

---

# Fast Approximate Dynamic Programming for Infinite-Horizon Markov Decision Processes

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We consider the infinite-horizon, discounted cost, optimal control of discrete-  
2 time systems with separable cost and constraint in the state and input variables.  
3 Specifically, we introduce a novel numerical scheme for implementation of the  
4 value iteration (VI) algorithm in the conjugate domain, using the linear-time  
5 Legendre transform algorithm. Detailed analyses of the convergence, complexity,  
6 and error of the proposed algorithm are provided. In particular, with a discretization  
7 of size  $X$  and  $U$  for the state and input spaces, respectively, the time complexity of  
8 each iteration in the VI algorithm reduces from  $\mathcal{O}(XU)$  to  $\mathcal{O}(X + U)$ .

## 9 1 Introduction

10 Value iteration (VI) is one of the most basic and wide-spread algorithms employed for tackling  
11 problems in reinforcement learning and optimal control [8, 25] formulated as Markov decision  
12 processes (MDPs). The VI algorithm simply involves the consecutive applications of the dynamic  
13 programming (DP) operator  $\mathcal{T}V(x_t) = \min_{u_t} \{C(x_t, u_t) + \gamma \mathbb{E}V(x_{t+1})\}$ , where  $C(x_t, u_t)$  is the  
14 cost of taking the control action  $u_t$  at the state  $x_t$ . This fixed point iteration is known to converge  
15 to the optimal value function for discount factors  $\gamma \in (0, 1)$ . However, this algorithm suffers from  
16 a high computational cost for large scale finite state spaces. For problems with a continuous state  
17 space, the DP operation becomes an infinite-dimensional optimization problem, rendering the exact  
18 implementation of the VI algorithm impossible in most cases. A common approach is to incorporate  
19 function approximation techniques and compute the output of the DP operator for a finite sample (i.e.,  
20 a discretization) of the underlying continuous state space. This approximation again suffers from a  
21 high computational cost for fine discretizations of the state space, particularly in high-dimensional  
22 problems. We refer the reader to [8, 23] for various schemes of approximate implementation of DP.

23 For some problems, however, it is possible to partially address this issue by using duality theory,  
24 i.e., approaching the minimization problem in the conjugate domain. In particular, as we will see  
25 in Section 3, the minimization in the primal domain in DP can be transformed to a simple addition  
26 in the dual domain, at the expense of three conjugate transforms. However, proper application of  
27 this transformation relies on efficient numerical algorithms for conjugation. Fortunately, such an  
28 algorithm, known as linear-time Legendre transform (LLT) has been developed in late 90s [20].  
29 Other than the classical application of LLT (and other fast algorithms for conjugate transform) in  
30 solving Hamilton-Jacobi equation [1, 12, 13], these algorithms are used in image processing [21],  
31 thermodynamics [11], and optimal transport [16].

32 The application of conjugate duality for the DP problem is not new and actually goes back to  
33 Bellman [4]. Further applications of this idea for reducing the computational complexity were  
34 later explored in [14, 17]. However, surprisingly, the application of LLT for solving discrete-time  
35 optimal control problems has been limited. In particular, in [10], the authors propose the “fast  
36 value iteration” algorithm (without a rigorous analysis of the complexity and error of the proposed

37 algorithm) for a particular class of infinite-horizon optimal control problems with state-independent  
38 stage cost  $C(x, u) = C(x)$  and linear dynamics  $x_{t+1} = Ax_t + Bu_t$ , where  $A$  is a non-negative,  
39 monotone, invertible matrix. More recently, in [18], the authors also considered the application of  
40 LLT for solving the DP operation in finite-horizon optimal control problems with a more general  
41 input-affine dynamics. In particular, they introduced the “discrete conjugate DP (d-CDP) operator,”  
42 and provided a detailed analysis of its complexity and error. As we will discuss shortly, the current  
43 study is an extension of the corresponding d-CDP algorithm that, among other things, considers  
44 infinite horizon, discounted cost problems. We also note that the algorithms developed in [15, 21]  
45 for distance transform can also potentially tackle the optimal control problems similar to the ones  
46 of interest in the current study. However, these algorithms require the stage cost to be reformulated  
47 as a convex function of the “distance” between the current and next states, which can be restrictive.  
48 Another line of work that is closely related to ours involves utilizing max-plus algebra in solving  
49 deterministic, continuous-state, continuous-time optimal control problems; see, e.g., [2, 22]. These  
50 works exploit the compatibility of the DP operation with max-plus operations, and approximate the  
51 value function as a max-plus linear combination. Recently, in [3, 5], the authors used this idea to  
52 propose an approximate value iteration algorithm for deterministic MDPs. In this regard, we note  
53 that the proposed approach in the current study also involves approximating the value function as a  
54 max-plus linear combination, namely, the maximum of a finite number of affine functions. The key  
55 difference is however that by choosing a grid-like (factorized) set of slopes for the linear terms (i.e.,  
56 the basis of the max-plus linear combination), we take advantage of linear time complexity of LLT in  
57 computing the constant terms (i.e., the coefficients of the max-plus linear combination). Moreover,  
58 unlike previous works, we use the result of our error analysis to provide concrete guidelines on how  
59 to construct the basis dynamically at each iteration of the proposed VI algorithm.

60 **Main contribution.** In this study, we focus on the optimal control of discrete-time systems, with  
61 continuous state-input space. In particular, we incorporate an approximation of VI algorithm involving  
62 discretization of the state-input space. Building upon the earlier work [18], we employ conjugate  
63 duality to speed-up the VI algorithm for problems with separable stage cost (in state and input)  
64 and input-affine dynamics. We introduce the conjugate VI (CVI) algorithm based on a modified  
65 version of the discrete conjugate DP operator introduced in [18], and extend the existing results  
66 for *infinite-horizon, discounted cost* problems, while taking into account *stochastic dynamics* and  
67 *numerical approximation of the conjugate of input-dependent stage cost* in our analysis. In particular,  
68 (i) we provide sufficient conditions for the convergence of the CVI algorithm (Theorem 3.9); (ii) we  
69 show that the CVI algorithm can achieve a linear time complexity of  $\mathcal{O}(X + U)$  in each iteration  
70 (Theorem 3.10), where  $X$  and  $U$  are the cardinality the discrete state and input spaces, respectively;  
71 (iii) we analyze the error of the CVI algorithm (Theorem 3.11), and use that result to provide specific  
72 guidelines on the dynamic construction of discrete dual domain (Section 3.5).

