

Set-based value operators for non-stationary and uncertain Markov decision processes

Journal Article**Author(s):**

Li, Sarah H.Q.; Adjé, Assalé; Garoche, Pierre-Loïc; Açıkmeşe, Behçet

Publication date:

2025-01

Permanent link:

<https://doi.org/10.3929/ethz-b-000703814>

Rights / license:

[Creative Commons Attribution 4.0 International](#)

Originally published in:

Automatica 171, <https://doi.org/10.1016/j.automatica.2024.111970>



Set-based value operators for non-stationary and uncertain Markov decision processes[☆]



Sarah H.Q. Li^{a,*}, Assalé Adjé^c, Pierre-Loïc Garoche^d, Behçet Açıkmeşe^b

^a Automatic Control Laboratory, ETH Zürich, Zürich, Switzerland

^b Department of Aeronautics and Astronautics, University of Washington, Seattle, USA

^c LAMPS, Université de Perpignan Via Domitia, Perpignan, France

^d École Nationale de l'Aviation Civile, Université de Toulouse, Toulouse, France

ARTICLE INFO

Article history:

Received 27 July 2022

Received in revised form 5 September 2023

Accepted 6 September 2024

Available online xxxx

Keywords:

Markov decision process

Contraction operator

Stochastic control

Decision making

Autonomy

ABSTRACT

This paper analyzes finite-state Markov Decision Processes (MDPs) with nonstationary and uncertain parameters via set-based fixed point theory. Given compact parameter ambiguity sets, we demonstrate that a family of contraction operators, including the Bellman operator and the policy evaluation operator, can be extended to set-based contraction operators with a unique fixed point—a compact value function set. For non-stationary MDPs, we show that while the value function trajectory diverges, its Hausdorff distance from this fixed point converges to zero. In parameter uncertain MDPs, the fixed point's extremum value functions are equivalent to the min-max value function in robust dynamic programming under the rectangularity condition. Furthermore, we show that the rectangularity condition is a sufficient condition for the fixed point to contain its own extremum value functions. Finally, we derive novel guarantees for probabilistic path planning in capricious wind fields and stratospheric station-keeping.

© 2024 The Authors. Published by Elsevier Ltd. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

1. Introduction

Markov decision process (MDP) is a versatile model for decision making in stochastic environments and is widely used in trajectory planning (Al-Sabban, Gonzalez, & Smith, 2013), robotics (Van Hoof, Hermans, Neumann, & Peters, 2015), and operations research (Doshi, Goodwin, Akkiraju, & Verma, 2005). Given state-action costs and transition probabilities, finding an optimal policy of the MDP is equivalent to solving for the fixed point value function of the corresponding Bellman operator. However, finding the fixed point value function may not be straight forward when the MDP model does not accurately describe its operation environment. In MDP applications such as traffic light control and dexterous manipulation, the MDP model is prone to two sources of inaccuracy: *environmental non-stationarity* and *parameter uncertainty*.

[☆] This research is partly funded by National Science Foundation grant CMMI-210563 and the University of Washington Aero&Astro Condit fellowship. The material in this paper was not presented at any conference. This paper was recommended for publication in revised form by Associate Editor Subhrakanti Dey under the direction of Editor Florian Dorfler.

* Corresponding author.

E-mail addresses: sarahli@control.ee.ethz.ch (S.H.Q. Li), assale.adje@univ-perp.fr (A. Adjé), Pierre-Loic.Garoche@enac.fr (P.-L. Garoche), behcet@uw.edu (B. Açıkmeşe).

Different from an MDP's internal stochasticity, environmental non-stationarity refers to the time-varying changes in the MDP parameters, in particular as induced by external factors or the presence of interfering decision makers in an unpredictable way. Standard dynamic programming under environmental non-stationarity has no convergence guarantees and can be shown to diverge.

Example 1 (Navigating in Non-Stationary Wind). An autonomous aircraft navigates a non-stationary wind field to reach a non-stationary target. The non-stationary wind field alternates between N known wind patterns over time and is predictable for the next time step but unpredictable in the long term. This creates environmental non-stationarity in both the aircraft's transition dynamics as well as the target position. Should its policy optimization use the average wind pattern, the worst-case wind pattern, or the short-term wind forecast?

On the other hand, parameter uncertainty refers to the discrepancy between the modeling parameters used in computation vs. the parameters that accurately model a time-invariant system. Minimizing the risk or worst-case failure for a parameter-uncertain MDP can be tackled from the min-max approach via robust MDPs (Iyengar, 2005; Mannor, Mebel, & Xu, 2016) and distributionally robust MDPs (Yang, 2017) under rectangularity assumptions on the parameter uncertainty's state-action space structure (Iyengar, 2005; Mannor et al., 2016).

In this paper, we develop a set-theoretic extension to the Bellman operator that is applicable to MDPs experiencing both environmental non-stationarity and parameter uncertainty. We use this approach to derive convergence guarantees for environmentally non-stationary MDPs under the Hausdorff distance and to extend the rectangularity assumptions for parameter-uncertain MDPs.

Contributions. We propose the set-extensions of *value operators*: a general class of contraction operators that includes the Bellman operator and the policy evaluation operator. For these value operators, we extend their input–output domain to sets of value functions, show the existence of a compact *fixed point* and prove that value iteration over sets of value functions converges. Under non-stationarity assumptions, we show that the Hausdorff distance between the value iteration trajectory and the fixed point of the value operator always converges to zero. Under MDP parameter uncertainty, we show that the rectangularity assumption in the min–max MDP model is sufficient for the fixed point of the value operator to contain its own extremal value functions. Given a value operator and a compact parameter uncertainty set, we present an algorithm that computes the bounds of the corresponding fixed point set. Finally, we apply our results to the wind-assisted navigation of high altitude platform systems relevant to space exploration (Wolf et al., 2010) and show that our algorithms can be used to derive policies with better guarantees.

Related research. MDP with parameter uncertainty is well-studied in robust control and reinforcement learning. In control theory, the worst-case cost-to-go with respect to state-decoupled parameter uncertainties is derived via a minmax variation of the Bellman operator in Givan, Leach, and Dean (2000), Iyengar (2005), Nilim and El Ghaoui (2005) and Wiesemann, Kuhn, and Rustem (2013). The cost-to-go under parameter uncertainty with coupling between states and time steps is similarly bounded in Goyal and Grand-Clement (2022) and Mannor et al. (2016). The effect of statistical uncertainty on the optimal cost-to-go is studied in Mannor et al. (2016), Nilim and El Ghaoui (2005), Wiesemann et al. (2013) and Yang (2017). Recently, parameter-uncertain MDPs gained traction in the reinforcement learning community due to the presence of uncertainty in real world problems such as traffic signal control and multi-agent coordination (Kumar, Zhou, Tucker, & Levine, 2020; Lecarpentier & Rachelson, 2019; Padakandla, KJ, & Bhatnagar, 2020). Most RL research extends the worst-case analysis to methods such as Q-learning and SARSA. Recently, methods for value-based RL using non-contracting operators have been investigated in Bellemare, Ostrovski, Guez, Thomas, and Munos (2016).

As an alternative to optimizing for the worst-case, distributionally robust MDP optimizes a risk metric over the value functions of an MDP whose parameters are described by a known set of probability distributions. Distributionally robust MDPs boil down to solving a min–max MDP formulation (Mannor et al., 2016; Xu & Mannor, 2010; Yang, 2017; Yu & Xu, 2015) over an ambiguity set, and require a rectangularity condition to be solvable.

We do not optimize the worst-case cost-to-go by assuming adversarial MDP parameter selection. Instead, we derive a set of cost-to-go values that is invariant with respect to a compact parameter uncertainty set. We continue from our previous work (Li, Adjé, Garoche, & Açıkmeşe, 2021), in which we analyzed the set-based Bellman operator for cost uncertainty only.

Notation: A set of N elements is given by $[N] = \{0, \dots, N - 1\}$. We denote the set of matrices of i rows and j columns with real (non-negative) entries as $\mathbb{R}^{i \times j}$ ($\mathbb{R}_+^{i \times j}$), respectively. Matrices and some integers are denoted by capital letters, X , while sets are denoted by cursive typeset \mathcal{X} . The set of all compact subsets of

\mathbb{R}^d is denoted by $\mathcal{K}(\mathbb{R}^d)$. The column vector of ones of size $N \in \mathbb{N}$ is denoted by $\mathbf{1}_N = [1, \dots, 1]^T \in \mathbb{R}^{N \times 1}$. The identity matrix of size S is denoted by I_S . The simplex of dimension S is denoted by

$$\Delta_S = \{p \in \mathbb{R}^S \mid \mathbf{1}_S^\top p = 1, p \geq 0\}. \quad (1)$$

A vector $h \in \mathbb{R}^S$ has equivalent notation (h_1, \dots, h_s) , where h_s is the value of h in the s th coordinate, $s \in [S]$.

2. Discounted infinite-horizon MDP

A *discounted infinite-horizon finite state MDP* is given by $([S], [A], P, C, \gamma)$, where $\gamma \in (0, 1)$ is a discount factor, $[S] = \{1, \dots, S\}$ is a finite set of states and $[A] = \{1, \dots, A\}$ is a finite set of actions. Without loss of generality, assume that every action is admissible from every state.

MDP Costs. $C \in \mathbb{R}^{S \times A}$ is the matrix encoding the MDP cost. Each $C_{sa} \in \mathbb{R}$ is the cost of taking action $a \in [A]$ from state $s \in [S]$. We also denote the cost of all actions at state s by $c_s = [C_{s1}, \dots, C_{sA}] \in \mathbb{R}^A$, such that $C = [c_1, \dots, c_S]^\top$.

MDP Transition Dynamics. The transition probabilities when action a is taken from state s are given by $p_{sa} \in \Delta_S$. Collectively, all possible transition probabilities from state $s \in [S]$ are given by the matrix $P_s = [p_{s1}, \dots, p_{sA}] \in \Delta_S^A \subset \mathbb{R}^{S \times A}$, and all possible transition probabilities in the MDP are given by the matrix $P = [P_1, \dots, P_S] \in \Delta_S^A \subset \mathbb{R}^{S \times S \times A}$.

MDP Policy. The policy is denoted as $\pi = [\pi_1, \dots, \pi_S] \in \Delta_A^S$, where the a th element of $\pi_s \in \Delta_A$ is the conditional probability of action a being chosen from state s . We also use $\pi(s)$ to denote the corresponding discrete random variable with probability density π_s .

MDP Objective. Under policy π , the decision maker's expected cost-to-go is given per state by

$$V_s^\pi := \mathbb{E}_s \left\{ \sum_{t=0}^{\infty} \gamma^t C_{st a^t} \mid s^0 = s, a^t \sim \pi(s^t) \right\}, \quad \forall s \in [S], \quad (2)$$

where $\mathbb{E}_s\{\cdot\}$ is the expected value of the input with respect to initial state s , and (s^t, a^t) are the state and action at time t . The decision maker's objective is to minimize V_s for all $s \in [S]$. We denote the minimum as V_s^* .

$$V_s^* := \min_{\pi \in \Delta_A^S} \mathbb{E}_s \left\{ \sum_{t=0}^{\infty} \gamma^t C_{st a^t} \mid s^0 = s, a^t \sim \pi(s^t) \right\}, \quad (3)$$

Under policy π_s , the expected immediate cost at s is given by $c_s^\top \pi_s \in \mathbb{R}$ and the expected transition probabilities from s is given by $P_s \pi_s \in \Delta_S$.

Remark 1. Although *value function* is the standard term for the expected cost-to-go, we use value vector in this paper to emphasize that the cost-to-go values of finite MDPs belong in a finite-dimensional space.

2.1. Value operators

We can find the optimal policy by finding the value vector that minimizes the objective (3). Typical solution methods utilize *order preserving* (Schröder, 2003, Def. 3.1), *α -contractive operators* whose fixed points are the MDP's optimal value vectors (e.g. Bellman operator (Puterman, 2014, Thm. 6.2.3), Q-value operator (Melo, 2001)).

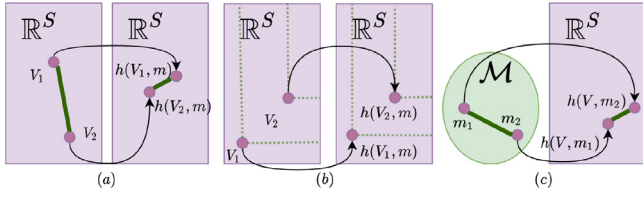


Fig. 1. Visualization of the three value operator properties. (a) α -contraction on \mathbb{R}^S , (b) Order preservation on \mathbb{R}^S , and (c) $K(V)$ -Lipschitz in input space \mathcal{M} .

Definition 1 (α -Contraction). Let (\mathcal{X}, d) be a metric space with metric d . The operator $H : \mathcal{X} \mapsto \mathcal{X}$ is an α -contraction if and only if there exists $\alpha \in [0, 1)$ such that

$$d(H(V), H(V')) \leq \alpha d(V, V'), \quad \forall V, V' \in \mathcal{X}. \quad (4)$$

Definition 2 (Order Preservation). Let (\mathcal{X}, \leq) be an ordered space with partial order \leq . The operator $H : \mathcal{X} \mapsto \mathcal{X}$ is order preserving if for all $V, V' \in \mathcal{X}$ such that $V \leq V'$, $H(V) \leq H(V')$.

