, V., Kallus, N.: Data-driven robust optimization. *Mathematical Programming* **167**(2), 235–292 (2018)
- Birge, J.R.: Decomposition and partitioning methods for multistage stochastic linear programs. *Operations Research* **33**(5), 989–1007 (1985)
- Bonnans, J.F., Cen, Z., Christel, T.: Energy contracts management by stochastic programming techniques. *Annals of Operations Research* **200**(1), 199–222 (2012)
- Bradley, S.P., Crane, D.B.: A dynamic model for bond portfolio management. *Management Science* **19**(2), 139–151 (1972)
- Carino, D.R., Kent, T., Myers, D.H., Stacy, C., Sylvanus, M., Turner, A.L., Watanabe, K., Ziemba, W.T.: The russell-yasuda kasai model: An asset/liability model for a Japanese insurance company using multistage stochastic programming. *Interfaces* **24**(1), 29–49 (1994)
- Castro, M.P., Bodur, M., Song, Y.: Markov chain-based policies for multi-stage stochastic integer linear programming with an application to disaster relief logistics. *arXiv preprint arXiv:2207.14779* (2022)

14. Cerisola, S., Baíllo, Á., Fernández-López, J.M., Ramos, A., Gollmer, R.: Stochastic power generation unit commitment in electricity markets: A novel formulation and a comparison of solution methods. *Operations Research* **57**(1), 32–46 (2009)
15. Chen, Z.L., Powell, W.B.: Convergent cutting-plane and partial-sampling algorithm for multistage stochastic linear programs with recourse. *Journal of Optimization Theory and Applications* **102**(3), 497–524 (1999)
16. Dantzig, G.B., Infanger, G.: Multi-stage stochastic linear programs for portfolio optimization. *Annals of Operations Research* **45**(1), 59–76 (1993)
17. De Matos, V.L., Philpott, A.B., Finardi, E.C.: Improving the performance of stochastic dual dynamic programming. *Journal of Computational and Applied Mathematics* **290**, 196–208 (2015)
18. Delage, E., Ye, Y.: Distributionally robust optimization under moment uncertainty with application to data-driven problems. *Operations Research* **58**(3), 595–612 (2010)
19. DeMiguel, V., Garlappi, L., Uppal, R.: Optimal versus naive diversification: How inefficient is the 1/n portfolio strategy? *The review of Financial studies* **22**(5), 1915–1953 (2009)
20. Downward, A., Dowson, O., Baucke, R.: Stochastic dual dynamic programming with stagewise-dependent objective uncertainty. *Operations Research Letters* **48**(1), 33–39 (2020)
21. Dowson, O.: Applying stochastic optimisation to the new zealand dairy industry. Ph.D. thesis, University of Auckland (2018)
22. Duque, D., Morton, D.P.: Distributionally robust stochastic dual dynamic programming. *SIAM Journal on Optimization* **30**(4), 2841–2865 (2020)
23. Escudero, L.F., Kamesam, P.V., King, A.J., Wets, R.J.: Production planning via scenario modelling. *Annals of Operations Research* **43**(6), 309–335 (1993)
24. Georghiou, A., Tsoukalas, A., Wiesemann, W.: Robust dual dynamic programming. *Operations Research* **67**(3), 813–830 (2019)
25. Girardeau, P., Leclere, V., Philpott, A.B.: On the convergence of decomposition methods for multistage stochastic convex programs. *Mathematics of Operations Research* **40**(1), 130–145 (2015)
26. Gjelsvik, A., Belsnes, M.M., Haugstad, A.: An algorithm for stochastic medium-term hydrothermal scheduling under spot price uncertainty. In: *Proceedings of 13th Power Systems Computation Conference* (1999)
27. Golub, B., Holmer, M., McKendall, R., Pohlman, L., Zenios, S.A.: A stochastic programming model for money management. *European Journal of Operational Research* **85**(2), 282–296 (1995)
28. Györfi, L., Kohler, M., Krzyzak, A., Walk, H., et al.: *A Distribution-Free Theory of Nonparametric Regression*, vol. 1. Springer (2002)
29. Hanasusanto, G.A., Kuhn, D.: Robust data-driven dynamic programming. In *Advances in Neural Information Processing Systems* **26** (2013)
30. Hannah, L., Dunson, D.B.: Approximate dynamic programming for storage problems. In: *International Conference on Machine Learning* (2011)
31. Huang, J., Zhou, K., Guan, Y.: A study of distributionally robust multistage stochastic optimization. arXiv preprint arXiv:1708.07930 (2017)
32. Infanger, G., Morton, D.P.: Cut sharing for multistage stochastic linear programs with interstage dependency. *Mathematical Programming* **75**(2), 241–256 (1996)
33. Kelley Jr, J.E.: The cutting-plane method for solving convex programs. *Journal of the society for Industrial and Applied Mathematics* **8**(4), 703–712 (1960)
34. Kim, J.H., Powell, W.B.: Optimal energy commitments with storage and intermittent supply. *Operations Research* **59**(6), 1347–1360 (2011)
35. Klabjan, D., Simchi-Levi, D., Song, M.: Robust stochastic lot-sizing by means of histograms. *Production and Operations Management* **22**(3), 691–710 (2013)
36. Kuhn, D., Wiesemann, W., Georghiou, A.: Primal and dual linear decision rules in stochastic and robust optimization. *Mathematical Programming* **130**(1), 177–209 (2011)
37. Lan, G.: Complexity of stochastic dual dynamic programming. *Mathematical Programming* **191**(2), 717–754 (2022)
38. Löhdorf, N., Shapiro, A.: Modeling time-dependent randomness in stochastic dual dynamic programming. *European Journal of Operational Research* **273**(2), 650–661 (2019)
39. Löhdorf, N., Wozabal, D., Minner, S.: Optimizing trading decisions for hydro storage systems using approximate dual dynamic programming. *Operations Research* **61**(4), 810–823 (2013)
40. Love, D., Bayraksan, G.: Phi-divergence constrained ambiguous stochastic programs for data-driven optimization. Technical report, Department of Integrated Systems Engineering, The Ohio State University, Columbus, Ohio (2015)
41. Mehrotra, S., Zhang, H.: Models and algorithms for distributionally robust least squares problems. *Mathematical Programming* **146**(1), 123–141 (2014)
42. Mulvey, J.M., Vladimirou, H.: Stochastic network programming for financial planning problems. *Management Science* **38**(11), 1642–1664 (1992)
43. Nadaraya, E.A.: On estimating regression. *Theory of Probability & Its Applications* **9**(1), 141–142 (1964)
44. Namkoong, H., Duchi, J.C.: Variance-based regularization with convex objectives. *Advances in neural information processing systems* **30** (2017)
45. Pereira, M.V., Pinto, L.M.: Stochastic optimization of a multireservoir hydroelectric system: A decomposition approach. *Water Resources Research* **21**(6), 779–792 (1985)
46. Pereira, M.V., Pinto, L.M.: Multi-stage stochastic optimization applied to energy planning. *Mathematical Programming* **52**(1), 359–375 (1991)
47. Philpott, A., de Matos, V., Finardi, E.: On solving multistage stochastic programs with coherent risk measures. *Operations Research* **61**(4), 957–970 (2013)
48. Philpott, A.B., De Matos, V.L.: Dynamic sampling algorithms for multi-stage stochastic programs with risk aversion. *European Journal of Operational Research* **218**(2), 470–483 (2012)
49. Philpott, A.B., De Matos, V.L., Kapelevich, L.: Distributionally robust SDDP. *Computational Management Science* **15**(3), 431–454 (2018)
50. Philpott, A.B., Guan, Z.: On the convergence of stochastic dual dynamic programming and related methods. *Operations Research Letters* **36**(4), 450–455 (2008)
51. Scarf, H.: A min-max solution of an inventory problem. *Studies in the Mathematical Theory of Inventory and Production* (1958)