73 **Notations.** The standard inner product in  $\mathbb{R}^n$  and the corresponding induced 2-norm are denoted  
74 by  $\langle \cdot, \cdot \rangle$  and  $\|\cdot\|_2$ , respectively.  $\|\cdot\|_\infty$  denotes the infinity norm. We use the superscript  $d$  to denote  
75 *finite (discrete)* sets (as in  $\mathbb{X}^d$ ) and *discrete* functions (as in  $h^d : \mathbb{X}^d \rightarrow \mathbb{R}$ ). We use the superscript  $g$   
76 to denote *grid-like* finite sets (as in  $\mathbb{X}^g = \prod_{i=1}^n \mathbb{X}_i^g$ ). We also use  $\mathbb{X}_{\text{sub}}^g$  to denote the *sub-grid* of  $\mathbb{X}^g$   
77 derived by omitting the smallest and the largest elements of  $\mathbb{X}^g$  in each dimension. The cardinality  
78 of the finite set  $\mathbb{X}^d$  ( $\mathbb{X}^g$ ) is denoted by  $X$ . We use  $\tilde{h}^d : \mathbb{R}^n \rightarrow \bar{\mathbb{R}} = \mathbb{R} \cup \{\infty\}$  to denote a generic  
79 *extension* of a discrete function  $h^d$ , and  $\bar{h}^d : \mathbb{R}^n \rightarrow \bar{\mathbb{R}}$  to denote *multilinear interpolation and*  
80 *extrapolation (LERP)* extension of a discrete function with a *grid-like* domain. Let  $\mathbb{X}, \mathbb{Y}$  be two  
81 arbitrary sets in  $\mathbb{R}^n$ .  $\text{co}(\mathbb{X})$  is the convex hull of  $\mathbb{X}$ . We use  $d(\mathbb{X}, \mathbb{Y}) := \inf_{x \in \mathbb{X}, y \in \mathbb{Y}} \|x - y\|_2$  to  
82 denote the distance between  $\mathbb{X}$  and  $\mathbb{Y}$ . The one-sided Hausdorff distance *from*  $\mathbb{X}$  *to*  $\mathbb{Y}$  is defined  
83 as  $d_H(\mathbb{X}, \mathbb{Y}) := \sup_{x \in \mathbb{X}} \inf_{y \in \mathbb{Y}} \|x - y\|_2$ . Let  $h : \mathbb{R}^n \rightarrow \bar{\mathbb{R}}$  be an extended real-valued function  
84 with a non-empty effective domain  $\text{dom}(h) = \mathbb{X} := \{x \in \mathbb{R}^n : h(x) < \infty\}$ . The range of  $h$  is  
85 denoted by  $\text{Rng}(h) := \max_{x \in \mathbb{X}} h(x) - \min_{x \in \mathbb{X}} h(x)$ , and the subdifferential of  $h$  at a point  $x \in \mathbb{X}$   
86 is defined as  $\partial h(x) := \{y \in \mathbb{R}^n : h(\tilde{x}) \geq h(x) + \langle y, \tilde{x} - x \rangle, \forall \tilde{x} \in \mathbb{X}\}$ . We define  $\mathbb{L}(h) :=$   
87  $\prod_{i=1}^n [\mathbb{L}_i^-(h), \mathbb{L}_i^+(h)]$ , where  $\mathbb{L}_i^+(h)$  (resp.  $\mathbb{L}_i^-(h)$ ) is the maximum (resp. minimum) slope of the  
88 function  $h$  along the  $i$ -th dimension. We report the complexities using the standard big O notations  $\mathcal{O}$   
89 and  $\tilde{\mathcal{O}}$ , where the latter hides the logarithmic factors. In this study, we are mainly concerned with the  
90 dependence of the computational complexities on *the size of the finite sets* involved (discretization of  
91 the primal and dual domains). In particular, we ignore the possible dependence of the computational  
92 complexities on the dimension of the variables, unless they appear in the power of the size of those  
93 discrete sets.

94 We note that the extended version of this article, including the technical proofs, is available in the  
 95 supplementary material.

## 96 2 VI in primal domain

97 We are concerned with the infinite-horizon, discounted cost, optimal control problems of the form

$$V_*(x) = \min \mathbb{E}_{w_t} \left[ \sum_{t=0}^{\infty} \gamma^t C(x_t, u_t) \middle| x_0 = x \right]$$

s.t.  $x_{t+1} = g(x_t, u_t, w_t)$ ,  $x_t \in \mathbb{X}$ ,  $u_t \in \mathbb{U}$ ,  $w_t \sim \mathbb{P}(\mathbb{W})$ ,  $\forall t \in \{0, 1, \dots\}$ ,

98 where  $x_t \in \mathbb{R}^n$ ,  $u_t \in \mathbb{R}^m$ , and  $w_t \in \mathbb{R}^l$  are the state, input and disturbance variables at time  $t$ ,  
 99 respectively;  $\gamma \in (0, 1)$  is the discount factor;  $C : \mathbb{X} \times \mathbb{U} \rightarrow \mathbb{R}$  is the stage cost;  $g : \mathbb{R}^n \times \mathbb{R}^m \times$   
 100  $\mathbb{R}^l \rightarrow \mathbb{R}^n$  describes the dynamics;  $\mathbb{X} \subset \mathbb{R}^n$  and  $\mathbb{U} \subset \mathbb{R}^m$  describe the state and input constraints,  
 101 respectively; and,  $\mathbb{P}(\cdot)$  is the distribution of the disturbance over the support  $\mathbb{W} \subseteq \mathbb{R}^l$ . Assuming  
 102 the stage cost  $C$  is bounded, the optimal value function solves the Bellman equation  $V_* = \mathcal{T}V_*$ ,  
 103 where  $\mathcal{T}$  is the corresponding DP operator ( $C$  and  $V$  are extended to infinity outside their effective  
 104 domains) [6, Prop. 1.2.2]

$$\mathcal{T}V(x) := \min_u \{C(x, u) + \gamma \cdot \mathbb{E}_w V(g(x, u, w))\}, \quad \forall x \in \mathbb{X}. \quad (1)$$

105 Indeed, the DP operator is  $\gamma$ -contractive in the infinity-norm [6, Prop. 1.2.4]. This property then gives  
 106 rise to the VI algorithm  $V_{k+1} = \mathcal{T}V_k$ ,  $k = 0, 1, \dots$ , which converges to  $V_*$  as  $k \rightarrow \infty$ , for arbitrary  
 107 initialization  $V_0$ . Moreover, assuming that the composition  $V \circ g$  (for each  $w$ ) and the cost  $C$  are  
 108 jointly convex in state and input variables, the DP operator also preserves convexity [7, Prop. 3.3.1].

109 For numerical implementation of VI algorithm, we need to address three issues. First, we need to  
 110 compute the expectation of the value function in (1). In order to simplify the exposition and include  
 111 the computational cost of this operation explicitly, we focus on disturbances with finite support:

112 **Assumption 2.1** (Disturbance with finite support). *The disturbance  $w$  has a finite support  $\mathbb{W}^d \subset \mathbb{R}^l$*   
 113 *with a given probability mass function (p.m.f.)  $p : \mathbb{W}^d \rightarrow [0, 1]$ .*

114 Under the preceding assumption, we have  $\mathbb{E}_w V(g(x, u, w)) = \sum_{w \in \mathbb{W}^d} p(w) \cdot V(g(x, u, w))$ . The  
 115 second and more important issue is that the optimization problem (1) is infinite-dimensional for the  
 116 continuous state space  $\mathbb{X}$ . This renders the exact implementation of the VI algorithm impossible,  
 117 except for a few cases with available closed form solution. A common solution to this problem is to  
 118 deploy a sample-based approach, accompanied by a function approximation scheme. To be precise,  
 119 for a finite subset  $\mathbb{X}^d$  of  $\mathbb{X}$ , at each iteration  $k \geq 0$ , we take the discrete function  $V_k^d : \mathbb{X}^d \rightarrow \mathbb{R}$   
 120 as the input, and compute the discrete function  $V_{k+1}^d = [\mathcal{T}V_k^d]^d : \mathbb{X}^d \rightarrow \mathbb{R}$ , by employing an  
 121 extension  $\widetilde{V}_k^d : \mathbb{X} \rightarrow \mathbb{R}$  of  $V_k^d$ , as an approximation of  $V_k$  in the DP operation (1). Finally, we have  
 122 to actually solve the minimization problem over the control input in (1). Here, again, a common  
 123 approximation involves enumeration over a finite subset  $\mathbb{U}^d$  of the inputs space  $\mathbb{U}$ . Incorporating  
 124 these approximations, we end up with the *approximate* VI algorithm  $V_{k+1}^d = \mathcal{T}^d V_k^d$ , characterized  
 125 by the *discrete* DP (d-DP) operator

$$\mathcal{T}^d V^d(x) := \min_{u \in \mathbb{U}^d} \left\{ C(x, u) + \gamma \cdot \sum_{w \in \mathbb{W}^d} p(w) \cdot \widetilde{V}^d(g(x, u, w)) \right\}, \quad \forall x \in \mathbb{X}^d. \quad (2)$$

126 The convergence of the approximate VI algorithm described above depends on the properties of the  
 127 extension operation  $\widetilde{[\cdot]}$ . In particular, if  $\widetilde{[\cdot]}$  is non-expansive (in the infinity-norm), then  $\mathcal{T}^d$  is also  
 128  $\gamma$ -contractive. The error of the approximate VI algorithm also depends on the extension operation  $\widetilde{[\cdot]}$   
 129 and its representative power. We refer the interested reader to [6, 9, 23] for detailed discussions on  
 130 the convergence and error of approximate VI algorithm using different approximation schemes.