These operators are typically locally Lipschitz in the MDP parameter space.

Definition 3 ($K(V)$ -Lipschitz). Let $(\mathcal{X}, d_{\mathcal{X}})$ be a metric space with metric $d_{\mathcal{X}}$ and $(\mathcal{Y}, d_{\mathcal{Y}})$ be a metric space with metric $d_{\mathcal{Y}}$. The operator $H : \mathcal{X} \times \mathcal{Y} \mapsto \mathcal{X}$ is $K(V)$ -Lipschitz with respect to $\mathcal{M} \subset \mathcal{Y}$ if for all $V \in \mathcal{X}$, there exists $K(V) \in \mathbb{R}_+$ such that

$$d_{\mathcal{X}}(H(V, m), H(V, m')) \leq K(V) d_{\mathcal{Y}}(m, m'), \quad \forall m, m' \in \mathcal{M}. \quad (5)$$

Remark 2. The α -contraction property is a special case of Lipschitz continuity when the input and output spaces are identical and the Lipschitz constant is less than 1.

To capture operators with these properties, we define a value operator that takes inputs: the value vector, the cost, and the transition dynamics. The cost and transition dynamics are selected from a parameter set \mathcal{M} .

Definition 4 (Value Operator). Consider the operator h , given by

$$h : \mathbb{R}^S \times \mathcal{M} \mapsto \mathbb{R}^S, \quad \mathcal{M} \subseteq \mathbb{R}^{S \times A} \times \Delta_S^{SA}. \quad (6)$$

We say that h (6) is a **value operator** (see Fig. 1) on $\mathbb{R}^S \times \mathcal{M}$ if

- (1) For all $m \in \mathcal{M}$, $h(\cdot, m)$ is α -contractive on \mathbb{R}^S .
- (2) For all $m \in \mathcal{M}$, $h(\cdot, m)$ is order preserving in \mathbb{R}^S .
- (3) For all $V \in \mathbb{R}^S$, $h(V, m)$ is $K(V)$ -Lipschitz on \mathcal{M} .

Remark 3. Definition 4 and the subsequent results can be extended to the space of Q-value functions by swapping \mathbb{R}^S for $\mathbb{R}^{S \times A}$ in Definition 4 (Melo, 2001).

An immediate consequence an operator being α -contractive and order-preserving on \mathbb{R}^S is that it is continuous on $\mathbb{R}^S \times \mathcal{M}$.

Lemma 1 (Continuity). If h (6) is a value operator on $\mathbb{R}^S \times \mathcal{M}$, h is continuous on $\mathbb{R}^S \times \mathcal{M}$.

See Appendix B for proof. Examples of value operators include the Bellman operator and the policy evaluation operators when the cost and transition dynamics are input parameters rather than fixed parameters.

Definition 5 (Policy Evaluation Operator). Given a policy $\pi \in \Delta_A^S$, the vector-valued operator $g^\pi = (g_1^\pi, \dots, g_S^\pi) : \mathbb{R}^S \times \mathbb{R}^{S \times A} \times \Delta_S^{SA} \mapsto \mathbb{R}^S$ is defined per state as

$$g_s^\pi(V, C, P) := c_s^\top \pi_s + \gamma (P_s \pi_s)^\top V, \quad \forall s \in [S]. \quad (7)$$

Given (C, P) , $g^\pi(\cdot, C, P) : \mathbb{R}^S \mapsto \mathbb{R}^S$ is a vector-valued operator whose fixed point is the expected cost-to-go of the MDP $([S], [A], C, P, \gamma)$ under π , denoted as $V^\pi(C, P)$ (Puterman, 2014, Thm. 6.2.5).

$$V^\pi(C, P) = g^\pi(V^\pi, C, P), \quad V^\pi(C, P) \in \mathbb{R}^S. \quad (8)$$

When the context is clear, we denote $V^\pi(C, P)$ as V^π .

Definition 6 (Bellman Operator). The vector-valued operator $f = (f_1, \dots, f_S) : \mathbb{R}^S \times \mathbb{R}^{S \times A} \times \Delta_S^{SA} \mapsto \mathbb{R}^S$ is defined per each state as

$$f_s(V, C, P) := \inf_{\pi_s \in \Delta_A} g_s^\pi(V, C, P), \quad \forall s \in [S]. \quad (9)$$

The corresponding optimal policy $\pi^* = (\pi_1^*, \dots, \pi_S^*)$ is defined per state as $\pi_s^* \in \operatorname{argmin}_{\pi_s \in \Delta_A} g_s^\pi(V, C, P)$ (9). One such policy is defined for all $(s, a) \in [S] \times [A]$ by

$$\pi_{sa}^* := \begin{cases} > 0 & a \in \operatorname{argmin}_{a \in [A]} C_{sa} + \gamma p_{sa}^\top V \\ 0 & \text{otherwise,} \end{cases} \quad \sum_{a \in [A]} \pi_{sa}^* = 1. \quad (10)$$

where $\operatorname{argmin}_{a \in [A]}(h)$ is the set of minimizing actions for the function h . An optimal policy in the form (10) always exists for a discounted infinite horizon MDP (Puterman, 2014, Thm 6.2.10). For a parameters (C, P) , $f(\cdot, C, P) : \mathbb{R}^S \mapsto \mathbb{R}^S$ is a vector operator whose fixed point is the optimal cost-to-go for the MDP $([S], [A], P, C, \gamma)$, denoted as $V^B(C, P)$.

$$V^B(C, P) = f(V^B, C, P), \quad V^B(C, P) \in \mathbb{R}^S. \quad (11)$$

When the context is clear, we denote $V^B(C, P)$ as V^B .

We show that both (7) and (9) are value operators.

Lemma 2. The Bellman operator (9) and the policy evaluation operators (7) for all $\pi \in \Delta_A^S$ are value operators on $\mathbb{R}^S \times \mathcal{M}$ where $\mathcal{M} \subseteq \mathbb{R}^{S \times A} \times \Delta_S^{SA}$ (6).

See Appendix C for proof.

Remark 4. Beyond the policy evaluation operator and the Bellman operator, many algorithms in reinforcement learning can be cast as value operators. For example, the Q-learning operator (Melo, 2001) and the off-policy temporal difference operator (Chen, Maguluri, Shakkottai, & Shanmugam, 2021) are both value operators on \mathbb{R}^{SA} .

3. Set-based value operators

We now define set-based value operators with respect to a compact set of MDP parameters, beginning with Hausdorff-type set distances.

Definition 7 (Point-to-Set Distance). The distance between a value vector and a set $\mathcal{V} \subseteq \mathbb{R}^S$ is given by

$$W \mapsto d(W, \mathcal{V}) := \inf_{V \in \mathcal{V}} \|W - V\|_\infty. \quad (12)$$

On the space of compact subsets of \mathbb{R}^S , given by $\mathcal{K}(\mathbb{R}^S)$, the distance between value vector sets extends (12) and is given by the Hausdorff distance (Henrikson, 1999).

Definition 8 (Set-to-Set Distance). The Hausdorff distance between two compact value vector sets $\mathcal{V}, \mathcal{W} \subseteq \mathbb{R}^S$ is given by

$$d_{\mathcal{K}}(\mathcal{V}, \mathcal{W}) := \max \left\{ \sup_{V \in \mathcal{V}} d(V, \mathcal{W}), \sup_{W \in \mathcal{W}} d(W, \mathcal{V}) \right\}. \quad (13)$$

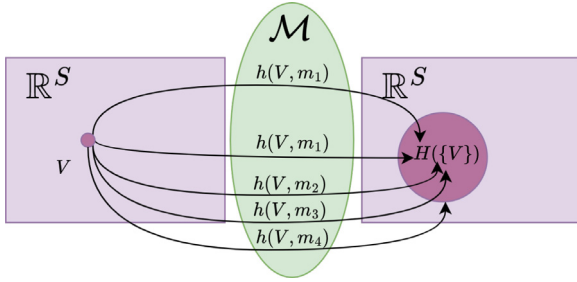


Fig. 2. Visualization of the set-based operator $H(\mathcal{V})$ applied to a singleton set $\mathcal{V} = \{V\} \subset \mathbb{R}^S$: $H(\mathcal{V}) = \cup_{m \in \mathcal{M}} h(V, m)$.

We use $(\mathcal{K}(\mathbb{R}^S), d_{\mathcal{K}})$ to denote the metric space formed by the set of all compact subsets of \mathbb{R}^S under the Hausdorff distance $d_{\mathcal{K}}$. The induced Hausdorff space is complete if and only if the original metric space is complete (Henrikson, 1999, Thm 3.3). Therefore, $(\mathcal{K}(\mathbb{R}^S), d_{\mathcal{K}})$ is a complete metric space.

Given a value operator h (6), we ask: what is the set of possible value vectors when the MDP exhibits environmental non-stationarity captured by \mathcal{M} ? To resolve this, we define the set-based value operator H .

Definition 9 (Set-Based Value Operator). The set-valued operator H is induced by h on $\mathbb{R}^S \times \mathcal{M}$ (6) and is defined as

$$H(\mathcal{V}) := \{h(V, m) \mid (V, m) \in \mathcal{V} \times \mathcal{M}\} \subseteq \mathbb{R}^S, \quad (14)$$

where $\mathcal{V} \subseteq \mathbb{R}^S$ is a subset of the value vector space (see Fig. 2).

We denote the set-based value operator induced by the Bellman operator (9) and the policy evaluation operators (7) as F and G^π , respectively, such that for any value vector set $\mathcal{V} \subseteq \mathbb{R}^S$,

$$F(\mathcal{V}) := \{f(V, C, P) \mid (V, C, P) \in \mathcal{V} \times \mathcal{M}\}, \quad (15)$$

$$G^\pi(\mathcal{V}) := \{g^\pi(V, C, P) \mid (V, C, P) \in \mathcal{V} \times \mathcal{M}\}, \forall \pi \in \Delta_A^S. \quad (16)$$

The set-based Bellman F is the union over all one-step optimal value vectors, which may result from different policies, while G^π is the union over all value vectors from the same policy π .

We ask the following: is there a set of value vectors that is invariant with respect to H ? Similar to the value operators h from Definition 4, we affirmatively answer this by demonstrating that H is α -contractive on $\mathcal{K}(\mathbb{R}^S)$.

Theorem 1. If h is a value operator on $\mathbb{R}^S \times \mathcal{M}$ (6) and \mathcal{M} is compact, then the induced set value operator H (14) satisfies

- (1) For all $\mathcal{V} \in \mathcal{K}(\mathbb{R}^S)$, $H(\mathcal{V}) \in \mathcal{K}(\mathbb{R}^S)$;
- (2) H is α -contractive on $(\mathcal{K}(\mathbb{R}^S), d_{\mathcal{K}})$ (13) with a unique fixed point set \mathcal{V}^* given by

$$H(\mathcal{V}^*) = \mathcal{V}^*, \quad \mathcal{V}^* \in \mathcal{K}(\mathbb{R}^S); \quad (17)$$

- (3) The sequence $\{\mathcal{V}^k\}_{k \in \mathbb{N}}$ where $\mathcal{V}^{k+1} = H(\mathcal{V}^k)$ converges to \mathcal{V}^* for any $\mathcal{V}^0 \in \mathcal{K}(\mathbb{R}^S)$.

In particular, these hold for F (15) and G^π (16), whose fixed point sets are denoted as \mathcal{V}^B and \mathcal{V}^π , respectively.

$$F(\mathcal{V}^B) = \mathcal{V}^B \in \mathcal{K}(\mathbb{R}^S), \quad G^\pi(\mathcal{V}^\pi) = \mathcal{V}^\pi \in \mathcal{K}(\mathbb{R}^S), \quad \forall \pi \in \Delta_A^S. \quad (18)$$

Proof. The first statement follows from Lemma 1, since the image of a compact set by a continuous function is compact (Rudin et al., 1964). Let us prove the second statement: for some $\beta \in (0, 1)$, for

all, $\mathcal{V}, \mathcal{V}' \in \mathcal{K}(\mathbb{R}^S)$:

$$\begin{aligned} & d_{\mathcal{K}}(H(\mathcal{V}), H(\mathcal{V}')) \\ &= \max \left\{ \sup_{\substack{V \in \mathcal{V} \\ m \in \mathcal{M}}} d(h(V, m), H(\mathcal{V}')), \sup_{\substack{V' \in \mathcal{V}' \\ m' \in \mathcal{M}}} d(h(V', m'), H(\mathcal{V})) \right\} \\ &\leq \beta d_{\mathcal{K}}(\mathcal{V}, \mathcal{V}') \end{aligned}$$

Take $(V, m) \in \mathcal{V} \times \mathcal{M}$, then $d(h(V, m), H(\mathcal{V}')) \leq \inf_{V' \in \mathcal{V}'} \|h(V, m) - h(V', m)\|_\infty \leq \alpha \inf_{V' \in \mathcal{V}'} \|V - V'\|_\infty$ holds from the α -contractive property of h . Finally,

$$\begin{aligned} \sup_{\substack{V \in \mathcal{V} \\ m \in \mathcal{M}}} d(h(V, m), H(\mathcal{V}')) &\leq \alpha \sup_{V \in \mathcal{V}} \inf_{V' \in \mathcal{V}'} \|V - V'\|_\infty \\ &\leq \alpha d_{\mathcal{K}}(\mathcal{V}, \mathcal{V}') \end{aligned}$$

We use the same technique to prove that

$$\sup_{\substack{V' \in \mathcal{V}' \\ m' \in \mathcal{M}}} d(h(V', m'), H(\mathcal{V})) \leq \alpha d_{\mathcal{K}}(\mathcal{V}, \mathcal{V}'). \quad (19)$$

Finally, $d_{\mathcal{K}}(H(\mathcal{V}), H(\mathcal{V}')) \leq \alpha d_{\mathcal{K}}(\mathcal{V}, \mathcal{V}')$. From the Banach fixed point theorem and the completeness of $(\mathcal{K}(\mathbb{R}^S), d_{\mathcal{K}})$ (Henrikson, 1999, Thm 3.3), H has a unique fixed point \mathcal{V}^* in $\mathcal{K}(\mathbb{R}^S)$.

