# Robust Average-Reward Markov Decision Processes

**Yue Wang,**[1] **Alvaro Velasquez,**[2] **George Atia,**[3] **Ashley Prater-Bennette,**[4] **Shaofeng Zou**[1]

[1] University at Buffalo, The State University of New York
[2] Information Innovation Office, Defense Advanced Research Projects Agency
[3] University of Central Florida
[4] Air Force Research Laboratory
ywang294@buffalo.com, alvaro.velasquez@darpa.mil, george.atia@ucf.edu, ashley.prater-bennette@us.af.mil,
szou3@buffalo.edu

## Abstract

In robust Markov decision processes (MDPs), the uncertainty in the transition kernel is addressed by finding a policy that optimizes the worst-case performance over an uncertainty set of MDPs. While much of the literature has focused on discounted MDPs, robust average-reward MDPs remain largely unexplored. In this paper, we focus on robust average-reward MDPs, where the goal is to find a policy that optimizes the worst-case average reward over an uncertainty set. We first take an approach that approximates average-reward MDPs using discounted MDPs. We prove that the robust discounted value function converges to the robust average-reward as the discount factor $\gamma$ goes to 1, and moreover, when $\gamma$ is large, any optimal policy of the robust discounted MDP is also an optimal policy of the robust average-reward. We further design a robust dynamic programming approach, and theoretically characterize its convergence to the optimum. Then, we investigate robust average-reward MDPs directly without using discounted MDPs as an intermediate step. We derive the robust Bellman equation for robust average-reward MDPs, prove that the optimal policy can be derived from its solution, and further design a robust relative value iteration algorithm that provably find its solution, or equivalently, the optimal robust policy.

## Introduction

A Markov decision process (MDP) is an effective mathematical tool for sequential decision-making in stochastic environments (Derman 1970; Puterman 1994). Solving an MDP problem entails finding an optimal policy that maximizes a cumulative reward according to a given criterion. However, in practice there could exist a mismatch between the assumed MDP model and the underlying environment due to various factors, such as non-stationarity of the environment, modeling error, exogenous perturbation, partial observability, and adversarial attacks. The ensuing model mismatch could result in solution policies with poor performance.

This challenge spurred noteworthy efforts on developing and analyzing a framework of robust MDPs e.g., (Bagnell, Ng, and Schneider 2001; Nilim and El Ghaoui 2004; Iyengar 2005). Rather than adopting a fixed MDP model, in the robust MDP setting, one seeks to optimize the worst-case performance over an uncertainty set of possible MDP models. The

solution to the robust MDP problem provides performance guarantee for all uncertain MDP models, and is thus robust to the model mismatch.

Robust MDP problems falling under different reward optimality criteria are fundamentally different. In robust discounted MDPs, the goal is to find a policy that maximizes the discounted cumulative reward in the worst case. In this setting, as the agent interacts with the environment, the reward received diminishes exponentially over time. Much of the prior work in the robust setting has focused on the discounted reward formulation. The model-based method, e.g., (Iyengar 2005; Nilim and El Ghaoui 2004; Bagnell, Ng, and Schneider 2001; Satia and Lave Jr 1973; Wiesemann, Kuhn, and Rustem 2013; Tamar, Mannor, and Xu 2014; Lim and Autef 2019; Xu and Mannor 2010; Yu and Xu 2015; Lim, Xu, and Mannor 2013), where information about the uncertainty set is assumed to be known to the learner, unveiled several fundamental characterizations of robust discounted MDPs. This was further extended to the more practical model-free setting in which only samples from a simulator (the centroid of the uncertainty set) are available to the learner. For example, the value-based method (Roy, Xu, and Pokutta 2017; Badrinath and Kalathil 2021; Wang and Zou 2021; Tessler, Efroni, and Mannor 2019; Zhou et al. 2021; Yang, Zhang, and Zhang 2021; Panaganti and Kalathil 2021; Goyal and Grand-Clement 2018; Kaufman and Schaefer 2013; Ho, Petrik, and Wiesemann 2018, 2021; Si et al. 2020) optimizes the worst-case performance using the robust value function as an intermediate step; on the other hand, the model-free policy-based method (Russel, Benosman, and Van Baar 2020; Derman, Geist, and Mannor 2021; Eysenbach and Levine 2021; Wang and Zou 2022) directly optimizes the policy and is thus scalable to large/continuous state and action spaces.

Although discounted MDPs induce an elegant Bellman operator that is a contraction, and have been studied extensively, the policy obtained usually has poor long-term performance when a system operates for an extended period of time. When the discount factor is very close to 1, the agent may prefer to compare policies on the basis of their average expected reward instead of their expected total discounted reward, e.g., queueing control, inventory management in supply chains, scheduling automatic guided vehicles and applications in communication networks (Kober, Bagnell, and Peters 2013). Therefore, it is also important to optimize the long-term aver-

age performance of a system.

However, robust MDPs under the average-reward criterion are largely understudied. Compared to the discounted setting, the average-reward setting depends on the limiting behavior of the underlying stochastic process, and hence is markedly more intricate. A recognized instance of such intricacy concerns the one-to-one correspondence between the stationary policies and the limit points of state-action frequencies, which while true for discounted MDPs, breaks down under the average-reward criterion even in the non-robust setting except in some very special cases (Puterman 1994; Atia et al. 2021). This is largely due to dependence of the necessary conditions for establishing a contraction in average-reward settings on the graph structure of the MDP, versus the discounted-reward setting where it simply suffices to have a discount factor that is strictly less than one. Heretofore, only a handful of studies have considered average-reward MDPs in the robust setting. The first work by (Tewari and Bartlett 2007) considers robust average-reward MDPs under a specific finite interval uncertainty set, but their method is not easily applicable to other uncertainty sets. More recently, (Lim, Xu, and Mannor 2013) proposed an algorithm for robust average-reward MDPs under the $\ell_1$ uncertainty set. However, obtaining fundamental characterizations of the problem and convergence guarantee remains elusive.

## Challenges and Contributions

In this paper, we derive characterizations of robust average-reward MDPs with general uncertainty sets, and develop model-based approaches with provable theoretical guarantee. Our approach is fundamentally different from previous work on robust discounted MDPs, robust and non-robust average-reward MDPs. In particular, the key challenges and the main contributions are summarized below.

- **We characterize the limiting behavior of robust discounted value function as the discount factor $\gamma \to 1$.** For the standard *non-robust* setting and for a specific transition kernel, the discounted non-robust value function converges to the average-reward non-robust value function as $\gamma \to 1$ (Puterman 1994). However, in the robust setting, we need to consider the worst-case limiting behavior under all possible transition kernels in the uncertainty set. Hence, the previous point-wise convergence result (Puterman 1994) cannot be directly applied. In (Tewari and Bartlett 2007), a finite interval uncertainty set is studied, where due to its special structure, the number of possible worst-case transition kernels of robust discounted MDPs is finite, and hence the order of $\min$ (over transition kernel) and $\lim_{\gamma \to 1}$ can be exchanged, and therefore, the robust discounted value function converges to the robust average-reward value function. This result, however, does not for general uncertainty sets investigated in this paper. We first prove the *uniform* convergence of discounted non-robust value function to average-reward w.r.t. the transition kernels and policies. Based on this uniform convergence, we show the convergence of the robust discounted value function to the robust average-reward. This uniform convergence result is the first in the literature and is of key importance to motivate

our algorithm design and to guarantee convergence to the optimal robust policy in the average-reward setting.

- **We design algorithms for robust policy evaluation and optimal control based on the limit method.** Based on the uniform convergence, we then use robust discounted MDPs to approximate robust average-reward MDPs. We show that when $\gamma$ is large, any optimal policy of the robust discounted MDP is also an optimal policy of the robust average-reward, and hence solves the robust optimal control problem in the average reward setting. This result is similar to the Blackwell optimality (Blackwell 1962; Hordijk and Yushkevich 2002) for the non-robust setting, however, our proof is fundamentally different. Technically, the proof in (Blackwell 1962; Hordijk and Yushkevich 2002) is based on the fact that the difference between the discounted value functions of two policies is a rational function of the discount factor, which has a finite number of zeros. However, in the robust setting with a general uncertainty set, the difference is no longer a rational function due to the min over the transition kernel. We construct a novel proof based on the limiting behavior of robust discounted MDPs, and show that the (optimal) robust discounted value function converges to the (optimal) robust average-reward as $\gamma \to 1$. Motivated by these insights, we then design our algorithms by applying a sequence of robust discounted Bellman operators while increasing the discount factor at a certain rate. We prove that our method can (i) evaluate the robust average-reward for a given policy and; (ii) find the optimal robust value function and, in turn, the optimal robust policy for general uncertainty sets.

- **We design a robust relative value iteration method without using the discounted MDPs as an intermediate step.** We further pursue a direct approach that solves the robust average-reward MDPs without using the limit method, i.e., without using discounted MDPs as an intermediate step. We derive a robust Bellman equation for robust average-reward MDPs, and show that the pair of robust relative value function and robust average-reward is a solution to the robust Bellman equation under the average-reward setting. We further prove that if we can find any solution to the robust Bellman equation, then the optimal policy can be derived by a greedy approach. The problem hence can be equivalently solved by solving the robust Bellman equation. We then design a robust value iteration method which provably converges to the solution of the robust Bellman equation, i.e., solve the optimal policy for the robust average-reward MDP problem.

## Related Work

**Robust discounted MDPs.** Model-based methods for robust discounted MDPs were studied in (Iyengar 2005; Nilim and El Ghaoui 2004; Bagnell, Ng, and Schneider 2001; Satia and Lave Jr 1973; Wiesemann, Kuhn, and Rustem 2013; Lim and Autef 2019; Xu and Mannor 2010; Yu and Xu 2015; Lim, Xu, and Mannor 2013; Tamar, Mannor, and Xu 2014), where the uncertainty set is assumed to be known, and the problem can be solved using robust dynamic programming. Later, the studies were generalized to the model-free setting where stochas-

tic samples from the centroid MDP of the uncertainty set are available in an online fashion (Roy, Xu, and Pokutta 2017; Badrinath and Kalathil 2021; Wang and Zou 2021, 2022; Tessler, Efroni, and Mannor 2019) and an offline fashion (Zhou et al. 2021; Yang, Zhang, and Zhang 2021; Panaganti and Kalathil 2021; Goyal and Grand-Clement 2018; Kaufman and Schaefer 2013; Ho, Petrik, and Wiesemann 2018, 2021; Si et al. 2020). There are also empirical studies on robust RL, e.g., (Vinitsky et al. 2020; Pinto et al. 2017; Abdullah et al. 2019; Hou et al. 2020; Rajeswaran et al. 2017; Huang et al. 2017; Kos and Song 2017; Lin et al. 2017; Pattanaik et al. 2018; Mandlekar et al. 2017). For discounted MDPs, the robust Bellman operator is a contraction, based on which robust dynamic programming and value-based methods can be designed. In this paper, we focus on robust average-reward MDPs. However, the robust Bellman operator for average-reward MDPs is not a contraction, and its fixed point may not be unique. Moreover, the average-reward setting depends on the limiting behavior of the underlying stochastic process, which is thus more intricate.

**Robust average-reward MDPs.** Studies on robust average-reward MDPs are quite limited in the literature. Robust average-reward MDPs under a specific finite interval uncertainty set was studied in (Tewari and Bartlett 2007), where the authors showed the existence of a Blackwell optimal policy, i.e., there exists some $\delta \in [0, 1)$, such that the optimal robust policy exists and remains unchanged for any discount factor $\gamma \in [\delta, 1)$. However, this result depends on the structure of the uncertainty set. For general uncertainty sets, the existence of a Blackwell optimal policy may not be guaranteed. More recently, (Lim, Xu, and Mannor 2013) designed a model-free algorithm for a specific $\ell_1$-norm uncertainty set and characterized its regret bound. However, their method also relies on the structure of the $\ell_1$-norm uncertainty set, and may not be generalizable to other types of uncertainty sets. In this paper, our results can be applied to various types of uncertainty sets, and thus is more general.

## Preliminaries and Problem Model

In this section, we introduce some preliminaries on discounted MDPs, average-reward MDPs, and robust MDPs.

**Discounted MDPs.** A discounted MDP $(\mathcal{S}, \mathcal{A}, \mathsf{P}, r, \gamma)$ is specified by: a state space $\mathcal{S}$, an action space $\mathcal{A}$, a transition kernel $\mathsf{P} = \{p_s^a \in \Delta(\mathcal{S}), a \in \mathcal{A}, s \in \mathcal{S}\}$[1], where $p_s^a$ is the distribution of the next state over $\mathcal{S}$ upon taking action $a$ in state $s$ (with $p_{s,s'}^a$ denoting the probability of transitioning to $s'$), a reward function $r : \mathcal{S} \times \mathcal{A} \to [0, 1]$, and a discount factor $\gamma \in [0, 1)$. At each time step $t$, the agent at state $s_t$ takes an action $a_t$, the environment then transitions to the next state $s_{t+1}$ according to $p_{s_t}^{a_t}$, and produces a reward signal $r(s_t, a_t) \in [0, 1]$ to the agent. In this paper, we also write $r_t = r(s_t, a_t)$ for convenience.

A stationary policy $\pi : \mathcal{S} \to \Delta(\mathcal{A})$ is a distribution over $\mathcal{A}$ for any given state $s$, and the agent takes action $a$ at state $s$ with probability $\pi(a|s)$. The discounted value function of a stationary policy $\pi$ starting from $s \in \mathcal{S}$ is defined as the

expected discounted cumulative reward by following policy $\pi$: $V_{\mathsf{P},\gamma}^\pi(s) \triangleq \mathbb{E}_{\pi,\mathsf{P}}\left[\sum_{t=0}^\infty \gamma^t r_t | S_0 = s\right]$.

**Average-Reward MDPs.** Different from discounted MDPs, average-reward MDPs do not discount the reward over time, and consider the behavior of the underlying Markov process under the steady-state distribution. More specifically, under a specific transition kernel $\mathsf{P}$, the average-reward of a policy $\pi$ starting from $s \in \mathcal{S}$ is defined as

$$g_{\mathsf{P}}^\pi(s) \triangleq \lim_{n \to \infty} \mathbb{E}_{\pi,\mathsf{P}}\left[\frac{1}{n}\sum_{t=0}^{n-1} r_t | S_0 = s\right], \qquad (1)$$

which we also refer to in this paper as the average-reward value function for convenience.

The average-reward value function can also be equivalently written as follows: $g_{\mathsf{P}}^\pi = \lim_{n \to \infty} \frac{1}{n}\sum_{t=0}^{n-1}(\mathsf{P}^\pi)^t r_\pi \triangleq \mathsf{P}_*^\pi r_\pi$, where $(\mathsf{P}^\pi)_{s,s'} \triangleq \sum_a \pi(a|s)p_{s,s'}^a$ and $r_\pi(s) \triangleq \sum_a \pi(a|s)r(s,a)$ are the transition matrix and reward function induced by $\pi$, and $\mathsf{P}_*^\pi \triangleq \lim_{n \to \infty} \frac{1}{n}\sum_{t=0}^{n-1}(\mathsf{P}^\pi)^t$ is the limit matrix of $\mathsf{P}^\pi$.

In the average-reward setting, we also define the following relative value function

$$V_{\mathsf{P}}^\pi(s) \triangleq \mathbb{E}_{\pi,\mathsf{P}}\left[\sum_{t=0}^\infty (r_t - g_{\mathsf{P}}^\pi)|S_0 = s\right], \qquad (2)$$

which is the cumulative difference over time between the reward and the average value $g_{\mathsf{P}}^\pi$. It has been shown that (Puterman 1994): $V_{\mathsf{P}}^\pi = H_{\mathsf{P}}^\pi r_\pi$, where $H_{\mathsf{P}}^\pi \triangleq (I - \mathsf{P}^\pi + \mathsf{P}_*^\pi)^{-1}(I - \mathsf{P}_*^\pi)$ is defined as the deviation matrix of $\mathsf{P}^\pi$.

The relationship between the average-reward and the relative value functions can be characterized by the following Bellman equation (Puterman 1994):

$$V_{\mathsf{P}}^\pi(s) = \mathbb{E}_\pi\left[r(s, A) - g_{\mathsf{P}}^\pi(s) + \sum_{s' \in \mathcal{S}} p_{s,s'}^A V_{\mathsf{P}}^\pi(s')\right]. \quad (3)$$

**Robust discounted and average-reward MDPs.** For robust MDPs, the transition kernel is not fixed but belongs to some uncertainty set $\mathcal{P}$. After the agent takes an action, the environment transits to the next state according to an arbitrary transition kernel $\mathsf{P} \in \mathcal{P}$. In this paper, we focus on the $(s, a)$-rectangular uncertainty set (Nilim and El Ghaoui 2004; Iyengar 2005), i.e., $\mathcal{P} = \bigotimes_{s,a} \mathcal{P}_s^a$, where $\mathcal{P}_s^a \subseteq \Delta(\mathcal{S})$. We note that there are also studies on relaxing the $(s, a)$-rectangular uncertainty set to $s$-rectangular uncertainty set, which is not the focus of this paper.