131 The d-DP operator and the corresponding approximate VI algorithm will be our benchmark for  
 132 evaluating the performance of the alternative algorithm developed in this study. So, we finish this  
 133 section with some remarks on the time complexity of the d-DP operation. Let the time complexity of  
 134 a single evaluation of the extension operation  $\widetilde{[\cdot]}$  be of  $\mathcal{O}(E)$ . Then, the time complexity of the d-DP  
 135 operation (2) is of  $\mathcal{O}(XUWE)$ . In this regard, note that the scheme described above essentially

136 involves approximating a continuous-state/action MDP with a finite-state/action MDP, and then  
 137 applying the VI algorithm. This, in turn, implies the lower bound  $\Omega(XU)$  for the time complexity  
 138 (corresponding to enumeration over  $u \in \mathbb{U}^d$  for each  $x \in \mathbb{X}^d$ ). This lower bound is also compatible  
 139 with the best existing time complexities in the literature for VI for finite MDPs; see, e.g., [3, 24].  
 140 However, as we will see shortly, for a particular class of problems, it is possible to exploit the structure  
 141 of the underlying continuous system in order to achieve a better time complexity.

### 142 3 Reducing complexity via conjugate duality

143 In this section, we introduce a class of problems that allow us to employ conjugate duality for the  
 144 VI problem and propose an alternative path for implementing the corresponding DP operator. We  
 145 also present the numerical scheme for implementing the proposed alternative path, and analyze its  
 146 convergence, complexity, and error. We note that the proposed algorithm and its analysis are based on  
 147 the d-CDP algorithm in [18, Sec. 5] for finite-horizon optimal control of deterministic systems. Here,  
 148 we extend those results for infinite-horizon, discounted cost, optimal control of stochastic systems.  
 149 Moreover, unlike [18], our analysis includes the case where the conjugate of input-dependent stage  
 150 cost is not analytically available and has to be computed numerically; see also [18, Assump. 5.1].

#### 151 3.1 Problem class

152 Throughout this section, we assume that the problem data satisfy the following conditions.

153 **Assumption 3.1** (Problem class). *The problem data has the following properties: (i) The disturbance*  
 154 *is additive; that is,  $g(x, u, w) = f(x, u) + w$ , where  $f : \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$  describes the deterministic*  
 155 *dynamics. (ii) The deterministic dynamics is of the form  $f(x, u) = f_s(x) + Bu$ , where  $f_s : \mathbb{R}^n \rightarrow \mathbb{R}^n$*   
 156 *is a Lipschitz continuous, possibly nonlinear map describing the “state” dynamics, and  $B \in \mathbb{R}^{n \times m}$ .*  
 157 *(iii) The constraint sets  $\mathbb{X} \subset \mathbb{R}^n$  and  $\mathbb{U} \subset \mathbb{R}^m$  are compact and convex. Moreover, for each  $x \in \mathbb{X}$ ,*  
 158 *the set of admissible inputs  $\mathbb{U}(x) := \{u \in \mathbb{U} : g(x, u, w) \in \mathbb{X}, \forall w \in \mathbb{W}^d\}$  is nonempty. (iv) The*  
 159 *stage cost  $C$  is separable in state and input, that is,  $C(x, u) = C_s(x) + C_i(u)$ , where  $C_s : \mathbb{X} \rightarrow \mathbb{R}$*   
 160 *and  $C_i : \mathbb{U} \rightarrow \mathbb{R}$  are Lipschitz continuous and convex.*

161 Some remarks are in order regarding the preceding assumptions. First, note that the setting of As-  
 162 sumption 3.1 goes beyond the classical LQR. In particular, it includes nonlinear state dynamics, state  
 163 and input constraints, and non-quadratic stage costs. Second, the properties laid out in Assumption 3.1  
 164 imply that the set of admissible inputs  $\mathbb{U}(x)$  is a compact set for each  $x \in \mathbb{X}$ . This, in turn, implies  
 165 that the optimal value in (1) is achieved if  $V : \mathbb{X} \rightarrow \mathbb{R}$  is also assumed to be lower semi-continuous.  
 166 Hereafter, we also assume that the joint discretization of the state-input space is “proper” in the sense  
 167 that the feasibility condition of Assumption 3.1-(iii) also holds for the discrete state-input space:

168 **Assumption 3.2** (Feasible discretization). *The discrete state space  $\mathbb{X}^d \subset \mathbb{X}$  and input space  $\mathbb{U}^d \subset \mathbb{U}$*   
 169 *are such that  $\mathbb{U}^d(x) := \mathbb{U}(x) \cap \mathbb{U}^d \neq \emptyset$  for all  $x \in \mathbb{X}^d$ .*

#### 170 3.2 VI in conjugate domain

171 In this subsection, we use duality theory to present an alternative path for computing the output of the  
 172 DP operator through the conjugate domain. This path forms the basis for the algorithm proposed in  
 173 this study. Fix  $x \in \mathbb{X}$  and consider the following reformulation of the optimization problem (1)

$$\mathcal{T}V(x) = C_s(x) + \min_{u,z} \{C_i(u) + \gamma \cdot \mathbb{E}_w V(z + w) : z = f(x, u)\},$$

174 where we used additivity of disturbance and separability of stage cost in Assumptions 3.1-(i,iv). The  
 175 corresponding dual problem then reads as

$$\widehat{\mathcal{T}}V(x) := C_s(x) + \max_y \min_{u,z} \{C_i(u) + \gamma \cdot \mathbb{E}_w V(z + w) + \langle y, f(x, u) - z \rangle\}, \quad (3)$$

176 where  $y \in \mathbb{R}^n$  is dual variable corresponding to the equality constraint. For the dynamics of  
 177 Assumption 3.1-(ii), we can obtain the following representation for the dual problem.

178 **Proposition 3.3** (CDP operator). *The dual problem (3) equivalently reads as*

$$\epsilon(x) := \gamma \cdot \mathbb{E}_w V(x + w), \quad x \in \mathbb{X}, \quad (4a)$$

$$\phi(y) := C_i^*(-B^\top y) + \epsilon^*(y), \quad y \in \mathbb{R}^n, \quad (4b)$$

$$\widehat{\mathcal{T}}V(x) = C_s(x) + \phi^*(f_s(x)), \quad x \in \mathbb{X}, \quad (4c)$$

179 where  $h^*(y) = \sup_x \{ \langle y, x \rangle - h(x) \}$  is the (convex) conjugate of  $h$ .

180 Following [18], we call the operator  $\widehat{\mathcal{T}}$  defined by (4) the *conjugate DP (CDP) operator*. We next  
181 provide an alternative representation of the CDP operator that captures the essence of this operation.

182 **Proposition 3.4** (CDP reformulation). *The CDP operator  $\widehat{\mathcal{T}}$  equivalently reads as*

$$\widehat{\mathcal{T}}V(x) = C_s(x) + \min_u \{ C_i(u) + \gamma \cdot [\mathbb{E}_w V(\cdot + w)]^{**}(f(x, u)) \}, \quad (5)$$

183 where  $[\cdot]^{**} = [[\cdot]^*]^*$  denotes the biconjugate operation.

184 The preceding result implies that the indirect path through the conjugate domain essentially involves  
185 substituting the (expectation of the) value function by its biconjugate. In particular, it points to a  
186 sufficient condition for zero duality gap.