The third point is a consequence of the Banach fixed point theorem. Finally, (18) holds because f and g^π are value operators (6) on $\mathbb{R}^S \times \mathcal{M}$. \square

Remark 5. Theorem 1's results can be extended to continuous state-action domains if the operator h satisfies Definition 4.

A consequence of Theorem 1 is the existence of a set-based value iteration, given by

$$\mathcal{V}^{k+1} = H(\mathcal{V}^k), \quad \mathcal{V}^0 \in \mathcal{K}(\mathbb{R}^S). \quad (20)$$

Analogous to the standard value iteration, (20) defines a sequence of value vector sets in $\mathcal{K}(\mathbb{R}^S)$ that converges to the fixed point set $\mathcal{V}^* \in \mathcal{K}(\mathbb{R}^S)$. In the next section, we demonstrate how this set-based approach leads to convergence results for non-stationary value iteration.

4. Properties of the fixed point set

We analyze the properties of the fixed point set \mathcal{V}^* under non-stationary value iteration in this section and derive contraction operators that bound \mathcal{V}^* .

4.1. Non-stationary value iteration

Given a value operator h on $\mathbb{R}^S \times \mathcal{M}$, we consider value iteration under a dynamic parameter uncertainty model, as discussed in Nilim and El Ghaoui (2005), where at every iteration, a new set of MDP parameters m^k is chosen from \mathcal{M} such that

$$\mathcal{V}^{k+1} = h(\mathcal{V}^k, m^k), \quad \mathcal{V}^0 \in \mathbb{R}^S, \quad m^k \in \mathcal{M}, \quad \forall k \in \mathbb{N}. \quad (21)$$

In the robust MDP approach (Iyengar, 2005; Nilim & El Ghaoui, 2005), m^k is adversarially modified such that (21). We consider a different scenario in which m^k is chosen from the closed and bounded set \mathcal{M} . In this scenario, convergence of \mathcal{V}^k in \mathbb{R}^S will not occur for all possible sequences of $\{m^k\}_{k \in \mathbb{N}}$. However, we can show convergence results on the set domain by leveraging our fixed point analysis of the set-based operator H (14).

Proposition 1. Let \mathcal{V}^* be the fixed point set of the set-based value operator H (14) induced by h on $\mathbb{R}^S \times \mathcal{M}$ (6). If the non-stationary value iteration (21) satisfies $\{m^k\}_{k \in \mathbb{N}} \subset \mathcal{M}$, then the sequence $\{\mathcal{V}^k\}_{k \in \mathbb{N}}$ defined by (21) satisfies

- (1) $\lim_{k \rightarrow +\infty} d(V^k, \mathcal{V}^*) = 0$,
- (2) there exists a sub-sequence $\{V^{\varphi(k)}\}_{k \in \mathbb{N}}$ that converges to a point in \mathcal{V}^* as $\lim_{k \rightarrow \infty} V^{\varphi(k)} \in \mathcal{V}^*$.

Proof. Let $\{\mathcal{V}^k\}_{k \in \mathbb{N}}$ be a set sequence defined by $\mathcal{V}^0 = \{V^0\}$ and $\mathcal{V}^{k+1} = H(\mathcal{V}^k)$, where H (14) is the set operator induced by h on $\mathbb{R}^S \times \mathcal{M}$. We first show statement (1). From Theorem 1, $\lim_{k \rightarrow \infty} \mathcal{V}^k$ converges to \mathcal{V}^* in $d_{\mathcal{K}}$. Therefore, $0 \leq d(V^k, \mathcal{V}^*) = \inf_{y \in \mathcal{V}^*} \|V^k - y\|_{\infty} \leq \sup_{x \in \mathcal{V}^k} \inf_{y \in \mathcal{V}^*} \|x - y\|_{\infty} \leq d_{\mathcal{K}}(\mathcal{V}^k, \mathcal{V}^*) \rightarrow 0$ as k tends to $+\infty$.

Next, for all $k \in \mathbb{N}$, there exists $N \in \mathbb{N}$ such that for all $n \geq N$, $d(V^n, \mathcal{V}^*) \leq (k+1)^{-1}$. We define the strictly increasing function $\psi_1 : \mathbb{N} \rightarrow \mathbb{N}$, such that $\psi_1(0) = 0$ and for all $k \neq 0$, $\psi_1(k) := \min\{N > \psi_1(k-1) : \forall n \geq N, d(V^n, \mathcal{V}^*) < (k+1)^{-1}\}$. Then, for all $k \in \{1, 2, \dots\}$, there exists $y^{\psi_1(k)} \in \mathcal{V}^*$ such that $\|V^{\psi_1(k)} - y^{\psi_1(k)}\|_{\infty} < (k+1)^{-1}$. As \mathcal{V}^* is compact, there exists $\psi_2 : \mathbb{N} \rightarrow \mathbb{N}$ strictly increasing such that $(y^{\psi_1(\psi_2(k))})_k$ converges to some $y^* \in \mathcal{V}^*$ (Rudin et al., 1964, Thm 3.6). Finally, let $\varepsilon > 0$, there exist $K_1, K_2 \in \mathbb{N}$ such that for all $l \geq K_1$, $(\psi_2(l))^{-1} < \varepsilon/2$ and for all $l' \geq K_2$, $\|y^{\psi_1(\psi_2(l'))} - y^*\|_{\infty} < \varepsilon/2$. So, taking $k \geq \max\{K_1, K_2\}$, we have $\|V^{\psi_1(\psi_2(k))} - y^*\|_{\infty} \leq \|V^{\psi_1(\psi_2(k))} - y^{\psi_1(\psi_2(k))}\|_{\infty} + \|y^{\psi_1(\psi_2(k))} - y^*\|_{\infty} \leq \varepsilon$ and $(V^{\psi_1(\psi_2(k))})_k$ converges to $y^* \in \mathcal{V}^*$. \square

Beyond bounding the asymptotic behavior of value vector trajectories under non-stationary parameters, the fixed point set \mathcal{V} also contains all fixed points of the value operator $h(\cdot, m)$ when $m \in \mathcal{M}$ (6) is fixed.

Corollary 1. Let h (6) be a value operator on $\mathbb{R}^S \times \mathcal{M}$ where \mathcal{M} is compact. For all $m \in \mathcal{M}$, if $V = h(V, m) \in \mathbb{R}^S$ and \mathcal{V}^* is the fixed point set of the induced set-based value operator H (14), $V \in \mathcal{V}^*$.

Proof. We construct sequence $\{V^k\}$ where $V^{k+1} = h(V^k, m)$ and $V^0 = V$. Then $V^k = V$ for all $k \in \mathbb{N}$. From the second point of Proposition 1, $V \in \mathcal{V}^*$ follows. \square

Furthermore, we can bound the transient behavior of (21) when V^0 is an element of the fixed point set \mathcal{V}^* .

Corollary 2 (Transient Behavior). Let \mathcal{V}^* be the fixed point of the set-based value operator H (14) induced by h on $\mathbb{R}^S \times \mathcal{M}$. If \mathcal{M} is compact and $V^0 \in \mathcal{V}^*$, then the sequence generated by (21) satisfies $\{V^k\}_{k \in \mathbb{N}} \subseteq \mathcal{V}^*$.

Proof. As the fixed point of H (14), \mathcal{V}^* (17) satisfies $\mathcal{V}^* = H(\mathcal{V}^*)$, then the following is true by definition of H : if $V^k \in \mathcal{V}^*$, then $V^{k+1} = h(V^k, m^k) \in \mathcal{V}^*$. If $V^0 \in \mathcal{V}^*$, then $\{V^k\}_{k \in \mathbb{N}} \subseteq \mathcal{V}^*$ follows by induction. \square

Proposition 1 and Corollary 2 bound the asymptotic and transient behavior of the sequence $\{h(V^k, m^k)\}_{k \in \mathbb{N}}$ from (21), irrespective of the convergence of the value vector sequence. This is a more general result than the classic convergence results for MDPs and robust MDPs.

4.2. Bounds of the fixed point set

Since \mathcal{V}^* is compact, it must have finite supremum and infimum value vectors. We show that it is possible to define contraction operators whose fixed points correspond to these extremal elements.

Greatest and least elements. We define the supremum and infimum elements of a value vector set $\mathcal{V} \in \mathcal{K}(\mathbb{R}^S)$ element-wise as follows,

$$\bar{V}_s := \sup_{V \in \mathcal{V}} V_s, \quad \underline{V}_s := \inf_{V \in \mathcal{V}} V_s, \quad \forall s \in [S]. \quad (22)$$

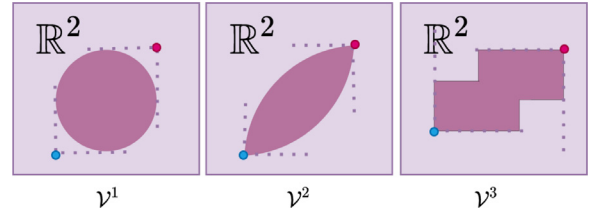


Fig. 3. The greatest least bounds of three different value function sets $\mathcal{V}^i \in \mathbb{R}^2$, where $(0, 0)$ the origin is located on the lower left. Sets \mathcal{V}^2 and \mathcal{V}^3 contain their own greatest and least elements, but \mathcal{V}^1 does not.

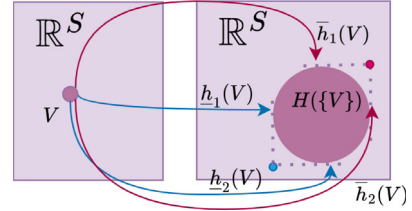


Fig. 4. We visualize \underline{h}/\bar{h} for $H(\mathcal{V})$ on $\mathbb{R}^S \times \mathcal{M}$. The input $\mathcal{V} = \{V\} \in \mathbb{R}^2$. Because \underline{h}_1 and \underline{h}_2 are achieved for two different $m \in \mathcal{M}$, the resulting $\underline{h}(V)$ lies outside of $H(\mathcal{V})$.

If a set $\mathcal{V} \subseteq \mathbb{R}^S$ is compact, the coordinate-wise supremum and infimum values are achieved by elements of \mathcal{V} . However, it is not obvious that a single element in \mathcal{V} can achieve the minimum over all states—i.e., $\bar{V}(\mathcal{V})$ may not be an element of \mathcal{V} . This is illustrated in Fig. 3. We introduce the following bound operators (see Fig. 4).

Definition 10 (Bound Operators). The bound operators induced by the value operator h on $\mathbb{R}^S \times \mathcal{M}$ are coordinate-wise defined at each $s \in [S]$ as

$$\underline{h}_s(V) = \inf_{m \in \mathcal{M}} h_s(V, m), \quad \bar{h}_s(V) = \sup_{m \in \mathcal{M}} h_s(V, m). \quad (23)$$

Our goal is to bound the fixed point set \mathcal{V} of the set-based value operator H (14) using the bound operators \underline{h}/\bar{h} (23). First we show that \underline{h}/\bar{h} are α -contractive and order preserving on \mathbb{R}^S .

Lemma 3 (α -Contraction). If h (6) is a value operator on $\mathbb{R}^S \times \mathcal{M}$ and \mathcal{M} is compact, then \underline{h} and \bar{h} (23) are α -contractions with fixed points \underline{X}, \bar{X} , respectively.

$$\bar{h}(\bar{X}) = \bar{X}, \quad \underline{h}(\underline{X}) = \underline{X}, \quad \underline{X}, \bar{X} \in \mathbb{R}^S. \quad (24)$$

Proof. From Lemma 1, h is continuous and \mathcal{M} is compact, then for all $X, Y \in \mathbb{R}^S$, there exists $\hat{m}(s) \in \mathcal{M}$ such that $\underline{h}_s(Y) = h_s(Y, \hat{m}(s))$ and $\underline{h}_s(X) \leq h_s(X, \hat{m}(s))$. We upper-bound $\underline{h}_s(X) - \underline{h}_s(Y)$ by $h_s(X, \hat{m}(s)) - h_s(Y, \hat{m}(s))$, and use the α -contraction property of h to derive

$$\begin{aligned} \underline{h}_s(X) - \underline{h}_s(Y) &\leq |h_s(X, \hat{m}(s)) - h_s(Y, \hat{m}(s))| \\ &\leq \alpha \|X - Y\|_{\infty}. \end{aligned}$$

Since X and Y are arbitrarily ordered, we conclude that $\|\underline{h}(X) - \underline{h}(Y)\|_{\infty} \leq \alpha \|X - Y\|_{\infty}$. The proof for \bar{h} follows a similar reasoning and takes $\hat{m}(s) = \operatorname{argmax}_{m \in \mathcal{M}} h_s(X, m)$. The existence of $\underline{X}(\bar{X})$ follows from applying Banach's fixed point theorem. \square

Lemma 4 (Order Preservation). The bound operators \underline{h} and \bar{h} (23) are order-preserving on \mathbb{R}^S (Definition 2), i.e.,

$$\forall U, V \in \mathbb{R}^S, U \leq V \Rightarrow \underline{h}(U) \leq \underline{h}(V), \bar{h}(U) \leq \bar{h}(V).$$

Proof. The lemma statement follows directly from the fact that order preservation is conserved through composition with inf and

sup. If $h(U, m) \leq h(V, m)$, then $\inf_{m \in \mathcal{M}} h(U, m) \leq \inf_{m \in \mathcal{M}} h(V, m)$. A similar argument follows for $\bar{h}(\cdot) = \sup_{m \in \mathcal{M}} h(\cdot, m)$. \square

We show that the fixed points \underline{X} and \bar{X} bound the fixed point set \mathcal{V}^* of the set-based value operator H (14).