Under the robust setting, we consider the worst-case performance over the uncertainty set of MDPs. More specifically, the robust discounted value function of a policy $\pi$ for a discounted MDP is defined as

$$V_{\mathcal{P},\gamma}^\pi(s) \triangleq \min_{\kappa \in \bigotimes_{t \geq 0} \mathcal{P}} \mathbb{E}_{\pi,\kappa}\left[\sum_{t=0}^\infty \gamma^t r_t | S_0 = s\right], \quad (4)$$

where $\kappa = (\mathsf{P}_0, \mathsf{P}_1...) \in \bigotimes_{t \geq 0} \mathcal{P}$.

---

[1] $\Delta(\mathcal{S})$: the $(|\mathcal{S}| - 1)$-dimensional probability simplex on $\mathcal{S}$.

In this paper, we focus on the following worst-case average-reward for a policy $\pi$:

$$g_{\mathcal{P}}^{\pi}(s) \triangleq \min_{\kappa \in \bigotimes_{t \geq 0} \mathcal{P}} \lim_{n \to \infty} \mathbb{E}_{\pi,\kappa} \left[ \frac{1}{n} \sum_{t=0}^{n-1} r_t | S_0 = s \right], \quad (5)$$

to which, for convenience, we refer as the robust average-reward value function.

For robust discounted MDPs, it has been shown that the robust discounted value function is the unique fixed-point of the robust discounted Bellman operator (Nilim and El Ghaoui 2004; Iyengar 2005; Puterman 1994):

$$\mathbf{T}_{\pi} V(s) \triangleq \sum_{a \in \mathcal{A}} \pi(a|s) \left( r(s,a) + \gamma \sigma_{\mathcal{P}_s^a}(V) \right), \quad (6)$$

where $\sigma_{\mathcal{P}_s^a}(V) \triangleq \min_{p \in \mathcal{P}_s^a} p^\top V$ is the support function of $V$ on $\mathcal{P}_s^a$. Based on the contraction of $\mathbf{T}_{\pi}$, robust dynamic programming approaches, e.g., robust value iteration, can be designed (Nilim and El Ghaoui 2004; Iyengar 2005) (see Appendix for a review of these methods). However, there is no such contraction result for robust average-reward MDPs. In this paper, our goal is to find a policy that optimizes the robust average-reward value function:

$$\max_{\pi \in \Pi} g_{\mathcal{P}}^{\pi}(s), \text{ for any } s \in \mathcal{S}, \quad (7)$$

where $\Pi$ is the set of all stationary policies, and we denote by $g_{\mathcal{P}}^{*}(s) \triangleq \max_{\pi} g_{\mathcal{P}}^{\pi}(s)$ the optimal robust average-reward.

## Limit Approach for Robust Average-Reward MDPs

We first take a limit approach to solve the problem of robust average-reward MDPs in eq. (7). It is known that under the non-robust setting, for any fixed $\pi$ and P, the discounted value function converges to the average-reward value function as the discount factor $\gamma$ approaches 1 (Puterman 1994), i.e.,

$$\lim_{\gamma \to 1} (1 - \gamma) V_{\mathsf{P},\gamma}^{\pi} = g_{\mathsf{P}}^{\pi}. \quad (8)$$

We take a similar idea, and show that the same result holds in the robust case: $\lim_{\gamma \to 1} (1 - \gamma) V_{\mathcal{P},\gamma}^{\pi} = g_{\mathcal{P}}^{\pi}$. Based on this result, we further design algorithms (Algorithms 1 and 2) that apply a sequence of robust discounted Bellman operators while increasing the discount factor at a certain rate. We then theoretically prove that our algorithms converge to the optimal solutions.

In the following, we first show that the convergence $\lim_{\gamma \to 1} (1 - \gamma) V_{\mathsf{P},\gamma}^{\pi} = g_{\mathsf{P}}^{\pi}$ is uniform on the set $\Pi \times \mathcal{P}$. We make a mild assumption as follows.

**Assumption 1.** *For any $s \in \mathcal{S}, a \in \mathcal{A}$, the uncertainty set $\mathcal{P}_s^a$ is a compact subset of $\Delta(\mathcal{S})$.*

The set $\mathcal{P}_s^a$ is compact if and only if it is bounded and closed. Since $\mathcal{P}_s^a \subseteq \Delta(\mathcal{S})$, it is clearly bounded. Hence, Assumption 1 amounts to assuming that the uncertainty set is closed. We remark that many standard uncertainty sets satisfy this assumption, e.g., those defined by $\epsilon$-contamination (Huber 1965), finite interval (Tewari and Bartlett 2007), total-variation (Rahimian, Bayraksan, and De-Mello 2022) and KL-divergence (Hu and Hong 2013).

In (Puterman 1994), the convergence $\lim_{\gamma \to 1} (1 - \gamma) V_{\mathsf{P},\gamma}^{\pi} = g_{\mathsf{P}}^{\pi}$ for a fixed policy $\pi$ and a fixed transition kernel P (non-robust setting) is point-wise. However, such point-wise convergence does not provide any convergence guarantee on the robust discounted value function, as the robust value function measures the worst-case performance over the uncertainty set and the order of $\lim$ and $\min$ may not be exchanged in general. In the following theorem, we prove the uniform convergence of the discounted value function under the foregoing assumption.

**Theorem 1** (Uniform convergence). *Under Assumption 1, the discounted value function converges uniformly to the average-reward value function on $\Pi \times \mathcal{P}$ as $\gamma \to 1$, i.e.,*

$$\lim_{\gamma \to 1} (1 - \gamma) V_{\mathsf{P},\gamma}^{\pi} = g_{\mathsf{P}}^{\pi}, \text{ uniformly.} \quad (9)$$

With uniform convergence in Theorem 1, the order of the limit $\gamma \to 1$ and min over P can be interchanged, then the following convergence of the robust discounted value function can be established.

**Theorem 2.** *The robust discounted value function in eq.* (4) *converges to the robust average-reward uniformly on $\Pi$:*

$$\lim_{\gamma \to 1} (1 - \gamma) V_{\mathcal{P},\gamma}^{\pi} = g_{\mathcal{P}}^{\pi}, \text{ uniformly.} \quad (10)$$

We note that a similar convergence result is shown in (Tewari and Bartlett 2007), but only for a special uncertainty set of finite interval. Our Theorem 2 holds for general compact uncertainty sets. Moreover, it is worth highlighting that our proof technique is fundamentally different from the one in (Tewari and Bartlett 2007). Specifically, under the finite interval uncertainty set, the worst-case transition kernels are from a finite set, i.e., $V_{\mathcal{P},\gamma}^{\pi} = \min_{\mathsf{P} \in \mathcal{M}} V_{\mathsf{P},\gamma}^{\pi}$ for a finite set $\mathcal{M} \subseteq \mathcal{P}$. This hence implies the interchangeability of $\lim$ and $\min$. However, for general uncertainty sets, the number of worst-case transition kernels may not be finite. We demonstrate the interchangeability via our uniform convergence result in Theorem 1.

The convergence result in Theorem 2 is of key importance to motivate the design of the following two algorithms, the basic idea of which is to apply a sequence of robust discounted Bellman operators on an arbitrary initialization while increasing the discount factor at a certain rate.

We first consider the robust policy evaluation problem, which aims to estimate the robust average-reward $g_{\mathcal{P}}^{\pi}$ for a fxied policy $\pi$. This problem for robust discounted MDPs is well studied in the literature, however, results for robust average-reward MDPs are quite limited except for the one in (Tewari and Bartlett 2007) for a specific finite interval uncertainty set. We present the a robust value iteration (robust VI) algorithm for evaluating the robust average-reward with general compact uncertainty sets in Algorithm 1.

At each time step $t$, the discount factor $\gamma_t$ is set to $\frac{t+1}{t+2}$, which converges to 1 as $t \to \infty$. Subsequently, a robust Bellman operator w.r.t discount factor $\gamma_t$ is applied on the current estimate $V_t$ of the robust discounted value function $(1 - \gamma_t) V_{\mathcal{P},\gamma_t}^{\pi}$. As the discount factor approaches 1, the estimated robust discounted value function converges to the robust average-reward $g_{\mathcal{P}}^{\pi}$ by Theorem 2.

---

**Algorithm 1: Robust VI: Policy Evaluation**

**Input**: $\pi, V_0(s) = 0, \forall s, T$
1: **for** $t = 0, 1, ..., T - 1$ **do**
2:     $\gamma_t \leftarrow \frac{t+1}{t+2}$
3:     **for all** $s \in \mathcal{S}$ **do**
4:         $V_{t+1}(s) \leftarrow \mathbb{E}_\pi[(1 - \gamma_t)r(s, A) + \gamma_t \sigma_{\mathcal{P}_s^A}(V_t)]$
5:     **end for**
6: **end for**
7: **return** $V_T$

---

**Theorem 3.** *Algorithm 1 converges to robust average reward, i.e.,* $\lim_{T\to\infty} V_T \to g_{\mathcal{P}}^\pi$.

Theorem 3 shows that the output of Algorithm 1 converges to the robust average-reward.

Besides the robust policy evaluation problem, it is also of great practical importance to find an optimal policy that maximizes the worst-case average-reward, i.e., to solve eq. (7). Based on a similar idea as the one of Algorithm 1, we extend our limit approach to solve the robust optimal control problem in Algorithm 2.

---

**Algorithm 2: Robust VI: Optimal Control**

**Input**: $V_0(s) = 0, \forall s, T$
1: **for** $t = 0, 1, ..., T - 1$ **do**
2:     $\gamma_t \leftarrow \frac{t+1}{t+2}$
3:     **for all** $s \in \mathcal{S}$ **do**
4:         $V_{t+1}(s) \leftarrow \max_{a \in \mathcal{A}} \left\{ (1 - \gamma_t)r(s, a) + \gamma_t \sigma_{\mathcal{P}_s^a}(V_t) \right\}$
5:     **end for**
6: **end for**
7: **for** $s \in \mathcal{S}$ **do**
8:     $\pi_T(s) \leftarrow \arg\max_{a \in \mathcal{A}} \left\{ (1 - \gamma_t)r(s, a) + \gamma_t \sigma_{\mathcal{P}_s^a}(V_t) \right\}$
9: **end for**
10: **return** $V_T, \pi_T$

---

Similar to Algorithm 1, at each time step, the discount factor $\gamma_t$ is set to be closer to 1, and a one-step robust discounted Bellman operator (for optimal control) w.r.t. $\gamma_t$ is applied to the current estimate $V_t$. The following theorem establishes that $V_T$ in Algorithm 2 converges to the optimal robust value function, hence can find the optimal robust policy.

**Theorem 4.** *The output $V_T$ in Algorithm 2 converges to the optimal robust average-reward $g_{\mathcal{P}}^*$: $V_T \to g_{\mathcal{P}}^*$ as $T \to \infty$.*

As discussed in (Blackwell 1962; Hordijk and Yushkevich 2002), the average-reward criterion is insensitive and under selective since it is only interested in the performance under the steady-state distribution. For example, two policies providing rewards: $100 + 0 + 0 + \cdots$ and $0 + 0 + 0 + \cdots$ are equally good/bad. Towards this issue, for the non-robust setting, a more sensitive term of optimality was introduced by Blackwell (Blackwell 1962). More specifically, a policy is said to be Blackwell optimal if it optimizes the discounted value function for all discount factor $\gamma \in (\delta, 1)$ for some $\delta \in (0, 1)$. Together with eq. (8), the optimal policy obtained by taking $\gamma \to 1$ is optimal not only for the average-reward

criterion, but also for the discounted criterion with large $\gamma$. Intuitively, it is optimal under the average-reward setting, and is sensitive to early rewards.

Following a similar idea, we justify that the obtained policy from Algorithm 2 is not only optimal in the robust average-reward setting, but also sensitive to early rewards.

Denote by $\Pi^*$ the set of all the optimal policies for robust average-reward, i.e. $\Pi^* = \{\pi : g_{\mathcal{P}}^\pi = g_{\mathcal{P}}^*\}$.

**Theorem 5** (Blackwell optimality). *There exists $0 < \delta < 1$ and a finite policy set $\Pi^*$, such that for any $\gamma > \delta$, the optimal robust policy for robust discounted value function $V_{\mathcal{P},\gamma}^*$ belongs to $\Pi^*$, i.e., for any $\delta < \gamma < 1$, $\exists \pi^* \in \Pi^*$, s.t. $V_{\mathcal{P},\gamma}^* = V_{\mathcal{P},\gamma}^{\pi^*}$.*

This result implies that using the limit method in this section to find the optimal robust policy for average-reward MDPs has an additional advantage that the policy it finds not only optimizes the average reward in steady state, but also is sensitive to early rewards.

It is worth highlighting the distinction of our results from the technique used in the proof of Blackwell optimality (Blackwell 1962). In the non-robust setting, the existence of a stationary Blackwell optimal policy is proved via contradiction, where a difference function of two policies $\pi$ and $\nu$: $f_{\pi,\nu}(\gamma) \triangleq V_{\mathsf{P},\gamma}^\pi - V_{\mathsf{P},\gamma}^\mu$ is used in the proof. It was shown by contradiction that $f$ has infinitely many zeros, which however contradicts with the fact that $f$ is a rational function of $\gamma$ with a finite number of zeros. A similar technique was also used in (Tewari and Bartlett 2007) for the finite interval uncertainty set. Specifically, in (Tewari and Bartlett 2007), it was shown that the worst-case transition kernels for any $\pi, \gamma$ are from a finite set $\mathcal{M}$, hence $f_{\pi,\nu}(\gamma) \triangleq \min_{\mathsf{P} \in \mathcal{M}} V_{\mathsf{P},\gamma}^\pi - \min_{\mathsf{P} \in \mathcal{M}} V_{\mathsf{P},\gamma}^\mu$ can also be shown to be a rational function with a finite number of zeroes. For a general uncertainty set $\mathcal{P}$, the difference function $f_{\pi,\nu}(\gamma)$, however, may not be rational. This makes the method in (Blackwell 1962; Tewari and Bartlett 2007) inapplicable to our problem.

## Direct Approach for Robust Average-Reward MDPs

The limit approach in Section is based on the uniform convergence of the discounted value function, and uses discounted MDPs to approximate average-reward MDPs. In this section, we develop a direct approach to solving the robust average-reward MDPs that does not adopt discounted MDPs as intermediate steps.

For average-reward MDPs, the relative value iteration (RVI) approach (Puterman 1994) is commonly used since it is numerically stable and has convergence guarantee. In the following, we generalize the RVI algorithm to the robust setting, and design the robust RVI algorithm in Algorithm 3.

We first generalize the relative value function in eq. (2) to the robust relative value function. The robust relative value function measures the difference between the worst-case cumulative reward and the worst-case average-reward for a policy $\pi$.

**Definition 1.** *The robust relative value function is defined as*

$$V_{\mathcal{P}}^{\pi}(s) \triangleq \min_{\kappa \in \bigotimes_{t \geq 0} \mathcal{P}} \mathbb{E}_{\kappa,\pi} \left[ \sum_{t=0}^{\infty} (r_t - g_{\mathcal{P}}^{\pi}) | S_0 = s \right], \quad (11)$$

*where $g_{\mathcal{P}}^{\pi}$ is the worst-case average-reward defined in eq. (5).*

The following theorem presents a robust Bellman equation for robust average-reward MDPs.

**Theorem 6.** *For any $s$ and $\pi$, $(V_{\mathcal{P}}^{\pi}, g_{\mathcal{P}}^{\pi})$ is a solution to the following robust Bellman equation:*

$$V(s) + g = \sum_{a} \pi(a|s) \left( r(s,a) + \sigma_{\mathcal{P}_s^a}(V) \right). \quad (12)$$

It can be seen that the robust Bellman equation for average-reward MDPs has a similar structure to the one for discounted MDPs in eq. (6) except for a discount factor. This actually reveals a fundamental difference between the robust Bellman operator of the discounted MDPs and the average-reward ones. For a discounted MDP, its robust Bellman operator is a contraction with constant $\gamma$ (Nilim and El Ghaoui 2004; Iyengar 2005), and hence the fixed point is unique. Based on this, the robust value function can be found by recursively applying the robust Bellman operator (see Appendix ). In sharp contrast, in the average-reward setting, the robust Bellman is not necessarily a contraction, and the fixed point may not be unique. Therefore, repeatedly applying the robust Bellman operator in the average-reward setting may not even converge, which underscores that the two problem settings are fundamentally different.

Using the robust Bellman equation in Theorem 6, we derive the following equivalent optimality condition for robust average-reward MDPs.

**Theorem 7.** *For any $(g, V)$ that is a solution to*

$$\max_{a} \left\{ r(s,a) - g + \sigma_{\mathcal{P}_s^a}(V) - V(s) \right\} = 0, \forall s, \quad (13)$$

*$g = g_{\mathcal{P}}^*$. If we further set*

$$\pi^*(s) = \arg\max_{a} \left\{ r(s,a) + \sigma_{\mathcal{P}_s^a}(V) \right\} \quad (14)$$

*for any $s \in \mathcal{S}$, then $\pi^*$ is an optimal robust policy.*

Theorem 7 suggests that as long as we find a solution $(g, V)$ to eq. (13), which though may not be unique, then $g$ is the optimal robust average-reward $g_{\mathcal{P}}^*$, and the greedy policy $\pi^*$ is the optimal policy to our robust average-reward MDP problem in eq. (7). Based on Theorem 7, our problem in eq. (7) can be equivalently solved by finding a solution to eq. (13). We note that eq. (12) holds for any $\pi$ and if we let the $\pi$ in eq. (12) be the greedy policy, then eq. (12) and eq. (13) are equivalent.