187 **Corollary 3.5** (Equivalence of  $\mathcal{T}$  and  $\widehat{\mathcal{T}}$ ). *If  $V : \mathbb{X} \rightarrow \mathbb{R}$  is convex, then  $\widehat{\mathcal{T}}V = \mathcal{T}V$ .*

188 Hence,  $\widehat{\mathcal{T}}$  has the same properties as  $\mathcal{T}$  if  $V$  is convex. In particular,  $\widehat{\mathcal{T}}$  is  $\gamma$ -contractive, and hence,  
189 the *conjugate VI (CVI) algorithm*  $V_{k+1} = \widehat{\mathcal{T}}V_k$ ,  $k = 0, 1, \dots$ , is expected to converge to the optimal  
190 value function  $V_*$  with arbitrary convex initialization  $V_0$ , if  $\widehat{\mathcal{T}}$  preserves convexity (so that the output  
191 of each iteration of the CVI algorithm is convex). As before, this last condition holds if  $V \circ f$  is jointly  
192 convex in  $x$  and  $u$ . This is for example the case for linear state dynamics  $f_s(x) = Ax$ ,  $A \in \mathbb{R}^{n \times n}$ .  
193 We note that, if  $\widehat{\mathcal{T}}$  does not preserve convexity, then the alternative path suffers from duality gap.

### 194 3.3 CVI algorithm

195 The approximate CVI algorithm involves consecutive applications of an approximate implementation  
196 of the CDP operator (4) until some termination condition is satisfied. Algorithm 1 provides the  
197 pseudo-code of this procedure. In particular, we consider solving (4) for a finite set  $\mathbb{X}^d$  of points in  
198 the state space, and terminate the iterations when the difference between two consecutive discrete  
199 value functions (in the infinity-norm) is less than a given constant  $e_t > 0$ ; see Algorithm 1:4. In what  
200 follows, we describe the main steps within the initialization and the iterations of Algorithm 1. In  
201 particular, the three conjugate operations in (4) are handled numerically via the linear-time Legendre  
202 transform (LLT) algorithm [20]. LLT is an efficient algorithm for computing the *discrete* conjugate  
203 over a finite *grid-like* dual domain. Precisely, to compute the conjugate of the function  $h : \mathbb{X} \rightarrow \mathbb{R}$ ,  
204 LLT takes its discretization  $h^d : \mathbb{X}^d \rightarrow \mathbb{R}$  as an input, and outputs  $h^{d*d}(y) = [[h^d]^*]^d(y) =$   
205  $\max_{x \in \mathbb{X}^d} \{ \langle y, x \rangle - h^d(x) \}$  for dual variables  $y$  belonging to a grid-like dual domain  $\mathbb{Y}^g$ . We refer  
206 the reader to [20] for a detailed description of the LLT algorithm.

207 The main steps of the proposed CVI algorithm are as follows: **(i)** For the expectation operation in (4a),  
208 by Assumption 2.1, we again have  $\mathbb{E}_w V(\cdot + w) = \sum_{w \in \mathbb{W}^d} p(w) \cdot V(\cdot + w)$ . Notice, however, that  
209 at each iteration, we only have the discrete value function  $V^d : \mathbb{X}^d \rightarrow \mathbb{R}$  at our disposal. Hence, we  
210 also need to employ some form of extension  $\widetilde{V}^d : \mathbb{X} \rightarrow \mathbb{R}$  of  $V^d$ . Therefore, we first pass the value  
211 function  $V^d : \mathbb{X}^d \rightarrow \mathbb{R}$  through the “scaled expectation filter” to derive  $\varepsilon^d : \mathbb{X}^d \rightarrow \mathbb{R}$  in (6a), as an  
212 approximation of  $\varepsilon$ . **(ii)** In order to compute  $\phi$  in (4b), we need access to two conjugate functions.  
213 First, for  $\varepsilon^*$ , we use the approximation  $\varepsilon^{d*d} : \mathbb{Y}^g \rightarrow \mathbb{R}$  in (6b), by applying LLT to the data points  
214  $\varepsilon^d : \mathbb{X}^d \rightarrow \mathbb{R}$  for a properly chosen state dual grid  $\mathbb{Y}^g \subset \mathbb{R}^n$ . We also need the conjugate  $C_i^*$  of the  
215 input-dependent stage cost. If this function is not analytically available, we approximate it as follows:  
216 For a properly chosen input dual grid  $\mathbb{V}^g \subset \mathbb{R}^m$ , we employ LLT to compute  $C_i^{d*d} : \mathbb{V}^g \rightarrow \mathbb{R}$  using  
217 the data points  $C_i^d : \mathbb{U}^d \rightarrow \mathbb{R}$ ; see (6c). With these conjugate functions at hand, we can now compute  
218  $\varphi^d : \mathbb{Y}^g \rightarrow \mathbb{R}$  in (6d), as an approximation of  $\phi$ . In particular, notice that we use the LERP extension  
219  $\overline{C}_i^{d*d}$  of  $C_i^{d*d}$  to approximate  $C_i^{d*}$  at the required point  $-B^\top y$  for each  $y \in \mathbb{Y}^g$ . **(iii)** To be able to  
220 compute the output according to (4c), we need to perform another conjugate transform. In particular,  
221 we need the value of  $\phi^*$  at  $f_s(x)$  for  $x \in \mathbb{X}^d$ . Here, we use the approximation  $\varphi^{d*d} : \mathbb{Z}^g \rightarrow \mathbb{R}$  in  
222 (6e), by applying LLT to the data points  $\varphi^d : \mathbb{Y}^g \rightarrow \mathbb{R}$  for a properly chosen grid  $\mathbb{Z}^g \subset \mathbb{R}^n$ . Finally,  
223 we use the LERP extension  $\overline{\varphi}^{d*d}$  of  $\varphi^{d*d}$  to approximate  $\varphi^{d*}$  at the required point  $f_s(x)$  for each  
224  $x \in \mathbb{X}^d$ , and compute  $\widehat{\mathcal{T}}^d V^d$  in (6f) as an approximation of  $\widehat{\mathcal{T}}V$ . With these approximations, we can

---

**Algorithm 1** CVI: Approximate VI in conjugate domain
 

---

**Input:** dynamics  $f_s : \mathbb{R}^n \rightarrow \mathbb{R}^n$ ,  $B \in \mathbb{R}^{n \times m}$ ; finite state space  $\mathbb{X}^d \subset \mathbb{X}$ ; finite input space  $\mathbb{U}^d \subset \mathbb{U}$ ; state-dependent cost function  $C_s^d : \mathbb{X}^d \rightarrow \mathbb{R}$ ; input-dependent cost function  $C_i^d : \mathbb{U}^d \rightarrow \mathbb{R}$ ; finite disturbance space  $\mathbb{W}^d$  and its p.m.f.  $p : \mathbb{W}^d \rightarrow [0, 1]$ ; discount factor  $\gamma$ ; termination constant  $e_t$ .

**Output:** discrete value function  $\widehat{V}^d : \mathbb{X}^d \rightarrow \mathbb{R}$ .

*initialization:*

- 1: construct the grid  $\mathbb{V}^g$ ; use LLT to compute  $C_i^{d*} : \mathbb{V}^g \rightarrow \mathbb{R}$  from  $C_i^d : \mathbb{U}^d \rightarrow \mathbb{R}$ ;
  - 2: construct the grid  $\mathbb{Z}^g$ ;
  - 3:  $V^d(x) \leftarrow 0$  for  $x \in \mathbb{X}^d$ ;  $V_+^d(x) \leftarrow C_s^d(x) - \min C_i^d$  for  $x \in \mathbb{X}^d$ ;
  - iteration:*
  - 4: **while**  $\|V_+^d - V^d\|_\infty \geq e_t$  **do**
  - 5:    $V^d \leftarrow V_+^d$ ;
  - 6:   construct the grid  $\mathbb{Y}^g$ ;
  - d-CDP operation:*
  - 7:    $\varepsilon^d(x) \leftarrow \gamma \cdot \sum_{w \in \mathbb{W}^d} p(w) \widetilde{V}^d(x+w)$  for  $x \in \mathbb{X}^d$ ; use LLT to compute  $\varepsilon^{d*} : \mathbb{Y}^g \rightarrow \mathbb{R}$  from  $\varepsilon^d : \mathbb{X}^d \rightarrow \mathbb{R}$ ;
  - 8:   **for each**  $y \in \mathbb{Y}^g$  **do**
  - 9:     use LERP to compute  $\overline{C_i^{d*}}(-B^\top y)$  from  $C_i^{d*} : \mathbb{V}^g \rightarrow \mathbb{R}$ ;  $\varphi^d(y) \leftarrow \overline{C_i^{d*}}(-B^\top y) + \varepsilon^{d*}(y)$ ;
  - 10:   **end for**
  - 11:   use LLT to compute  $\varphi^{d*} : \mathbb{Z}^g \rightarrow \mathbb{R}$  from  $\varphi^d : \mathbb{Y}^g \rightarrow \mathbb{R}$ ;
  - 12:   **for each**  $x \in \mathbb{X}^d$  **do**
  - 13:     use LERP to compute  $\overline{\varphi^{d*}}(f_s(x))$  from  $\varphi^{d*} : \mathbb{Z}^g \rightarrow \mathbb{R}$ ;  $V_+^d(x) \leftarrow C_s(x) + \overline{\varphi^{d*}}(f_s(x))$ ;
  - 14:   **end for**
  - 15: **end while**
  - 16: output  $\widehat{V}^d \leftarrow V_+^d$ .
- 