52. Sen, S., Yu, L., Genc, T.: A stochastic programming approach to power portfolio optimization. *Operations Research* **54**(1), 55–72 (2006)
53. Shapiro, A.: Inference of statistical bounds for multistage stochastic programming problems. *Mathematical Methods of Operations Research* **58**(1), 57–68 (2003)
54. Shapiro, A.: On complexity of multistage stochastic programs. *Operations Research Letters* **34**(1), 1–8 (2006)
55. Shapiro, A.: Analysis of stochastic dual dynamic programming method. *European Journal of Operational Research* **209**(1), 63–72 (2011)
56. Shapiro, A.: Tutorial on risk neutral, distributionally robust and risk averse multistage stochastic programming. *European Journal of Operational Research* **288**(1), 1–13 (2021)
57. Shapiro, A., Ahmed, S.: On a class of minimax stochastic programs. *SIAM Journal on Optimization* **14**(4), 1237–1249 (2004)
58. Shapiro, A., Dentcheva, D., Ruszczyński, A.: *Lectures on Stochastic Programming: Modeling and Theory*. SIAM (2021)
59. Shapiro, A., Nemirovski, A.: On complexity of stochastic programming problems. In: V. Jeyakumar, A. Rubinov (eds.) *Continuous Optimization: Current Trends and Modern Applications*, pp. 111–146. Springer (2005)
60. Shapiro, A., Tekaya, W., da Costa, J.P., Soares, M.P.: Risk neutral and risk averse stochastic dual dynamic programming method. *European Journal of Operational Research* **224**(2), 375–391 (2013)
61. Silva, T., Valladão, D., Homem-de Mello, T.: A data-driven approach for a class of stochastic dynamic optimization problems. *Computational Optimization and Applications* **80**(3), 687–729 (2021)
62. Singh, K.J., Philpott, A.B., Wood, R.K.: Dantzig-wolfe decomposition for solving multistage stochastic capacity-planning problems. *Operations Research* **57**(5), 1271–1286 (2009)
63. Srivastava, P.R., Wang, Y., Hanasusanto, G.A., Ho, C.P.: On data-driven prescriptive analytics with side information: A regularized Nadaraya-Watson approach. *arXiv preprint arXiv:2110.04855* (2021)
64. Watson, G.S.: Smooth regression analysis. *Sankhyā: The Indian Journal of Statistics, Series A* pp. 359–372 (1964)
65. Wiesemann, W., Kuhn, D., Sim, M.: Distributionally robust convex optimization. *Operations Research* **62**(6), 1358–1376 (2014)
66. Zhang, S., Sun, X.A.: On distributionally robust multistage convex optimization: new algorithms and complexity analysis. *arXiv preprint arXiv:2010.06759* (2020)
67. Zhang, S., Sun, X.A.: On distributionally robust multistage convex optimization: Data-driven models and performance. *arXiv preprint arXiv:2210.08433* (2022)
68. Zhang, S., Sun, X.A.: Stochastic dual dynamic programming for multistage stochastic mixed-integer nonlinear optimization. *Mathematical Programming* pp. 1–51 (2022)
69. Zou, J., Ahmed, S., Sun, X.A.: Multistage stochastic unit commitment using stochastic dual dynamic integer programming. *IEEE transactions on Power Systems* **34**(3), 1814–1823 (2018)
70. Zou, J., Ahmed, S., Sun, X.A.: Stochastic dual dynamic integer programming. *Mathematical Programming* **175**(1), 461–502 (2019)
71. Zymler, S., Kuhn, D., Rustem, B.: Distributionally robust joint chance constraints with second-order moment information. *Mathematical Programming* **137**(1), 167–198 (2013)

A Proofs in Section 2

A.1 Proof of Lemma 1

The proof of Lemma 1 relies on the Hoffman's Lemma shown below.

Lemma 6 (Hoffman Lemma (Theorem 9.14 in [58])) *Let $\mathcal{M}(\mathbf{u}) = \{\mathbf{x} \in \mathbb{R}_+^n \mid \mathbf{A}\mathbf{x} = \mathbf{u}\}$ be a nonempty polyhedron parameterized by the right-hand side vector $\mathbf{u} \in \mathbb{R}^m$. Consider $\mathbf{u}' \in \mathbb{R}^m$ such that $\mathcal{M}(\mathbf{u}') \neq \emptyset$. Then, there exists a positive constant $r = \max_{\lambda_1 \in S_0} \|\lambda_1\|_1$, such that for any $\mathbf{x} \in \mathcal{M}(\mathbf{u})$,*

$$\text{dist}(\mathbf{x}, \mathcal{M}(\mathbf{u}')) \leq r \|\mathbf{u} - \mathbf{u}'\|.$$

Here, S_0 is a bounded polyhedral set satisfying $S = S_0 + C$, where $S = \{(\lambda_1, \lambda_2) : \|\mathbf{A}^\top \lambda_1 + \lambda_2\|_1 \leq 1\}$ is a polyhedron set that depends only on \mathbf{A} , and C is a polyhedral cone.

Proof (Proof of Lemma 1) We proceed by backward induction from the terminal stage. At stage T , we have $Q_T(\mathbf{x}_{T-1}, \boldsymbol{\xi}_T) = \hat{Q}_T(\mathbf{x}_{T-1}, \boldsymbol{\xi}_T) = \min \left\{ \mathbf{c}_T^\top \mathbf{x}_T : \mathbf{A}_T \mathbf{x}_T + \mathbf{B}_T \mathbf{x}_{T-1} = \mathbf{b}_T, \mathbf{x}_T \in \mathbb{R}_+^{d_T} \right\}$. Hence, it is sufficient to show the result for $Q_T(\mathbf{x}_{T-1}, \boldsymbol{\xi}_T)$. For any $\boldsymbol{\xi}_T \in \Xi_T$, let \mathbf{x}_{T-1} and \mathbf{x}'_{T-1} be two points in the domain of $Q_T(\mathbf{x}_{T-1}, \boldsymbol{\xi}_T)$. Fix any feasible point $\mathbf{x}_T \in \mathcal{X}_T(\mathbf{x}_{T-1}, \boldsymbol{\xi}_T)$. Applying Lemma 6 with right-hand side vectors $\mathbf{u}_T = \mathbf{b}_T - \mathbf{B}_T \mathbf{x}_{T-1}$ and $\mathbf{u}'_T = \mathbf{b}_T - \mathbf{B}_T \mathbf{x}'_{T-1}$, we have that there exists a feasible point $\mathbf{x}'_T \in \mathcal{X}_T(\mathbf{x}'_{T-1}, \boldsymbol{\xi}_T)$ such that

$$\|\mathbf{x}_T - \mathbf{x}'_T\| \leq r_T \|\mathbf{u}_T - \mathbf{u}'_T\| \leq r_T \|(\mathbf{b}_T - \mathbf{B}_T \mathbf{x}_{T-1}) - (\mathbf{b}_T - \mathbf{B}_T \mathbf{x}'_{T-1})\| \leq r_T \|\mathbf{B}_T\| \cdot \|\mathbf{x}_{T-1} - \mathbf{x}'_{T-1}\|,$$

where $r_T > 0$ is a constant that depends only on \mathbf{A}_T . Therefore, the Lipschitz continuity of $\mathbf{c}_T^\top \mathbf{x}_T$ implies

$$Q_T(\mathbf{x}'_{T-1}, \boldsymbol{\xi}_T) \leq \mathbf{c}_T^\top \mathbf{x}'_T \leq \mathbf{c}_T^\top \mathbf{x}_T + r_T \|\mathbf{c}_T\| \cdot \|\mathbf{B}_T\| \cdot \|\mathbf{x}_{T-1} - \mathbf{x}'_{T-1}\|.$$

Taking minimization over $\mathbf{x}_T \in \mathcal{X}_T(\mathbf{x}_{T-1}, \boldsymbol{\xi}_T)$, we have

$$Q_T(\mathbf{x}'_{T-1}, \boldsymbol{\xi}_T) \leq Q_T(\mathbf{x}_{T-1}, \boldsymbol{\xi}_T) + r_T \|\mathbf{c}_T\| \cdot \|\mathbf{B}_T\| \cdot \|\mathbf{x}_{T-1} - \mathbf{x}'_{T-1}\|.$$

By symmetry, we have

$$|Q_T(\mathbf{x}_{T-1}, \boldsymbol{\xi}_T) - Q_T(\mathbf{x}'_{T-1}, \boldsymbol{\xi}_T)| \leq r_T \|\mathbf{c}_T\| \cdot \|\mathbf{B}_T\| \cdot \|\mathbf{x}_{T-1} - \mathbf{x}'_{T-1}\|.$$

By the compactness of the uncertainty set in the assumption (A4), we know there exists $\bar{r}_T, \bar{\mathbf{c}}_T$, and $\bar{\mathbf{B}}_T$ such that the right-hand side is upper bounded by $\bar{r}_T \|\bar{\mathbf{c}}_T\| \cdot \|\bar{\mathbf{B}}_T\| \cdot \|\mathbf{x}_{T-1} - \mathbf{x}'_{T-1}\|$. Setting $L_T = \bar{r}_T \|\bar{\mathbf{c}}_T\| \cdot \|\bar{\mathbf{B}}_T\|$, we conclude that

$$|Q_T(\mathbf{x}_{T-1}, \boldsymbol{\xi}_T) - Q_T(\mathbf{x}'_{T-1}, \boldsymbol{\xi}_T)| \leq L_T \|\mathbf{x}_{T-1} - \mathbf{x}'_{T-1}\|, \quad \forall \mathbf{x}_{T-1}, \mathbf{x}'_{T-1} \in \mathcal{X}_{T-1}(\mathbf{x}_{T-2}, \boldsymbol{\xi}_{T-1}), \forall \boldsymbol{\xi}_T \in \Xi_T.$$