In the following, we generalize the RVI approach to the robust setting, and design a robust RVI algorithm in Algorithm 3. We will further show that the output of this algorithm converges to a solution to eq. (13), and further the optimal policy could be obtained by eq. (14). Here $\mathbb{1}$ denotes the all-ones vector, and $sp$ denotes the span semi-norm: $sp(w) = \max_s w(s) - \min_s w(s)$. Different from Algorithm 2, in Algorithm 3, we do not need to apply the robust discounted Bellman operator. The method directly solves the

---

**Algorithm 3: Robust RVI**

**Input**: $V_0$, $\epsilon$ and arbitrary $s^* \in \mathcal{S}$
1: $w_0 \leftarrow V_0 - V_0(s^*)\mathbb{1}$
2: **while** $sp(w_t - w_{t+1}) \geq \epsilon$ **do**
3:    **for** all $s \in \mathcal{S}$ **do**
4:       $V_{t+1}(s) \leftarrow \max_a(r(s,a) + \sigma_{\mathcal{P}_s^a}(w_t))$
5:       $w_{t+1}(s) \leftarrow V_{t+1} - V_{t+1}(s^*)\mathbb{1}$
6:    **end for**
7: **end while**
8: **return** $w_t, V_t$

---

robust optimal control problem for average-reward robust MDPs.

In studies of average-reward MDPs, it is usually the case that a certain class of MDPs are considered, e.g., unichain and communicating (Wei et al. 2020; Zhang and Ross 2021; Chen, Jain, and Luo 2022; Wan, Naik, and Sutton 2021). In this paper, we focus on the unichain setting to highlight the major technical novelty to achieve robustness.

**Assumption 2.** *For any $\mathsf{P} = \{p_s^a \in \Delta(\mathcal{S})\} \in \mathcal{P}$ and any $a \in \mathcal{A}, s, s' \in \mathcal{S}$, $p_{s,s'}^a > 0$, and the induced Markov process is a unichain.*

In the following theorem, we show that our Algorithm 3 converges to a solution of eq. (13), hence according to Theorem 7 if we set $\pi$ according to (14), then $\pi$ is the optimal robust policy.

**Theorem 8.** *$(w_t, V_t)$ converges to a solution $(w, V)$ to eq. (13) as $\epsilon \to 0$, which satisfies*

$$w(s) + \max_{a}\{r(s^*, a) + \sigma_{\mathcal{P}_{s^*}^a}(w)\}$$
$$= \max_{a}\{r(s,a) + \sigma_{\mathcal{P}_s^a}(w)\}. \quad (15)$$

**Remark 1.** *In this section, we mainly present the robust RVI algorithm for the robust optimal control problem, and its convergence and optimality guarantee. A robust RVI algorithm for robust policy evaluation can be similarly designed by replacing the $\max$ in line 4, Algorithm 3 with an expectation w.r.t. $\pi$. The convergence results in Theorem 8 can also be similarly derived.*

*Assumption 2 can be also replaced using some weaker ones, e.g., Proposition 4.3.2 of (Bertsekas 2011), or be removed by designing a variant of RVI, e.g., Proposition 4.3.4 of (Bertsekas 2011).*

## Examples and Numerical Results

In this section, we study several commonly used uncertainty set models, including contamination model, Kullback-Lerbler (KL) divergence defined model and total-variation defined model.

As can be observed from Algorithms 1 to 3, for different uncertainty sets, the only difference lies in how the support function $\sigma_{\mathcal{P}_s^a}(V)$ is calculated. In the sequel, we discuss how to efficiently calculate the support function for various uncertainty sets.

We numerically compare our robust (relative) value iteration methods v.s. non-robust (relative) value iteration method

on different uncertainty sets. Our experiments are based on the Garnet problem $\mathcal{G}(20, 40)$ (Archibald, McKinnon, and Thomas 1995). More specifically, there are 20 states and 30 actions; the nominal transition kernel $\mathsf{P} = \{p_s^a \in \Delta(\mathcal{S})\}$ is randomly generated according to the uniform distribution, and the reward functions $r(s, a) \sim \mathcal{N}(0, \sigma_{s,a})$, where $\sigma_{s,a} \sim \text{Uniform}[0, 1]$. In our experiments, the uncertainty sets are designed to be centered at the nominal transition kernel. We run different algorithms, i.e., (robust) value iteration and (robust) relative value iteration, and obtain the greedy policies at each time step. Then, we use robust average-reward policy evaluation (Algorithm 1) to evaluate the robust average-reward of these policies. We plot the robust average-reward against the number of iterations.

**Contamination model.** For any $(s, a)$ the uncertainty set $\mathcal{P}_s^a$ is defined as $\mathcal{P}_s^a = \{q : q = (1 - R)p_s^a + Rp', p' \in \Delta(\mathcal{S})\}$, where $p_s^a$ is the nominal transition kernel. It can be viewed as an adversarial model, where at each time-step, the environment transits according to the nominal transition kernel $p$ with probability $1 - R$, and according to an arbitrary kernel $p'$ with probability $R$.

It can be easily shown that the value of the problem $\sigma_{\mathcal{P}_s^a}(V) = (1-R)(p_s^a)^\top V + R \min_s V(s)$. Our experimental results under the contamination model are shown in Figure 1.
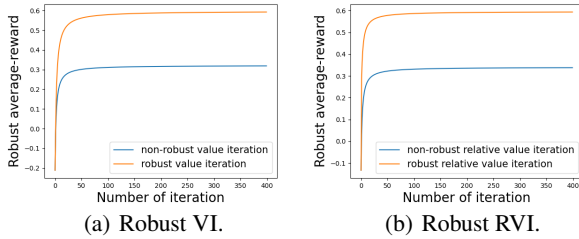


Figure 1: Comparison on contamination model with $R = 0.4$.

**Total variation.** The total variation distance is another commonly used distance metric to measure the difference between two distributions. Specifically, the total variation between two distributions $p$ and $q$ is defined as $D_{TV}(p, q) = \frac{1}{2}\|p - q\|_1$. Consider an uncertainty set defined via total variation: $\mathcal{P}_s^a = \{q : D_{TV}(q\|p_s^a) \le R\}$. Then, its support function can be efficiently solved as follows (Iyengar 2005): $\sigma_{\mathcal{P}_s^a}(V) = p^\top V - R \min_{\mu \ge 0} \{\max_s(V(s) - \mu(s)) - \min_s(V(s) - \mu(s))\}$.

Our experimental results under the total variation model are shown in Figure 2.

**Kullback-Lerbler (KL) divergence.** The Kullback–Leibler divergence is widely used to measure the distance between two probability distributions. The KL-divergence of two distributions $p, q$ is defined as $D_{KL}(q\|p) = \sum_s q(s) \log \frac{q(s)}{p(s)}$. Consider an uncertainty set defined via KL divergence: $\mathcal{P}_s^a = \{q : D_{KL}(q\|p_s^a) \le R\}$. Then, its support function can be efficiently solved using the duality result in (Hu and Hong 2013): $\sigma_{\mathcal{P}_s^a}(V) = -\min_{\alpha \ge 0}\left\{R\alpha + \alpha \log\left(p^\top e^{\frac{-V}{\alpha}}\right)\right\}$.

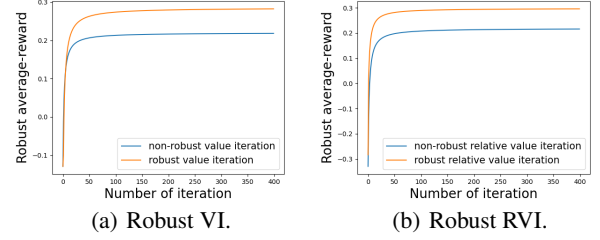Our experimental results under the KL-divergence model are shown in Figure 3.



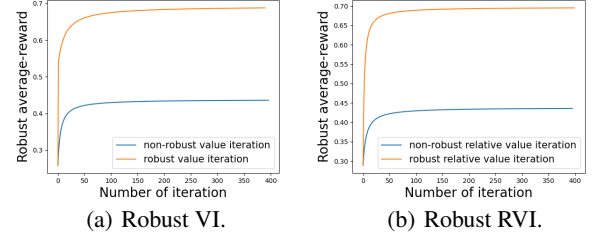Figure 2: Comparison on total variation model with $R = 0.6$.



Figure 3: Comparison on KL-divergence model with $R = 0.8$.

It can be seen that our robust methods can obtain policies that achieve higher worst-case reward. Also, both our limit-based robust value iteration and our direct method of robust relative value iteration converge to the optimal robust policies, which validates our theoretical results.

## Conclusion

In this paper, we investigated the problem of robust MDPs under the average-reward setting. We established *uniform* convergence of the discounted value function to average-reward, which further implies the uniform convergence of the *robust* discounted value function to *robust* average-reward. Based on this insight, we designed a robust dynamic programming approach using the robust discounted MDPs as an approximation (the limit method). We theoretically proved their convergence and optimality and proved a robust version of the Blackwell optimality (Blackwell 1962), i.e., any optimal policy of the robust discounted MDP when $\gamma$ is large enough is also an optimal policy of the robust average-reward. We then designed a direct approach for robust average-reward MDPs, where we derived the robust Bellman equation for robust average-reward MDPs. We further designed a robust RVI method, which was proven to converge to the optimal robust solution. Technically, our proof techniques are fundamentally different from existing studies on average-reward robust MDPs, e.g., those in (Blackwell 1962; Tewari and Bartlett 2007).

## Acknowledgment

# References

Abdullah, M. A.; Ren, H.; Ammar, H. B.; Milenkovic, V.; Luo, R.; Zhang, M.; and Wang, J. 2019. Wasserstein robust reinforcement learning. *arXiv preprint arXiv:1907.13196*.

Archibald, T.; McKinnon, K.; and Thomas, L. 1995. On the generation of Markov decision processes. *Journal of the Operational Research Society*, 46(3): 354–361.

Atia, G. K.; Beckus, A.; Alkhouri, I.; and Velasquez, A. 2021. Steady-State Planning in Expected Reward Multichain MDPs. *Journal of Artificial Intelligence Research*, 72: 1029–1082.

Badrinath, K. P.; and Kalathil, D. 2021. Robust Reinforcement Learning using Least Squares Policy Iteration with Provable Performance Guarantees. In *Proc. International Conference on Machine Learning (ICML)*, 511–520. PMLR.

Bagnell, J. A.; Ng, A. Y.; and Schneider, J. G. 2001. Solving uncertain Markov decision processes.

Bertsekas, D. P. 2011. Dynamic Programming and Optimal Control 3rd edition, volume II. *Belmont, MA: Athena Scientific*.

Blackwell, D. 1962. Discrete dynamic programming. *The Annals of Mathematical Statistics*, 719–726.

Chen, L.; Jain, R.; and Luo, H. 2022. Learning Infinite-Horizon Average-Reward Markov Decision Processes with Constraints. *arXiv preprint arXiv:2202.00150*.

Derman, C. 1970. *Finite state Markovian decision processes*. Academic Press, Inc.

Derman, E.; Geist, M.; and Mannor, S. 2021. Twice regularized MDPs and the equivalence between robustness and regularization. In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*.

Eysenbach, B.; and Levine, S. 2021. Maximum entropy RL (provably) solves some robust RL problems. *arXiv preprint arXiv:2103.06257*.

Goyal, V.; and Grand-Clement, J. 2018. Robust Markov decision process: Beyond rectangularity. *arXiv preprint arXiv:1811.00215*.

Ho, C. P.; Petrik, M.; and Wiesemann, W. 2018. Fast Bellman updates for robust MDPs. In *Proc. International Conference on Machine Learning (ICML)*, 1979–1988. PMLR.

Ho, C. P.; Petrik, M.; and Wiesemann, W. 2021. Partial policy iteration for L1-robust Markov decision processes. *Journal of Machine Learning Research*, 22(275): 1–46.

Hordijk, A.; and Yushkevich, A. A. 2002. Blackwell optimality. In *Handbook of Markov decision processes*, 231–267. Springer.

Hou, L.; Pang, L.; Hong, X.; Lan, Y.; Ma, Z.; and Yin, D. 2020. Robust Reinforcement Learning with Wasserstein Constraint. *arXiv preprint arXiv:2006.00945*.

Hu, Z.; and Hong, L. J. 2013. Kullback-Leibler divergence constrained distributionally robust optimization. *Available at Optimization Online*, 1695–1724.

Huang, S.; Papernot, N.; Goodfellow, I.; Duan, Y.; and Abbeel, P. 2017. Adversarial attacks on neural network policies. In *Proc. International Conference on Learning Representations (ICLR)*.

Huber, P. J. 1965. A Robust Version of the Probability Ratio Test. *Ann. Math. Statist.*, 36: 1753–1758.

Iyengar, G. N. 2005. Robust dynamic programming. *Mathematics of Operations Research*, 30(2): 257–280.

Kaufman, D. L.; and Schaefer, A. J. 2013. Robust modified policy iteration. *INFORMS Journal on Computing*, 25(3): 396–410.

Kober, J.; Bagnell, J. A.; and Peters, J. 2013. Reinforcement Learning in Robotics: A Survey. *The International Journal of Robotics Research*, 32(11): 1238–1274.

Kos, J.; and Song, D. 2017. Delving into adversarial attacks on deep policies. In *Proc. International Conference on Learning Representations (ICLR)*.

Lim, S. H.; and Autef, A. 2019. Kernel-based reinforcement learning in robust Markov decision processes. In *Proc. International Conference on Machine Learning (ICML)*, 3973–3981. PMLR.

Lim, S. H.; Xu, H.; and Mannor, S. 2013. Reinforcement learning in robust Markov decision processes. In *Proc. Advances in Neural Information Processing Systems (NIPS)*, 701–709.

Lin, Y.-C.; Hong, Z.-W.; Liao, Y.-H.; Shih, M.-L.; Liu, M.-Y.; and Sun, M. 2017. Tactics of adversarial attack on deep reinforcement learning agents. In *Proc. International Joint Conferences on Artificial Intelligence (IJCAI)*, 3756–3762.

Mandlekar, A.; Zhu, Y.; Garg, A.; Fei-Fei, L.; and Savarese, S. 2017. Adversarially robust policy learning: Active construction of physically-plausible perturbations. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 3932–3939. IEEE.

Nilim, A.; and El Ghaoui, L. 2004. Robustness in Markov decision problems with uncertain transition matrices. In *Proc. Advances in Neural Information Processing Systems (NIPS)*, 839–846.

Panaganti, K.; and Kalathil, D. 2021. Sample Complexity of Robust Reinforcement Learning with a Generative Model. *arXiv preprint arXiv:2112.01506*.

Pattanaik, A.; Tang, Z.; Liu, S.; Bommannan, G.; and Chowdhary, G. 2018. Robust Deep Reinforcement Learning with Adversarial Attacks. In *Proc. International Conference on Autonomous Agents and MultiAgent Systems*, 2040–2042.

Pinto, L.; Davidson, J.; Sukthankar, R.; and Gupta, A. 2017. Robust adversarial reinforcement learning. In *Proc. International Conference on Machine Learning (ICML)*, 2817–2826. PMLR.

Puterman, M. L. 1994. Markov Decision Processes: Discrete Stochastic Dynamic Programming.

Rahimian, H.; Bayraksan, G.; and De-Mello, T. H. 2022. Effective scenarios in multistage distributionally robust optimization with a focus on total variation distance. *SIAM Journal on Optimization*, 32(3): 1698–1727.

Rajeswaran, A.; Ghotra, S.; Ravindran, B.; and Levine, S. 2017. Epopt: Learning robust neural network policies using model ensembles. In *Proc. International Conference on Learning Representations (ICLR)*.

Roy, A.; Xu, H.; and Pokutta, S. 2017. Reinforcement learning under model mismatch. In *Proc. Advances in Neural Information Processing Systems (NIPS)*, 3046–3055.

Rudin, W. 2022. *Functional Analysis*. McGraw-Hill Science &Engineering &Math, 2nd edition.

Russel, R. H.; Benosman, M.; and Van Baar, J. 2020. Robust Constrained-MDPs: Soft-Constrained Robust Policy Optimization under Model Uncertainty. *arXiv preprint arXiv:2010.04870*.

Satia, J. K.; and Lave Jr, R. E. 1973. Markovian decision processes with uncertain transition probabilities. *Operations Research*, 21(3): 728–740.

Si, N.; Zhang, F.; Zhou, Z.; and Blanchet, J. 2020. Distributionally robust policy evaluation and learning in offline contextual bandits. In *Proc. International Conference on Machine Learning (ICML)*, 8884–8894. PMLR.