225 introduce the *discrete* CDP (d-CDP) operator as follows

$$\varepsilon^d(x) := \gamma \cdot \sum_{w \in \mathbb{W}^d} p(w) \cdot \widetilde{V}^d(x+w), \quad x \in \mathbb{X}^d, \quad (6a)$$

$$\varepsilon^{d*}(y) = \max_{x \in \mathbb{X}^d} \{ \langle x, y \rangle - \varepsilon^d(x) \}, \quad y \in \mathbb{Y}^g, \quad (6b)$$

$$C_i^{d*}(v) = \max_{u \in \mathbb{U}^d} \{ \langle u, v \rangle - C^d(u) \}, \quad v \in \mathbb{V}^g, \quad (6c)$$

$$\varphi^d(y) := \overline{C_i^{d*}}(-B^\top y) + \varepsilon^{d*}(y), \quad y \in \mathbb{Y}^g, \quad (6d)$$

$$\varphi^{d*}(z) = \max_{y \in \mathbb{Y}^g} \{ \langle y, z \rangle - \varphi^d(y) \}, \quad z \in \mathbb{Z}^g, \quad (6e)$$

$$\widehat{\mathcal{T}}^d V^d(x) := C_s(x) + \overline{\varphi^{d*}}(f_s(x)), \quad x \in \mathbb{X}^d. \quad (6f)$$

226 The proper construction of the grids  $\mathbb{Y}^g$ ,  $\mathbb{V}^g$ , and  $\mathbb{Z}^g$  will be discussed in Section 3.5.

### 227 3.4 Analysis of CVI algorithm

228 We now provide our main theoretical results concerning the convergence, complexity, and error of  
 229 the proposed algorithm. Before that, we present the assumptions to be called in this subsections.

230 **Assumption 3.6** (Grids). *Consider the following properties for the grids in Algorithm 1: (i) The*  
 231 *grid  $\mathbb{V}^g$  is constructed such that  $\text{co}(\mathbb{V}_{\text{sub}}^g) \supseteq \mathbb{L}(C_i^d)$ ; (ii) The grid  $\mathbb{Z}^g$  is constructed such that*  
 232  *$\text{co}(\mathbb{Z}^g) \supseteq f_s(\mathbb{X}^d)$ ; (iii) The construction of  $\mathbb{Y}^g$ ,  $\mathbb{V}^g$ , and  $\mathbb{Z}^g$  requires at most  $\mathcal{O}(X + U)$  operations.*  
 233  *$\mathbb{Y}^g$  and  $\mathbb{Z}^g$  are of (approximately) the same size as  $\mathbb{X}^d$  in each dimension.  $\mathbb{V}^g$  is of (approximately)*  
 234 *the same size as  $\mathbb{U}^d$  in each dimension.*

235 **Assumption 3.7** (The extension operator). *Consider the following properties for the extension*  
 236 *operator  $\widetilde{[\cdot]}$  in (6a): (i)  $\widetilde{[\cdot]}$  is non-expansive w.r.t. infinity norm; (ii) Given a function  $V : \mathbb{X} \rightarrow \mathbb{R}$  and*  
 237 *its discretization  $V^d : \mathbb{X}^d \rightarrow \mathbb{R}$ , we have  $\|V - \widetilde{V}^d\|_\infty \leq e_e$  for some constant  $e_e \geq 0$ . (iii) Each*  
 238 *evaluation of  $\widetilde{[\cdot]}$  has a time complexity of  $\mathcal{O}(E)$ .*

239 **Assumption 3.8** (The dynamics). *Given a convex function  $V : \mathbb{X} \rightarrow \mathbb{R}$ , the composition  $V \circ f$  is*  
 240 *jointly convex in the state and input variables.*

241 We begin with the following result on the contractiveness of the d-CDP operator.

242 **Theorem 3.9** (Convergence). *Let Assumptions 3.6-(ii) and 3.7-(i) hold. Then, the d-CDP operator (6)*  
 243 *is  $\gamma$ -contractive in the infinity-norm.*

244 The preceding theorem implies that the approximate CVI Algorithm 1 is indeed convergent given  
 245 that the required conditions are satisfied. We next consider the time complexity of our algorithm.

246 **Theorem 3.10** (Complexity). *Let Assumptions 3.6-(iii) and 3.7-(iii) hold. The time complexity of*  
 247 *initialization and each iteration in Algorithm 1 are of  $\mathcal{O}(X + U)$  and  $\tilde{\mathcal{O}}(XWE + U)$ , respectively.*

248 As a result, each iteration of the CVI algorithm requires at most  $\tilde{\mathcal{O}}(XWE + U)$  operations, where  $E$   
 249 denotes the complexity of the extension operator  $[\cdot]$  in (6a). On the other hand, recall that the time  
 250 complexity of each iteration of the VI algorithm (in primal domain) is of  $\mathcal{O}(XUWE)$ , where  $E$   
 251 denotes the complexity of extension operation used in (2). In particular, for a deterministic dynamics  
 252 ( $W = 1$ ), if  $\mathbb{X}^d = \mathbb{X}^g$  is grid-like, and all the extension operations are handled using LERP (so  
 253 that  $E = \log X$ ), then the complexity of the CVI and VI algorithms are of  $\tilde{\mathcal{O}}(X + U)$  and  $\tilde{\mathcal{O}}(XU)$ ,  
 254 respectively. This is a reduction from quadratic complexity to linear complexity. However, the CVI  
 255 algorithm, like VI and other approximation schemes that utilize discretization/abstraction of the  
 256 continuous state space, still suffers from the so-called ‘‘curse of dimensionality.’’ This is because the  
 257 size  $X$  of the discrete state space increases exponentially with the dimension of the state space. We  
 258 finish with the following result on the error of the proposed CVI algorithm.

259 **Theorem 3.11** (Error). *Let Assumptions 3.6-(i,ii), 3.7-(i), and 3.8 hold. Consider the true optimal*  
 260 *value function  $V_\star = TV_\star : \mathbb{X} \rightarrow \mathbb{R}$  and its discretization  $V_\star^d : \mathbb{X}^d \rightarrow \mathbb{R}$ . Let Assumption 3.7-(ii) hold*  
 261 *for  $V_\star$ . Also, let  $\hat{V}^d : \mathbb{X}^d \rightarrow \mathbb{R}$  be the output of Algorithm 1 with fixed grids  $\mathbb{Y}^g, \mathbb{V}^g$ , and  $\mathbb{Z}^g$ . Then,*

$$\|\hat{V}^d - V_\star^d\|_\infty \leq \frac{\gamma(e_e + e_t) + e_d}{1 - \gamma}, \quad (7)$$

262 where  $e_d = e_u + e_v + e_x + e_y + e_z$ , and

$$\begin{aligned} e_u &= c_u \cdot d_H(\mathbb{U}, \mathbb{U}^d), \quad e_v = c_v \cdot d_H(\text{co}(\mathbb{V}^g), \mathbb{V}^g), \quad e_x = c_x \cdot d_H(\mathbb{X}, \mathbb{X}^d), \\ e_y &= c_y \cdot \max_{x \in \mathbb{X}^d} d(\partial(V_\star - C_s)(x), \mathbb{Y}^g), \quad e_z = c_z \cdot d_H(f_s(\mathbb{X}^d), \mathbb{Z}^g), \end{aligned} \quad (8)$$

263 with constants  $c_u, c_v, c_x, c_y, c_z > 0$  depending on the problem data.