Hence the desired result holds for stage T . By induction, suppose the statement in the lemma holds for stage $t+1 \leq T$. That is, there exists a constant $L_{t+1} > 0$ such that

$$\begin{aligned} |Q_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}) - Q_{t+1}(\mathbf{x}'_t, \boldsymbol{\xi}_{t+1})| &\leq L_{t+1} \|\mathbf{x}_t - \mathbf{x}'_t\|, \quad \forall \mathbf{x}_t, \mathbf{x}'_t \in \mathcal{X}_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) \forall \boldsymbol{\xi}_{t+1} \in \Xi_{t+1}, \text{ and} \\ |\hat{Q}_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}) - \hat{Q}_{t+1}(\mathbf{x}'_t, \boldsymbol{\xi}_{t+1})| &\leq L_{t+1} \|\mathbf{x}_t - \mathbf{x}'_t\|, \quad \forall \mathbf{x}_t, \mathbf{x}'_t \in \mathcal{X}_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) \forall \boldsymbol{\xi}_{t+1} \in \Xi_{t+1}. \end{aligned}$$

We show the Lipschitz continuity of $Q_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$, the Lipschitz continuity of $\hat{Q}_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$ can be proved in a similar way. For any feasible points $\mathbf{x}_t, \mathbf{x}'_t \in \mathcal{X}_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$ and any $\boldsymbol{\xi}_t \in \Xi_t$, the following chain of inequalities holds for the true expected cost-to-go function.

$$\begin{aligned} |Q_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_t) - Q_{t+1}(\mathbf{x}'_t, \boldsymbol{\xi}_t)| &= \left| \int_{\boldsymbol{\xi}_{t+1} \in \Xi_{t+1}} Q_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}) f(\boldsymbol{\xi}_{t+1} \mid \boldsymbol{\xi}_t) d\boldsymbol{\xi}_{t+1} \right. \\ &\quad \left. - \int_{\boldsymbol{\xi}_{t+1} \in \Xi_{t+1}} Q_{t+1}(\mathbf{x}'_t, \boldsymbol{\xi}_{t+1}) f(\boldsymbol{\xi}_{t+1} \mid \boldsymbol{\xi}_t) d\boldsymbol{\xi}_{t+1} \right| \\ &\leq \int_{\boldsymbol{\xi}_{t+1} \in \Xi_{t+1}} |Q_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}) - Q_{t+1}(\mathbf{x}'_t, \boldsymbol{\xi}_{t+1})| f(\boldsymbol{\xi}_{t+1} \mid \boldsymbol{\xi}_t) d\boldsymbol{\xi}_{t+1} \\ &\leq L_{t+1} \|\mathbf{x}_t - \mathbf{x}'_t\|. \end{aligned}$$

Therefore, $\mathbf{c}_t^\top \mathbf{x}_t + Q_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_t)$ is a Lipschitz function with constant $L_{t+1} + \|\mathbf{c}_t\|$. For any $\boldsymbol{\xi}_t \in \Xi_t$, let \mathbf{x}_{t-1} and \mathbf{x}'_{t-1} be two points in the domain of $Q_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$. Fix any feasible point $\mathbf{x}_t \in \mathcal{X}_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$. Applying Lemma 6 with right-hand side vectors $\mathbf{u}_t = \mathbf{b}_t - \mathbf{B}_t \mathbf{x}_{t-1}$ and $\mathbf{u}'_t = \mathbf{b}_t - \mathbf{B}_t \mathbf{x}'_{t-1}$, we have that there exists a feasible point $\mathbf{x}'_t \in \mathcal{X}_t(\mathbf{x}'_{t-1}, \boldsymbol{\xi}_t)$ such that

$$\|\mathbf{x}_t - \mathbf{x}'_t\| \leq r_t \|\mathbf{u}_t - \mathbf{u}'_t\| \leq r_t \|(\mathbf{b}_t - \mathbf{B}_t \mathbf{x}_{t-1}) - (\mathbf{b}_t - \mathbf{B}_t \mathbf{x}'_{t-1})\| \leq r_t \|\mathbf{B}_t\| \cdot \|\mathbf{x}_{t-1} - \mathbf{x}'_{t-1}\|$$

where $r_t > 0$ is a constant that depends only on \mathbf{A}_t . Therefore, the Lipschitz continuity of $\mathbf{c}_t^\top \mathbf{x}_t + Q_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_t)$ implies

$$\begin{aligned} Q_t(\mathbf{x}'_{t-1}, \boldsymbol{\xi}_t) &\leq \mathbf{c}_t^\top \mathbf{x}'_t + Q_{t+1}(\mathbf{x}'_t, \boldsymbol{\xi}_t) \\ &\leq \mathbf{c}_t^\top \mathbf{x}_t + Q_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_t) + r_t (L_{t+1} + \|\mathbf{c}_t\|) \cdot \|\mathbf{B}_t\| \cdot \|\mathbf{x}_{t-1} - \mathbf{x}'_{t-1}\|. \end{aligned}$$

Taking minimization over $\mathbf{x}_t \in \mathcal{X}_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$, we have

$$Q_t(\mathbf{x}'_{t-1}, \boldsymbol{\xi}_t) \leq Q_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) + r_t (L_{t+1} + \|\mathbf{c}_t\|) \cdot \|\mathbf{B}_t\| \cdot \|\mathbf{x}_{t-1} - \mathbf{x}'_{t-1}\|.$$

By symmetry, we have

$$|Q_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) - Q_t(\mathbf{x}'_{t-1}, \boldsymbol{\xi}_t)| \leq r_t(L_{t+1} + \|\mathbf{c}_t\|) \cdot \|\mathbf{B}_t\| \cdot \|\mathbf{x}_{t-1} - \mathbf{x}'_{t-1}\|.$$

By the compactness of the uncertainty set in the assumption (A4), we know there exists $\bar{r}_t, \bar{\mathbf{c}}_t$, and $\bar{\mathbf{B}}_t$ such that the right-hand side is upper bounded by $\bar{r}_t(L_{t+1} + \|\bar{\mathbf{c}}_t\|) \cdot \|\bar{\mathbf{B}}_t\| \cdot \|\mathbf{x}_{t-1} - \mathbf{x}'_{t-1}\|$. Setting $L_t = \bar{r}_t(L_{t+1} + \|\bar{\mathbf{c}}_t\|) \cdot \|\bar{\mathbf{B}}_t\|$, we conclude that

$$|Q_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) - Q_t(\mathbf{x}'_{t-1}, \boldsymbol{\xi}_t)| \leq L_t \|\mathbf{x}_{t-1} - \mathbf{x}'_{t-1}\|, \quad \forall \mathbf{x}_{t-1}, \mathbf{x}'_{t-1} \in \mathcal{X}_{t-1}(\mathbf{x}_{t-2}, \boldsymbol{\xi}_{t-1}), \forall \boldsymbol{\xi}_t \in \Xi_T.$$

Therefore, the result holds for $Q_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$. We omit the proof for $\hat{Q}_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$ since it is similar to $Q_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$. Hence, the desired result holds for stage t and this completes the induction. \blacksquare

A.2 Proof of Theorem 2

Proof Based on the assumptions (A3) and (A4), there exists a positive constant D_t satisfying

$$\sup_{\mathbf{x}_t, \mathbf{x}'_t \in \mathcal{X}_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)} \|\mathbf{x}_t - \mathbf{x}'_t\| \leq D_t, \quad \forall \mathbf{x}_{t-1}, \forall \boldsymbol{\xi}_t,$$

i.e., $\mathcal{X}_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) \subseteq \mathbb{R}^{d_t}$ has a finite diameter less than or equal to D_t for any \mathbf{x}_{t-1} and $\boldsymbol{\xi}_t$.