Sigaud, O.; and Buffet, O. 2013. *Markov decision processes in artificial intelligence*. John Wiley & Sons.

Sutton, R. S.; and Barto, A. G. 2018. *Reinforcement Learning: An Introduction*. Cambridge, Massachusetts: The MIT Press.

Tamar, A.; Mannor, S.; and Xu, H. 2014. Scaling up robust MDPs using function approximation. In *Proc. International Conference on Machine Learning (ICML)*, 181–189. PMLR.

Tessler, C.; Efroni, Y.; and Mannor, S. 2019. Action robust reinforcement learning and applications in continuous control. In *International Conference on Machine Learning*, 6215–6224. PMLR.

Tewari, A.; and Bartlett, P. L. 2007. Bounded parameter Markov decision processes with average reward criterion. In *International Conference on Computational Learning Theory*, 263–277. Springer.

Vinitsky, E.; Du, Y.; Parvate, K.; Jang, K.; Abbeel, P.; and Bayen, A. 2020. Robust Reinforcement Learning using Adversarial Populations. *arXiv preprint arXiv:2008.01825*.

Wan, Y.; Naik, A.; and Sutton, R. S. 2021. Learning and planning in average-reward markov decision processes. In *International Conference on Machine Learning*, 10653–10662. PMLR.

Wang, Y.; and Zou, S. 2021. Online Robust Reinforcement Learning with Model Uncertainty. In *Proc. Advances in Neural Information Processing Systems (NeurIPS)*.

Wang, Y.; and Zou, S. 2022. Policy Gradient Method For Robust Reinforcement Learning. In *Proc. International Conference on Machine Learning (ICML)*, volume 162, 23484–23526. PMLR.

Wei, C.-Y.; Jahromi, M. J.; Luo, H.; Sharma, H.; and Jain, R. 2020. Model-free reinforcement learning in infinite-horizon average-reward markov decision processes. In *International conference on machine learning*, 10170–10180. PMLR.

Wiesemann, W.; Kuhn, D.; and Rustem, B. 2013. Robust Markov decision processes. *Mathematics of Operations Research*, 38(1): 153–183.

Xu, H.; and Mannor, S. 2010. Distributionally Robust Markov Decision Processes. In *Proc. Advances in Neural Information Processing Systems (NIPS)*, 2505–2513.

Yang, W.; Zhang, L.; and Zhang, Z. 2021. Towards Theoretical Understandings of Robust Markov Decision Processes: Sample Complexity and Asymptotics. *arXiv preprint arXiv:2105.03863*.

Yu, P.; and Xu, H. 2015. Distributionally robust counterpart in Markov decision processes. *IEEE Transactions on Automatic Control*, 61(9): 2538–2543.

Zhang, Y.; and Ross, K. W. 2021. On-policy deep reinforcement learning for the average-reward criterion. In *Proc. International Conference on Machine Learning (ICML)*, 12535–12545. PMLR.

Zhou, Z.; Bai, Q.; Zhou, Z.; Qiu, L.; Blanchet, J.; and Glynn, P. 2021. Finite-Sample Regret Bound for Distributionally Robust Offline Tabular Reinforcement Learning. In *Proc. International Conference on Artifical Intelligence and Statistics (AISTATS)*, 3331–3339. PMLR.

# Review of Robust Discounted MDPs

In this section, we provide a brief review on the existing methods and results for robust discounted MDPs.

## Robust Policy Evaluation

We first consider the robust policy evaluation problem, where we aim to estimate the robust value function $V^\pi_{\mathcal{P},\gamma}$ for any policy $\pi$. It has been shown that the robust Bellman operator $\mathbf{T}_\pi$ is a $\gamma$-contraction, and the robust value function $V^\pi_{\mathcal{P},\gamma}$ is its unique fixed-point. Hence by recursively applying the robust Bellman operator, we can find the robust discounted value function (Nilim and El Ghaoui 2004; Iyengar 2005).

---

**Algorithm 4:** Policy evaluation for robust discounted MDPs

---

**Input**: $\pi, V_0, T$

1: **for** $t = 0, 1, ..., T - 1$ **do**
2:     **for** all $s \in \mathcal{S}$ **do**
3:         $V_{t+1}(s) \leftarrow \mathbb{E}_\pi[r(s, A) + \gamma \sigma_{\mathcal{P}^A_s}(V_t)]$
4:     **end for**
5: **end for**
6: **return** $V_T$

---

## Robust Optimal Control

Another important problem in robust MDP is to find the optimal policy which maximizes the robust discounted value function:

$$\pi^* = \arg\max_\pi V^\pi_{\mathcal{P},\gamma}. \tag{16}$$

A robust value iteration approach is developed in (Nilim and El Ghaoui 2004; Iyengar 2005) as follows.

---

**Algorithm 5:** Optimal Control for robust discounted MDPs

---

**Input**: $V_0, T$

1: **for** $t = 0, 1, ..., T - 1$ **do**
2:     **for** all $s \in \mathcal{S}$ **do**
3:         $V_{t+1}(s) \leftarrow \max_a \left\{ r(s, a) + \gamma \sigma_{\mathcal{P}^a_s}(V_t) \right\}$
4:     **end for**
5: **end for**
6: $\pi^*(s) \leftarrow \arg\max_a \left\{ r(s, a) + \gamma \sigma_{\mathcal{P}^a_s}(V_T) \right\}, \forall s$
7: **return** $\pi^*$

---

# Equivalence between Time-Varying and Stationary Models

We first provide an equivalence result between time-varying and stationary transition kernel models under stationary policies, which is an analog result to the one for robust discounted MDPs (Iyengar 2005; Nilim and El Ghaoui 2004). This result will be used in our following proofs.

Recall the definitions of robust discounted value function and worst-case average reward in eqs. (4) and (5), the worst-case is taken w.r.t. $\kappa = (\mathsf{P}_0, \mathsf{P}_1...) \in \bigotimes_{t\geq 0} \mathcal{P}$, therefore, the transition kernel at each time step could be different. This model is referred to as time-varying transition kernel model (as in (Iyengar 2005; Nilim and El Ghaoui 2004)). Another commonly used setting is that the transition kernels at different time step are the same, which is referred to as the stationary model (Iyengar 2005; Nilim and El Ghaoui 2004). In this paper, we use the following notations to distinguish the two models. By $\mathbb{E}_\mathsf{P}[\cdot]$, we denote the expectation when the transition kernels at all time steps are the same, $\mathsf{P}$, i.e., the stationary model. We also denote by $g^\pi_\mathsf{P}(s) \triangleq \lim_{n\to\infty} \mathbb{E}_{\mathsf{P},\pi} \left[ \frac{1}{n} \sum_{t=0}^{n-1} r_t \big| S_0 = s \right]$ and $V^\pi_{\mathsf{P},\gamma}(s) \triangleq \mathbb{E}_{\mathsf{P},\pi} \left[ \sum_{t=0}^\infty \gamma^t r_t \big| S_0 = s \right]$ being the expected average-reward and expected discounted value function under the stationary model $\mathsf{P}$. By $\mathbb{E}_\kappa[\cdot]$, we denote the expectation when the transition kernel at time $t$ is $\mathsf{P}_t$, i.e., the time-varying model.

For the discounted setting, it has been shown in (Nilim and El Ghaoui 2004) that for a stationary policy $\pi$, any $\gamma \in [0, 1)$, and any $s \in \mathcal{S}$,

$$V^\pi_{\mathcal{P},\gamma}(s) = \min_{\kappa \in \bigotimes_{t\geq 0} \mathcal{P}} \mathbb{E}_{\pi,\kappa} \left[ \sum_{t=0}^\infty \gamma^t r_t \big| S_0 = s \right]$$

$$= \min_{\mathsf{P} \in \mathcal{P}} \mathbb{E}_{\pi,\mathsf{P}} \left[ \sum_{t=0}^{\infty} \gamma^t r_t | S_0 = s \right]. \tag{17}$$

In the following theorem, we prove an analog of eq. (17) for robust-average reward MDPs that if we consider stationary policies, then the robust average-reward problem with the time-varying model can be equivalently solved by a stationary model.

Specifically, we define the worst-case average reward for the stationary transition kernel model as follows:

$$\min_{\mathsf{P} \in \mathcal{P}} \lim_{n \to \infty} \mathbb{E}_{\pi,\mathsf{P}} \left[ \frac{1}{n} \sum_{t=0}^{n-1} r_t | S_0 = s \right]. \tag{18}$$

Recall the worst-case average reward for the time-varying model in eq. (5). We will show that for any stationary policy, eq. (5) can be equivalently solved by solving eq. (18).

**Theorem 9.** *Consider an arbitrary stationary policy $\pi$. Then, the worst-case average-reward under the time-varying model is the same as the one under the stationary model:*

$$g_{\mathcal{P}}^\pi(s) \triangleq \min_{\kappa \in \otimes_{t \geq 0} \mathcal{P}} \lim_{n \to \infty} \mathbb{E}_{\kappa,\pi} \left[ \frac{1}{n} \sum_{t=0}^{n-1} r_t | S_0 = s \right]$$

$$= \min_{\mathsf{P} \in \mathcal{P}} \lim_{n \to \infty} \mathbb{E}_{\mathsf{P},\pi} \left[ \frac{1}{n} \sum_{t=0}^{n-1} r_t | S_0 = s \right]. \tag{19}$$

*Similar result also holds for the robust relative value function:*

$$V_{\mathcal{P}}^\pi(s) \triangleq \min_{\kappa \in \otimes_{t \geq 0} \mathcal{P}} \mathbb{E}_{\kappa,\pi} \left[ \sum_{t=0}^{\infty} (r_t - g_{\mathcal{P}}^\pi) | S_0 = s \right]$$

$$= \min_{\mathsf{P} \in \mathcal{P}} \mathbb{E}_{\mathsf{P},\pi} \left[ \sum_{t=0}^{\infty} (r_t - g_{\mathcal{P}}^\pi) | S_0 = s \right]. \tag{20}$$

*Proof.* From the robust Bellman equation in Theorem 6 [2], we have that

$$V_{\mathcal{P}}^\pi(s) + g_{\mathcal{P}}^\pi = \sum_a \pi(a|s) \left( r(s,a) + \sigma_{\mathcal{P}_s^a}(V_{\mathcal{P}}^\pi) \right). \tag{21}$$

Denote by $\arg\min_{p \in \mathcal{P}_s^a} (p)^\top V_{\mathcal{P}}^\pi \triangleq p_s^a$ [3], and denote by $\mathsf{P}^\pi \triangleq \{p_s^a : s \in \mathcal{S}, a \in \mathcal{A}\}$. It then follows that

$$V_{\mathcal{P}}^\pi(s) = \sum_a \pi(a|s) \left( r(s,a) - g_{\mathcal{P}}^\pi + \sigma_{\mathcal{P}_s^a}(V_{\mathcal{P}}^\pi) \right)$$

$$= \sum_a \pi(a|s)(r(s,a) - g_{\mathcal{P}}^\pi) + \sum_a \pi(a|s)\mathbb{E}_{\mathsf{P}^\pi}[V_{\mathcal{P}}^\pi(S_1)|S_0 = s, A_0 = a]$$

$$= \sum_a \pi(a|s)(r(s,a) - g_{\mathcal{P}}^\pi) + \mathbb{E}_{\mathsf{P}^\pi,\pi}[V_{\mathcal{P}}^\pi(S_1)|S_0 = s]$$

$$= \sum_a \pi(a|s)(r(s,a) - g_{\mathcal{P}}^\pi) + \mathbb{E}_{\mathsf{P}^\pi,\pi}\left[ \sum_a \pi(a|S_1)(r(S_1,a) - g_{\mathcal{P}}^\pi)|S_0 = s \right] + \mathbb{E}_{\mathsf{P}^\pi,\pi}\left[ \sum_a \pi(a|S_1)\sigma_{\mathcal{P}_{S_1}^a}(V_{\mathcal{P}}^\pi)|S_0 = s \right]$$

$$= \sum_a \pi(a|s)(r(s,a) - g_{\mathcal{P}}^\pi) + \mathbb{E}_{\mathsf{P}^\pi,\pi}[r_1 - g_{\mathcal{P}}^\pi|S_0 = s] + \mathbb{E}_{\mathsf{P}^\pi,\pi}\left[ \sigma_{\mathcal{P}_{S_1}^{A_1}}(V_{\mathcal{P}}^\pi)|S_0 = s \right]$$

$$= \sum_a \pi(a|s)(r(s,a) - g_{\mathcal{P}}^\pi) + \mathbb{E}_{\mathsf{P}^\pi,\pi}\left[ r_1 - g_{\mathcal{P}}^\pi|S_0 = s \right] + \mathbb{E}_{\mathsf{P}^\pi,\pi}\left[ (p_{S_1}^{A_1})^\top V_{\mathcal{P}}^\pi|S_0 = s \right]$$

$$= \mathbb{E}_{\mathsf{P}^\pi,\pi}\left[ r_0 - g_{\mathcal{P}}^\pi + r_1 - g_{\mathcal{P}}^\pi|S_0 = s \right] + \mathbb{E}_{\mathsf{P}^\pi,\pi}[V_{\mathcal{P}}^\pi(S_2)|S_0 = s]$$

$$\ldots\ldots$$

---

[2] The proof of Theorem 6 is independent of theorem 9 and does not relay on the results to be showed here.

[3] We pick one arbitrarily, if there are multiple minimizers.

$$= \mathbb{E}_{\mathsf{P}^\pi, \pi}\left[\sum_{t=0}^\infty (r_t - g_{\mathcal{P}}^\pi)|s\right]. \tag{22}$$

By the definition, the following always hold:

$$\min_{\kappa \in \bigotimes_{t \geq 0} \mathcal{P}} \mathbb{E}_{\kappa, \pi}\left[\sum_{t=0}^\infty (r_t - g_{\mathcal{P}}^\pi)|S_0 = s\right] \leq \min_{\mathsf{P} \in \mathcal{P}} \mathbb{E}_{\mathsf{P}, \pi}\left[\sum_{t=0}^\infty (r_t - g_{\mathcal{P}}^\pi)|S_0 = s\right]. \tag{23}$$

This hence implies that a stationary transition kernel sequence $\kappa = (\mathsf{P}^\pi, \mathsf{P}^\pi, ...)$ is one of the worst-case transition kernels for $V_{\mathcal{P}}^\pi$. Therefore, eq. (20) can be proved.

Consider the transition kernel $\mathsf{P}^\pi$. We denote its non-robust average-reward and the non-robust relative value function by $g_{\mathsf{P}^\pi}^\pi$ and $V_{\mathsf{P}^\pi}^\pi$. By the non-robust Bellman equation (Sutton and Barto 2018), we have that

$$V_{\mathsf{P}^\pi}^\pi(s) = \sum_a \pi(a|s)(r(s,a) - g_{\mathsf{P}^\pi}^\pi) + \mathbb{E}_{\mathsf{P}^\pi, \pi}[V_{\mathsf{P}^\pi}^\pi(S_1)|s]. \tag{24}$$

On the other hand, the robust Bellman equation shows that

$$V_{\mathcal{P}}^\pi(s) = V_{\mathsf{P}^\pi}^\pi(s) = \sum_a \pi(a|s)(r(s,a) - g_{\mathcal{P}}^\pi) + \mathbb{E}_{\mathsf{P}^\pi, \pi}[V_{\mathsf{P}^\pi}^\pi(S_1)|s]. \tag{25}$$

These two equations hence implies that $g_{\mathcal{P}}^\pi = g_{\mathsf{P}^\pi}^\pi$, and hence the stationary kernel $(\mathsf{P}^\pi, \mathsf{P}^\pi, ...)$ is also a worst-case kernel of robust average-reward in the time-varying setting. This proves eq. (19). $\square$

## Proof of Theorem 1

In the proof, unless otherwise specified, we denote by $\|v\|$ the $l_\infty$ norm of a vector $v$, and for a matrix $A$, we denote by $\|A\|$ its matrix norm induced by $l_\infty$ norm, i.e., $\|A\| = \sup_{x \in \mathbb{R}^d} \frac{\|Ax\|_\infty}{\|x\|_\infty}$.

**Lemma 1.** *[Theorem 8.2.3 in (Puterman 1994)] For any* $\mathsf{P}$, $\gamma$, $\pi$,

$$V_{\mathsf{P}, \gamma}^\pi = \frac{1}{1-\gamma} g_{\mathsf{P}}^\pi + h_{\mathsf{P}}^\pi + f_{\mathsf{P}}^\pi(\gamma), \tag{26}$$

*where* $h_{\mathsf{P}}^\pi = H_{\mathsf{P}}^\pi r_\pi$, *and* $f_{\mathsf{P}}^\pi(\gamma) = \frac{1}{\gamma} \sum_{n=1}^\infty (-1)^n \left(\frac{1-\gamma}{\gamma}\right)^n (H_{\mathsf{P}}^\pi)^{n+1} r_\pi$.

Following Proposition 8.4.6 in (Puterman 1994), we can show the following lemma.