264 The error terms in (7) relate to the main sources of error in the proposed approximate implementation  
 265 of the CVI algorithm: (i)  $e_e$  is due to the approximation of the value function using the extension  
 266 operator  $[\cdot]$ ; (ii)  $e_t$  corresponds to the termination of the algorithm after a finite number of iterations;  
 267 (iii)  $e_d$  captures the error due to the discretization of the primal and dual state and input domains. We  
 268 again note that Assumption 3.8 implies that the CDP operator preserves convexity and hence the  
 269 duality gap is zero. Otherwise, the proposed scheme can suffer from large errors due to dualization.

### 270 3.5 Construction of the grids

271 In this final subsection, we provide specific guidelines for the construction of the grids  $\mathbb{Y}^g, \mathbb{V}^g$   
 272 and  $\mathbb{Z}^g$ . The presented guidelines are based on our error analysis (Theorem 3.11) and correspond  
 273 to the properties laid out in Assumption 3.6. We refer the reader to [18] for further details about  
 274 Assumption 3.6-(i,ii) on the grids  $\mathbb{V}^g$  and  $\mathbb{Z}^g$ .

275 **Static construction of  $\mathbb{V}^g$ .** By Assumption 3.6-(i), we need to find the smallest grid  $\mathbb{V}^g$  for  
 276 which  $\text{co}(\mathbb{V}_{\text{sub}}^g) \supseteq \mathbb{L}(C_i^d)$ , so that  $e_v$  in (8) is controlled by the resolution of  $\mathbb{V}^g$ . This essen-  
 277 tially means that  $\mathbb{V}^g$  must ‘‘more than cover the range of slopes’’ of the function  $C_i^d$ . Since  $C_i$   
 278 is convex (Assumption 3.1-(iv)), and hence  $C_i^d$  is convex-extensible, we can compute the re-  
 279 quired range of slopes with an acceptable computational cost. In particular, assuming the do-  
 280 main  $\mathbb{U}^d = \mathbb{U}^g = \prod_{i=1}^m \mathbb{U}_i^g$  of  $C_i^d$  is grid-like, we can take  $L_i^-(C_i^d)$  (resp.  $L_i^+(C_i^d)$ ) to be the  
 281 minimum first forward (resp. maximum last backward) difference of  $C_i^d$  along the  $i$ -th dimension.  
 282 We then construct  $\mathbb{V}_{\text{sub}}^g = \prod_{i=1}^m \mathbb{V}_{\text{sub}i}^g$  such that, in each dimension  $i$ ,  $\mathbb{V}_{\text{sub}i}^g$  is uniform with the same  
 283 cardinality as  $\mathbb{U}_i^g$ , and  $\text{co}(\mathbb{V}_{\text{sub}i}^g) = [L_i^-(C_i^d), L_i^+(C_i^d)]$ . Finally, we construct  $\mathbb{V}^g$  by adding one  
 284 smaller element and one greater element to  $\mathbb{V}_{\text{sub}}^g$ , while preserving the resolution, in each dimension.  
 285 This construction satisfies the condition of Assumption 3.6-(iii).

286 **Static construction of  $\mathbb{Z}^g$ .** According to Assumption 3.6-(ii), the grid  $\mathbb{Z}^g$  must be constructed such  
 287 that  $\text{co}(\mathbb{Z}^g) \supseteq f_s(\mathbb{X}^d)$ . This can be simply done by finding the vertices of the smallest box that  
 288 contains the set  $f_s(\mathbb{X}^d)$ . Those vertices give the range of  $\mathbb{Z}^g$  in each dimension. We can then, for  
 289 example, take  $\mathbb{Z}^g$  to be the uniform grid with the same cardinality as  $\mathbb{X}^d$  in each dimension. This  
 290 way,  $d_H(f_s(\mathbb{X}^d), \mathbb{Z}^g) \leq d_H(\text{co}(\mathbb{Z}^g), \mathbb{Z}^g)$ , and hence  $e_z$  in (8) is controlled by the resolution of  $\mathbb{Z}^g$ .  
 291 One again, this construction satisfies the complexity requirement of Assumption 3.6-(iii).

292 **Static construction of  $\mathbb{Y}^g$ .** Construction of the dual state grid  $\mathbb{Y}^g$  is more involved. In particular,  
 293 according to Theorem 3.11, we need to choose a *fixed* grid that minimizes  $e_y$  in (8). This can be  
 294 done by choosing  $\mathbb{Y}^g$  such that  $\mathbb{Y}^g \cap \partial(V_* - C_s) \neq \emptyset$  for all  $x \in \mathbb{X}^g$  so that  $e_y = 0$ . Even if we had  
 295 access to the optimal value function  $V_*$ , satisfying such a condition can lead to dual grids  $\mathbb{Y}^g \subset \mathbb{R}^n$   
 296 of size  $\mathcal{O}(X^n)$ . Such a large size violates Assumption 3.6-(iii) on the size of  $\mathbb{Y}^g$ , and essentially  
 297 renders the proposed algorithm impractical for dimensions  $n \geq 2$ . A more practical condition is  
 298  $\text{co}(\mathbb{Y}^g) \cap \partial(V_* - C_s) \neq \emptyset$  for all  $x \in \mathbb{X}^d$  so that  $\max_{x \in \mathbb{X}^d} d(\partial(V_* - C_s)(x), \mathbb{Y}^g) \leq d_H(\text{co}(\mathbb{Y}^g), \mathbb{Y}^g)$ ,  
 299 and hence  $e_y$  is controlled by the resolution of the grid  $\mathbb{Y}^g$ . In order to satisfy the latter condition, we  
 300 need to approximate the range of slopes of  $V_* - C_s$  (over  $\mathbb{X}^d$ ). To this end, we can use the fact that  
 301  $V_*$  is the fixed point of DP operator (1) to approximate  $\text{Rng}(V_* - C_s)$  by  $R = \frac{\text{Rng}(C_1^d) + \gamma \cdot \text{Rng}(C_s^d)}{1 - \gamma}$ .  
 302 We then construct the grid  $\mathbb{Y}^g = \prod_{i=1}^n \mathbb{Y}_i^g$  such that for each dimension  $i = 1, \dots, n$ , we have  
 303  $\pm \alpha R / \Delta_{x_i^d} \in \text{co}(\mathbb{Y}_i^g)$ , where  $\Delta_{x_i^d}$  (with some abuse of notation) denotes the diameter of  $\mathbb{X}^d$  along  $i$ -th  
 304 dimension, and  $\alpha > 0$  is a scaling factor mainly depending on the dimension of the state space. This  
 305 construction requires  $\mathcal{O}(X + U)$  operations which is in consistence with Assumption 3.6-(iii).