Next, we define a finite set of points $\mathcal{X}_t^\eta(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) \subseteq \mathcal{X}_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$. For any $\mathbf{x}_t \in \mathcal{X}_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$, there exists $\mathbf{x}'_t \in \mathcal{X}_t^\eta(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$ such that $\|\mathbf{x}_t - \mathbf{x}'_t\| \leq \eta$, i.e., for a fixed tolerance level η , we have the cardinality $|\mathcal{X}_t^\eta(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)| = O(1)(D_t/\eta)^{d_t}$. Since the cost-to-go function $Q_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1})$ is L_{t+1} -Lipschitz continuous in \mathbf{x}_t , Lemma 2 implies that for any $\mathbf{x}_t \in \mathcal{X}_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$, there exists $\mathbf{x}'_t \in \mathcal{X}_t^\eta(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$, $\|\mathbf{x}_t - \mathbf{x}'_t\| \leq \eta$, such that:

$$\left| \mathbb{E} \left[Q_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t \right] - \mathbb{E} \left[Q_{t+1}(\mathbf{x}'_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t \right] \right| \leq L_{t+1}\eta. \quad (\text{A.1})$$

Furthermore, from Theorem 1, we have that for a fixed $\mathbf{x}'_t \in \mathcal{X}_t^\eta(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$,

$$\left| \mathbb{E} \left[Q_{t+1}(\mathbf{x}'_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t \right] - \hat{\mathbb{E}} \left[Q_{t+1}(\mathbf{x}'_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t \right] \right| \leq \sqrt{\frac{\mathbb{V} \left[Q_{t+1}(\mathbf{x}'_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t \right] \log \left(\frac{1}{\delta_{t+1}} \right)}{O(N^{\frac{4}{p+4}})(1+o(1))g(\boldsymbol{\xi}_t)}} \quad (\text{A.2})$$

with probability at least $1 - \delta_{t+1}$. Applying union bound, we get

$$\begin{aligned} \left| \mathbb{E} \left[Q_{t+1}(\mathbf{x}'_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t \right] - \hat{\mathbb{E}} \left[Q_{t+1}(\mathbf{x}'_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t \right] \right| &\leq \sqrt{\frac{\mathbb{V} \left[Q_{t+1}(\mathbf{x}'_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t \right] \log \left(\frac{|\mathcal{X}_t^\eta(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)|}{\delta_{t+1}} \right)}{O(N^{\frac{4}{p+4}})(1+o(1))g(\boldsymbol{\xi}_t)}} \\ &\leq \sqrt{\frac{\sigma_{t+1}^2 \frac{\log \left(\frac{O(1)(D_t/\eta)^{d_t}}{\delta_{t+1}} \right)}{O(N^{\frac{4}{p+4}})(1+o(1))g(\boldsymbol{\xi}_t)}}{O(N^{\frac{4}{p+4}})(1+o(1))g(\boldsymbol{\xi}_t)}} \quad \forall \mathbf{x}'_t \in \mathcal{X}_t^\eta(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) \end{aligned} \quad (\text{A.3})$$

with probability at least $1 - \delta_{t+1}$. Using the Lipschitz continuity of $Q_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1})$, from Lemma 2, we get

$$\left| \mathbb{E} \left[Q_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t \right] - \hat{\mathbb{E}} \left[Q_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t \right] \right| \leq \sqrt{\frac{\sigma_{t+1}^2 \frac{\log \left(\frac{O(1)(D_t/\eta)^{d_t}}{\delta_{t+1}} \right)}{O(N^{\frac{4}{p+4}})(1+o(1))g(\boldsymbol{\xi}_t)}}{O(N^{\frac{4}{p+4}})(1+o(1))g(\boldsymbol{\xi}_t)}} + 2L_{t+1}\eta$$

$\forall \mathbf{x}_t \in \mathcal{X}_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$

with probability at least $1 - \delta_{t+1}$. This completes the proof. \blacksquare

B Proofs in Section 3

B.1 Proof of Proposition 1

Proof The proof of this proposition follows from the approach discussed in [44]. To simplify the notation, we define a random variable $\tilde{z} = \hat{Q}_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1})$ and a vector $\mathbf{z} \in \mathbb{R}^N$ where $z_i = \hat{Q}_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i)$. We denote

$$\begin{aligned} \bar{z} &= \hat{\mathbb{E}} \left[\hat{Q}_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t \right] = \sum_{i \in [N]} \hat{w}_t^i(\boldsymbol{\xi}_t) \cdot \hat{Q}_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i) \\ s &= \hat{\mathbb{V}} \left[\hat{Q}_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t \right] = \sum_{i \in [N]} \hat{w}_t^i(\boldsymbol{\xi}_t) \cdot \left[\hat{Q}_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i) - \bar{z} \right]^2. \end{aligned}$$

The DRO problem $\max_{\mathbb{P}_t \in \mathcal{P}_t^\lambda(\tilde{\mathbb{P}}_t)} \mathbb{E}_{\mathbb{P}_t} [\hat{Q}_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t]$ can be equivalently written as

$$\max_{\mathbf{w}_t} \left\{ \mathbf{w}_t^\top \mathbf{z} : \sum_{i \in [N]} \left| \frac{w_t^i - \hat{w}_t^i(\boldsymbol{\xi}_t)}{\sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)}} \right| \leq \sqrt{N}\lambda, \max_{i \in [N]} \left| \frac{w_t^i - \hat{w}_t^i(\boldsymbol{\xi}_t)}{\sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)}} \right| \leq \lambda, \mathbf{e}^\top \mathbf{w}_t = 1, \mathbf{w}_t \in \mathbb{R}_+^N \right\},$$

which is lower bounded by

$$\max_{\mathbf{w}_t} \left\{ \mathbf{w}_t^\top \mathbf{z} : \sum_{i \in [N]} \frac{[w_t^i - \hat{w}_t^i(\boldsymbol{\xi}_t)]^2}{\hat{w}_t^i(\boldsymbol{\xi}_t)} \leq \lambda^2, \mathbf{e}^\top \mathbf{w}_t = 1, \mathbf{w}_t \in \mathbb{R}_+^N \right\}.$$

By change of variable $\mathbf{u}_t = \mathbf{w}_t - \hat{\mathbf{w}}_t(\boldsymbol{\xi}_t)$, the above problem is equivalent to

$$\max_{\mathbf{u}_t} \left\{ \bar{z} + \mathbf{u}_t^\top (\mathbf{z} - \bar{z} \cdot \mathbf{e}) : \|\mathbf{u}_t\|_W \leq \lambda, \mathbf{e}^\top \mathbf{u}_t = 0, \mathbf{u}_t + \hat{\mathbf{w}}_t(\boldsymbol{\xi}_t) \geq \mathbf{0} \right\},$$

where $\|\mathbf{u}_t\|_W := \sqrt{\sum_{i \in [N]} \frac{1}{\hat{w}_t^i(\boldsymbol{\xi}_t)} (u_t^i)^2}$ is defined to be a weighted norm. We further define its dual norm $\|\mathbf{u}_t\|_{W^{-1}} := \sqrt{\sum_{i \in [N]} \hat{w}_t^i(\boldsymbol{\xi}_t) \cdot (u_t^i)^2}$, and the upper bound of the above optimization problem is

$$\bar{z} + \mathbf{u}_t^\top (\mathbf{z} - \bar{z} \cdot \mathbf{e}) \leq \bar{z} + \|\mathbf{u}_t\|_W \|\mathbf{z} - \bar{z} \cdot \mathbf{e}\|_{W^{-1}} \leq \bar{z} + \lambda \|\mathbf{z} - \bar{z} \cdot \mathbf{e}\|_{W^{-1}} = \bar{z} + \lambda \sqrt{s},$$

where the last equality holds because

$$\|\mathbf{z} - \bar{z} \cdot \mathbf{e}\|_{W^{-1}} = \sqrt{\sum_{n \in [N]} \hat{w}_t^i(\boldsymbol{\xi}_t) \cdot (z_i - \bar{z})^2} = \sqrt{\hat{\mathbb{V}} [\hat{Q}_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t]}.$$

The above upper bound can be achieved by selecting

$$u_t^i = \frac{\lambda \hat{w}_t^i(\boldsymbol{\xi}_t) \cdot (z_i - \bar{z})}{\sqrt{s}}.$$

The above choice of \mathbf{u}_t satisfies the constraints $\|\mathbf{u}_t\|_W^2 \leq \lambda$ and $\mathbf{u}_t^\top \mathbf{e} = 0$. Therefore, such \mathbf{u}_t is feasible as long as

$$u_t^i = \frac{\lambda \hat{w}_t^i(\boldsymbol{\xi}_t) \cdot (z_i - \bar{z})}{\sqrt{s}} \geq -\hat{w}_t^i(\boldsymbol{\xi}_t) \iff \frac{\lambda(z_i - \bar{z})}{\sqrt{s}} \geq -1.$$