**Lemma 2.** $H_{\mathsf{P}}^\pi$ *is continuous on* $\Pi \times \mathcal{P}$. *If* $\Pi$ *and* $\mathcal{P}$ *are compact,* $\|H_{\mathsf{P}}^\pi\|$ *is uniformly bounded on* $\Pi \times \mathcal{P}$, *i.e., there exists a constant* $h$, *such that* $\|H_{\mathsf{P}}^\pi\| \leq h$ *for any* $\pi, \mathsf{P}$.

For simplicity, denote by

$$S_\infty^\pi(\mathsf{P}, \gamma) \triangleq \frac{1}{\gamma} \sum_{n=1}^\infty (-1)^n \left(\frac{1-\gamma}{\gamma}\right)^n (H_{\mathsf{P}}^\pi)^{n+1} r_\pi,$$

$$S_N^\pi(\mathsf{P}, \gamma) \triangleq \frac{1}{\gamma} \sum_{n=1}^N (-1)^n \left(\frac{1-\gamma}{\gamma}\right)^n (H_{\mathsf{P}}^\pi)^{n+1} r_\pi. \tag{27}$$

Clearly $S_\infty^\pi(\mathsf{P}, \gamma) = f_{\mathsf{P}}^\pi(\gamma)$ and $\lim_{N \to \infty} S_N^\pi(\mathsf{P}, \gamma) = S_\infty^\pi(\mathsf{P}, \gamma)$ for any specific $\pi, \mathsf{P}$.

**Lemma 3.** *There exists* $\delta \in (0, 1)$, *such that*

$$\lim_{N \to \infty} S_N^\pi(\mathsf{P}, \gamma) = S_\infty^\pi(\mathsf{P}, \gamma) \tag{28}$$

*uniformly on* $\Pi \times \mathcal{P} \times [\delta, 1]$.

*Proof.* Note that $\|H_{\mathsf{P}}^\pi\| \leq h$, hence there exists $\delta$, s.t.

$$\frac{1-\delta}{\delta} h \leq k < 1 \tag{29}$$

for some constant $k$. Then for any $\gamma \geq \delta$,

$$\frac{1-\gamma}{\gamma} h \leq \frac{1-\delta}{\delta} h \leq k. \tag{30}$$

Moreover, note that

$$\left\|\frac{1}{\gamma}(-1)^n\left(\frac{1-\gamma}{\gamma}\right)^n (H_{\mathsf{P}}^\pi)^{n+1}r\right\| \leq \frac{1}{\gamma}\left(\frac{1-\gamma}{\gamma}\right)^n h^{n+1} \leq \frac{hk^n}{\delta} \triangleq M_n, \tag{31}$$

which is because $\|A+B\| \leq \|A\| + \|B\|$ for induced $l_\infty$ norm, $\|Ax\| \leq \|A\|\|x\|$ and $\|r_\pi\|_\infty \leq 1$.

Note that

$$\sum_{n=1}^\infty M_n = \frac{h}{\delta}\frac{k}{1-k}, \tag{32}$$

hence by Weierstrass $M$-test (Rudin 2022), $S_N^\pi(\mathsf{P}, \gamma)$ uniformly converges to $S_\infty^\pi(\mathsf{P}, \gamma)$ on $\Pi \times \mathcal{P} \times [\delta, 1]$. $\qquad\square$

**Lemma 4.** *There exists a uniform constant L, such that*

$$\|S_N^\pi(\mathsf{P}, \gamma_1) - S_N^\pi(\mathsf{P}, \gamma_2)\| \leq L|\gamma_1 - \gamma_2|, \tag{33}$$

*for any $N$, $\pi$, $\mathsf{P}$, $\gamma_1, \gamma_2 \in [\delta, 1]$.*

*Proof.* We first show that $\gamma S_N^\pi(\mathsf{P}, \gamma) = \sum_{n=1}^N (-1)^n \left(\frac{1-\gamma}{\gamma}\right)^n (H_{\mathsf{P}}^\pi)^{n+1}r_\pi \triangleq T_N^\pi(\mathsf{P}, \gamma)$ is uniformly Lipschitz w.r.t. the $l_\infty$ norm, i.e.,

$$\|T_N^\pi(\mathsf{P}, \gamma_1) - T_N^\pi(\mathsf{P}, \gamma_2)\| \leq l|\gamma_1 - \gamma_2|, \tag{34}$$

for any $N$, $\pi$, $\mathsf{P}$, $\gamma_1, \gamma_2 \in [\delta, 1]$ and some constant $l$.

Clearly, it can be shown by verifying $\nabla T_N^\pi(\mathsf{P}, \gamma)$ is uniformly bounded for any $\pi, N, \mathsf{P}$ or $\gamma$.

First, it can be shown that

$$\nabla T_N^\pi(\mathsf{P}, \gamma) = \sum_{n=1}^N (-1)^n n \left(\frac{1-\gamma}{\gamma}\right)^{n-1} \frac{-1}{\gamma^2}(H_{\mathsf{P}}^\pi)^{n+1}r_\pi, \tag{35}$$

and moreover

$$\|\nabla T_N^\pi(\mathsf{P}, \gamma)\| \leq \sum_{n=1}^N n \left(\frac{1-\gamma}{\gamma}\right)^{n-1} \frac{1}{\gamma^2}h^{n+1} \triangleq l_N(\gamma). \tag{36}$$

Note that

$$h\frac{1-\gamma}{\gamma}l_N(\gamma) = \sum_{n=1}^N n \left(\frac{1-\gamma}{\gamma}\right)^n \frac{1}{\gamma^2}h^{n+2}, \tag{37}$$

then, we can show that

$$\left(1 - h\frac{1-\gamma}{\gamma}\right)l_N(\gamma)$$

$$= \sum_{n=1}^N n \left(\frac{1-\gamma}{\gamma}\right)^{n-1} \frac{1}{\gamma^2}h^{n+1} - \sum_{n=1}^N n \left(\frac{1-\gamma}{\gamma}\right)^n \frac{1}{\gamma^2}h^{n+2}$$

$$= \frac{1}{\gamma^2}h^2 - N\left(\frac{1-\gamma}{\gamma}\right)^N \frac{1}{\gamma^2}h^{N+2} + \sum_{n=2}^N \left(\frac{1-\gamma}{\gamma}\right)^{n-1} \frac{1}{\gamma^2}h^{n+1}$$

$$\leq \frac{1}{\gamma^2}h^2 + \frac{h^2}{\gamma^2}\frac{1-\gamma}{\gamma}h\frac{1}{1-\frac{1-\gamma}{\gamma}h}$$

$$= \frac{h^2}{\gamma^2} + \frac{h^2}{\gamma^2}\frac{1-\gamma}{\gamma}h\frac{1}{1-\frac{1-\gamma}{\gamma}h}. \tag{38}$$

Hence, we have that

$$\|\nabla T_N^\pi(\mathsf{P}, \gamma)\| \leq l_N(\gamma) \leq \frac{1}{1-h\frac{1-\gamma}{\gamma}}\left(\frac{h^2}{\gamma^2} + \frac{h^2}{\gamma^2}\frac{1-\gamma}{\gamma}h\frac{1}{1-\frac{1-\gamma}{\gamma}h}\right)$$

$$\leq \frac{1}{1-k}\left(\frac{h^2}{\delta^2} + \frac{h^2}{\delta^2}\frac{k}{1-k}\right),\tag{39}$$

which implies a uniform bound on $\|\nabla T_N^\pi(\mathsf{P},\gamma)\|$.

Now, we have that

$$|S_N^\pi(\mathsf{P},\gamma_1) - S_N^\pi(\mathsf{P},\gamma_2)|$$
$$\leq \frac{|\gamma_2 - \gamma_1|}{\gamma_1\gamma_2}\|T_N^\pi(\mathsf{P},\gamma_1)\| + \frac{\|T_N^\pi(\mathsf{P},\gamma_1) - T_N^\pi(\mathsf{P},\gamma_2)\|}{\gamma_2}.\tag{40}$$

To show $\|T_N^\pi(\mathsf{P},\gamma)\|$ is uniformly bounded, we have that

$$\|T_N^\pi(\mathsf{P},\gamma)\| \leq \sum_{n=1}^N \left\|\left(\frac{1-\gamma}{\gamma}\right)^n (H_\mathsf{P}^\pi)^{n+1}r\right\|$$
$$\leq \sum_{n=1}^N \left(\frac{1-\gamma}{\gamma}\right)^n h^{n+1}$$
$$\leq \sum_{n=1}^N k^n h$$
$$\leq h\frac{k}{1-k}.\tag{41}$$

Then, it follows that

$$\|S_N^\pi(\mathsf{P},\gamma_1) - S_N^\pi(\mathsf{P},\gamma_2)\|$$
$$= \left\|\frac{\gamma_2 - \gamma_1}{\gamma_1\gamma_2}T_N^\pi(\mathsf{P},\gamma_1) + \frac{T_N^\pi(\mathsf{P},\gamma_1) - T_N^\pi(\mathsf{P},\gamma_2)}{\gamma_2}\right\|$$
$$\leq \left(\frac{1}{\delta^2}h\frac{k}{1-k} + \frac{1}{\delta}\frac{1}{1-k}\left(\frac{h^2}{\delta^2} + \frac{h^2}{\delta^2}\frac{k}{1-k}\right)\right)|\gamma_1 - \gamma_2|$$
$$\triangleq L|\gamma_1 - \gamma_2|,\tag{42}$$

where $L = \left(\frac{1}{\delta^2}h\frac{k}{1-k} + \frac{1}{\delta}\frac{1}{1-k}\left(\frac{h^2}{\delta^2} + \frac{h^2}{\delta^2}\frac{k}{1-k}\right)\right)$ is a universal constant that does not depend on $N, \mathsf{P}, \pi$ or $\gamma$. $\square$

**Lemma 5.** $S_\infty^\pi(\mathsf{P},\gamma)$ *uniformly converges as $\gamma \to 1$ on $\Pi \times \mathcal{P}$. Also, $S_\infty^\pi(\mathsf{P},\gamma)$ is L-Lipschitz for any $\gamma > \delta$: for any $\pi, \mathsf{P}$ and any $\gamma_1, \gamma_2 \in (\delta, 1]$.*

$$\|S_\infty^\pi(\mathsf{P},\gamma_1) - S_\infty^\pi(\mathsf{P},\gamma_2)\| \leq L|\gamma_1 - \gamma_2|.\tag{43}$$

*Proof.* From Lemma 3, for any $\epsilon$, there exists $N_\epsilon$, such that for any $n \geq N_\epsilon, \pi, \mathsf{P}, \gamma > \delta$,

$$\|S_\infty^\pi(\mathsf{P},\gamma) - S_n^\pi(\mathsf{P},\gamma)\| < \epsilon.\tag{44}$$

Thus for any $\gamma_1, \gamma_2 \in (\delta, 1]$,

$$\|S_\infty^\pi(\mathsf{P},\gamma_1) - S_\infty^\pi(\mathsf{P},\gamma_2)\|$$
$$\leq \|S_\infty^\pi(\mathsf{P},\gamma_1) - S_n^\pi(\mathsf{P},\gamma_1)\| + \|S_n^\pi(\mathsf{P},\gamma_1) - S_n^\pi(\mathsf{P},\gamma_2)\| + \|S_n^\pi(\mathsf{P},\gamma_2) - S_\infty^\pi(\mathsf{P},\gamma_2)\|$$
$$\leq 2\epsilon + \|S_n^\pi(\mathsf{P},\gamma_1) - S_n^\pi(\mathsf{P},\gamma_2)\|$$
$$\leq 2\epsilon + L|\gamma_1 - \gamma_2|,\tag{45}$$

where the last step is from Lemma 4.

Thus, for any $\epsilon$, there exists $\omega = \max\{\delta, 1-\epsilon\}$, such that for any $\gamma_1, \gamma_2 > \omega$,

$$\|S_\infty^\pi(\mathsf{P},\gamma_1) - S_\infty^\pi(\mathsf{P},\gamma_2)\| < (2+L)\epsilon,\tag{46}$$

and hence by Cauchy's criterion we conclude that $S_\infty^\pi(\mathsf{P},\gamma)$ converges uniformly on $\Pi \times \mathcal{P}$.

On the other hand, since eq. (45) holds for any $\epsilon$, it implies that

$$\|S_\infty^\pi(\mathsf{P},\gamma_1) - S_\infty^\pi(\mathsf{P},\gamma_2)\| \leq L|\gamma_1 - \gamma_2|,\tag{47}$$

which completes the proof. $\square$

We now prove Theorem 1. For any $\mathsf{P}, \pi$, we have that

$$V^\pi_{\mathsf{P},\gamma} = \frac{1}{1-\gamma}g^\pi_\mathsf{P} + h^\pi_\mathsf{P} + f^\pi_\mathsf{P}(\gamma). \tag{48}$$

It then follows that

$$(1-\gamma)V^\pi_{\mathsf{P},\gamma} = g^\pi_\mathsf{P} + (1-\gamma)h^\pi_\mathsf{P} + (1-\gamma)f^\pi_\mathsf{P}(\gamma). \tag{49}$$

Clearly $(1-\gamma)h^\pi_\mathsf{P} \to 0$ uniformly on $\Pi \times \mathcal{P}$ because $\|h^\pi_\mathsf{P}\| = \|H^\pi_\mathsf{P} r_\pi\| \le h$ is uniformly bounded. Then,

$$\begin{aligned}
&\|(1-\gamma_1)f^\pi_\mathsf{P}(\gamma_1) - (1-\gamma_2)f^\pi_\mathsf{P}(\gamma_2)\| \\
&\le \|(1-\gamma_1)f^\pi_\mathsf{P}(\gamma_1) - (1-\gamma_1)f^\pi_\mathsf{P}(\gamma_2)\| + \|(1-\gamma_1)f^\pi_\mathsf{P}(\gamma_2) - (1-\gamma_2)f^\pi_\mathsf{P}(\gamma_2)\| \\
&\le (1-\gamma_1)L|\gamma_1 - \gamma_2| + \|f^\pi_\mathsf{P}(\gamma_2)\||\gamma_1 - \gamma_2|.
\end{aligned} \tag{50}$$

For any $\pi, \mathsf{P}, \gamma > \delta$,

$$\begin{aligned}
\|f^\pi_\mathsf{P}(\gamma)\| &= \left\| \frac{1}{\gamma}\sum_{n=1}^{\infty}(-1)^n\left(\frac{1-\gamma}{\gamma}\right)^n (H^\pi_\mathsf{P})^{n+1}r_\pi \right\| \\
&\le \left| \frac{1}{\gamma}\sum_{n=1}^{\infty}\left(\frac{1-\gamma}{\gamma}\right)^n h^{n+1} \right| \\
&\le \frac{h}{\delta}\frac{1-\gamma}{\gamma}h\frac{1}{1-\frac{1-\gamma}{\gamma}h} \\
&\le \frac{h}{\delta}\frac{k}{1-k} \\
&\triangleq c_f.
\end{aligned} \tag{51}$$

Hence, $(1-\gamma)f^\pi_\mathsf{P}(\gamma) \to 0$ uniformly on $\Pi \times \mathcal{P}$ due to the fact that $\|f^\pi_\mathsf{P}(\gamma)\|$ is uniformly bounded for any $\pi, \gamma > \delta, \mathsf{P}$.
Then we have that $\lim_{\gamma \to 1}(1-\gamma)V^\pi_{\mathsf{P},\gamma} = g^\pi_\mathsf{P}$ uniformly on $\mathcal{P} \times \Pi$. This completes the proof of Theorem 1.

## Proof of Theorem 2

We first show a lemma which allows us to interchange the order of $\lim$ and $\max$.

**Lemma 6.** *If a function $f(x,y)$ converges uniformly to $F(x)$ on $\mathcal{X}$ as $y \to y_0$, then*

$$\max_x \lim_{y \to y_0} f(x,y) = \lim_{y \to y_0}\max_x f(x,y). \tag{52}$$

*Proof.* For each $f(x,y)$, denote by $\arg\max_x f(x,y) = x_y$, and hence $f(x_y, y) \ge f(x,y)$ for any $x, y$. Also denote by $\arg\max_x F(x) = x'$. Now because $f(x,y)$ uniformly converges to $F(x)$, then for any $\epsilon$, there exists $\delta'$, such that $\forall|y-y_0| < \delta'$,

$$|f(x,y) - F(x)| \le \epsilon \tag{53}$$

for any $x$. Now consider $|f(x_y, y) - F(x')|$ for $|y - y_0| < \delta'$. If $f(x_y, y) - F(x') > 0$, then

$$|f(x_y, y) - F(x')| = f(x_y, y) - F(x') = f(x_y, y) - F(x_y) + F(x_y) - F(x') \le \epsilon; \tag{54}$$

On the other hand if $f(x_y, y) - F(x') < 0$, then

$$|f(x_y, y) - F(x')| = F(x') - f(x_y, y) = F(x') - f(x', y) + f(x', y) - f(x_y, y) \le \epsilon. \tag{55}$$

Hence, we showed that for any $\epsilon$, there exists $\delta'$, such that $\forall|y-y_0| < \delta'$,

$$|f(x_y, y) - F(x')| = |\max_x f(x,y) - \max_x F(x)| \le \epsilon, \tag{56}$$

and hence

$$\lim_{y \to y_0}\max_x f(x,y) = \max_x F(x) = \max_x \lim_{y \to y_0} f(x,y), \tag{57}$$

and this completes the proof. $\qquad\square$

Then, we show that the robust discounted value function converges uniformly to the robust average-reward as the discounted factor approaches 1.