Theorem 2 (Bounding Fixed Point Sets). If h (6) is a value operator on $\mathbb{R}^S \times \mathcal{M}$ and \mathcal{M} is compact,

$$\underline{X} \leq V \leq \bar{X}, \quad \forall V \in \mathcal{V}^*, \quad (25)$$

where \underline{X} and \bar{X} (24) are the fixed points of the bound operators \underline{h} and \bar{h} (23). Here, \mathcal{V}^* is the fixed point set of the set-based value operator H (14) induced by h (6) on $\mathbb{R}^S \times \mathcal{M}$.

Proof. For $\mathcal{V}^0 = \{\underline{X}, \bar{X}\}$ and $\mathcal{V}^{k+1} = H(\mathcal{V}^k)$ (20), we first show

$$\underline{X} \leq V \leq \bar{X}, \quad \forall V \in \mathcal{V}^k, \quad (26)$$

via induction. Suppose that (26) is satisfied for \mathcal{V}^k . The order preserving property of $h(\cdot, m)$ implies that $h(\underline{X}, m) \leq h(V, m) \leq h(\bar{X}, m)$ holds for all $(V, m) \in \mathcal{V}^k \times \mathcal{M}$. We take the infimum and supremum over $h(\underline{X}, m)$ and $h(\bar{X}, m)$, respectively, to show that for all $(V, m) \in \mathcal{V}^k \times \mathcal{M}$ and $s \in [S]$,

$$\inf_{m' \in \mathcal{M}} h_s(\underline{X}, m') \leq h_s(V, m) \leq \sup_{m' \in \mathcal{M}} h_s(\bar{X}, m').$$

Since \underline{X} and \bar{X} are the fixed points of $\inf_{m' \in \mathcal{M}} h_s(\cdot, m')$ and $\sup_{m' \in \mathcal{M}} h_s(\cdot, m')$ for all $s \in [S]$, respectively, we conclude that (26) holds for \mathcal{V}^{k+1} .

Next, we show that \underline{X} and \bar{X} bound the fixed point set \mathcal{V}^* for the h -induced operator H (14). From Lemma 5, we know that for all $V \in \mathcal{V}^*$, there exists a strictly increasing sequence $\phi : \mathbb{N} \mapsto \mathbb{N}$ and corresponding value vectors $\{W^{\phi(n)}\}$ such that $\lim_{n \rightarrow \infty} W^{\phi(n)} = V$ and $W^{\phi(n)} \in \mathcal{V}^{\phi(n)}$ for the sequence of value vector sets generated from $\mathcal{V}^0 = \{\underline{X}, \bar{X}\}$. Since $\underline{X} \leq W^{\phi(n)} \leq \bar{X}$ holds for all n , we conclude (25) holds. \square

5. Revisiting min-max MDP

The extremum of the set-based Bellman operator's fixed point set is equivalent to the min-max value vector under the rectangularity condition, though it exists under fewer restrictions. We formally prove this relationship by re-examining robust MDPs using a set-theoretic approach. Recall the optimistic value vector $W^o \in \mathbb{R}^S$ and robust value vector $W^r \in \mathbb{R}^S$ of a discounted MDP $([S], [A], C, P, \gamma)$ from Iyengar (2005) and Nilim and El Ghaoui (2005) as the fixed points of the following operators.

$$W_s^o = \min_{\pi_s \in \Delta_A} \min_{(C,P) \in \mathcal{M}} g_s^\pi(W^o, C, P), \quad \forall s \in [S], \quad (27)$$

$$W_s^r = \min_{\pi_s \in \Delta_A} \max_{(C,P) \in \mathcal{M}} g_s^\pi(W^r, C, P), \quad \forall s \in [S]. \quad (28)$$

The optimistic policy π^o and robust policy π^r are the optimal policies corresponding to (27) and (28), respectively.

$$\pi_s^o \in \operatorname{argmin}_{\pi_s \in \Delta_A} \min_{(C,P) \in \mathcal{M}} g_s^\pi(W^o, C, P), \quad \forall s \in [S] \quad (29)$$

$$\pi_s^r \in \operatorname{argmin}_{\pi_s \in \Delta_A} \max_{(C,P) \in \mathcal{M}} g_s^\pi(W^r, C, P), \quad \forall s \in [S] \quad (30)$$

For clarity, we denote the policy evaluation operator (7) under π^o as g^o and the policy evaluation operator (7) under π^r as g^r .

When \mathcal{M} is (s, a) -rectangular (35), the set of policies satisfying (29) and (30) is non-empty and includes deterministic policies (Iyengar, 2005, Thm 3.1). When \mathcal{M} is s -rectangular and convex, the set of policies satisfying (30) is non-empty but may be mixed (Wiesemann et al., 2013, Thm 4). When \mathcal{M} is convex, we show that policies (29) and (30) exist.

Proposition 2. If the MDP parameter set \mathcal{M} is compact and convex, then

- (1) W^o (27) and W^r (28) exist and satisfy $\bar{f}(W^r) = W^r$, $\underline{f}(W^o) = W^o$, where \bar{f} and \underline{f} (23) are the bound operators of the Bellman operator (9).
- (2) π^o (29) and π^r (30) exist.

Proof. Recall the Bellman operator f (9). When $\mathcal{M} \times \Delta_A$ is compact, the formulation of the fixed point of \underline{f} (23) is equivalently given by

$$\underline{f}_s(\underline{X}) = \min_{(C,P) \in \mathcal{M}} \min_{\pi_s \in \Delta_A} g_s^\pi(\underline{X}, C, P), \quad \forall s \in [S]. \quad (31)$$

We note that (31) is identical to the formulation of W^o (27). Therefore, $W^o = \underline{X}$ is the fixed point of \underline{f} . When \mathcal{M} is compact, W^o exists due to Lemma 3. From (29), π_s^o is the optimal argument of $g_s^\pi(W^o, C, P)$, a continuous function in π_s, C, P minimized over compact sets $\Delta_A \times \mathcal{M}$ for all $s \in [S]$. Therefore π_s^o exists. Since $\pi^o = (\pi_1^o, \dots, \pi_S^o)$, the optimal $\pi^o \in \Delta_A^S$ exists.

For the robust scenario: when \mathcal{M} is compact, the fixed point of \bar{f} (23), \bar{X} , exists from Lemma 3 and is given by

$$\bar{X}_s = \max_{(C,P) \in \mathcal{M}} \min_{\pi_s \in \Delta_A} g_s^\pi(\bar{X}, C, P), \quad \forall s \in [S]. \quad (32)$$

The function $g_s^\pi(\bar{X}, C, P)$ is concave in (C, P) and convex in π . If \mathcal{M} is convex, then we apply the minimax theorem (Neumann, 1928) to switch the order of min and max in (32) to derive

$$\bar{X}_s = \min_{\pi_s \in \Delta_A} \max_{(C,P) \in \mathcal{M}} g_s^\pi(\bar{X}, C, P), \quad \forall s \in [S]. \quad (33)$$

Eq. (33) is identical to (28), therefore $W^r = \bar{X}$ and exists by Lemma 3. In (33), $\max_{(C,P) \in \mathcal{M}} g_s^\pi(\bar{X}, C, P)$ is piece-wise linear in π_s and Δ_A is compact for all $s \in [S]$, thus $\operatorname{argmin}_{\pi_s \in \Delta_A} \max_{(C,P) \in \mathcal{M}} g_s^\pi(\bar{X}, C, P)$ is non-empty. Finally, since $\pi^r = (\pi_1^r, \dots, \pi_S^r)$, π^r exists. \square

Remark 6. Since $\max_{(C,P) \in \mathcal{M}} g_s^\pi(\bar{X}, C, P)$ is piecewise linear in π_s , the optimal π_s^r is mixed in general. This is consistent with the results in Wiesemann et al. (2013).

Proposition 2 generalizes the results from Wiesemann et al. (2013) to show that (28) exists when \mathcal{M} is compact and convex instead of s -rectangular and convex. In particular, if we construct $\hat{\mathcal{M}} = \prod_{s \in [S]} \operatorname{proj}_s(\mathcal{M})$, where $\operatorname{proj}_s(\mathcal{M})$ is the projection of the elements of \mathcal{M} onto the s coordinate, then the policies π_o and π_r can be computed using the robust policy iteration from Wiesemann et al. (2013). Diverging from the min-max approach, W^o and W^r always exist and do not require the rectangularity conditions (Iyengar, 2005).

5.1. Containment and rectangularity

Since the fixed point set \mathcal{V}^* is a subset of a multi-dimensional vector space \mathbb{R}^S , it is possible that \mathcal{V}^* does not contain its own extremal elements. We illustrate this in Fig. 5. We also show that interestingly, MDP rectangularity is a sufficient condition for the fixed point set to contain its own extrema.

Fig. 5 implies that the structure of \mathcal{M} may dictate whether the fixed point set contains its own extrema points. We formally prove this in the following section.

Assumption 1 (Containment Condition). The parameter uncertainty set \mathcal{M} satisfies the containment condition with respect to h if \mathcal{M} is compact and for all $V \in \mathbb{R}^S$,

$$\bigcap_{s \in [S]} \operatorname{argmin}_{m \in \mathcal{M}} h_s(V, m) \neq \emptyset, \quad \bigcap_{s \in [S]} \operatorname{argmax}_{m \in \mathcal{M}} h_s(V, m) \neq \emptyset. \quad (34)$$

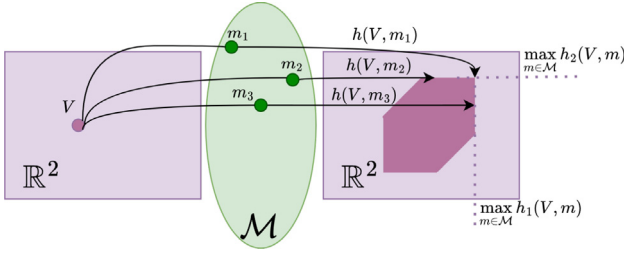


Fig. 5. The $\operatorname{argmax}_{m \in \mathcal{M}} h_s(V, m)$ for a value operator h given \mathcal{M} when $S = 2$. Here, $\operatorname{argmax}_{m \in \mathcal{M}} h_1(V, m) = \{m_1, m_3\}$, $\operatorname{argmax}_{m \in \mathcal{M}} h_2(V, m) = \{m_1, m_2\}$. Therefore, \mathcal{M} satisfies [Assumption 1](#).

[Assumption 1](#) states that (1) there exists a parameter m such that $h(\cdot, m)$'s fixed point is the extremum of \mathcal{V}^* and (2) $m \in \mathcal{M}$. Intuitively, if there exists $m \in \mathcal{M}$ such that $V = h(V, m)$, then from [Corollary 1](#), $V \in \mathcal{V}^*$.

Remark 7. [Assumption 1](#) is an h -dependent condition imposed on the structure of \mathcal{M} , and is independent of \mathcal{M} 's convexity and connectivity.

With respect to the Bellman operator f (9) and the policy evaluation operators g^π (7), the following conditions in min-max MDP are sufficient to satisfy [Assumption 1](#).

Definition 11 (*(s, a)-Rectangular Sets* [Iyengar, 2005](#); [Nilim & El Ghaoui, 2005](#)). The uncertainty set $\mathcal{M} \subset \mathbb{R}^{S \times A} \times \Delta_S^{SA}$ is (s, a) -rectangular if

$$\mathcal{M} = \bigtimes_{(s,a) \in [S] \times [A]} \mathcal{M}_{sa}, \quad \mathcal{M}_{sa} \subset \mathbb{R} \times \Delta_S, \quad \forall (s, a) \in [S] \times [A]. \quad (35)$$

Intuitively, (s, a) -rectangularity implies that the MDP parameter uncertainty is *decoupled* between each state-action.

Definition 12 (*s-Rectangular Sets*). The uncertainty set $\mathcal{M} \subset \mathbb{R}^{S \times A} \times \Delta_S^{SA}$ is s -rectangular if

$$\mathcal{M} = \bigtimes_{s \in [S]} \mathcal{M}_s, \quad \mathcal{M}_s \subset \mathbb{R}^A \times \Delta_S^A, \quad \forall s \in [S]. \quad (36)$$

Remark 8. [Definition 12](#) applies to ambiguity sets that arise in distributionally robust MDPs ([Xu & Mannor, 2010](#); [Yang, 2017](#); [Yu & Xu, 2015](#)).