By Lipschitz continuity of the cost-to-go function and compactness of the feasible region according to the assumption **(A3)**, $\hat{Q}_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i)$ is bounded. Denote \bar{U}_{t+1} to be the upper bound, that is, $\hat{Q}_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i) \leq \bar{U}_{t+1}, \forall \boldsymbol{\xi}_{t+1}^i \in \tilde{\Xi}_{t+1}$. Hence, $|z_i - \bar{z}| = |\hat{Q}_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i) - \sum_{n \in [N]} \hat{w}_t^n \cdot \hat{Q}_{t+1}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^n)| \leq \bar{U}_{t+1}$, then a sufficient condition of the above is

$$\frac{\lambda^2 \bar{U}_{t+1}^2}{s} \leq 1 \iff s \geq \lambda^2 \bar{U}_{t+1}^2 \iff \lambda \sqrt{s} \geq \lambda^2 \bar{U}_{t+1}.$$

if $s - \lambda^2 \bar{U}_{t+1}^2 \geq 0$, \mathbf{u}_t is a feasible solution. On the other hand, $\mathbf{u}_t = \mathbf{0}$ is another feasible solution for this problem. Thus

$$\begin{aligned} \bar{z} + \left(\lambda \sqrt{\hat{\mathbb{V}} [\hat{Q}_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t]} - \lambda^2 \bar{U}_{t+1} \right) &\leq \max_{\mathbf{w}_t} \left\{ \mathbf{w}_t^\top \mathbf{z} : \sum_{i \in [N]} \frac{(w_t^i - \hat{w}_t^i(\boldsymbol{\xi}_t))^2}{\hat{w}_t^i(\boldsymbol{\xi}_t)} \leq \lambda^2, \mathbf{e}^\top \mathbf{w}_t = 1, \mathbf{w}_t \in \mathbb{R}_+^N \right\} \\ &\leq \max_{\mathbb{P}_t \in \mathcal{P}_t^\lambda(\tilde{\mathbb{P}}_t)} \mathbb{E}_{\mathbb{P}_t} [\hat{Q}_{t+1}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t]. \end{aligned}$$

Rearranging the terms, we complete the proof. \blacksquare

B.2 Proof of Proposition 2

Proof We consider the inner maximization problem, which given a feasible solution $\mathbf{x}_t \in \mathcal{X}_t(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$, yields the distribution with the worst-case expected loss as given below

$$\max_{\mathbb{P}_t \in \mathcal{P}_t^\lambda(\tilde{\mathbb{P}}_t)} \hat{Q}_{t+1}^{\text{DRO}}(\mathbf{x}_t, \boldsymbol{\xi}_t).$$

To simplify the notation, we define the vector $\mathbf{z} \in \mathbb{R}^i$ whose n -th component $z_i = \hat{Q}_{t+1}^{\text{DRO}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i)$ denotes the loss function evaluated for the n -th data point. We then have the following linear programming formulation for the worst-case expected

loss:

$$\begin{aligned}
\max \quad & \sum_{n \in [N]} z_i \hat{w}_t^i(\boldsymbol{\xi}_t) \\
\text{s.t.} \quad & \mathbf{w}_t \in \mathbb{R}_+^N, \quad \mathbf{f} \in \mathbb{R}^N, \\
& \mathbf{e}^\top \mathbf{w}_t = 1 && : (\gamma), \\
& \frac{1}{\sqrt{N}} \sum_{n \in [N]} f_i \leq \lambda && : (\beta), \\
& \frac{w_t^i - \hat{w}_t^i(\boldsymbol{\xi}_t)}{\sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)}} \leq f_i \quad \forall i \in [N] && : (\mu), \\
& \frac{-w_t^i + \hat{w}_t^i(\boldsymbol{\xi}_t)}{\sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)}} \leq f_i \quad \forall i \in [N] && : (\zeta), \\
& f_i \leq \lambda \quad \forall i \in [N] && : (\psi).
\end{aligned} \tag{B.1}$$

Strong linear programming duality holds because the ambiguity set (4.5) is always nonempty with the nominal distribution as the trivial solution. Therefore, the dual problem can be written as:

$$\begin{aligned}
\min \quad & \gamma + \lambda \left(\beta + \sum_{n \in [N]} \psi_i \right) + \sum_{n \in [N]} \sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)} (\mu_i - \zeta_i) \\
\text{s.t.} \quad & \gamma \in \mathbb{R}, \quad \beta \in \mathbb{R}_+, \quad \boldsymbol{\mu}, \boldsymbol{\zeta}, \boldsymbol{\psi} \in \mathbb{R}_+^N, \\
& z_i \leq \gamma + \frac{\mu_i - \zeta_i}{\sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)}} \quad \forall i \in [N], \\
& \mu_i + \zeta_i = \psi_i + \frac{\beta}{\sqrt{N}} \quad \forall i \in [N].
\end{aligned} \tag{B.2}$$

Combining with the outer minimization problem, we have the following desired result:

$$\begin{aligned}
\min \quad & \mathbf{c}_t^\top \mathbf{x}_t + \gamma + \lambda \left(\beta + \sum_{n \in [N]} \psi_i \right) + \sum_{n \in [N]} \sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)} (\mu_i - \zeta_i) \\
\text{s.t.} \quad & \mathbf{x}_t \in \mathbb{R}_+^{d_t}, \quad \gamma \in \mathbb{R}, \quad \beta \in \mathbb{R}_+, \quad \boldsymbol{\mu}, \boldsymbol{\zeta}, \boldsymbol{\psi} \in \mathbb{R}_+^N, \\
& \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{x}_{t-1} = \mathbf{b}_t, \\
& \hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}) \leq \gamma + \frac{\mu_i - \zeta_i}{\sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)}} \quad \forall i \in [N], \\
& \mu_i + \zeta_i = \psi_i + \frac{\beta}{\sqrt{N}} \quad \forall i \in [N].
\end{aligned} \tag{B.3}$$

■

B.3 Proof of Lemma 4

Proof We proceed by backward induction in the stages.

Stage $t = T$: For any $\boldsymbol{\xi}_T \in \hat{\Xi}_T$, using the assumption (A3) and strong duality, we have

$$\begin{aligned}
\hat{Q}_T^{\text{DR}\mathcal{O}}(\mathbf{x}_{T-1}, \boldsymbol{\xi}_T) &= \min_{\mathbf{x}_T} \left\{ \mathbf{c}_T^\top \mathbf{x}_T : \mathbf{A}_T \mathbf{x}_T + \mathbf{B}_T \mathbf{x}_{T-1} = \mathbf{b}_T, \mathbf{x}_T \in \mathbb{R}_+^{d_T} \right\} \\
&= \max_{\boldsymbol{\pi}_T} \left\{ \boldsymbol{\pi}_T^\top (\mathbf{b}_T - \mathbf{B}_T \mathbf{x}_{T-1}) : \mathbf{A}_T^\top \boldsymbol{\pi}_T \leq \mathbf{c}_T, \boldsymbol{\pi}_T \in \mathbb{R}^{d_T} \right\} \\
&= \max \left\{ \boldsymbol{\pi}_{T,j}^\top (\mathbf{b}_T - \mathbf{B}_T \mathbf{x}_{T-1}) : j \in [\mathcal{S}_T] \right\}
\end{aligned}$$

where $\boldsymbol{\pi}_{T,j}$, $j \in [\mathcal{S}_T]$ are the extreme points of the dual problem. Note that the number of dual extreme points in \mathcal{S}_T is finite. Therefore the claim holds for the stage T .