**Theorem 10** (Restatement of Theorem 2)**.** *The robust discounted value function converges uniformly to the robust average-reward on $\Pi$:*

$$\lim_{\gamma \to 1} (1-\gamma) V^{\pi}_{\mathcal{P},\gamma} = g^{\pi}_{\mathcal{P}}. \tag{58}$$

*Proof.* Due to Theorem 9, for any stationary policy $\pi$, $g^{\pi}_{\mathcal{P}}(s) = \min_{\mathsf{P} \in \mathcal{P}} g^{\pi}_{\mathsf{P}}(s)$ under the stationary model. Hence from the uniform convergence in Theorem 1, we first show the following:

$$
\begin{aligned}
g^{\pi}_{\mathcal{P}} &= \min_{\mathsf{P} \in \mathcal{P}} g^{\pi}_{\mathsf{P}} \\
&= \min_{\mathsf{P} \in \mathcal{P}} \lim_{\gamma \to 1} (1-\gamma) V^{\pi}_{\mathsf{P},\gamma} \\
&\overset{(a)}{=} \lim_{\gamma \to 1} \min_{\mathsf{P} \in \mathcal{P}} (1-\gamma) V^{\pi}_{\mathsf{P},\gamma} \\
&= \lim_{\gamma \to 1} (1-\gamma) V^{\pi}_{\mathcal{P},\gamma},
\end{aligned} \tag{59}
$$

where $(a)$ is because Lemma 6. Moreover, note that $\lim_{\gamma \to 1} (1-\gamma) V^{\pi}_{\mathsf{P},\gamma} = g^{\pi}_{\mathsf{P}}$ uniformly on $\Pi \times \mathcal{P}$, hence the convergence in (59) is also uniform on $\Pi$. Thus, we complete the proof. $\qquad\square$

## Proof of Theorem 3

**Theorem 11** (Restatement of Theorem 3)**.** *$V_T$ generated by Algorithm 1 converges to the robust average-reward $g^{\pi}_{\mathcal{P}}$ as $T \to \infty$.*

*Proof.* From discounted robust Bellman equation (Nilim and El Ghaoui 2004), it can be shown that

$$(1-\gamma_t) V^{\pi}_{\mathcal{P},\gamma_t} = (1-\gamma_t) \sum_a \pi(a|s)(r(s,a) + \gamma_t \sigma_{\mathcal{P}^a_s}(V^{\pi}_{\mathcal{P},\gamma_t})). \tag{60}$$

Then we can show that for any $s \in \mathcal{S}$,

$$
\begin{aligned}
&|V_{t+1}(s) - (1-\gamma_{t+1}) V^{\pi}_{\mathcal{P},\gamma_{t+1}}(s)| \\
&= |V_{t+1}(s) - (1-\gamma_t) V^{\pi}_{\mathcal{P},\gamma_t}(s) + (1-\gamma_t) V^{\pi}_{\mathcal{P},\gamma_t}(s) - (1-\gamma_{t+1}) V^{\pi}_{\mathcal{P},\gamma_{t+1}}(s)| \\
&\leq |(1-\gamma_t) V^{\pi}_{\mathcal{P},\gamma_t}(s) - (1-\gamma_{t+1}) V^{\pi}_{\mathcal{P},\gamma_{t+1}}(s)| + |V_{t+1}(s) - (1-\gamma_t) V^{\pi}_{\mathcal{P},\gamma_t}(s)| \\
&= |(1-\gamma_t) V^{\pi}_{\mathcal{P},\gamma_t}(s) - (1-\gamma_{t+1}) V^{\pi}_{\mathcal{P},\gamma_{t+1}}(s)| \\
&\quad + \left| \sum_a \pi(a|s)\Big( (1-\gamma_t) r(s,a) + \gamma_t \sigma_{\mathcal{P}^a_s}(V_t) - ((1-\gamma_t) r(s,a) + \gamma_t \sigma_{\mathcal{P}^a_s}((1-\gamma_t) V^{\pi}_{\mathcal{P},\gamma_t})) \Big) \right| \\
&= |(1-\gamma_t) V^{\pi}_{\mathcal{P},\gamma_t}(s) - (1-\gamma_{t+1}) V^{\pi}_{\mathcal{P},\gamma_{t+1}}(s)| + \left| \sum_a \pi(a|s)\Big( \gamma_t \sigma_{\mathcal{P}^a_s}(V_t) - \gamma_t \sigma_{\mathcal{P}^a_s}((1-\gamma_t) V^{\pi}_{\mathcal{P},\gamma_t}) \Big) \right| \\
&= |(1-\gamma_t) V^{\pi}_{\mathcal{P},\gamma_t}(s) - (1-\gamma_{t+1}) V^{\pi}_{\mathcal{P},\gamma_{t+1}}(s)| + \gamma_t \left| \sum_a \pi(a|s)\Big( \sigma_{\mathcal{P}^a_s}(V_t) - \sigma_{\mathcal{P}^a_s}((1-\gamma_t) V^{\pi}_{\mathcal{P},\gamma_t}) \Big) \right|.
\end{aligned} \tag{61, 62}
$$

If we denote by $\Delta_t \triangleq \|V_t - (1-\gamma_t) V^{\pi}_{\mathcal{P},\gamma_t}\|_\infty$, then

$$\Delta_{t+1} \leq \|(1-\gamma_t) V^{\pi}_{\mathcal{P},\gamma_t} - (1-\gamma_{t+1}) V^{\pi}_{\mathcal{P},\gamma_{t+1}}\|_\infty + \gamma_t \max_s \left\{ \sum_a \pi(a|s) \left| \sigma_{\mathcal{P}^a_s}(V_t) - \sigma_{\mathcal{P}^a_s}((1-\gamma_t) V^{\pi}_{\mathcal{P},\gamma_t}) \right| \right\}. \tag{63}$$

It can be easily verified that $\sigma_{\mathcal{P}^a_s}(V)$ is a 1-Lipschitz function, thus the second term in (63) can be further bounded as

$$
\begin{aligned}
&\sum_a \pi(a|s) \left| \sigma_{\mathcal{P}^a_s}(V_t) - \sigma_{\mathcal{P}^a_s}((1-\gamma_t) V^{\pi}_{\mathcal{P},\gamma_t}) \right| \\
&\leq \sum_a \pi(a|s) \|V_t - (1-\gamma_t) V^{\pi}_{\mathcal{P},\gamma_t}\|_\infty \\
&= \|V_t - (1-\gamma_t) V^{\pi}_{\mathcal{P},\gamma_t}\|_\infty,
\end{aligned} \tag{64}
$$

and hence

$$\Delta_{t+1} \leq \|(1-\gamma_t) V^{\pi}_{\mathcal{P},\gamma_t} - (1-\gamma_{t+1}) V^{\pi}_{\mathcal{P},\gamma_{t+1}}\|_\infty + \gamma_t \Delta_t. \tag{65}$$

Recall that

$$(1 - \gamma_t)V^{\pi}_{\mathcal{P},\gamma_t} = (1 - \gamma_t)\min_{\mathsf{P}} V^{\pi}_{\mathsf{P},\gamma_t}. \tag{66}$$

Let $s_t^* \triangleq \arg\max_s |(1 - \gamma_t)V^{\pi}_{\mathcal{P},\gamma_t}(s) - (1 - \gamma_{t+1})V^{\pi}_{\mathcal{P},\gamma_{t+1}}(s)|$. Then it follows that

$$\|(1 - \gamma_t)V^{\pi}_{\mathcal{P},\gamma_t} - (1 - \gamma_{t+1})V^{\pi}_{\mathcal{P},\gamma_{t+1}}\|_{\infty} = |(1 - \gamma_t)V^{\pi}_{\mathcal{P},\gamma_t}(s_t^*) - (1 - \gamma_{t+1})V^{\pi}_{\mathcal{P},\gamma_{t+1}}(s_t^*)|. \tag{67}$$

Note that from (Nilim and El Ghaoui 2004; Iyengar 2005), for any stationary policy $\pi$, there exists a stationary model P such that $V^{\pi}_{\mathcal{P},\gamma}(s) = \mathbb{E}_{\mathsf{P},\pi}\left[\sum_{t=0}^{\infty}\gamma^t r_t | S_0 = s\right] \triangleq V^{\pi}_{\mathsf{P},\gamma}$. Hence in the following, for each $\gamma_t$, we denote the worst-case transition kernel of $V^{\pi}_{\mathcal{P},\gamma_t}$ by $\mathsf{P}_t$.

If $(1 - \gamma_t)V^{\pi}_{\mathcal{P},\gamma_t}(s_t^*) \geq (1 - \gamma_{t+1})V^{\pi}_{\mathcal{P},\gamma_{t+1}}(s_t^*)$, then

$$\begin{aligned}
&|(1 - \gamma_t)V^{\pi}_{\mathcal{P},\gamma_t}(s_t^*) - (1 - \gamma_{t+1})V^{\pi}_{\mathcal{P},\gamma_{t+1}}(s_t^*)| \\
&= \min_{\mathsf{P}}(1 - \gamma_t)V^{\pi}_{\mathsf{P},\gamma_t}(s_t^*) - \min_{\mathsf{P}}(1 - \gamma_{t+1})V^{\pi}_{\mathsf{P},\gamma_{t+1}}(s_t^*) \\
&= (1 - \gamma_t)V^{\pi}_{\mathsf{P}_t,\gamma_t}(s_t^*) - (1 - \gamma_{t+1})V^{\pi}_{\mathsf{P}_{t+1},\gamma_{t+1}}(s_t^*) \\
&= (1 - \gamma_t)V^{\pi}_{\mathsf{P}_t,\gamma_t}(s_t^*) - (1 - \gamma_t)V^{\pi}_{\mathsf{P}_{t+1},\gamma_t}(s_t^*) + (1 - \gamma_t)V^{\pi}_{\mathsf{P}_{t+1},\gamma_t}(s_t^*) - (1 - \gamma_{t+1})V^{\pi}_{\mathsf{P}_{t+1},\gamma_{t+1}}(s_t^*) \\
&\overset{(a)}{\leq} (1 - \gamma_t)V^{\pi}_{\mathsf{P}_{t+1},\gamma_t}(s_t^*) - (1 - \gamma_{t+1})V^{\pi}_{\mathsf{P}_{t+1},\gamma_{t+1}}(s_t^*) \\
&\leq \|(1 - \gamma_t)V^{\pi}_{\mathsf{P}_{t+1},\gamma_t} - (1 - \gamma_{t+1})V^{\pi}_{\mathsf{P}_{t+1},\gamma_{t+1}}\|_{\infty},
\end{aligned} \tag{68}$$

where $(a)$ is due to $(1 - \gamma_t)V^{\pi}_{\mathsf{P}_t,\gamma_t}(s_t^*) = \min_{\mathsf{P}}(1 - \gamma_t)V^{\pi}_{\mathsf{P},\gamma_t}(s_t^*) \leq (1 - \gamma_t)V^{\pi}_{\mathsf{P}_{t+1},\gamma_t}(s_t^*)$.

Now, according to Lemma 1,

$$(1 - \gamma_t)V^{\pi}_{\mathsf{P}_{t+1},\gamma_t} = g^{\pi}_{\mathsf{P}_{t+1}} + (1 - \gamma_t)h^{\pi}_{\mathsf{P}_{t+1}} + (1 - \gamma_t)f^{\pi}_{\mathsf{P}_{t+1}}(\gamma_t), \tag{69}$$

$$(1 - \gamma_{t+1})V^{\pi}_{\mathsf{P}_{t+1},\gamma_{t+1}} = g^{\pi}_{\mathsf{P}_{t+1}} + (1 - \gamma_{t+1})h^{\pi}_{\mathsf{P}_{t+1}} + (1 - \gamma_{t+1})f^{\pi}_{\mathsf{P}_{t+1}}(\gamma_{t+1}). \tag{70}$$

Hence, for any $\gamma_t > \delta$, eq. (68) can be further bounded as

$$\begin{aligned}
&\|(1 - \gamma_t)V^{\pi}_{\mathsf{P}_{t+1},\gamma_t} - (1 - \gamma_{t+1})V^{\pi}_{\mathsf{P}_{t+1},\gamma_{t+1}}\|_{\infty} \\
&= \|(\gamma_{t+1} - \gamma_t)h^{\pi}_{\mathsf{P}_{t+1}} + (1 - \gamma_t)f^{\pi}_{\mathsf{P}_{t+1}}(\gamma_t) - (1 - \gamma_{t+1})f^{\pi}_{\mathsf{P}_{t+1}}(\gamma_{t+1})\|_{\infty} \\
&\leq (\gamma_{t+1} - \gamma_t)\|h^{\pi}_{\mathsf{P}_{t+1}}\|_{\infty} + \|f^{\pi}_{\mathsf{P}_{t+1}}(\gamma_t) - f^{\pi}_{\mathsf{P}_{t+1}}(\gamma_{t+1})\|_{\infty} + \|\gamma_{t+1}f^{\pi}_{\mathsf{P}_{t+1}}(\gamma_{t+1}) - \gamma_t f^{\pi}_{\mathsf{P}_{t+1}}(\gamma_t)\|_{\infty} \\
&\overset{(a)}{\leq} h(\gamma_{t+1} - \gamma_t) + L(\gamma_{t+1} - \gamma_t) + \|\gamma_{t+1}f^{\pi}_{\mathsf{P}_{t+1}}(\gamma_{t+1}) - \gamma_t f^{\pi}_{\mathsf{P}_{t+1}}(\gamma_t)\|_{\infty} \\
&\leq h(\gamma_{t+1} - \gamma_t) + L(\gamma_{t+1} - \gamma_t) + \|\gamma_{t+1}f^{\pi}_{\mathsf{P}_{t+1}}(\gamma_{t+1}) - \gamma_{t+1}f^{\pi}_{\mathsf{P}_{t+1}}(\gamma_t)\|_{\infty} + \|\gamma_{t+1}f^{\pi}_{\mathsf{P}_{t+1}}(\gamma_t) - \gamma_t f^{\pi}_{\mathsf{P}_{t+1}}(\gamma_t)\|_{\infty} \\
&\leq h(\gamma_{t+1} - \gamma_t) + L(\gamma_{t+1} - \gamma_t) + \gamma_{t+1}\|f^{\pi}_{\mathsf{P}_{t+1}}(\gamma_{t+1}) - f^{\pi}_{\mathsf{P}_{t+1}}(\gamma_t)\|_{\infty} + \|f^{\pi}_{\mathsf{P}_{t+1}}(\gamma_t)\|_{\infty}(\gamma_{t+1} - \gamma_t) \\
&\overset{(b)}{\leq} (h + L + \gamma_{t+1}L + \sup_{\pi,\mathsf{P},\gamma}\|f^{\pi}_{\mathsf{P}}(\gamma)\|_{\infty})(\gamma_{t+1} - \gamma_t) \\
&\leq K(\gamma_{t+1} - \gamma_t),
\end{aligned} \tag{71}$$

where $(a)$ is from Lemma 5 for any $\gamma_t > \delta$, $c_f$ is defined in (51) and $K \triangleq h + 2L + c_f$ is a uniform constant; And $(b)$ is from Lemma 5.

Similarly, the inequality also holds for the case when $(1 - \gamma_t)V^{\pi}_{\mathcal{P},\gamma_t}(s_t^*) \leq (1 - \gamma_{t+1})V^{\pi}_{\mathcal{P},\gamma_{t+1}}(s_t^*)$. Thus we have that for any $t$ such that $\gamma_t > \delta$,

$$\Delta_{t+1} \leq K(\gamma_{t+1} - \gamma_t) + \gamma_t\Delta_t, \tag{72}$$

where $K$ is a uniform constant.

Following Lemma 8 from (Tewari and Bartlett 2007), we have that $\Delta_t \to 0$. Note that

$$\|V_t - g^{\pi}_{\mathcal{P}}\|_{\infty} \leq \|V_t - (1 - \gamma_t)V^{\pi}_{\mathcal{P},\gamma_t}\|_{\infty} + \|(1 - \gamma_t)V^{\pi}_{\mathcal{P},\gamma_t} - g^{\pi}_{\mathcal{P}}\|_{\infty} = \Delta_t + \|(1 - \gamma_t)V^{\pi}_{\mathcal{P},\gamma_t} - g^{\pi}_{\mathcal{P}}\|_{\infty}. \tag{73}$$

Together with Theorem 2, we further have that

$$\lim_{t\to\infty}\|V_t - g^{\pi}_{\mathcal{P}}\|_{\infty} = 0, \tag{74}$$

which completes the proof.