Example 2 (*Wind Uncertainty*). Consider the navigation problem presented in [Example 1](#). If the wind pattern strictly switches between N wind patterns, then the transition uncertainty at state $s \in [S]$ is given by $\mathcal{P}_s = \{P_s^1, \dots, P_s^N\}$. If the wind pattern is a mixture of N discrete wind trends, the transition uncertainty at state $s \in [S]$ is $\mathcal{P}_s = \{\sum_i \alpha_i P_s^i \mid \alpha \in \Delta_N\}$. Both wind patterns lead to s -rectangular uncertainty, given by $\mathcal{P} = \bigtimes_{s \in [S]} \mathcal{P}_s$.

We show that the rectangularity conditions indeed are sufficient for satisfying [Assumption 1](#) with respect to f (9) and g^π (7).

Proposition 3. If \mathcal{M} is compact and s -rectangular ([Definition 12](#)), \mathcal{M} satisfies [Assumption 1](#) with respect to f (9) and g^π (7) for all $\pi \in \Delta_A^S$.

Proof. We first show that \mathcal{M} satisfies [Assumption 1](#) with respect to the Bellman operator. Given $s \in [S]$, $f_s(V, C, P)$ only depends on the s component of C and P . From [Lemma 1](#), f_s is continuous

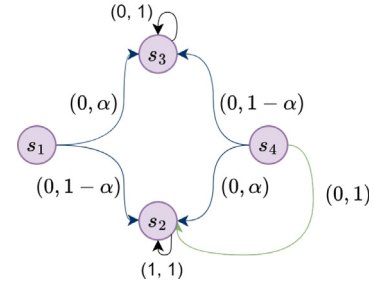


Fig. 6. MDP with parameter coupling in transition probability across different states.

in (c_s, P_s) . Let (c_s^*, P_s^*) be the solution to $\operatorname{argmin}_{(c_s, P_s) \in \mathcal{M}_s} f_s(V, C, P)$ for all $\forall s \in [S]$. If \mathcal{M}_s is compact, $(c_s^*, P_s^*) \in \mathcal{M}_s$. We can construct $C^* = [c_1^*, \dots, c_S^*]$ and $P^* = [P_1^*, \dots, P_S^*]$. If \mathcal{M} is s -rectangular, then $(C^*, P^*) \in \mathcal{M}$ and $(C^*, P^*) \in \operatorname{argmin}_{(C, P) \in \mathcal{M}} f_s(V, C, P)$ for all $s \in [S]$. We conclude that \mathcal{M} satisfies [Assumption 1](#).

Given $\pi \in \Delta_A^S$ and $s \in [S]$, g_s^π only depends on c_s and P_s as well. We can similarly show that there exists an optimal parameter $(C^*, P^*) \in \operatorname{argmin}_{(C, P) \in \mathcal{M}} g_s^\pi(V, C, P)$ for all $s \in [S]$ such that $(C^*, P^*) \in \mathcal{M}$. \square

Beyond s -rectangularity, there are sets that satisfy [Assumption 1](#) with respect to specific value operators.

Example 3 (*Beyond Rectangularity*). Consider a four state MDP with α -parametrized transition uncertainty \mathcal{M} in [Fig. 6](#), where states are the nodes and actions are the multi-headed arrows. Each head has an associated tuple $(C_{sa}, P_{sa, s'})$ denoting its cost and transition probability. The states s_2 and s_3 have values $V_2 = \frac{1}{1-\gamma}$ and $V_3 = 0$ for both f and g^π for all $\pi \in \Delta_A^S$.

The states s_1 and s_4 have transition uncertainty jointly parametrized by $\alpha \in [0, 1]$, therefore violating s -rectangularity ([Definition 12](#)). The optimal cost-to-go values V_1 and V_4 occur at different α 's, therefore violating [Assumption 1](#) with respect to f . However, suppose that at s_4 , we only choose the action with cost-transition $(0, 1)$ to s_2 in [Fig. 6](#). Then V_4 is independent of α . The minimum and maximum V_1 occur at $\alpha = 1$ and $\alpha = 0$, respectively. Therefore, \mathcal{M} satisfies [Assumption 1](#) with respect to operator g^π for all $\pi = [\pi_{s_1}, \dots, \pi_{s_4}]$ where $\pi_{s_4} = [1, 0]$.

When [Assumption 1](#) is satisfied, the fixed point of H (14) contains its own supremum and infimum value vectors.

Theorem 3. If h (6) on $\mathbb{R}^S \times \mathcal{M}$ satisfies [Assumption 1](#), then there exists $\underline{m}, \bar{m} \in \mathcal{M}$ such that \bar{h} and \underline{h} (23) and their fixed points \underline{X} and \bar{X} (24) satisfies

$$\underline{h}(\underline{X}) = h(\underline{X}, \underline{m}) = \underline{X}, \quad \bar{h}(\bar{X}) = h(\bar{X}, \bar{m}) = \bar{X}. \quad (37)$$

Additionally, \underline{X} and \bar{X} are the least and the greatest elements of H 's fixed point set \mathcal{V}^* , $\underline{V}^*, \bar{V}^*$ (22) respectively, and both belong to \mathcal{V}^* (17).

$$\underline{X} = \underline{V}^*, \quad \bar{X} = \bar{V}^*, \quad \underline{X}, \bar{X} \in \mathcal{V}^*.$$

Proof. From [Theorem 2](#), \underline{X} and \bar{X} are the lower and upper bounds on the fixed point set \mathcal{V}^* . We show that these are the infimum and supremum elements of \mathcal{V}^* by showing that they are also elements of \mathcal{V}^* . From [Assumption 1](#), there exists $\underline{m}, \bar{m} \in \mathcal{M}$ such that $h_s(\underline{X}, \underline{m}) = \min_{m \in \mathcal{M}} h_s(\underline{X}, m)$ and $h_s(\bar{X}, \bar{m}) = \max_{m \in \mathcal{M}} h_s(\bar{X}, m)$ for all $s \in [S]$. Since \underline{X} and \bar{X} are fixed points of $h(\cdot, \underline{m})$ and $h(\cdot, \bar{m})$, we apply [Corollary 1](#) to conclude that $\underline{X}, \bar{X} \in \mathcal{V}^*$. \square

5.2. Performance of non-stationary Bellman update

Consider the wind navigation problem in [Example 1](#), one question posed is whether it is better to use the robust policy, optimistic policy, or the next step optimal policy. In this section, we use the set-theoretic tools to compare the performance of these approaches. First, we introduce some notations: let $G^o = G^{\pi^o}$ ([7](#)), the fixed point of G^o be \mathcal{V}^o , $G^r = G^{\pi^r}$, and the fixed point of G^r be \mathcal{V}^r .

$$\mathcal{V}^o = \{g^o(V, C, P) \mid (C, P) \in \mathcal{M}, V \in \mathcal{V}^o\}, \quad (38)$$

$$\mathcal{V}^r = \{g^r(V, C, P) \mid (C, P) \in \mathcal{M}, V \in \mathcal{V}^r\}. \quad (39)$$

Additionally, the supremum value vectors of \mathcal{V}^o and \mathcal{V}^r are \bar{V}^o and \bar{V}^r respectively and the infimum value vectors are \underline{V}^o and \underline{V}^r , respectively.

$$\underline{V}_s^r = \min_{V \in \mathcal{V}^r} V_s, \quad \bar{V}_s^r = \max_{V \in \mathcal{V}^r} V_s, \quad \forall s \in [S]. \quad (40)$$

$$\underline{V}_s^o = \min_{V \in \mathcal{V}^o} V_s, \quad \bar{V}_s^o = \max_{V \in \mathcal{V}^o} V_s, \quad \forall s \in [S]. \quad (41)$$

We compare these with the fixed point set of the Bellman operator, $\mathcal{V}^B = \{\min_{\pi} g^{\pi}(V, C, P) \mid (C, P) \in \mathcal{M}, V \in \mathcal{V}^B\}$ ([17](#)), denoted by \bar{V}^B and \underline{V}^B as

$$\underline{V}_s^B = \min_{V \in \mathcal{V}^B} V_s, \quad \bar{V}_s^B = \max_{V \in \mathcal{V}^B} V_s, \quad \forall s \in [S]. \quad (42)$$

Theorem 4. *If f, g^o, g^r satisfy [Assumption 1](#) on $\mathbb{R}^S \times \mathcal{M}$, then the bounding value vectors [\(42\)](#) [\(41\)](#) [\(40\)](#) of the corresponding fixed point sets \mathcal{V}^B , \mathcal{V}^o [\(38\)](#) and \mathcal{V}^r [\(39\)](#) are ordered as*

$$\underline{V}^B = \underline{V}^o \leq \underline{V}^r, \quad \bar{V}^B = \bar{V}^r \leq \bar{V}^o. \quad (43)$$

Proof. Since \underline{V}^o is the infimum element for the fixed point set \mathcal{V}^o [\(41\)](#), we can apply [Theorem 3](#) to derive

$$\underline{V}^o = \min_{(C,P) \in \mathcal{M}} g^o(\underline{V}^o, C, P). \quad (44)$$

By definition of π^o [\(29\)](#), $\min_{(C,P) \in \mathcal{M}} g^o(\underline{V}^o, C, P) = \min_{(C,P) \in \mathcal{M}} \min_{\pi \in \Delta_A^S} g^{\pi}(\underline{V}^o, C, P)$. As the two minima commute,

$$\min_{(C,P) \in \mathcal{M}} g^o(\underline{V}^o, C, P) = \min_{(C,P) \in \mathcal{M}} \min_{\pi \in \Delta_A^S} g^{\pi}(\underline{V}^o, C, P). \quad (45)$$

Combining [\(44\)](#) and [\(45\)](#), \underline{V}^o is exactly the unique fixed point of $\min_{(C,P) \in \mathcal{M}} \min_{\pi \in \Delta_A^S} g^{\pi}(\cdot, C, P)$. However, by applying [Theorem 3](#) to f on $\mathbb{R}^S \times \mathcal{M}$, \underline{V}^B is also the unique fixed point of $\min_{(C,P) \in \mathcal{M}} \min_{\pi \in \Delta_A^S} g^{\pi}(\cdot, C, P)$. Therefore $\underline{V}^o = \underline{V}^B$.

From [\(40\)](#), $\underline{V}^r = \min_{(C,P) \in \mathcal{M}} g^r(\underline{V}^r, C, P)$, we can minimize over the policy space to lower bound \underline{V}^r as

$$\underline{V}^r \geq \min_{\pi \in \Delta_A^S} \min_{(C,P) \in \mathcal{M}} g^{\pi}(\underline{V}^r, C, P). \quad (46)$$

Since the right hand side of [\(46\)](#) is equivalent to $f(\underline{V}^r)$, [\(46\)](#) is equivalent to $\underline{V}^r \geq f(\underline{V}^r)$. From [Lemma 4](#), f is order-preserving in V , we conclude that $\underline{V}^o = \underline{V}^* \leq \underline{V}^r$.

From [Theorem 3](#), \bar{V}^r is the fixed point of \bar{g}^r , such that

$$\bar{V}^r = \max_{(C,P) \in \mathcal{M}} g^r(\bar{V}^r, C, P). \quad (47)$$

We apply \min_{π} to both sides of [\(47\)](#) and use the definition of π^r to derive that \bar{V}^r is the fixed point of $\min_{\pi \in \Delta_A^S} \max_{(C,P) \in \mathcal{M}} g^{\pi}(V^r, C, P)$. From [Assumption 1](#), there exists $(\bar{C}, \bar{P}) \in \mathcal{M}$ that maximizes $g^{\pi}(\bar{V}, C, P)$, so \bar{V}^r equivalently satisfies

$$\bar{V}^r = \min_{\pi \in \Delta_A^S} g^{\pi}(\bar{V}^r, \bar{C}, \bar{P}).$$

From [Corollary 1](#), $\bar{V}^r \in \mathcal{V}^B$ and therefore $\bar{V}^r \leq \bar{V}^B$. Next we show $\bar{V}^B \leq \bar{V}^r$. From [Theorem 3](#), \bar{V}^B is the fixed point of \bar{f} , such that

$$\bar{V}^B = \max_{(C,P) \in \mathcal{M}} \min_{\pi} g^{\pi}(\bar{V}^B, C, P),$$

From the min-max inequality,

$$\bar{V}^B \leq \min_{\pi \in \Delta_A^S} \max_{(C,P) \in \mathcal{M}} g^{\pi}(\bar{V}^B, C, P).$$

Since $\pi^r \in \Delta_A^S$,

$$\bar{V}^B \leq \max_{(C,P) \in \mathcal{M}} g^r(\bar{V}^B, C, P). \quad (48)$$

The right-hand side of [\(48\)](#) is $\bar{g}^r(\bar{V}^B)$ [\(23\)](#), such that [\(48\)](#) is equivalent to $\bar{V}^B \leq \bar{g}^r(\bar{V}^B)$. Consider the sequence $V^{k+1} = \bar{g}^r(V^k)$ where $V^1 = \bar{V}^B$. Since \bar{g}^r is a contraction, $\lim_{k \rightarrow \infty} V^k = V^r$, the fixed point of \bar{g}^r . From [Lemma 4](#), \bar{g}^r is order preserving. Therefore $\bar{V}^B = V^1 \leq V^r$.