Suppose $\hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1})$ is piecewise linear convex with finite pieces in \mathbf{x}_t for any $\boldsymbol{\xi}_{t+1} \in \hat{\Xi}_{t+1}$. For a specific $\boldsymbol{\xi}_{t+1}^i \in \hat{\Xi}_{t+1}$, $i \in [N]$, let \mathcal{P}_{t+1}^i denote the set of all finite pieces of $\hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i)$, i.e.,

$$\mathcal{P}_{t+1}^i = \max \left\{ \mathcal{G}_{t,j}^i \top \mathbf{x}_t + \mathcal{V}_{t,j}^i : j \in [\mathcal{S}_{t+1}^i] \right\}$$

where $\mathcal{S}_{t+1}^i < \infty$ is the total number of piecewise linear functions that describe $\hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i)$. Now for any $\boldsymbol{\xi}_t \in \hat{\Xi}_t$, we can rewrite (4.6) as

$$\begin{aligned} \hat{Q}_t^{\text{DR}\mathcal{O}}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) = \min & \mathbf{c}_t^\top \mathbf{x}_t + \gamma + \lambda \left(\beta + \sum_{n \in [N]} \psi_n \right) + \sum_{n \in [N]} \sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)} (\mu_i - \zeta_i) \\ \text{s.t. } & \mathbf{x}_t \in \mathbb{R}_+^{d_t}, \quad \gamma \in \mathbb{R}, \quad \beta \in \mathbb{R}_+, \quad \boldsymbol{\mu}, \boldsymbol{\zeta}, \boldsymbol{\psi} \in \mathbb{R}_+^N, \\ & \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{x}_{t-1} = \mathbf{b}_t, & : (\boldsymbol{\pi}_t), \\ & \mathcal{G}_{t,j}^i{}^\top \mathbf{x}_t + \mathcal{V}_{t,j}^i \leq \gamma + \frac{\mu_i - \zeta_i}{\sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)}} \quad \forall i \in [N], \forall j \in [\mathcal{S}_{t+1}^i] & : (y_{t,j}^i), \\ & \mu_i + \zeta_i = \psi_i + \frac{\beta}{\sqrt{N}} \quad \forall i \in [N] & : (r_t^i). \end{aligned} \quad (\text{B.4})$$

Under assumption (A3), problem (B.4) has a finite optimal solution. The corresponding dual problem is

$$\begin{aligned} \hat{Q}_t^{\text{DR}\mathcal{O}}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) = \max & \boldsymbol{\pi}_t^\top (\mathbf{b}_t - \mathbf{B}_t \mathbf{x}_{t-1}) + \sum_{i \in [N]} \sum_{j \in [\mathcal{S}_{t+1}^i]} y_{t,j}^i \mathcal{V}_{t,j}^i \\ \text{s.t. } & \boldsymbol{\pi}_t \in \mathbb{R}^M, \quad \mathbf{r}_t \in \mathbb{R}^N, \quad y_{t,j}^i \in \mathbb{R}_+ \quad \forall i \in [N], \forall j \in [\mathcal{S}_{t+1}^i], \\ & \mathbf{A}_t^\top \boldsymbol{\pi}_t \leq \mathbf{c}_t + \sum_{i \in [N]} \sum_{j \in [\mathcal{S}_{t+1}^i]} y_{t,j}^i \mathcal{G}_{t,j}^i, \\ & \sum_{i \in [N]} \sum_{j \in [\mathcal{S}_{t+1}^i]} y_{t,j}^i = 1, \\ & \frac{\sum_{i \in [N]} r_t^i}{\sqrt{N}} \leq \lambda, \\ & \frac{\sum_{j \in [\mathcal{S}_{t+1}^i]} y_{t,j}^i}{\sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)}} \leq \sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)} + r_t^i, \\ & \sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)} \leq \frac{\sum_{j \in [\mathcal{S}_{t+1}^i]} y_{t,j}^i}{\sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)}} + r_t^i, \\ & r_t^i \leq \lambda. \end{aligned} \quad (\text{B.5})$$

Note that here $\{\sum_{j \in [\mathcal{S}_{t+1}^i]} y_{t,j}^i, i \in [N]\}$ is actually the worst-case distribution. Using a finite number, \mathcal{S}_t , of extreme points of the feasible region in (B.5), we can write the above problem, as follows:

$$\hat{Q}_t^{\text{DR}\mathcal{O}}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) = \max \left\{ \boldsymbol{\pi}_{t,l}^\top (\mathbf{b}_t - \mathbf{B}_t \mathbf{x}_{t-1}) + \sum_{i \in [N]} \sum_{j \in [\mathcal{S}_{t+1}^i]} (y_{t,j}^i)_l \mathcal{V}_{t,j}^i : l \in [\mathcal{S}_t] \right\} \quad (\text{B.6})$$

which implies $\hat{Q}_t^{\text{DR}\mathcal{O}}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$ is piecewise linear convex with finite pieces in \mathbf{x}_{t-1} for any $\boldsymbol{\xi}_t \in \hat{\Xi}_t$. This completes the proof. \blacksquare

B.4 Proof of Lemma 5

Proof We proceed by backward induction in the stages. Before we start, we present the lower bound problem at the k -th iteration:

$$\begin{aligned} \underline{Q}_{t,k}^{\text{DR}\mathcal{O}}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) = \min & \mathbf{c}_t^\top \mathbf{x}_t + \gamma + \lambda \left(\beta + \sum_{n \in [N]} \psi_n \right) + \sum_{n \in [N]} \sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)} (\mu_i - \zeta_i) \\ \text{s.t. } & \mathbf{x}_t \in \mathbb{R}_+^{d_t}, \quad \gamma \in \mathbb{R}, \quad \beta \in \mathbb{R}_+, \quad \boldsymbol{\mu}, \boldsymbol{\zeta}, \boldsymbol{\psi} \in \mathbb{R}_+^N \\ & \mathbf{A}_t \mathbf{x}_t + \mathbf{B}_t \mathbf{x}_{t-1} = \mathbf{b}_t, & : (\boldsymbol{\pi}_t), \\ & \mathcal{G}_{t,j}^i{}^\top \mathbf{x}_t + \mathcal{V}_{t,j}^i \leq \gamma + \frac{\mu_i - \zeta_i}{\sqrt{\hat{w}_t^i(\boldsymbol{\xi}_t)}} \quad \forall i \in [N], \forall j \in [\mathcal{S}_{t+1}^{i,k}] & : (y_{t,j}^i), \\ & \mu_i + \zeta_i = \psi_i + \frac{\beta}{\sqrt{N}} \quad \forall i \in [N] & : (r_t^i). \end{aligned} \quad (\text{B.7})$$

Here, we use $\mathcal{S}_{t+1}^{i,k}$ to denote the total number of cuts generated in the first k iterations for $\hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i)$. Note in (B.4), we use \mathcal{S}_{t+1}^i to the total number of piecewise linear functions that describe $\hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i)$. Now we proceed to the main proof.

Stage $t = T$: Since the cost-to-go function $\hat{Q}_{T+1}^{\text{DR}\mathcal{O}}(\cdot) \equiv 0$, there are no cuts. For any feasible \mathbf{x}_{T-1} and $\boldsymbol{\xi}_T \in \hat{\Xi}_T$, let $\hat{\boldsymbol{\pi}}_T$ be a dual feasible extreme point (here $\hat{\boldsymbol{\pi}}_T$ is the only dual variable) to problem (B.7), we have

$$\begin{aligned} \mathcal{G}_{T-1}^\top \mathbf{x}_{T-1} + \mathcal{V}_{T-1} &= \hat{\boldsymbol{\pi}}_T^\top (\mathbf{b}_T - \mathbf{B}_T \mathbf{x}_{T-1}) \\ &\leq \max \left\{ \boldsymbol{\pi}_T^\top (\mathbf{b}_T - \mathbf{B}_T \mathbf{x}_{T-1}) : \mathbf{A}_T^\top \boldsymbol{\pi}_T \leq \mathbf{c}_T, \boldsymbol{\pi}_T \in \mathbb{R}^{d_T} \right\} \\ &= \min \left\{ \mathbf{c}_T^\top \mathbf{x}_T : \mathbf{A}_T \mathbf{x}_T + \mathbf{B}_T \mathbf{x}_{T-1} = \mathbf{b}_T, \mathbf{x}_T \in \mathbb{R}_+^{d_T} \right\} \\ &= \hat{Q}_T^{\text{DR}\mathcal{O}}(\mathbf{x}_{T-1}, \boldsymbol{\xi}_T) \end{aligned}$$

where \mathcal{G}_{T-1} is the slope and \mathcal{V}_{T-1} is the intercept. The second equality holds by strong duality. This proves the validity of the cut in stage $T-1$. Hence,

$$\underline{Q}_{T-1,k}^{\text{DR}\mathcal{O}}(\mathbf{x}_{T-2}, \boldsymbol{\xi}_{T-1}) \leq \hat{Q}_{T-1}^{\text{DR}\mathcal{O}}(\mathbf{x}_{T-2}, \boldsymbol{\xi}_{T-1}). \quad (\text{B.8})$$