$\square$

# Proof of Theorem 4

Note that the optimal robust average-reward is defined as

$$g_{\mathcal{P}}^*(s) \triangleq \max_\pi g_{\mathcal{P}}^\pi(s). \tag{75}$$

We further define

$$V_{\mathcal{P},\gamma}^*(s) \triangleq \max_\pi V_{\mathcal{P},\gamma}^\pi(s). \tag{76}$$

**Theorem 12** (Restatement of Theorem 4). *$V_T$ generated by Algorithm 2 converges to the optimal robust average-reward $g_{\mathcal{P}}^*$ as $T \to \infty$.*

*Proof.* Firstly, from the uniform convergence in Theorem 2, it can be shown that

$$\lim_{t\to\infty} (1-\gamma_t)V_{\mathcal{P},\gamma_t}^* = g_{\mathcal{P}}^*. \tag{77}$$

We then show that for any $s \in \mathcal{S}$,

$$\begin{aligned}
&|V_{t+1}(s) - (1-\gamma_{t+1})V_{\mathcal{P},\gamma_{t+1}}^*(s)| \\
&\leq |V_{t+1}(s) - (1-\gamma_t)V_{\mathcal{P},\gamma_t}^*(s)| + |(1-\gamma_t)V_{\mathcal{P},\gamma_t}^*(s) - (1-\gamma_{t+1})V_{\mathcal{P},\gamma_{t+1}}^*(s)| \\
&\overset{(a)}{=} |(1-\gamma_t)V_{\mathcal{P},\gamma_t}^*(s) - (1-\gamma_{t+1})V_{\mathcal{P},\gamma_{t+1}}^*(s)| \\
&\quad + \left| \max_a \left( (1-\gamma_t)r(s,a) + \gamma_t\sigma_{\mathcal{P}_s^a}(V_t) \right) - \max_a \left( ((1-\gamma_t)r(s,a) + \gamma_t\sigma_{\mathcal{P}_s^a}((1-\gamma_t)V_{\mathcal{P},\gamma_t}^*)) \right) \right| \\
&\leq |(1-\gamma_t)V_{\mathcal{P},\gamma_t}^*(s) - (1-\gamma_{t+1})V_{\mathcal{P},\gamma_{t+1}}^*(s)| \\
&\quad + \max_a \left| (1-\gamma_t)r(s,a) + \gamma_t\sigma_{\mathcal{P}_s^a}(V_t) - ((1-\gamma_t)r(s,a) + \gamma_t\sigma_{\mathcal{P}_s^a}((1-\gamma_t)V_{\mathcal{P},\gamma_t}^*)) \right|,
\end{aligned} \tag{78}$$

where $(a)$ is because the optimal robust Bellman equation, and the last inequality is from the fact that $|\max_x f(x) - \max_x g(x)| \leq \max_x |f(x) - g(x)|$.

Hence eq. (78) can be further bounded as

$$\begin{aligned}
&|V_{t+1}(s) - (1-\gamma_{t+1})V_{\mathcal{P},\gamma_{t+1}}^*(s)| \\
&\leq |(1-\gamma_t)V_{\mathcal{P},\gamma_t}^*(s) - (1-\gamma_{t+1})V_{\mathcal{P},\gamma_{t+1}}^*(s)| + \gamma_t \max_a \left| \sigma_{\mathcal{P}_s^a}(V_t) - \sigma_{\mathcal{P}_s^a}((1-\gamma_t)V_{\mathcal{P},\gamma_t}^*) \right|.
\end{aligned} \tag{79}$$

If we denote by $\Delta_t \triangleq \|V_t - (1-\gamma_t)V_{\mathcal{P},\gamma_t}^*\|_\infty$, then

$$\Delta_{t+1} \leq \|(1-\gamma_t)V_{\mathcal{P},\gamma_t}^* - (1-\gamma_{t+1})V_{\mathcal{P},\gamma_{t+1}}^*\|_\infty + \gamma_t \max_{s,a} \left| \sigma_{\mathcal{P}_s^a}(V_t) - \sigma_{\mathcal{P}_s^a}((1-\gamma_t)V_{\mathcal{P},\gamma_t}^*) \right|. \tag{80}$$

Since the support function $\sigma_{\mathcal{P}_s^a}(V)$ is 1-Lipschitz, then it can be shown that for any $s, a$,

$$\left| \sigma_{\mathcal{P}_s^a}(V_t) - \sigma_{\mathcal{P}_s^a}((1-\gamma_t)V_{\mathcal{P},\gamma_t}^*) \right| \leq \|V_t - (1-\gamma_t)V_{\mathcal{P},\gamma_t}^*\|_\infty. \tag{81}$$

Hence

$$\Delta_{t+1} \leq \|(1-\gamma_t)V_{\mathcal{P},\gamma_t}^* - (1-\gamma_{t+1})V_{\mathcal{P},\gamma_{t+1}}^*\|_\infty + \gamma_t\Delta_t. \tag{82}$$

Similar to (71) in Theorem 3, we can show that

$$\|(1-\gamma_t)V_{\mathcal{P},\gamma_t}^* - (1-\gamma_{t+1})V_{\mathcal{P},\gamma_{t+1}}^*\|_\infty \leq K|\gamma_t - \gamma_{t+1}|, \tag{83}$$

and similar to Lemma 8 from (Tewari and Bartlett 2007),

$$\lim_{t\to\infty} \Delta_t = 0. \tag{84}$$

Moreover, note that

$$\|V_t - g_{\mathcal{P}}^*\|_\infty \leq \|V_t - (1-\gamma_t)V_{\mathcal{P},\gamma_t}^*\|_\infty + \|(1-\gamma_t)V_{\mathcal{P},\gamma_t}^* - g_{\mathcal{P}}^*\|_\infty = \Delta_t + \|(1-\gamma_t)V_{\mathcal{P},\gamma_t}^* - g_{\mathcal{P}}^*\|_\infty, \tag{85}$$

which together with eq. (77) implies that

$$\|V_t - g_{\mathcal{P}}^*\|_\infty \to 0, \tag{86}$$

and hence it completes the proof.

$\square$

# Proof of Theorem 5

We denote the set of all stationary deterministic polices by $\Pi^D$ in this section.

**Theorem 13** (Restatement of Theorem 5). *There exists a Blackwell optimal policy set $\Pi^*$, i.e., there exists $0 < \delta < 1$ and a finite policy set $\Pi^* \subseteq \Pi^D$, such that for any $\gamma > \delta$, the optimal robust policy for robust discounted value function $V^*_{\mathcal{P},\gamma}$ belongs to $\Pi^*$, i.e.,*

$$V^*_{\mathcal{P},\gamma} = V^{\pi^*}_{\mathcal{P},\gamma}, \tag{87}$$

*for some $\pi^* \in \Pi^*$. Moreover, when $\arg\max_{\pi \in \Pi^D} g^\pi_{\mathcal{P}}$ is a singleton, there exists a unique Blackwell optimal policy.*

*Proof.* According to Theorem 7[4], there exists $\pi^* \in \Pi^D$ such that

$$g^*_{\mathcal{P}} = g^{\pi^*}_{\mathcal{P}}. \tag{88}$$

Assume the robust average-reward of all deterministic policies are sorted in a descending order:

$$g^*_{\mathcal{P}} = g^{\pi^*_1}_{\mathcal{P}} = g^{\pi^*_2}_{\mathcal{P}} = ... = g^{\pi^*_m}_{\mathcal{P}} > g^{\pi_1}_{\mathcal{P}} \geq ... \geq g^{\pi_n}_{\mathcal{P}} \tag{89}$$

for all $\pi^*_i, \pi_i \in \Pi^D$, and we define $\Pi^* = \{\pi^*_i : i = 1, ..., m\}$. Denote by $d = g^{\pi^*_i}_{\mathcal{P}} - g^{\pi_1}_{\mathcal{P}}$.

From Theorem 2, we know that for any $\pi \in \Pi^D$,

$$\lim_{\gamma \to 1} (1 - \gamma) V^\pi_{\mathcal{P},\gamma} = g^\pi_{\mathcal{P}}. \tag{90}$$

Because the set $\Pi^D$ is finite, for any $\epsilon < \frac{d}{2}$, there exists $\delta' < 1$, such that for any $\gamma > \delta'$, $\pi^*_i$ and $\pi_j$,

$$|(1 - \gamma) V^{\pi^*_i}_{\mathcal{P},\gamma} - g^*_{\mathcal{P}}| < \epsilon, \tag{91}$$

$$|(1 - \gamma) V^{\pi_j}_{\mathcal{P},\gamma} - g^{\pi_j}_{\mathcal{P}}| < \epsilon. \tag{92}$$

It hence implies that

$$(1 - \gamma) V^{\pi^*_i}_{\mathcal{P},\gamma} \geq (d - 2\epsilon) + (1 - \gamma) V^{\pi_j}_{\mathcal{P},\gamma} > (1 - \gamma) V^{\pi_j}_{\mathcal{P},\gamma}, \tag{93}$$

and

$$V^{\pi^*_i}_{\mathcal{P},\gamma} > V^{\pi_j}_{\mathcal{P},\gamma}. \tag{94}$$

Note that from Theorem 3.1 in (Iyengar 2005), i.e., $\max_{\pi \in \Pi^D} V^\pi_{\mathcal{P},\gamma} = V^*_{\mathcal{P},\gamma}$, we have that for any $\gamma$, there exists a deterministic policy $\pi \in \Pi^D$, such that $V^*_{\mathcal{P},\gamma} = V^\pi_{\mathcal{P},\gamma}$. Together with (94), it implies that all the possible optimal robust polices of $V^\pi_{\mathcal{P},\gamma}$ belong to $\{\pi^*_1, ...\pi^*_m\}$, i.e., the set $\Pi^*$. Hence, there exists $\pi^*_j \in \Pi^*$, such that

$$V^{\pi^*_j}_{\mathcal{P},\gamma} = \max_{\pi \in \Pi^D} V^\pi_{\mathcal{P},\gamma} = V^*_{\mathcal{P},\gamma}. \tag{95}$$

For the second part, when the optimal robust policy of robust average-reward is unique, i.e., $\Pi^* = \{\pi^*\}$. Then from the results above, there exists $\delta'$, such that for any $\gamma > \delta'$, $V^{\pi^*}_{\mathcal{P},\gamma} > V^\pi_{\mathcal{P},\gamma}$ for any $\pi^* \neq \pi \in \Pi^D$, and hence $\pi^*$ is the optimal policy for discounted robust MDPs, which is the unique Blackwell optimal policy.

$\square$

# Proof of Results for Direct Approach

Recall that

$$V^\pi_{\mathcal{P}}(s) \triangleq \min_{\kappa \in \bigotimes_{t \geq 0} \mathcal{P}} \mathbb{E}_{\kappa,\pi} \left[ \sum_{t=0}^{\infty} (r_t - g^\pi_{\mathcal{P}}) \big| S_0 = s \right], \tag{96}$$

where

$$g^\pi_{\mathcal{P}} = \min_{\kappa \in \bigotimes_{t \geq 0} \mathcal{P}} \lim_{n \to \infty} \mathbb{E}_{\kappa,\pi} \left[ \frac{1}{n} \sum_{t=0}^{n-1} r_t \big| S_0 = s \right]. \tag{97}$$

We first show that the robust relative function is always finite.

---

[4]The proof of Theorem 7 is independent of theorem 5 and does not relay on the results to be showed here.

**Lemma 7.** *For any $\pi$, $V_{\mathcal{P}}^{\pi}$ is finite.*

*Proof.* According to Theorem 9, $V_{\mathcal{P}}^{\pi} = \min_{\mathsf{P} \in \mathcal{P}} V_{\mathsf{P}}^{\pi} = \min_{\mathsf{P} \in \mathcal{P}} \mathbb{E}_{\mathsf{P},\pi} \left[ \sum_{t=0}^{\infty} (r_t - g_{\mathcal{P}}^{\pi}) \right]$. Note that $V_{\mathcal{P}}^{\pi}$ can be rewritten as

$$
\begin{aligned}
V_{\mathcal{P}}^{\pi} &= \min_{\mathsf{P} \in \mathcal{P}} \mathbb{E}_{\mathsf{P},\pi} \left[ \sum_{t=0}^{\infty} (r_t - g_{\mathcal{P}}^{\pi}) \right] \\
&= \min_{\mathsf{P} \in \mathcal{P}} \mathbb{E}_{\mathsf{P},\pi} \left[ \lim_{n \to \infty} \sum_{t=0}^{n} (r_t - g_{\mathcal{P}}^{\pi}) \right] \\
&= \min_{\mathsf{P} \in \mathcal{P}} \mathbb{E}_{\mathsf{P},\pi} \left[ \lim_{n \to \infty} \sum_{t=0}^{n} (r_t - g_{\mathsf{P}}^{\pi} + g_{\mathsf{P}}^{\pi} - g_{\mathcal{P}}^{\pi}) \right] \\
&= \min_{\mathsf{P} \in \mathcal{P}} \mathbb{E}_{\mathsf{P},\pi} \left[ \lim_{n \to \infty} (R_n - n g_{\mathsf{P}}^{\pi} + n g_{\mathsf{P}}^{\pi} - n g_{\mathcal{P}}^{\pi}) \right],
\end{aligned}
\tag{98}
$$

where $R_n = \sum_{t=0}^{n} r_t$. Note that for any $\mathsf{P} \in \mathcal{P}$ and $n$, $n g_{\mathsf{P}}^{\pi} \geq n g_{\mathcal{P}}^{\pi}$, hence

$$
\lim_{n \to \infty} (R_n - n g_{\mathsf{P}}^{\pi} + n g_{\mathsf{P}}^{\pi} - n g_{\mathcal{P}}^{\pi}) \geq \lim_{n \to \infty} (R_n - n g_{\mathsf{P}}^{\pi}),
\tag{99}
$$

and thus the lower bound of $V_{\mathcal{P}}^{\pi}$ can be derived as follows,

$$
\begin{aligned}
V_{\mathcal{P}}^{\pi} &\geq \min_{\mathsf{P} \in \mathcal{P}} \mathbb{E}_{\mathsf{P},\pi} \left[ \sum_{t=0}^{\infty} (r_t - g_{\mathsf{P}}^{\pi}) \right] \\
&= \min_{\mathsf{P} \in \mathcal{P}} V_{\mathsf{P}}^{\pi} \\
&= \min_{\mathsf{P} \in \mathcal{P}} H_{\mathsf{P}}^{\pi} r_{\pi}.
\end{aligned}
\tag{100}
$$

which is finite due to the fact that $H_{\mathsf{P}}^{\pi}$ is continuous on the compact set $\mathcal{P}$.

From Theorem 9, we denote the stationary worst-case transition kernel of $g_{\mathcal{P}}^{\pi}$ by $\mathsf{P}_g$. Then the upper bound of $V_{\mathcal{P}}^{\pi}$ can be bounded by noting that

$$
\begin{aligned}
V_{\mathcal{P}}^{\pi} &= \min_{\mathsf{P} \in \mathcal{P}} \mathbb{E}_{\mathsf{P},\pi} \left[ \sum_{t=0}^{\infty} (r_t - g_{\mathsf{P}_g}^{\pi}) \right] \\
&\leq \mathbb{E}_{\mathsf{P}_g,\pi} \left[ \sum_{t=0}^{\infty} (r_t - g_{\mathsf{P}_g}^{\pi}) \right] \\
&= V_{\mathsf{P}_g}^{\pi},
\end{aligned}
\tag{101}
$$

which is also finite and $\mathsf{P}_g$ denotes the worst-case transition kernel of $g_{\mathcal{P}}^{\pi}$. Hence we show that $V_{\mathcal{P}}^{\pi}$ is finite for any $\pi$ and hence complete the proof. $\square$

After showing that the robust relative value function is well-defined, we show the following robust Bellman equation for average-reward robust MDPs.