Finally, [Theorem 3](#) implies that \bar{V}^o is the fixed point of \bar{g}^o : $\bar{V}^o = \max_{(C,P) \in \mathcal{M}} g^o(\bar{V}^o, C, P)$. By construction, $\bar{V}^o \geq \min_{\pi \in \Delta_A^S} \max_{(C,P) \in \mathcal{M}} g^{\pi}(\bar{V}^o, C, P)$. From the min-max inequality,

$$\min_{\pi \in \Delta_A^S} \max_{(C,P) \in \mathcal{M}} g^{\pi}(\bar{V}^o, C, P) \geq \max_{(C,P) \in \mathcal{M}} \min_{\pi \in \Delta_A^S} g^{\pi}(\bar{V}^o, C, P),$$

such that the right hand side of the inequality is equivalent to $\bar{f}(\bar{V}^o)$. Following the monotonicity properties of the Bellman operator f ([Puterman, 2014, Thm. 6.2.2](#)), we conclude that $\bar{V}^o \geq \bar{V}^B$. \square

Remark 9. Our set-theoretic approach shows that in addition to having the best worst-case performance among $\{\mathcal{V}^o, \mathcal{V}^B, \mathcal{V}^r\}$, \mathcal{V}^r also has the smallest variation in performance for the same uncertainty set \mathcal{M} .

Recall [Example 1](#), iteratively updating policy using the next step cost and transition dynamics will result in a value vector trajectory that asymptotically converges to \mathcal{V}^B . In that case, [Theorem 4](#) implies this policy update scheme has comparable performance to the optimistic policy in non-adversarial wind fields but will perform no worse than the robust policy in adversarial winds.

Finally, we generalize the s -rectangularity condition by showing that the optimistic and robust policies exist when the MDP parameter set \mathcal{M} satisfies [Assumption 1](#).

Corollary 3 (*Robust MDP under [Assumption 1](#)*). *If \mathcal{M} is compact, convex, and f, g^o, g^r satisfy [Assumption 1](#) on $\mathbb{R}^S \times \mathcal{M}$, then W^o [\(27\)](#) and W^r [\(28\)](#) are the infimum and supremum value vectors of \mathcal{V}^o and \mathcal{V}^r ,*

$$W_s^o = \inf_{V \in \mathcal{V}^o} V_s, \quad W_s^r = \sup_{V \in \mathcal{V}^r} V_s, \quad \forall s \in [S], \quad (49)$$

where \mathcal{V}^o [\(38\)](#) and \mathcal{V}^r [\(39\)](#) are the fixed point sets of G^o and G^r , respectively (see [Fig. 7](#)).

Proof. When f satisfies [Assumption 1](#) on $\mathbb{R}^S \times \mathcal{M}$, [Theorem 3](#) shows that $\underline{V}^B = W^o$, $\bar{V}^B = W^r$. If g^o , and g^r also satisfies [Assumption 1](#) on $\mathbb{R}^S \times \mathcal{M}$, then we apply [Theorem 4](#) to derive $W^o = \underline{V}^o$ and $W^r = \bar{V}^r$. This proves the corollary statement. \square

Remark 10. When [Assumption 1](#) is not satisfied, W^o and W^r still bound \underline{V}^o and \bar{V}^r .

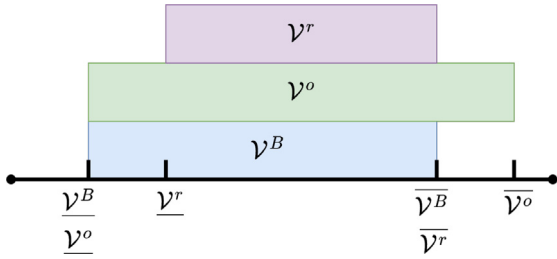


Fig. 7. Visualization of Theorem 4. The purple, green, blue regions indicate the ranges of ν^r , ν^o , and ν^B , respectively.

6. Value iteration for fixed point set computation

In the previous sections, we proved the existence of a fixed point set for value operators with compact parameter uncertainty sets and re-interpreted robust control. Next, we derive an iterative algorithm for computing the bounds of the fixed point set ν given a value operator h and parameter uncertainty set \mathcal{M} .

Algorithm Sketch. Based on the set-based value iteration (20), we iteratively find the one-step bounds of $H(\nu^k)$ to converge the bounds of the fixed point set.

For any compact set $\nu \in \mathcal{K}(\mathbb{R}^S)$, the one step bounds of $H(\nu)$ are the result of applying the one-step bound operators \underline{h} and \bar{h} (23) to the extremal points of ν .

Theorem 5 (One Step H Bounds). Consider the set operator H (14) and its bound operators \underline{h} and \bar{h} (23) induced by h on $\mathbb{R}^S \times \mathcal{M}$ (6). For a compact set $\nu \subset \mathbb{R}^S$, $H(\nu)$ is bounded by $\underline{h}(\underline{\nu})$ and $\bar{h}(\bar{\nu})$ (23) as

$$\underline{h}(\underline{\nu}) \leq V \leq \bar{h}(\bar{\nu}), \quad \forall V \in H(\nu). \quad (50)$$

where $\underline{\nu}$ and $\bar{\nu}$ (22) are the extremal elements of ν . If h satisfies Assumption 1 on $\mathbb{R}^S \times \mathcal{M}$ and $\underline{\nu}, \bar{\nu} \in \nu$, then $\underline{h}(\underline{\nu})$ and $\bar{h}(\bar{\nu})$ are the supremum and infimum elements of $H(\nu)$, respectively—for all $s \in [S]$, $\underline{h}_s(\underline{\nu})$ and $\bar{h}_s(\bar{\nu})$ satisfy

$$\underline{h}_s(\underline{\nu}) = \inf_{(V,m) \in \nu \times \mathcal{M}} h_s(V, m), \quad \bar{h}_s(\bar{\nu}) = \sup_{(V,m) \in \nu \times \mathcal{M}} h_s(V, m). \quad (51)$$

Proof. For all $s \in [S]$, $h_s(V, m) \leq \bar{h}_s(\bar{\nu})$ for all $m \in \mathcal{M}$. If h is $K(V)$ -Lipschitz and α -contractions in \mathcal{M} , then \bar{h} is order-preserving (Lemma 4) such that $\bar{h}_s(\underline{\nu}) \leq \bar{h}_s(\bar{\nu})$ for all $V \in \nu$. We conclude that

$$h(V, m) \leq \bar{h}(\bar{\nu}), \quad \forall (V, m) \in \nu \times \mathcal{M}. \quad (52)$$

Since \bar{h} is an upper bound, and \sup is the least upper bound, it holds that $\sup_{V,m} h_s(V, m) \leq \bar{h}(\bar{\nu})$. We use the definition of $H(\nu)$ (14) to conclude that $V \leq \bar{h}(\bar{\nu})$ for all $V \in H(\nu)$. The inequality $\underline{h}(\underline{\nu}) \leq V \forall V \in H(\nu)$ can be similarly proved.

If h satisfies Assumption 1 on $\mathbb{R}^S \times \mathcal{M}$ and $\underline{\nu}, \bar{\nu} \in \nu$, Assumption 1 states that there exists $\underline{m} \in \mathcal{M}$ such that $h(\underline{\nu}, \underline{m}) = \underline{h}(\underline{\nu})$. Therefore, $\underline{h}(\underline{\nu}) \in H(\nu)$. Since $\underline{h}(\underline{\nu})$ also lower bounds all the elements of $H(\nu)$, it is the infimum element of $H(\nu)$. The fact that the greatest element of $H(\nu)$ is $\bar{h}(\bar{\nu})$ can be similarly proved. \square

Based on Theorem 5, we propose the following bound approximation algorithm of the fixed point set ν^* (17) for a set-valued operator H (6).

6.1. Computing one-step optimal parameters

Algorithm 1 is stated for a general MDP parameter set \mathcal{M} and does not specify how to compute lines 4 and 5. Here are some solution methods for different types of \mathcal{M} .

Algorithm 1 Bounding the fixed point set ν

Input: $\mathcal{C}, \mathcal{P}, V^0, \epsilon$.

Output: \underline{V}, \bar{V}

- 1: $\underline{V}^0 := \bar{V}^0 := V^0$
- 2: $e^0 = \frac{1-\gamma}{\gamma} \epsilon$
- 3: **while** $\frac{\gamma}{1-\gamma} e^k \geq \epsilon$ **do**
- 4: $\underline{V}_s^{k+1} = \min_{m \in \mathcal{M}} h_s(\underline{V}^k, m), \quad \forall s \in [S]$
- 5: $\bar{V}_s^{k+1} = \max_{m \in \mathcal{M}} h_s(\bar{V}^k, m), \quad \forall s \in [S]$
- 6: $e^{k+1} = \max \left\{ \|\underline{V}^{k+1} - \underline{V}^k\|, \|\bar{V}^{k+1} - \bar{V}^k\| \right\}$
- 7: $k = k + 1$
- 8: **end while**

- (1) **Finite** \mathcal{M} . If $\mathcal{M} = \{m_1, \dots, m_N\}$ has a finite number of elements, we can directly compute line 4 as

$$\underline{V}^{k+1} = \min \left\{ h_s(\underline{V}^k, m_i) \mid i = \{1, \dots, N\} \right\}. \quad (53)$$

For line 5, we replace min with max in (53).

- (2) **Convex** \mathcal{M} . When \mathcal{M} is a convex set, the computation depends on h . If $h = g^\pi$ is the policy operator, lines 4 and 5 can be solved as convex optimization problems. If h is the Bellman operator f , lines 4 and 5 take on min-max formulation and is NP-hard to solve in the general form Wiesemann et al. (2013). When \mathcal{M} can be characterized by an ellipsoidal set of parameters, the solutions to lines 4 and 5 is given in Wiesemann et al. (2013).

6.2. Algorithm convergence rate

When lines 4 and 5 are solvable, Algorithm 1 asymptotically converges to approximations of the bounding elements of ν^* . If \mathcal{M} satisfies Assumption 1 with respect to h , Algorithm 1 derives the exact bounds of ν^* . Algorithm 1 has similar rates of convergence in Hausdorff distance as standard value iteration using h on \mathbb{R}^S .

Theorem 6. Consider the value operator h , compact uncertainty set \mathcal{M} , and the fixed point set ν^* of the set-based operator H (14) induced by h on $\mathbb{R}^S \times \mathcal{M}$. If \mathcal{M} satisfies Assumption 1 with respect to h , then at each iteration k ,

$$\begin{aligned} \|\underline{V}^{k+1} - \underline{V}^*\|_\infty &\leq \alpha \|\underline{V}^k - \underline{V}^*\|_\infty, \\ \|\bar{V}^{k+1} - \bar{V}^*\|_\infty &\leq \alpha \|\bar{V}^k - \bar{V}^*\|_\infty, \end{aligned} \quad (54)$$

where all norms are infinity norms, and $\underline{V}^*, \bar{V}^*$ are the infimum and supremum bounds of ν , respectively. At Algorithm 1's termination, $\underline{V}^k, \bar{V}^k$ satisfies

$$\max \left\{ \|\underline{V}^k - \underline{V}^*\|_\infty, \|\bar{V}^k - \bar{V}^*\|_\infty \right\} < \epsilon. \quad (55)$$

Proof. From Algorithm 1, $\bar{V}^{k+1} = \bar{h}(\bar{V}^k)$. From Lemma 3, \bar{h} is an α -contraction. We obtain

$$\|\bar{V}^{k+1} - \bar{V}^*\|_\infty \leq \alpha \|\bar{V}^k - \bar{V}^*\|_\infty$$

and note that (54) holds by induction. Next, we apply triangle inequality to $\|\bar{V}^k - \bar{V}^*\|_\infty$ as

$$\|\bar{V}^k - \bar{V}^*\|_\infty \leq \|\bar{V}^k - \bar{V}^{k+1}\|_\infty + \|\bar{V}^{k+1} - \bar{V}^*\|_\infty. \quad (56)$$

We can then use $\|\bar{V}^{k+1} - \bar{V}^*\|_\infty \leq \alpha \|\bar{V}^k - \bar{V}^*\|_\infty$ to bound (56) as $\|\bar{V}^k - \bar{V}^*\|_\infty \leq \frac{1}{1-\alpha} \|\bar{V}^k - \bar{V}^{k+1}\|_\infty$. A similar argument can show that $\|\underline{V}^k - \underline{V}^*\|_\infty \leq \frac{1}{1-\alpha} \|\underline{V}^k - \underline{V}^{k+1}\|_\infty$. When Algorithm 1's while

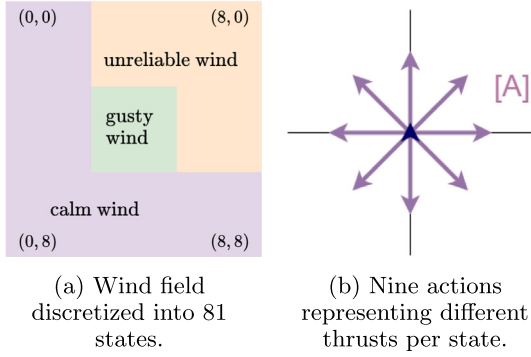


Fig. 8. Wind field MDP state space and action space.

condition is satisfied, $\max\{\|\bar{V}^k - \bar{V}^*\|_\infty, \|\underline{V}^k - \underline{V}^*\|_\infty\} \leq \epsilon$. This concludes our proof. \square

In particular, the Bellman operator f and policy operator g^π are γ -contractive on \mathbb{R}^S , where γ is the discount factor. Therefore, Theorem 5 applies with $\alpha = \gamma$.