306 **Dynamic construction of  $\mathbb{Y}^g$ .** Alternatively, we can construct  $\mathbb{Y}^g$  *dynamically* at each iteration  
 307 in order to minimize the corresponding error in each application of the d-CDP operator given by  
 308  $e_y = c_y \cdot \max_{x \in \mathbb{X}^d} d(\partial(\mathcal{T}V - C_s)(x), \mathbb{Y}^g)$ ; see Lemma 5.6 and Proposition 5.8 in the supplementary  
 309 material. Once again, the idea is to make sure  $\text{co}(\mathbb{Y}^g) \cap \partial(\mathcal{T}V - C_s) \neq \emptyset$  for all  $x \in \mathbb{X}^g$ . Since  
 310 we do not have  $[\mathcal{T}V]^d$  at our disposal (it is the output of the current iteration), we can again use the  
 311 definition of the DP operator (1) to approximate  $\text{Rng}(\mathcal{T}V - C_s)$  by  $R = \text{Rng}(C_1^d) + \gamma \cdot \text{Rng}(V^d)$ ,  
 312 where  $V^d$  is the output of the previous iteration. We can then construct  $\mathbb{Y}^g$  such that again for  
 313 each dimension  $i$ , we have  $\pm \alpha R / \Delta_{x_i^d} \in \text{co}(\mathbb{Y}_i^g)$ . This construction has a one-time computational  
 314 cost of  $\mathcal{O}(U)$  for computing  $\text{Rng}(C_1^d)$ , and per iteration computational cost of  $\mathcal{O}(X)$  for computing  
 315  $\text{Rng}(V^d)$ , which is again in consistence with Assumption 3.6-(iii).

## 316 4 Numerical simulations

317 We now showcase the application of the CVI algorithm in solving the optimal control problem of  
 318 a noisy inverted pendulum, and compare its performance with the VI algorithm. The details of  
 319 these simulations (and more numerical examples) are provided in Section 4 of the supplementary  
 320 material. For these simulations, we use modified versions of multiple routines in the d-CDP MATLAB  
 321 package [19]; see the supplementary material for the codes.

322 We use the setup (model and stage cost) of [18, App. C.2.2], with discount factor  $\gamma = 0.95$  and  
 323 termination constant  $e_t = 0.001$ . In particular, the stage cost is quadratic (in both state and input),  
 324 and the discrete-time dynamics is of the form  $x^+ = f_s(x) + Bu + w$ , where  $f_s(x_1, x_2) = [x_1 +$   
 325  $\alpha_{12}x_2, \alpha_{21} \sin x_1 + \alpha_{22}x_2]^T$  is the *nonlinear* state dynamics,  $B = [0, \beta]^T$ , and  $w \in \mathbb{R}^2$  is the  
 326 additive noise ( $\alpha_{12}, \alpha_{21}, \alpha_{22}, \beta \in \mathbb{R}$ ). We use uniform, grid-like discretizations of the state and input  
 327 spaces. This choice allows us to use LERP for the extension operators in (2) for VI and in (6a) for  
 328 CVI. The grids  $\mathbb{V}^g, \mathbb{Z}^g$ , and  $\mathbb{Y}^g$  are also constructed uniformly, following the guidelines of Section 3.5  
 329 (with  $\alpha = 1$ ). In particular, we consider both *static* and *dynamic* construction of the dual grid  $\mathbb{Y}^g$   
 330 in the CVI. Moreover, in each implementation of VI and CVI, all of the involved grids are chosen  
 331 to be of the same size  $N_i$  in each dimension, i.e.,  $X = Y = Z = N_i^2$  and  $U = V = N_i$ . We  
 332 are particularly interested in the performance of these algorithms, as the size of the discretizations  
 333 increases. We note that with these grid sizes, the time complexity of each iteration of VI and CVI is  
 334 of  $\mathcal{O}(XUW)$  and  $\mathcal{O}(XW + U)$ , respectively, where  $W = 5^2$  is the size of discrete disturbance set.

335 Let us also note that the optimal control problem at hand does not satisfy the feasibility condition  
 336 of Assumption 3.2. Indeed, there exist  $x \in \mathbb{X}^g$  for which there is no  $u \in \mathbb{U}^g$  such that  $x^+ =$

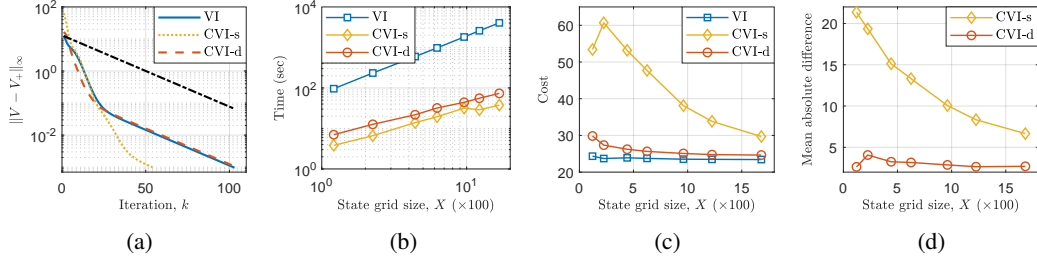


Figure 1: VI vs. CVI: (a) Convergence rate for grid size  $X = 41^2$ ; (b) Running time; (c) Average cost of fifty instances of the control problem; (d) Difference between outputs of VI and CVI algorithms. The black dashed-dotted line in (a) corresponds to exponential convergence with rate  $\gamma = 0.95$ . CVI-s and CVI-d correspond to *static* and *dynamic* construction of the dual grid  $\mathbb{Y}^g$ , respectively, in the CVI algorithm.

337  $f(x) + Bu \in \mathbb{X}$ . Nevertheless, both algorithms can be used for computing the value function.  
 338 However, the CVI algorithm, unlike the VI algorithm, assigns a finite value to the value function for  
 339 these problematic states. This does not contradict our theoretical results on the error of the proposed  
 340 algorithm, as the assumption on the optimal control problem to be feasible for all  $x \in \mathbb{X}^g$  is violated.  
 341 Indeed, for the feasible initial states, our theoretical results still hold true. And, practically, we  
 342 can simply impose the required state constraint while computing the control action in the forward  
 343 iterations using the computed value function.

344 The results of our numerical simulations are shown in Figure 1. As shown in Figure 1a, both  
 345 algorithms are indeed convergent with a rate greater than or equal to  $\gamma = 0.95$ . In particular, we see  
 346 that CVI-s converges in  $k_t = 54$  iterations, compared to  $k_t = 105$  iterations required for CVI-d (with  
 347 dynamic dual grid  $\mathbb{Y}^g$ ) and VI to reach the termination bound  $e_t = 0.001$ . Our simulations show that  
 348 the same increase in the convergence rate is also seen in CVI-d for *deterministic* systems; see Figure 2.  
 349 (These effects were also seen for other grid sizes). On the other hand, the lower time complexity per  
 350 iteration for CVI, leads to a significant reduction in the running time of this algorithm compared to  
 351 VI. This effect can be clearly seen in Figure 1b, where the run-time of CVI for  $X = 41^2$  is still less  
 352 than the that of VI for  $X = 11^2$ . Moreover, as expected, using a fixed grid  $\mathbb{Y}^g$  in CVI-s, as opposed  
 353 to dynamically constructing it at each iteration in CVI-d, leads to a reduced time requirement.

354 This reduction in the running time comes however with an  
 355 increase in the cost of underlying optimal control problem.  
 356 This effect is shown in Figure 1c, where we report the average  
 357 cost of fifty instances of the optimal control problem with  
 358 randomly chosen initial states, over  $T = 400$  time steps. For  
 359 these instances, at each time step, the control action is generated  
 360 using the greedy policy  $\mu(x) \in \arg\min_{u \in \mathbb{U}^g} C(x, u) + \gamma \cdot$   
 361  $\overline{V}^d(f(x, u))$ , w.r.t. the discrete value function  $V^d$  computed  
 362 using VI and CVI. Notice that, compared to VI, we see a small  
 363 increase in the cost of controlled trajectories in CVI-d, while  
 364 using a static grid  $\mathbb{Y}^g$  in CVI-s causes a huge increase in the  
 365 cost. This is because the dynamic construction of  $\mathbb{Y}^g$  in CVI  
 366 uses the available computational power (related to size of the  
 367 discretization) smartly by finding the smallest grid  $\mathbb{Y}^g$  in each  
 368 iteration, in order to minimize the error of that same iteration. This effect is also seen in Figure 1d,  
 369 where we report the mean absolute difference (over the discrete state space  $\mathbb{X}^g$ ) between the output  
 370 of the CVI and VI algorithms. Once again, using a dynamic dual grid  $\mathbb{Y}^g$  leads to less difference  
 371 between the two algorithm, which in turn implies smaller errors.