For the upper bound problem, first we have $\bar{Q}_T^{\text{DR}\mathcal{O}}(\cdot, \boldsymbol{\xi}_T) = \hat{Q}_T^{\text{DR}\mathcal{O}}(\cdot, \boldsymbol{\xi}_T)$ for $\boldsymbol{\xi}_T \in \hat{\Xi}_T$ due to the fact that there are no cost-to-go functions at stage T . By the convexity of $\hat{Q}_T^{\text{DR}\mathcal{O}}(\cdot, \boldsymbol{\xi}_T)$, we have

$$\hat{Q}_T^{\text{DR}\mathcal{O}}(\mathbf{x}_{T-1}, \boldsymbol{\xi}_T^i) - \mathcal{F}_T^\top \mathbf{y}_{T-1}^i \leq \hat{Q}_T^{\text{DR}\mathcal{O}}\left(\sum_{j \in [k]} \theta_j^i \bar{\mathbf{x}}_{T-1}^j, \boldsymbol{\xi}_T^i\right) \leq \sum_{j \in [k]} \theta_j^i \hat{Q}_T^{\text{DR}\mathcal{O}}(\bar{\mathbf{x}}_{T-1}^j, \boldsymbol{\xi}_T^i), \quad (\text{B.9})$$

where \mathcal{F}_T is any subgradient of $\hat{Q}_T^{\text{DR}\mathcal{O}}(\mathbf{x}_{T-1}, \boldsymbol{\xi}_T^i)$ at \mathbf{x}_{T-1} . Therefore,

$$\begin{aligned} \sum_{j \in [k]} \theta_j^i \bar{Q}_{T,j}^{\text{DR}\mathcal{O}}(\bar{\mathbf{x}}_{T-1}^j, \boldsymbol{\xi}_T^i) + M_T \|\mathbf{y}_{T-1}^i\|_1 &= \sum_{j \in [k]} \theta_j^i \hat{Q}_T^{\text{DR}\mathcal{O}}(\bar{\mathbf{x}}_{T-1}^j, \boldsymbol{\xi}_T^i) + M_T \|\mathbf{y}_{T-1}^i\|_1 \\ &\geq \sum_{j \in [k]} \theta_j^i \hat{Q}_T^{\text{DR}\mathcal{O}}(\bar{\mathbf{x}}_{T-1}^j, \boldsymbol{\xi}_T^i) + \|\mathcal{F}_T^\top\|_1 \cdot \|\mathbf{y}_{T-1}^i\|_1 \\ &\geq \sum_{j \in [k]} \theta_j^i \hat{Q}_T^{\text{DR}\mathcal{O}}(\bar{\mathbf{x}}_{T-1}^j, \boldsymbol{\xi}_T^i) + \|\mathcal{F}_T^\top\|_2 \cdot \|\mathbf{y}_{T-1}^i\|_2 \\ &\geq \hat{Q}_T^{\text{DR}\mathcal{O}}\left(\sum_{j \in [k]} \theta_j^i \bar{\mathbf{x}}_{T-1}^j, \boldsymbol{\xi}_T^i\right) + \mathcal{F}_T^\top \mathbf{y}_{T-1}^i \\ &\geq \hat{Q}_T^{\text{DR}\mathcal{O}}(\mathbf{x}_{T-1}, \boldsymbol{\xi}_T^i). \end{aligned}$$

The first inequality holds because $M_T \geq L_T^{\text{DR}\mathcal{O}}$ and thus any subgradient of $\hat{Q}_T^{\text{DR}\mathcal{O}}(\mathbf{x}_{T-1}, \boldsymbol{\xi}_T^i)$ under 1-norm is bounded by M_T . The second inequality comes from the fact that $\|\mathbf{y}\|_2 \leq \|\mathbf{y}\|_1$ for any \mathbf{y} . The convexity and the Cauchy-Schwarz inequality imply the third inequality, whereas the last inequality comes from (B.9). Hence,

$$\hat{Q}_{T-1}^{\text{DR}\mathcal{O}}(\mathbf{x}_{T-2}, \boldsymbol{\xi}_{T-1}) \leq \bar{Q}_{T-1,k}^{\text{DR}\mathcal{O}}(\mathbf{x}_{T-2}, \boldsymbol{\xi}_{T-1}).$$

Combining this with (B.8), we have

$$\underline{Q}_{T-1,k}^{\text{DR}\mathcal{O}}(\mathbf{x}_{T-2}, \boldsymbol{\xi}_{T-1}) \leq \hat{Q}_{T-1}^{\text{DR}\mathcal{O}}(\mathbf{x}_{T-2}, \boldsymbol{\xi}_{T-1}) \leq \bar{Q}_{T-1,k}^{\text{DR}\mathcal{O}}(\mathbf{x}_{T-2}, \boldsymbol{\xi}_{T-1}).$$

Therefore the claim holds for stage $T-1$.
Suppose the result holds for $t+1 \leq T-1$,

$$\underline{Q}_{t+1,k}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}) \leq \hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}) \leq \bar{Q}_{t+1,k}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}), \quad \forall \boldsymbol{\xi}_{t+1} \in \hat{\Xi}_{t+1}.$$

This implies the cuts generated for the lower bound problem are valid for $t+1$, i.e.,

$$\max \left\{ \mathcal{G}_{t,j}^i \top \mathbf{x}_t + \mathcal{V}_{t,j}^i : j \in [\mathcal{S}_{t+1}^{i,k}] \right\} \leq \hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}), \quad \forall \boldsymbol{\xi}_{t+1} \in \hat{\Xi}_{t+1}.$$

Hence the cuts in the feasible region of (B.7) are valid. Let \mathcal{D}_t denote the feasible region of $\hat{Q}_t^{\text{DR}\mathcal{O}}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$ in (B.4), and \mathcal{D}_t^k denote the feasible region of (B.7). Note the only difference between (B.4) and (B.7) is the set of cuts: in (B.4) \mathcal{S}_{t+1}^i represent the total number of cuts for $\hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1})$, while in (B.7) $\mathcal{S}_{t+1}^{i,k}$ represent the total number of cuts generated in the first k iterations for $\hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1})$. Let $(\hat{\boldsymbol{\pi}}_t, \hat{\mathbf{y}})$ be a dual feasible extreme point to problem (B.7). Using similar result as (B.6), we have

$$\begin{aligned} \mathcal{G}_{t-1}^\top \mathbf{x}_{t-1} + \mathcal{V}_{t-1} &= \hat{\boldsymbol{\pi}}_t^\top (\mathbf{b}_t - \mathbf{B}_t \mathbf{x}_{t-1}) + \sum_{i \in [N]} \sum_{j \in [\mathcal{S}_{t+1}^{i,k}]} \hat{\mathbf{y}}_{t,j}^i \mathcal{V}_{t,j}^i \\ &\leq \max \left\{ \boldsymbol{\pi}_t^\top (\mathbf{b}_t - \mathbf{B}_t \mathbf{x}_{t-1}) + \sum_{i \in [N]} \sum_{j \in [\mathcal{S}_{t+1}^{i,k}]} \mathbf{y}_{t,j}^i \mathcal{V}_{t,j}^i : (\boldsymbol{\pi}_t, \mathbf{y}) \in \mathcal{D}_t^k \right\} \\ &\leq \max \left\{ \boldsymbol{\pi}_t^\top (\mathbf{b}_t - \mathbf{B}_t \mathbf{x}_{t-1}) + \sum_{i \in [N]} \sum_{j \in [\mathcal{S}_{t+1}^i]} \mathbf{y}_{t,j}^i \mathcal{V}_{t,j}^i : (\boldsymbol{\pi}_t, \mathbf{y}) \in \mathcal{D}_t \right\} \\ &= \hat{Q}_t^{\text{DR}\mathcal{O}}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t). \end{aligned}$$

The first inequality holds because $(\hat{\boldsymbol{\pi}}_t, \hat{\mathbf{y}}) \in \mathcal{D}_t^k$ whereas the second inequality holds because \mathcal{D}_t includes all the cuts that define $\hat{Q}_t^{\text{DR}\mathcal{O}}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$ while the cuts indexed by \mathcal{S}_{t+1}^i in \mathcal{D}_t^k by assumption only includes parts of these valid cuts. The last equality comes from strong duality. This proves the validity of the cut in stage t . Hence,

$$\underline{Q}_{t,k}^{\text{DR}\mathcal{O}}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) \leq \hat{Q}_t^{\text{DR}\mathcal{O}}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t). \quad (\text{B.10})$$

For the upper bound problem, by hypothesis we have $\hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i) \leq \bar{Q}_{t+1,k}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i)$, $\forall \boldsymbol{\xi}_{t+1}^i \in \hat{\Xi}_{t+1}$. Now replace $\bar{Q}_{t+1,k}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i)$ with $\hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i)$ in the constraints of (4.8) and call the new problem $\bar{Q}_{t,k}^{\text{DR}\mathcal{O}L}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$, we have

$$\bar{Q}_{t,k}^{\text{DR}\mathcal{O}L}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) \leq \bar{Q}_{t,k}^{\text{DR}\mathcal{O}}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$$