Remark 11. Theorem 6 implies that at the termination of Algorithm 1, the fixed point set \mathcal{V}^* can be over-approximated by

$$\mathcal{V}^* \subseteq \mathcal{V}_{approx} := \prod_{s \in [S]} [\underline{V}_s^{k+1} - \epsilon, \bar{V}_s^{k+1} + \epsilon],$$

where k is the last iterate before Algorithm 1 terminates.

7. Path planning in non-stationary wind fields

We apply set-based value iteration to a wind-assisted probabilistic path planning problem in a non-stationary wind field (Wolf et al., 2010). MDP as a model for wind-assisted path planning of balloons in the stratosphere and exoplanets has recently gained traction (Bellemare et al., 2020; Wolf et al., 2010). Finite state-action MDPs have been shown to be a viable high-level path planning model (Wolf et al., 2010) for such applications.

Mission Objective. In a two-dimensional wind-field, the wind-assisted balloon aims to reach target state (8, 8) in Fig. 8 using minimum fuel.

Non-stationary Wind Fields. By collecting wind data on the environment's wind field, an MDP can be created and a policy that handles stochastic planning can be deployed. However, wind can be a time-varying factor that causes the *expected* optimal policy to have *worse-than-expected* worst-case performance. We built an ideal uncertain wind field to demonstrate how the set Bellman operator can be used to predict the best and worst-case behavior of a robust policy.

MDP Modeling Assumptions. Following Wolf et al. (2010), we model the path planning problem in a non-stationary wind field as an infinite horizon, discounted MDP with discrete state-actions. While balloons typically traverse in three dimensions, we assume that the wind is consistent in the vertical direction and that the final target is any vertical position along the given two-dimensional coordinates.

States. States are grouped into three different regions based on their wind variability as shown in Fig. 8: $[S] = [S_{calm}] \cup [S_{gusty}] \cup [S_{unreliable}]$.

- (1) $s \in [S_{calm}]$. The wind magnitude changes between $[0, 0.5]$, and the wind direction changes between $[0, 2\pi]$. $[S_{calm}] = \{(i, j) \mid (0, 0) \leq (i, j) \leq (2, 8), (0, 6) \leq (i, j) \leq (8, 8)\}$.

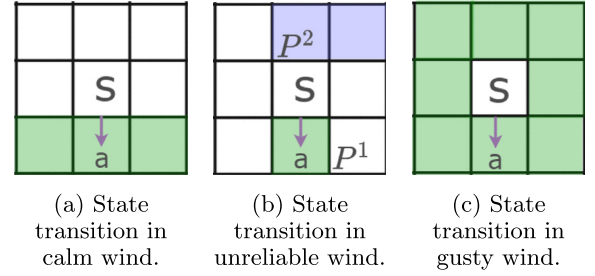


Fig. 9. Transition probabilities for the three different wind regions.

- (2) $s \in [S_{gusty}]$. The wind magnitude is 1, while the wind direction changes between $[0, 2\pi]$. $S_{gusty} = \{(i, j) \mid (3, 3) \leq (i, j) \leq (5, 5)\}$.
- (3) $s \in [S_{unreliable}]$. A wind front occasionally moves across an otherwise windless region. The wind magnitude is either 0 or 1 and the wind direction changes between $[\pi/4, \pi/2]$.

Actions. The balloon is equipped with an actuator that provides a constant thrust of 1 unit power in 8 discretized directions shown in Fig. 8b. We assume that this actuation force is enough to move the balloon across one state in wind with magnitude ≤ 0.5 , and is otherwise not strong enough to overcome the wind effects. In addition, each state has an action corresponding to no thrust.

Transition Probabilities. The transition probabilities are region-dependent. In the states $[S_{calm}]$ and $[S_{gusty}]$, the transition dynamics are stochastic but stationary in time. In the states $[S_{unreliable}]$, the transition dynamics are stochastic but change over time. We define the following neighboring states for each state $s \in [S]$.

- (1) $\mathcal{N}(s)$: all 8 neighboring states of state s .
- (2) $\mathcal{N}(s, a, 0)$: the neighboring state of s in the direction of a .
- (3) $\mathcal{N}(s, a, 1)$: the neighboring state of s in the direction of a plus the two states adjacent to the neighbor state, shown in green in Fig. 9a.
- (4) $\mathcal{N}(s, a, 2)$: state s 's neighboring state in the opposite direction of the action a plus its clockwise neighbor, shown in purple in Fig. 9b.

In the calm wind region, the transition probabilities are given by

$$P_{sa,s'} = \begin{cases} \frac{1}{\mathcal{N}(s,a,1)}, & s' \in \mathcal{N}(s, a, 1) \\ 0 & \text{otherwise,} \end{cases} \quad \forall s \in [S_{calm}]. \quad (57)$$

In the gusty wind region, the transition probabilities are given by

$$P_{sa,s'} = \begin{cases} \frac{1}{\mathcal{N}(s)}, & s' \in \mathcal{N}(s) \\ 0 & \text{otherwise,} \end{cases} \quad \forall s \in [S_{gusty}], \quad \forall a \in [A]. \quad (58)$$

In the unreliable wind region, the transition probabilities vary between transition dynamics P_s^1 and P_s^2 .

$$P_{sa,s'}^1 = \begin{cases} 1, & s' \in \mathcal{N}(s, a, 0) \\ 0 & \text{otherwise,} \end{cases} \quad \forall s \in [S_{unreliable}], \quad \forall a \in [A]. \quad (59)$$

$$P_{sa,s'}^2 = \begin{cases} 0.5, & s' \in \mathcal{N}(s, a, 2) \\ 0 & \text{otherwise,} \end{cases} \quad \forall s \in [S_{unreliable}], \quad \forall a \in [A]. \quad (60)$$

Collectively, P_s^1 and P_s^2 collectively form the uncertainty set $\mathcal{P}_s \subset \Delta_s^A$ defined at each state.

$$\mathcal{P}_s = \{P_{sa}^i \mid i \in \{1, 2\}, a \in [A]\}, \quad \forall s \in [S_{unreliable}]. \quad (61)$$

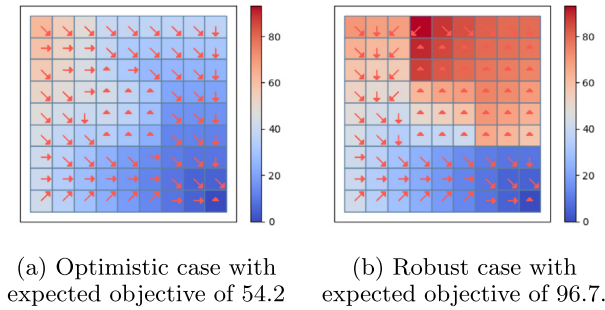


Fig. 10. Visualization of optimistic and robust policies.

Table 1

Bellman, optimistic policy, robust policy value bounds.

Set	Maximum value	Minimum value
\mathcal{V}^B	70.61	62.25
\mathcal{V}^o	101.58	62.25
\mathcal{V}^r	70.63	70.52

Cost. At each state–action, the cost is the sum of the current distance from target position $s_{targ} = (8, 8)$, as well as the fuel expended by the given action.

$$C((i, j), a) = \sqrt{(i - s_{targ}[0])^2 + (j - s_{targ}[1])^2} + \frac{1}{2} \|a\|_2.$$

We take $\|a\|_2 = 1$ for all actions except for the staying still action, where $\|a\|_2 = 0$.

7.1. Bellman, optimistic policy, and robust policy

We compute the optimistic and robust bounds with parameter uncertainty in \mathcal{P} when $s \in [S_{unreliable}]$ by running Algorithm 1. The results are shown in Fig. 10.

We denote the optimistic policy as π^o and the robust policy as π^r , and derive the bounds of their respective value vector sets \mathcal{V}^o (38) and \mathcal{V}^r (39) using Algorithm 1. The output is compared against the bounds of the set-based Bellman operator's fixed point set \mathcal{V}^* in Table 1.

Non-stationary wind field Next, we consider a non-stationary wind field: at each time step k , the transition probability P^k is chosen at random from \mathcal{P} (61). We compare three different policy update schemes: (1) stationary optimistic policy π^o as policy operator g^o (38), (2) stationary robust policy π^r as policy operator g^r (39), and (3) Bellman policy that is one-step optimal for the MDP $([S], [A], P^k, C, \gamma)$ as f (9). These three different policy schemes are given by

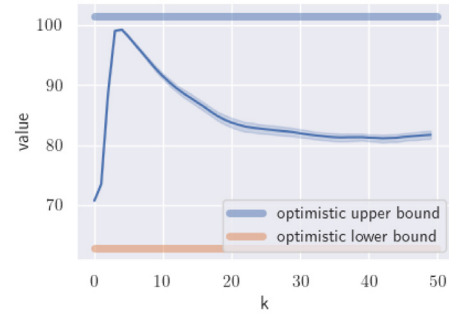
$$V^{k+1} = g^o(V^k, C, P^k), \quad (62)$$

$$V^{k+1} = g^r(V^k, C, P^k), \quad (63)$$

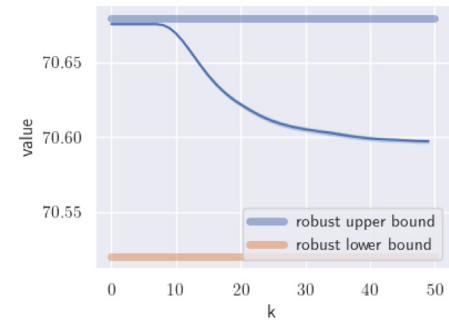
$$V^{k+1} = f(V^k, C, P^k). \quad (64)$$

The cost-to-go at state $s_{orig} = [0, 0]$ is plotted in Fig. 11. The optimistic policy (62) has the greatest variation in value over the course of 50 MDP time steps. Both the robust policy (63) and the Bellman policy (64) achieve better upper-bound at each MDP iteration. The Bellman policy (64) achieves less than 70 in cost-to-go on average, which is the best among all three policy schemes. As we discussed in Remark 9, the robust policy has the smallest variance in value in the presence of wind uncertainty, achieving a value difference of less than 0.1.

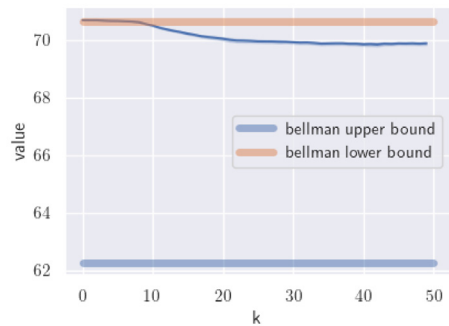
Sampled solutions. We can compute a sampled MDP model based on 50 samples of wind vectors for each state. Based on



(a) Optimistic Policy with \mathcal{V}^o 's bounding values.



(b) Robust Policy with \mathcal{V}^r 's bounding values.



(c) Bellman policy with \mathcal{V}^B 's bounding values.

Fig. 11. Comparison of robust policy, optimistic policy, and Bellman policy's value trajectories in non-stationary wind fields. Center blue line is the average over 50 trials. The shaded blue region denotes the standard variation. The top and bottom lines are the extrema values of the fixed points.

these samples, we add the action vector and compute the statistical distribution of state transitions. We then compute the value of these stationary sampled MDPs, and compare 9 randomly selected states' values. The resulting scatter plot is shown in Fig. 12.

8. Conclusion

In this paper, we lifted contraction operators that solve Markov decision processes to operate on compact sets of vectors. Using fixed point analysis, we showed that the set-based value operators have fixed point sets that are invariant to the parameter uncertainties. These sets were applied to non-stationary and parameter uncertain MDPs to derive novel results in these

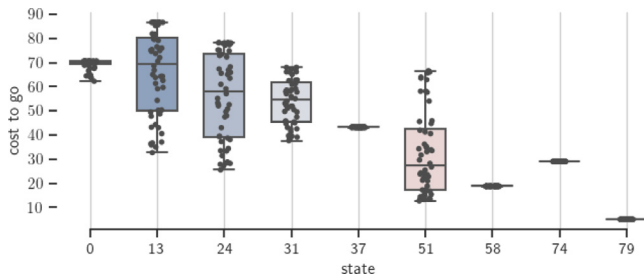


Fig. 12. Comparison of different optimal value vectors under the Bellman operator for 50 randomly sampled MDPs. On the x-axis, the state number is computed as $i \times 9 + j$.

settings. Finally, we demonstrated our results on path planning in a non-stationary wind field. In the future, we aim on apply set-based value operators to stochastic games in the presence of uncoordinated players.