**Theorem 14** (Restatement of Theorem 6). *For any $s$ and $\pi$, $(V_{\mathcal{P}}^{\pi}, g_{\mathcal{P}}^{\pi})$ is a solution to the following robust Bellman equation:*

$$
V(s) + g = \sum_{a} \pi(a|s) \left( r(s,a) + \sigma_{\mathcal{P}_s^a}(V) \right).
\tag{102}
$$

*Proof.* From the definition,

$$
V_{\mathcal{P}}^{\pi}(s) = \min_{\kappa \in \bigotimes_{t \geq 0} \mathcal{P}} \mathbb{E}_{\kappa,\pi} \left[ \sum_{t=0}^{\infty} (r_t - g_{\mathcal{P}}^{\pi}) \big| S_0 = s \right],
\tag{103}
$$

hence

$$
V_{\mathcal{P}}^{\pi}(s) = \min_{\kappa \in \bigotimes_{t \geq 0} \mathcal{P}} \mathbb{E}_{\kappa,\pi} \left[ \sum_{t=0}^{\infty} (r_t - g_{\mathcal{P}}^{\pi}) \big| S_0 = s \right]
$$

$$= \min_{\kappa \in \bigotimes_{t \geq 0} \mathcal{P}} \mathbb{E}_{\kappa,\pi} \left[ (r_0 - g_{\mathcal{P}}^\pi) + \sum_{t=1}^{\infty} (r_t - g_{\mathcal{P}}^\pi) \big| S_0 = s \right]$$

$$= \min_{\kappa \in \bigotimes_{t \geq 0} \mathcal{P}} \left\{ \sum_a \pi(a|s) r(s,a) - g_{\mathcal{P}}^\pi + \mathbb{E}_{\kappa,\pi} \left[ \sum_{t=1}^{\infty} (r_t - g_{\mathcal{P}}^\pi) \big| S_0 = s \right] \right\}$$

$$= \sum_a \pi(a|s) \left( r(s,a) - g_{\mathcal{P}}^\pi \right) + \min_{\kappa \in \bigotimes_{t \geq 0} \mathcal{P}} \left\{ \sum_{a,s'} \pi(a|s) \mathsf{P}_{s,s'}^a \mathbb{E}_{\kappa,\pi} \left[ \sum_{t=1}^{\infty} (r_t - g_{\mathcal{P}}^\pi) \big| S_1 = s' \right] \right\}$$

$$= \sum_a \pi(a|s) \left( r(s,a) - g_{\mathcal{P}}^\pi \right) + \min_{\mathsf{P}_0 \in \mathcal{P}} \min_{\kappa = (\mathsf{P}_1,\dots) \in \bigotimes_{t \geq 1} \mathcal{P}} \left\{ \sum_{a,s'} \pi(a|s)(\mathsf{P}_0)_{s,s'}^a \mathbb{E}_{\kappa,\pi} \left[ \sum_{t=1}^{\infty} (r_t - g_{\mathcal{P}}^\pi) \big| S_1 = s' \right] \right\}$$

$$= \sum_a \pi(a|s) \left( r(s,a) - g_{\mathcal{P}}^\pi \right) + \min_{\mathsf{P}_0 \in \mathcal{P}} \left\{ \sum_{a,s'} \pi(a|s)(\mathsf{P}_0)_{s,s'}^a \min_{\kappa = (\mathsf{P}_1,\dots) \in \bigotimes_{t \geq 1} \mathcal{P}} \left\{ \mathbb{E}_{\kappa,\pi} \left[ \sum_{t=1}^{\infty} (r_t - g_{\mathcal{P}}^\pi) \big| S_1 = s' \right] \right\} \right\}$$

$$= \sum_a \pi(a|s) \left( r(s,a) - g_{\mathcal{P}}^\pi \right) + \sum_a \pi(a|s) \sum_{s'} \min_{p_{s,s'}^a \in \mathcal{P}_s^a} p_{s,s'}^a V_{\mathcal{P}}^\pi(s')$$

$$= \sum_a \pi(a|s) \left( r(s,a) - g_{\mathcal{P}}^\pi \right) + \sum_a \pi(a|s) \sigma_{\mathcal{P}_s^a}(V_{\mathcal{P}}^\pi)$$

$$= \sum_a \pi(a|s) \left( r(s,a) - g_{\mathcal{P}}^\pi + \sigma_{\mathcal{P}_s^a}(V_{\mathcal{P}}^\pi) \right). \tag{104}$$

This hence completes the proof. $\qquad\square$

**Theorem 15.** *[Restatement of Theorem 7, Part 1] For any $(g, V)$ that is a solution to $\max_a \left\{ r(s,a) - g + \sigma_{\mathcal{P}_s^a}(V) - V(s) \right\} = 0, \forall s$, then $g = g_{\mathcal{P}}^*$.*

*Proof.* In this proof, for two vectors $v, w \in \mathbb{R}^n$, $v \geq w$ denotes that $v(s) \geq w(s)$ entry-wise.

Let $B(g, V)(s) \triangleq \max_a \left\{ r(s,a) - g + \sigma_{\mathcal{P}_s^a}(V) - V(s) \right\}$. Since $(g, V)$ is a solution to (13), hence for any $a \in \mathcal{A}$ and any $s \in \mathcal{S}$,

$$r(s,a) - g + \sigma_{\mathcal{P}_s^a}(V) - V(s) \leq 0, \tag{105}$$

from which it follows that for any policy $\pi$,

$$g(s) \geq r_\pi(s) + \sum_a \pi(a|s)\sigma_{\mathcal{P}_s^a}(V) - V(s) \triangleq r_\pi(s) + \sum_a \pi(a|s)(p_s^a)^\top V - V(s), \tag{106}$$

where $r_\pi(s) \triangleq \sum_a \pi(a|s)r(s,a)$, $p_s^a \triangleq \arg\min_{p \in \mathcal{P}_s^a} p^\top V$, and $\mathsf{P}_V = \{p_s^a : s \in \mathcal{S}, a \in \mathcal{A}\}$. We also denotes the state transition matrix induced by $\pi$ and $\mathsf{P}_V$ by $\mathsf{P}_V^\pi$.

Using these notations, and rewrite eq. (106), we have that

$$g\mathbb{1} \geq r_\pi + (\mathsf{P}_V^\pi - I)V. \tag{107}$$

Since the inequality in eq. (107) holds entry-wise, all entries of $\mathsf{P}_V^\pi$ are positive, then by multiplying both sides of eq. (107) by $\mathsf{P}_V^\pi$, we have that

$$g\mathbb{1} = g\mathsf{P}_V^\pi \mathbb{1} \geq \mathsf{P}_V^\pi r_\pi + \mathsf{P}_V^\pi (\mathsf{P}_V^\pi - I)V. \tag{108}$$

Multiplying the both sides of eq. (108) by $\mathsf{P}_V^\pi$, and repeatedly doing that, we have that

$$g\mathbb{1} \geq (\mathsf{P}_V^\pi)^2 r_\pi + (\mathsf{P}_V^\pi)^2 (\mathsf{P}_V^\pi - I)V, \tag{109}$$

$$\vdots \qquad\qquad \vdots \tag{110}$$

$$g\mathbb{1} \geq (\mathsf{P}_V^\pi)^{n-1} r_\pi + (\mathsf{P}_V^\pi)^{n-1}(\mathsf{P}_V^\pi - I)V. \tag{111}$$

Summing up these inequalities from eq. (107) to eq. (111), we have that

$$ng\mathbb{1} \geq (I + \mathsf{P}_V^\pi + \dots + (\mathsf{P}_V^\pi)^{n-1})r_\pi + (I + \mathsf{P}_V^\pi + \dots + (\mathsf{P}_V^\pi)^{n-1})(\mathsf{P}_V^\pi - I)V, \tag{112}$$

and from which, it follows that

$$g\mathbb{1} \geq \frac{1}{n}(I + \mathsf{P}_V^\pi + ... + (\mathsf{P}_V^\pi)^{n-1})r_\pi + \frac{1}{n}(I + \mathsf{P}_V^\pi + ... + (\mathsf{P}_V^\pi)^{n-1})(\mathsf{P}_V^\pi - I)V$$

$$= \frac{1}{n}(I + \mathsf{P}_V^\pi + ... + (\mathsf{P}_V^\pi)^{n-1})r_\pi + \frac{1}{n}((\mathsf{P}_V^\pi)^n - I)V. \tag{113}$$

It can be easily verified that $\lim_{n\to\infty} \frac{1}{n}((\mathsf{P}_V^\pi)^n - I)V = 0$, and hence it implies that

$$g\mathbb{1} \geq \lim_{n\to\infty} \frac{1}{n}(I + \mathsf{P}_V^\pi + ... + (\mathsf{P}_V^\pi)^{n-1})r_\pi$$

$$= \lim_{n\to\infty} \frac{1}{n}\mathbb{E}_{\mathsf{P}_V^\pi,\pi}\left[\sum_{t=0}^n r_t\right]$$

$$= g_{\mathsf{P}_V^\pi}^\pi \mathbb{1}$$

$$\geq g_{\mathcal{P}}^\pi \mathbb{1}. \tag{114}$$

Since eq. (114) holds for any policy $\pi$, it follows that $g \geq g_{\mathcal{P}}^*$. On the other hand, since $B(g, V) = 0$, there exists a policy $\tau$ such that

$$g\mathbb{1} = r_\tau + (\mathsf{P}_V^\tau - I)V, \tag{115}$$

where $r_\tau, \mathsf{P}_V^\tau$ are similarly defined as for $\pi$. From Theorem 9, there exists a stationary transition kernel $\mathsf{P}_{\text{ave}}^\tau$ such that $g_{\mathcal{P}}^\tau = g_{\mathsf{P}_{\text{ave}}^\tau}^\tau$. We denote the state transition matrix induced by $\tau$ and $\mathsf{P}_{\text{ave}}^\tau$ by $\mathsf{P}^\tau$. Then because $\mathsf{P}_V^\tau$ is the worst-case transition of $V$, it follows that

$$\mathsf{P}_V^\tau V \leq \mathsf{P}^\tau V. \tag{116}$$

Thus

$$g\mathbb{1} \leq r_\tau + (\mathsf{P}^\tau - I)V. \tag{117}$$

Similarly, we have that

$$g\mathbb{1} \leq (\mathsf{P}^\tau)^{j-1}r_\tau + (\mathsf{P}^\tau)^{j-1}(\mathsf{P}^\tau - I)V, \tag{118}$$

for $j = 2, ..., n$. Summing these inequalities together we have that

$$ng\mathbb{1} \leq (I + \mathsf{P}^\tau + ... + (\mathsf{P}^\tau)^{n-1})r_\tau + (I + \mathsf{P}^\tau + ... + (\mathsf{P}^\tau)^{n-1})(\mathsf{P}^\tau)^{n-1}(\mathsf{P}^\tau - I)V$$

$$= (I + \mathsf{P}^\tau + ... + (\mathsf{P}^\tau)^{n-1})r_\tau + ((\mathsf{P}^\tau)^n - I)V. \tag{119}$$

Hence

$$g\mathbb{1} \leq \lim_{n\to\infty} \frac{1}{n}\mathbb{E}_{\mathsf{P}_{\text{ave}}^\tau,\tau}\left[\sum_{t=0}^n r_t\right] = g_{\mathsf{P}_{\text{ave}}^\tau}^\tau \mathbb{1} = g_{\mathcal{P}}^\tau \mathbb{1} \leq g_{\mathcal{P}}^* \mathbb{1}. \tag{120}$$

Thus $g = g_{\mathcal{P}}^*$, and this concludes the proof. $\square$

**Theorem 16** (Restatement of Theorem 7, Part 2). *For any $(g, V)$ that is a solution to*

$$\max_a \left\{r(s, a) - g + \sigma_{\mathcal{P}_s^a}(V) - V(s)\right\} = 0, \forall s, \tag{121}$$

*if we set*

$$\pi^*(s) = \arg\max_a \left\{r(s, a) + \sigma_{\mathcal{P}_s^a}(V)\right\} \tag{122}$$

*for any $s \in \mathcal{S}$, then $\pi^*$ is an optimal robust policy.*

*Proof.* Note that for any stationary policy $\pi$, we denote by $\sigma_{\mathcal{P}^\pi}(V) \triangleq (\sum_a \pi(a|s_1)\sigma_{\mathcal{P}_{s_1}^a}(V), ..., \sum_a \pi(a|s_{|\mathcal{S}|})\sigma_{\mathcal{P}_{s_{|\mathcal{S}|}}^a}(V))$ being a vector in $\mathbb{R}^{|\mathcal{S}|}$. Then eq. (14) is equivalent to

$$r_{\pi^*} + \sigma_{\mathcal{P}^{\pi^*}}(V) = \max_\pi \left\{r_\pi + \sigma_{\mathcal{P}^\pi}(V)\right\}. \tag{123}$$

Hence,

$$r_{\pi^*} - g + \sigma_{\mathcal{P}^{\pi^*}}(V) - V = \max_\pi \left\{r_\pi - g + \sigma_{\mathcal{P}^\pi}(V) - V\right\}. \tag{124}$$

Since $(g, V)$ is a solution to (13), it follows that

$$r_{\pi^*} - g + \sigma_{\mathcal{P}^{\pi^*}}(V) - V = 0. \tag{125}$$

According to the robust Bellman equation eq. (12), $(g_{\mathcal{P}}^{\pi^*}, V_{\mathcal{P}}^{\pi^*})$ is a solution to eq. (125). Thus from Theorem 15, $g_{\mathcal{P}}^{\pi^*} = g_{\mathcal{P}}^*$, and hence $\pi^*$ is an optimal robust policy. $\square$

**Theorem 17** (Restatement of Theorem 8). $(w_T, V_t)$ *in Algorithm 3 converges to a solution of eq.* (13).

*Proof.* We first denote the update operator as

$$Lv(s) \triangleq \max_a(r(s,a) + \sigma_{\mathcal{P}_s^a}(v)). \tag{126}$$

Now, consider $sp(Lv - Lu)$. Denote by $\acute{s} \triangleq \arg\max_s(Lv(s) - Lu(s))$ and $\grave{s} \triangleq \arg\min_s(Lv(s) - Lu(s))$. Also denote by $a_v \triangleq \arg\max_a(r(\acute{s},a) + \sigma_{\mathcal{P}_{\acute{s}}^a}(v))$ and $a_u \triangleq \arg\max_a(r(\acute{s},a) + \sigma_{\mathcal{P}_{\acute{s}}^a}(u))$ Then

$$
\begin{aligned}
Lv(\acute{s}) - Lu(\acute{s}) &= \max_a(r(\acute{s},a) + \sigma_{\mathcal{P}_{\acute{s}}^a}(v)) - \max_a(r(\acute{s},a) + \sigma_{\mathcal{P}_{\acute{s}}^a}(u)) \\
&\triangleq r(\acute{s},a_v) + \sigma_{\mathcal{P}_{\acute{s}}^{a_v}}(v) - (r(\acute{s},a_u) + \sigma_{\mathcal{P}_{\acute{s}}^{a_u}}(u)) \\
&\leq r(\acute{s},a_v) + \sigma_{\mathcal{P}_{\acute{s}}^{a_v}}(v) - (r(\acute{s},a_v) + \sigma_{\mathcal{P}_{\acute{s}}^{a_v}}(u)) \\
&= \sigma_{\mathcal{P}_{\acute{s}}^{a_v}}(v) - \sigma_{\mathcal{P}_{\acute{s}}^{a_v}}(u) \\
&\triangleq (p_{\acute{s}}^{a_v,v})^\top v - (p_{\acute{s}}^{a_v,u})^\top u,
\end{aligned}
\tag{127}
$$

where $p_{\acute{s}}^{a_v,v} = \arg\min_{p\in\mathcal{P}_{\acute{s}}^{a_v}} p^\top v$ and $p_{\acute{s}}^{a_v,u} = \arg\min_{p\in\mathcal{P}_{\acute{s}}^{a_v}} p^\top u$. Thus eq. (127) can be further bounded as

$$
\begin{aligned}
Lv(\acute{s}) - Lu(\acute{s}) \\
\leq (p_{\acute{s}}^{a_v,v})^\top v - (p_{\acute{s}}^{a_v,u})^\top u \\
\leq (p_{\acute{s}}^{a_v,u})^\top (v - u).
\end{aligned}
\tag{128}
$$

Similarly,

$$Lv(\grave{s}) - Lu(\grave{s}) \geq (p_{\grave{s}}^{a_u,v})^\top(v-u). \tag{129}$$

Thus

$$sp(Lv - Lu) \leq (p_{\acute{s}}^{a_v,u})^\top(v-u) - (p_{\grave{s}}^{a_u,v})^\top(v-u). \tag{130}$$

Now denote by $v - u \triangleq (x_1, x_2, ..., x_n)$, $p_{\acute{s}}^{a_v,u} = (p_1, ..., p_n)$ and $p_{\grave{s}}^{a_u,v} = (q_1, ..., q_n)$. Further denote by $b_i \triangleq \min\{p_i, q_i\}$ Then

$$
\begin{aligned}
&\sum_{i=1}^n p_i x_i - \sum_{i=1}^n q_i x_i \\
&= \sum_{i=1}^n(p_i - b_i)x_i - \sum_{i=1}^n(q_i - b_i)x_i \\
&\leq \sum_{i=1}^n(p_i - b_i)\max\{x_i\} - \sum_{i=1}^n(q_i - b_i)\min\{x_i\} \\
&= \sum_{i=1}^n(p_i - b_i)sp(x) + \left(\sum_{i=1}^n(p_i - b_i) - \sum_{i=1}^n(q_i - b_i)\right)\min\{x_i\} \\
&= \left(1 - \sum_{i=1}^n b_i\right)sp(x).
\end{aligned}
\tag{131}
$$

Thus we showed that

$$sp(Lv - Lu) \leq \left(1 - \sum_{i=1}^n b_i\right)sp(v - u). \tag{132}$$

Now from Assumption 2, and following Theorem 8.5.3 from (Puterman 1994), it can be shown that there exists $1 > \lambda > 0$, such that for any $a, u, v$,

$$\sum_{i=1}^n b_i \geq \lambda. \tag{133}$$

Further, following Theorem 8.5.2 in (Puterman 1994), it can be shown that $L$ is a $J$-step contraction operator for some integer $J$, i.e.,

$$sp(L^J v - L^J u) \leq (1 - \lambda)sp(v - u). \tag{134}$$

Then, it can be shown that the relative value iteration converges to a solution of the optimal equation similar to the relative value iteration for non-robust MDPs under the average-reward criterion (Theorem 8.5.7 in (Puterman 1994), Section 1.6.4 in (Sigaud and Buffet 2013)), and hence $(w_t, V_t)$ converges to a solution to eq. (13) as $\epsilon \to 0$. $\square$