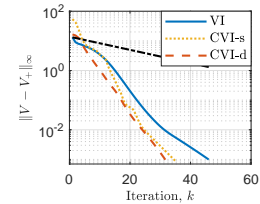


Figure 2: Convergence rate for grid size  $X = 41^2$  and *deterministic* dynamics  $x^+ = f_s(x) + Bu$ ; see also Figure 1.

## 372 References

- 373 [1] Y. Achdou, F. Camilli, and L. Corrias. On numerical approximation of the Hamilton-Jacobi-transport  
 374 system arising in high frequency approximations. *Discrete & Continuous Dynamical Systems-Series B*,  
 375 19(3), 2014.
- 376 [2] M. Akian, S. Gaubert, and A. Lakhoua. The max-plus finite element method for solving deterministic

- 377 optimal control problems: Basic properties and convergence analysis. *SIAM Journal on Control and*  
378 *Optimization*, 47(2):817–848, 2008.
- 379 [3] F. Bach. Max-plus matching pursuit for deterministic Markov decision processes. *arXiv preprint*  
380 *arXiv:1906.08524*, 2019.
- 381 [4] R. Bellman and W. Karush. Mathematical programming and the maximum transform. *Journal of the*  
382 *Society for Industrial and Applied Mathematics*, 10(3):550–567, 1962.
- 383 [5] E. Berthier and F. Bach. Max-plus linear approximations for deterministic continuous-state markov  
384 decision processes. *IEEE Control Systems Letters*, pages 1–1, 2020.
- 385 [6] D. P. Bertsekas. *Dynamic Programming and Optimal Control, Vol. II*. Athena Scientific, Belmont, MA,  
386 3rd edition, 2007.
- 387 [7] D. P. Bertsekas. *Convex Optimization Theory*. Athena Scientific, Belmont, MA, 2009.
- 388 [8] D. P. Bertsekas. *Reinforcement Learning and Optimal Control*. Athena Scientific, Belmont, MA, 2019.
- 389 [9] L. Busoniu, R. Babuska, B. De Schutter, and D. Ernst. *Reinforcement learning and dynamic programming*  
390 *using function approximators*. CRC press, 2017.
- 391 [10] R. Carpio and T. Kamihigashi. Fast value iteration: an application of Legendre-Fenchel duality to a class  
392 of deterministic dynamic programming problems in discrete time. *Journal of Difference Equations and*  
393 *Applications*, 26(2):209–222, 2020.
- 394 [11] L. Contento, A. Ern, and R. Vermiglio. A linear-time approximate convex envelope algorithm using the  
395 double Legendre–Fenchel transform with application to phase separation. *Computational Optimization*  
396 *and Applications*, 60(1):231–261, 2015.
- 397 [12] L. Corrias. Fast Legendre-Fenchel transform and applications to Hamilton-Jacobi equations and conserva-  
398 tion laws. *SIAM Journal on Numerical Analysis*, 33(4):1534–1558, 1996.
- 399 [13] G. Costeseque and J.-P. Lebacque. A variational formulation for higher order macroscopic traffic flow  
400 models: Numerical investigation. *Transportation Research Part B: Methodological*, 70:112 – 133, 2014.
- 401 [14] A. O. Esogbue and C. W. Ahn. Computational experiments with a class of dynamic programming  
402 algorithms of higher dimensions. *Computers & Mathematics with Applications*, 19(11):3 – 23, 1990.
- 403 [15] P. F. Felzenszwalb and D. P. Huttenlocher. Distance transforms of sampled functions. *Theory of computing*,  
404 8(1):415–428, 2012.
- 405 [16] M. Jacobs and F. Léger. A fast approach to optimal transport: the back-and-forth method. *arXiv preprint*  
406 *arXiv:1905.12154*, 2019.
- 407 [17] C. M. Klein and T. L. Morin. Conjugate duality and the curse of dimensionality. *European Journal of*  
408 *Operational Research*, 50(2):220 – 228, 1991.
- 409 [18] M. A. S. Kolarijani and P. Mohajerin Esfahani. Fast approximate dynamic programming for input-affine  
410 dynamics. *preprint arXiv:2008.10362*, 2020.
- 411 [19] M. A. S. Kolarijani and P. Mohajerin Esfahani. Discrete conjugate dynamic programming (d-CDP)  
412 MATLAB package. Available online at <https://github.com/AminKolarijani/d-CDP>, 2020.
- 413 [20] Y. Lucet. Faster than the fast Legendre transform, the linear-time Legendre transform. *Numerical*  
414 *Algorithms*, 16(2):171–185, 1997.
- 415 [21] Y. Lucet. New sequential exact Euclidean distance transform algorithms based on convex analysis. *Image*  
416 *and Vision Computing*, 27(1):37 – 44, 2009.
- 417 [22] W. M. McEneaney. Max-plus eigenvector representations for solution of nonlinear  $H_\infty$  problems: basic  
418 concepts. *IEEE Transactions on Automatic Control*, 48(7):1150–1163, 2003.
- 419 [23] W. B. Powell. *Approximate Dynamic Programming: Solving the Curses of Dimensionality*. John Wiley &  
420 Sons, Hoboken, NJ, 2nd edition, 2011.
- 421 [24] A. Sidford, M. Wang, X. Wu, and Y. Ye. Variance reduced value iteration and faster algorithms for  
422 solving Markov decision processes. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium*  
423 *on Discrete Algorithms*, pages 770–787. SIAM, 2018.
- 424 [25] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. MIT Press, 2018.

## 425 Checklist

- 426 1. For all authors...
- 427 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s  
428 contributions and scope? [Yes] See Theorems 3.9, 3.10 and 3.11.

- 429 (b) Did you describe the limitations of your work? [Yes] See the discussions following  
430 Theorems 3.10 and 3.11.
- 431 (c) Did you discuss any potential negative societal impacts of your work? [N/A]
- 432 (d) Have you read the ethics review guidelines and ensured that your paper conforms to  
433 them? [Yes]
- 434 2. If you are including theoretical results...
- 435 (a) Did you state the full set of assumptions of all theoretical results? [Yes]
- 436 (b) Did you include complete proofs of all theoretical results? [Yes] See Section 5 of the  
437 supplementary material.
- 438 3. If you ran experiments...
- 439 (a) Did you include the code, data, and instructions needed to reproduce the main experi-  
440 mental results (either in the supplemental material or as a URL)? [Yes]
- 441 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they  
442 were chosen)? [N/A]
- 443 (c) Did you report error bars (e.g., with respect to the random seed after running experi-  
444 ments multiple times)? [N/A]
- 445 (d) Did you include the total amount of compute and the type of resources used (e.g., type  
446 of GPUs, internal cluster, or cloud provider)? [Yes] See Section 4 of the supplementary  
447 material.
- 448 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 449 (a) If your work uses existing assets, did you cite the creators? [Yes] See [19].
- 450 (b) Did you mention the license of the assets? [N/A]
- 451 (c) Did you include any new assets either in the supplemental material or as a URL? [No]
- 452 (d) Did you discuss whether and how consent was obtained from people whose data you're  
453 using/curating? [N/A]
- 454 (e) Did you discuss whether the data you are using/curating contains personally identifiable  
455 information or offensive content? [N/A]
- 456 5. If you used crowdsourcing or conducted research with human subjects...
- 457 (a) Did you include the full text of instructions given to participants and screenshots, if  
458 applicable? [N/A]
- 459 (b) Did you describe any potential participant risks, with links to Institutional Review  
460 Board (IRB) approvals, if applicable? [N/A]
- 461 (c) Did you include the estimated hourly wage paid to participants and the total amount  
462 spent on participant compensation? [N/A]