Appendix A. Set sequence convergence

Lemma 5. Let $\{\mathcal{V}_n\} \subseteq \mathcal{K}(\mathbb{R}^S)$ be a converging sequence for $d_{\mathcal{K}}$ with $\mathcal{V}_n \rightarrow \mathcal{V}$ as $n \rightarrow \infty$. For all $V \in \mathcal{V}$, there exists a converging subsequence $\{V^{\varphi(n)}\}_{n \in \mathbb{N}}$ whose limit is V for $\|\cdot\|_{\infty}$.

Proof. Let $V \in \mathcal{V}$. We define the strictly increasing function φ on \mathbb{N} as follows: $\varphi(0) := 0$ and for all $n \in \mathbb{N}$, $\varphi(n+1) := \min\{j > \varphi(n) \mid \exists V^j \in \mathcal{V}^j, \|V - V^j\|_{\infty} = d(V, \mathcal{V}^j) \leq (n+1)^{-1}\}$. Finally, as for all $n \in \mathbb{N}^*$, $\|V - V^{\varphi(n)}\|_{\infty} \leq (\varphi(n)+1)^{-1}$, the result holds. \square

Appendix B. Proof of Lemma 1

Proof. Let $(V, m) \in \mathbb{R}^S \times \mathcal{M}$ and consider a sequence $\{(V_k, m_k)\}_{k \in \mathbb{N}} \subset \mathbb{R}^S \times \mathcal{M}$ that converges to (V, m) . It holds that $\|h(V_k, m_k) - h(V, m)\|_{\infty} \leq \|h(V_k, m_k) - h(V, m_k)\|_{\infty} + \|h(V, m_k) - h(V, m)\|_{\infty}$, where from the α -contractive property of $h(\cdot, m^k)$, $\|h(V_k, m_k) - h(V, m_k)\|_{\infty} \leq \alpha \|V_k - V\|_{\infty}$. From the $K(V)$ -Lipschitz property of $h(V, \cdot)$,

$$\|h(V, m_k) - h(V, m)\|_{\infty} \leq K(V) \|m_k - m\|_{\infty}.$$

As both $\lim_{k \rightarrow \infty} \|V_k - V\|_{\infty} \rightarrow 0$ and $\lim_{k \rightarrow \infty} \|m_k - m\|_{\infty} \rightarrow 0$, $\|h(V_k, m_k) - h(V, m)\|_{\infty} \rightarrow 0$ and h is continuous. \square

Appendix C. Proof of Lemma 2

Proof. We show that both the Bellman operator f and the policy evaluation operator g^{π} satisfy the contracting/order preserving/Lipschitz properties given in Definition 4. Contraction: given $(C, P) \in \mathcal{M}$, $g^{\pi}(\cdot, C, P)$ and $f(\cdot, C, P)$ are both γ -contractions (Puterman, 2014, Prop. 6.2.4) on the complete metric space $(\mathbb{R}^S, \|\cdot\|_{\infty})$, when $\gamma < 1$.

Order preservation: given $(C, P) \in \mathcal{M}$, the operator $g^{\pi}(\cdot, C, P)$ is order preserving (Puterman, 2014, Lem. 6.1.2). Consider $U, V \in \mathbb{R}^S$ where $U \leq V$. If $g^{\pi}(\cdot, C, P)$ is order-preserving, $g^{\pi}(U, C, P) \leq g^{\pi}(V, C, P)$ for all $\pi \in \Pi$. Taking the infimum over Π , we have $f(U, C, P) = \inf_{\pi \in \Pi} g^{\pi}(U, C, P) \leq \inf_{\pi \in \Pi} g^{\pi}(V, C, P) = f(V, C, P)$.

$K(V)$ -Lipschitz: given $(C, P), (C', P') \in \mathcal{M}$ and $V \in \mathbb{R}^S$, we prove the following for each $s \in [S]$,

$$\begin{aligned} & |f_s(V, C', P') - f_s(V, C, P)| \\ & \leq \|c'_s - c_s\|_{\infty} + \gamma \|P'_s - P_s\|_{\infty} \max\{\|\pi_s^*\|_{\infty}, \|\hat{\pi}_s\|_{\infty}\} \|V\|_{\infty}. \end{aligned} \quad (\text{C.1})$$

We prove (C.1) by case: (1) $f_s(V, C', P') \geq f_s(V, C, P)$, and (2) $f_s(V, C', P') \leq f_s(V, C, P)$. For case (1), let $\hat{\pi}$ (10) be the optimal policy for $f(V, C', P')$ and π^* be the optimal policy for $f(V, C, P)$. For $s \in [S]$, suppose $f_s(V, C', P') \geq f_s(V, C, P)$, then $0 \leq f_s(V, C', P') - f_s(V, C, P) \leq (c'_s)^{\top} \hat{\pi}_s - c_s^{\top} \pi_s^* + \gamma (P'_s \hat{\pi}_s)^{\top} V - \gamma (P_s \pi_s^*)^{\top} V$. Since π^* is sub-optimal for $f(V, C', P')$, we upper bound $|f_s(V, C', P') - f_s(V, C, P)| \leq (c'_s - c_s)^{\top} \pi_s^* + \gamma [(P'_s - P_s) \pi_s^*]^{\top} V$. Since $\pi_s^*, \hat{\pi}_s \in \Delta_A$, $\|\pi_s^*\|_{\infty} \leq 1$. We conclude that (C.1) holds when $f_s(V, C', P') \geq f_s(V, C, P)$. For case (2), $f_s(V, C', P') \leq f_s(V, C, P)$, (C.1) also holds by similar arguments.

Letting $m' = (C', P')$ and $m = (C, P)$, we can upper bound $f(V, m) - f(V, m') = f - f'$ as

$$\|f - f'\|_{\infty} \leq \max_{s \in [S]} \{\|c'_s - c_s\|_{\infty} + \gamma \|(P_s - P'_s)^{\top} V\|_{\infty}\} \quad (\text{C.2})$$

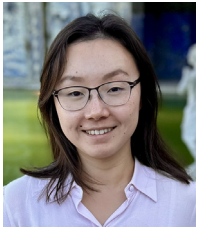
$$\leq \max(1, \gamma \|V\|_{\infty}) \|m - m'\|_{\infty}. \quad (\text{C.3})$$

The policy evaluation operator g^{π} satisfies (C.1) if $\max\{\|\pi_s^*\|_{\infty}, \|\hat{\pi}_s\|_{\infty}\}$ is replaced by $\|\pi_s\|_{\infty}$. Since $\|\pi_s\|_{\infty} \leq 1$, g^{π} is $K(V)$ -Lipschitz. \square

References

- Al-Sabban, Wesam H., Gonzalez, Luis F., & Smith, Ryan N. (2013). Wind-energy based path planning for unmanned aerial vehicles using markov decision processes. In *2013 IEEE international conference on robotics and automation* (pp. 784–789). IEEE.
- Bellemare, Marc G, Candido, Salvatore, Castro, Pablo Samuel, Gong, Jun, Machado, Marlos C, Moitra, Subhodeep, et al. (2020). Autonomous navigation of stratospheric balloons using reinforcement learning. *Nature*, 588, 77–82.
- Bellemare, Marc G, Ostrovski, Georg, Guez, Arthur, Thomas, Philip, & Munos, Rémi (2016). Increasing the action gap: New operators for reinforcement learning. Vol. 30, In *Proceedings of the AAAI conference on artificial intelligence*.
- Chen, Zaiwei, Maguluri, Siva Theja, Shakkottai, Sanjay, & Shanmugam, Karthikeyan (2021). Finite-sample analysis of off-policy TD-learning via generalized bellman operators. *Advances in Neural Information Processing Systems*, 34, 21440–21452.
- Doshi, Prashant, Goodwin, Richard, Akkiraju, Rama, & Verma, Kunal (2005). Dynamic workflow composition: Using markov decision processes. *International Journal of Web Services Research (IJWSR)*, 2(1), 1–17.
- Givan, Robert, Leach, Sonia, & Dean, Thomas (2000). Bounded-parameter Markov decision processes. *Artificial Intelligence*, 122(1–2), 71–109.
- Goyal, Vineet, & Grand-Clement, Julien (2022). Robust Markov decision processes: Beyond rectangularity. *Mathematics of Operations Research*.
- Henrikson, Jeff (1999). Completeness and total boundedness of the hausdorff metric. In *MIT undergraduate journal of mathematics*. Citeseer.
- Iyengar, Garud N. (2005). Robust dynamic programming. *Mathematics of Operations Research*, 30(2), 257–280.
- Kumar, Aviral, Zhou, Aurick, Tucker, George, & Levine, Sergey (2020). Conservative q-learning for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 33, 1179–1191.
- Lecarpentier, Erwan, & Rachelson, Emmanuel (2019). Non-stationary Markov decision processes, a worst-case approach using model-based reinforcement learning. *Advances in Neural Information Processing Systems*, 32.
- Li, Sarah H. Q., Adjé, Assalé, Garoche, Pierre-Loïc, & Açıkmeşe, Behçet (2021). Bounding fixed points of set-based bellman operator and Nash equilibria of stochastic games. *Automatica*, 130, Article 109685.
- Mannor, Shie, Mebel, Ofir, & Xu, Huan (2016). Robust MDPs with k-rectangular uncertainty. *Mathematics of Operations Research*, 41(4), 1484–1509.
- Melo, Francisco S. (2001). *Convergence of Q-learning: A simple proof*. Tech. Rep. (pp. 1–4). Institute of Systems and Robotics.
- Neumann, J. v (1928). Zur theorie der gesellschaftsspiele. *Mathematische Annalen*, 100(1), 295–320.
- Nilim, Arnab, & El Ghaoui, Laurent (2005). Robust control of Markov decision processes with uncertain transition matrices. *Operations Research*, 53(5), 780–798.
- Padakandla, Sindhu, KJ, Prabuchandran, & Bhatnagar, Shalabh (2020). Reinforcement learning algorithm for non-stationary environments. *Applied Intelligence: The International Journal of Artificial Intelligence, Neural Networks, and Complex Problem-Solving Technologies*, 50(11), 3590–3606.
- Puterman, Martin L. (2014). *Markov decision processes.: Discrete stochastic dynamic programming*. John Wiley & Sons.

- Rudin, Walter, et al. (1964). Vol. 3, *Principles of mathematical analysis*. McGraw-hill New York.
- Schröder, Bernd S. W. (2003). Ordered sets. *Springer*, 29, 30.
- Van Hoof, Herke, Hermans, Tucker, Neumann, Gerhard, & Peters, Jan (2015). Learning robot in-hand manipulation with tactile features. In *2015 IEEE-RAS 15th international conference on humanoid robots (humanoids)* (pp. 121–127). IEEE.
- Wiesemann, Wolfram, Kuhn, Daniel, & Rustem, Berç (2013). Robust Markov decision processes. *Mathematics of Operations Research*, 38(1), 153–183.
- Wolf, Michael T, Blackmore, Lars, Kuwata, Yoshiaki, Fathpour, Nanaz, Elfes, Alberto, & Newman, Claire (2010). Probabilistic motion planning of balloons in strong, uncertain wind fields. In *2010 IEEE international conference on robotics and automation* (pp. 1123–1129). IEEE.
- Xu, Huan, & Mannor, Shie (2010). Distributionally robust Markov decision processes. *Advances in Neural Information Processing Systems*, 23.
- Yang, Insoon (2017). A convex optimization approach to distributionally robust Markov decision processes with wasserstein distance. *IEEE Control Systems Letters*, 1(1), 164–169.
- Yu, Pengqian, & Xu, Huan (2015). Distributionally robust counterpart in Markov decision processes. *IEEE Transactions on Automatic Control*, 61(9), 2538–2543.



Sarah H.Q. Li is a postdoctoral scholar at ETH Zurich. She received her Ph.D. in Aeronautics and Astronautics Engineering at the University of Washington and her B.A.Sc degree in Engineering Physics at the University of British Columbia in 2017. She was a visiting researcher at Onera, the French Aerospace Lab in 2019 and has interned for Google Loon (2020), Deutsche Elektronen-Synchrotron (2015), and Macdonald Dettwiler and Associates (2014). Her research interests include game theory, Markov decision processes, and optimization.



Assalé Adjé obtained his PHD in applied mathematics/computer science from the École Polytechnique (Palaiseau, France) in 2011. Since 2016, he has been an associate professor in computer science at the University of Perpignan Via Domitia. His research interests include game theory and policy iteration algorithms. He is also interested in optimization theory and their applications in dynamical systems and programs verification.



Pierre-Loïc Garoche is a professor at ENAC, the French National School of Civil Aviation, and a contractor for NASA Ames in the Robust Software Engineering group. From 2008–2020, he worked at Onera, the French Aerospace Lab. His work is focused on the formal verification of control system software.



Behçet Açıkmeşe is a professor of Aerospace Optimization and Control at University of Washington, Seattle. He received his Ph.D. in Aerospace Engineering from Purdue University.

Previously, he was a senior technologist at JPL, and faculty member at the University of Texas at Austin. At JPL, he developed flyaway control algorithms that were successfully used in the landing of Curiosity and Perseverance rovers on Mars. His research interests include robust and nonlinear control, convex optimization and its applications control theory and its aerospace applications, and Markov decision processes. He is a recipient of many NASA and JPL achievement awards for his contributions to spacecraft control in planetary landing, formation flying, and asteroid and comet sample return missions. He is also a recipient of NSF CAREER Award and IEEE CSS Award for Technical Excellence in Aerospace Control. He is a fellow of IEEE and AIAA.