Therefore, it is sufficient to show $\bar{Q}_{t,k}^{\text{DR}\mathcal{O}L}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$ is an upper bound for $\hat{Q}_{t,k}^{\text{DR}\mathcal{O}}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t)$. Denote \mathcal{F}_t to be any subgradient of $\hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i)$ at \mathbf{x}_t . We have

$$\begin{aligned} \hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i) &\leq \sum_{j \in [k]} \theta_j^i \hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\bar{\mathbf{x}}_t^j, \boldsymbol{\xi}_{t+1}^i) + \mathcal{F}_t^\top \mathbf{y}_t^i \\ &\leq \sum_{j \in [k]} \theta_j^i \hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\bar{\mathbf{x}}_t^j, \boldsymbol{\xi}_{t+1}^i) + \|\mathcal{F}_t^\top\|_2 \cdot \|\mathbf{y}_t^i\|_2 \\ &\leq \sum_{j \in [k]} \theta_j^i \hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\bar{\mathbf{x}}_t^j, \boldsymbol{\xi}_{t+1}^i) + \|\mathcal{F}_t^\top\|_1 \cdot \|\mathbf{y}_t^i\|_1 \\ &\leq \sum_{j \in [k]} \theta_j^i \hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\bar{\mathbf{x}}_t^j, \boldsymbol{\xi}_{t+1}^i) + M_t \|\mathbf{y}_t^i\|_1 \\ &= \bar{Q}_{t+1,k}^{\text{DR}\mathcal{O}L}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i). \end{aligned}$$

The first inequality comes from the convexity of $\hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i)$, whereas the second inequality comes from the Cauchy-Schwarz inequality. The third inequality holds due to the fact that $\|\mathbf{y}\|_2 \leq \|\mathbf{y}\|_1$ for any \mathbf{y} . The last inequality holds because $M_t \geq L_t^{\text{DR}\mathcal{O}}$ and thus any subgradient of $\hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \boldsymbol{\xi}_{t+1}^i)$ under 1-norm is bounded by M_t . Hence,

$$\hat{Q}_{t,k}^{\text{DR}\mathcal{O}}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) \leq \bar{Q}_{t,k}^{\text{DR}\mathcal{O}L}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) \leq \bar{Q}_{t,k}^{\text{DR}\mathcal{O}}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t).$$

Combining this with (B.10), we have

$$\underline{Q}_{t,k}^{\text{DR}\mathcal{O}}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) \leq \hat{Q}_t^{\text{DR}\mathcal{O}}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t) \leq \bar{Q}_{t,k}^{\text{DR}\mathcal{O}}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t),$$

which completes the proof. \blacksquare

B.5 Proof of Theorem 4

Proof We first discuss the lower bound convergence, using the proof of the standard SDDP algorithm presented by [55].

As discussed in Lemma 4, $\hat{Q}_t^{\text{DR}\mathcal{O}}(\cdot, \boldsymbol{\xi}_t^i) \forall i \in [N]$ are piecewise linear convex with a finite number of pieces under the assumption (B1) since there exist only a finite number of basic feasible solutions of the dual problem of $\hat{Q}_t^{\text{DR}\mathcal{O}}(\cdot, \boldsymbol{\xi}_t^i) \forall i \in [N]$. In addition, the total number of possible scenarios is finite (i.e., $\prod_{t=2}^T |\hat{\Xi}_t| = N^{T-1}$) due to the discretization, and there is a nonzero probability for any possible scenario to occur in the forward step since scenarios are generated by Monte Carlo sampling. Hence, any possible scenario occurs infinitely many times in the forward step unless the algorithm terminates.

As shown in [55], the lower bound convergence holds if policy $\bar{\mathbf{x}}_t(\boldsymbol{\xi}_{[t]})$ for the current lower bound problem satisfies the following Bellman's optimality condition,

$$\bar{\mathbf{x}}_t(\boldsymbol{\xi}_{[t]}) \in \underset{\mathbf{x}_t \in \mathcal{X}_t(\bar{\mathbf{x}}_{t-1}(\boldsymbol{\xi}_{[t-1]}), \boldsymbol{\xi}_t)}{\text{argmin}} \quad \mathbf{c}_t^\top \mathbf{x}_t + \max_{\mathbb{P}_t \in \mathcal{P}_t^\lambda(\hat{\mathbb{P}}_t)} \mathbb{E}_{\mathbb{P}_t} \left[\hat{Q}_{t+1}^{\text{DR}\mathcal{O}}(\mathbf{x}_t, \tilde{\boldsymbol{\xi}}_{t+1}) \mid \boldsymbol{\xi}_t \right] \quad (\text{B.11})$$

for every stage and possible scenario. Note that the DRO formulation (B.11) is equivalent to the single-level DRO reformulation (4.6) (here, we use (B.11) to save space). Let $\bar{\mathbf{x}}_t^k(\boldsymbol{\xi}_{[t]})$ be a policy obtained by lower bound approximation $\underline{Q}_{t,k}^{\text{DR}\mathcal{O}}(\mathbf{x}_{t-1}, \boldsymbol{\xi}_t^i) \forall t = 2, \dots, T$ and $\forall i \in [N]$ at iteration k . Suppose that (B.11) does not hold for some (or all) stage $t \in \{2, \dots, T\}$ and some (or all) possible scenario, i.e., a policy for the current lower bound approximation is not optimal for the DRO problem. Let stage t' be the largest stage t such that the function $\bar{\mathbf{x}}_{t'}^k(\boldsymbol{\xi}_{[t']})$ does not satisfy (B.11), i.e., a candidate solution $\bar{\mathbf{x}}_{t'}^k = \bar{\mathbf{x}}_{t'}^k(\boldsymbol{\xi}_{[t']})$ is suboptimal for the current scenario. Then, at some iteration $k' > k$, we add a new cut corresponding to $\bar{\mathbf{x}}_{t'}^k$, updating the current approximation $\underline{Q}_{t',k'}^{\text{DR}\mathcal{O}}(\cdot, \cdot)$. Similarly, a new cut is added until (B.11) holds for stage t' . Such cut generations continue for some stage $t < t'$ until (B.11) holds. After a sufficiently large number of iterations, (B.11) holds for every stage and possible scenario. This completes the proof for the lower bound convergence.

Now we discuss the upper bound convergence. Let \bar{k}^* be the iteration after which the lower bound convergence holds. The lower bound convergence under the assumption (B1) implies that there only exist a finite number of optimal policies after the convergence, i.e., iteration $k \geq \bar{k}^*$. In other words, for each scenario, we have only a finite number of sequences of candidate solutions $\bar{\mathbf{x}}_{[T]}$ at which the upper bound approximation is evaluated (recall that candidate solutions are obtained by solving the lower bound problem in Algorithm 2). Hence, we will show that the gap between the lower and upper approximation is closed at those finite number of candidate solutions.

Without loss of generality, let us assume that there is only one optimal policy for each scenario. Let \mathcal{S} be a set of all possible scenarios and $\bar{\mathbf{x}}_{[T]}^s$ be the optimal sequence of solutions for scenario $s \in \mathcal{S}$ after the lower bound convergence. Here, we show

$$\bar{Q}_{t,k}^{\text{DR}\mathcal{O}}(\bar{\mathbf{x}}_{t-1}^s, \boldsymbol{\xi}_t^i) = \underline{Q}_{t,k}^{\text{DR}\mathcal{O}}(\bar{\mathbf{x}}_{t-1}^s, \boldsymbol{\xi}_t^i) \forall i \in [N], s \in \mathcal{S} \text{ for some iteration } k \geq \bar{k}^* \quad (\text{B.12})$$

for every stage.

Let \bar{k}_t be the iteration for which (B.12) holds at stage t . At the terminal stage T , (B.12) holds after the upper bound approximation at stage T is evaluated at all possible candidate solutions $\bar{\mathbf{x}}_T^s \forall s \in \mathcal{S}$ since $\bar{Q}_{T+1,k}^{\text{DR}\mathcal{O}}(\cdot, \cdot) = \underline{Q}_{T+1,k}^{\text{DR}\mathcal{O}}(\cdot, \cdot) \equiv 0$. Proceeding by induction, for $t = T-1, \dots, 2$, there exists iteration $k \geq \bar{k}_{t+1}$ such that (B.12) holds at stage t because $\bar{Q}_{t+1,k}^{\text{DR}\mathcal{O}}(\bar{\mathbf{x}}_t^s, \boldsymbol{\xi}_{t+1}^i) = \underline{Q}_{t+1,k}^{\text{DR}\mathcal{O}}(\bar{\mathbf{x}}_t^s, \boldsymbol{\xi}_{t+1}^i) \forall i \in [N], s \in \mathcal{S}$. This completes the upper bound convergence proof. \blacksquare