
Robust Reinforcement Learning using Least Squares Policy Iteration with Provable Performance Guarantees

Kishan Panaganti¹ Dileep Kalathil¹

Abstract

This paper addresses the problem of model-free reinforcement learning for Robust Markov Decision Process (RMDP) with large state spaces. The goal of the RMDP framework is to find a policy that is robust against the parameter uncertainties due to the mismatch between the simulator model and real-world settings. We first propose the Robust Least Squares Policy Evaluation algorithm, which is a multi-step online model-free learning algorithm for policy evaluation. We prove the convergence of this algorithm using stochastic approximation techniques. We then propose Robust Least Squares Policy Iteration (RLSPI) algorithm for learning the optimal robust policy. We also give a general weighted Euclidean norm bound on the error (closeness to optimality) of the resulting policy. Finally, we demonstrate the performance of our RLSPI algorithm on some standard benchmark problems.

1. Introduction

Model-free Reinforcement Learning (RL) algorithms typically learn a policy by training on a simulator. In the RL literature, it is nominally assumed that the testing environment is identical to the training environment (simulator model). However, in reality, the parameters of the simulator model can be different from the real-world setting. This can be due to the approximation errors incurred while modeling, due to the changes in the real-world parameters over time, and can even be due to possible adversarial disturbances in the real-world. For example, in many robotics applications, the standard simulator parameter settings (mass, friction, wind conditions, sensor noise, action delays) can be different from that of the actual robot in the real-world.

¹Department of Electrical and Computer Engineering, Texas A&M University, College Station, United States. Correspondence to: Kishan Panaganti <kpb@tamu.edu>, Dileep Kalathil <dileep.kalathil@tamu.edu>.

This mismatch between the training and testing environment parameters can significantly degrade the real-world performance of the model-free learning algorithms trained on a simulator model.

The RMDP framework (Iyengar, 2005; Nilim & El Ghaoui, 2005) addresses the *planning* problem of computing the optimal policy that is robust against parameter uncertainties that cause the mismatch between the training and testing environment parameters. The RMDP problem has been analyzed extensively in the tabular case (Iyengar, 2005; Nilim & El Ghaoui, 2005; Wiesemann et al., 2013; Xu & Mannor, 2010; Yu & Xu, 2015) and under the linear function approximation (Tamar et al., 2014). Algorithms for learning the optimal robust policy with provable guarantees have been proposed, both in the model-free (Roy et al., 2017) and model-based (Lim et al., 2013) reinforcement learning settings. However, the theoretical guarantees from these works are limited to the tabular RMDP settings. Learning policies for problems with large state spaces is computationally challenging. RL algorithms typically overcome this issue by using function approximation architectures, such as linear basis functions (Lagoudakis & Parr, 2003), reproducing kernel Hilbert spaces (RKHS) (Yang & Wang, 2020) and deep neural networks (Lillicrap et al., 2016). Recently, robust reinforcement learning problem has been addressed using deep RL methods (Pinto et al., 2017; Mankowitz et al., 2020; Zhang et al., 2020; Derman et al., 2018; Vinitzky et al., 2020). However, these works are empirical in nature and do not provide any theoretical guarantees for the learned policies. The problem of learning optimal robust policies with provable performance guarantees for RMDPs with large state spaces has not been well studied in the literature.

In this paper, we address the problem of learning a policy that is provably robust against the parameter uncertainties for RMDPs with large state spaces. In particular, we propose an online model-free reinforcement learning algorithm with linear function approximation for learning the optimal robust policy, and provide theoretical guarantees on the performance of the learned policy. Our choice of linear function approximation is motivated by its analytical tractability while providing the scaling to large state spaces. Indeed, linear function approximation based approaches have been successful in providing algorithms with provable

guarantees for many challenging problems in RL, including online model-free exploration (Jin et al., 2020; Yang & Wang, 2020), imitation learning (Abbeel & Ng, 2004; Arora et al., 2020), meta reinforcement learning (Wang et al., 2020; Kong et al., 2020), and offline reinforcement learning (Wang et al., 2021; Duan et al., 2020). Robust RL is much more challenging than the standard (non-robust) RL problems due to the inherent nonlinearity associated with the robust dynamic programming. We overcome this issue by a cleverly designed approximate dynamic programming approach. We then propose a model-free robust policy iteration using this approach with provable guarantees. Our algorithmic and technical contributions are as follows:

(i) Robust Least Squares Policy Evaluation (RLSPE(λ)) algorithm: A learning-based policy iteration algorithm needs to learn the value of a policy for performing (greedy) policy improvement. For this, we first propose RLSPE(λ) algorithm, a multi-step, online, model-free policy evaluation algorithm with linear function approximation. This can be thought as the robust version of classical least squares based RL algorithms for policy evaluation, like LSTD(λ) and LSPE(λ). We prove the convergence of this algorithm using stochastic approximation techniques, and also characterize its approximation error due to the linear architecture.

(ii) Robust Least Squares Policy Iteration (RLSPI) algorithm: We propose the RLSPI algorithm for learning the optimal robust policy. We also give a general L_2 -norm bound on the error (closeness to optimality) of the resulting policy at any iterate of the algorithm. To the best of our knowledge, this is the first work that presents a learning based policy iteration algorithm for robust reinforcement learning with such provable guarantees.

(iii) Finally, we demonstrate the performance of the RLSPI algorithm on various standard RL test environments.

1.1. Related Work

RMDP formulation to address the parameter uncertainty problem was first proposed by (Iyengar, 2005) and (Nilim & El Ghaoui, 2005). (Iyengar, 2005) showed that the optimal robust value function and policy can be computed using the robust counterparts of the standard value iteration and policy iteration. To tackle the parameter uncertainty problem, other works considered distributionally robust setting (Xu & Mannor, 2010), modified policy iteration (Kaufman & Schaefer, 2013), and more general uncertainty set (Wiesemann et al., 2013). We note that the focus of these works were mainly on the planning problem in the tabular setting. Linear function approximation method to solve large RMDPs was proposed in (Tamar et al., 2014). Though this work suggests a sampling based approach, a general model-free learning algorithm and analysis was not included. (Roy et al., 2017) proposed the robust versions of the classical model-free rein-

forcement learning algorithms such as Q-learning, SARSA, and TD-learning in the tabular setting. They also proposed function approximation based algorithms for the policy evaluation. However, this work does not have a policy iteration algorithm with provable guarantees for learning the optimal robust policy. (Derman et al., 2018) introduced soft-robust actor-critic algorithms using neural networks, but does not provide any global convergence guarantees for the learned policy. (Tessler et al., 2019) proposed a min-max game framework to address the robust learning problem focusing on the tabular setting. (Lim & Autef, 2019) proposed a kernel-based RL algorithm for finding the robust value function in a batch learning setting. (Mankowitz et al., 2020) employed an entropy-regularized policy optimization algorithm for continuous control using neural network, but does not provide any provable guarantees for the learned policy.

Our work differs from the above in two significant ways. Firstly, we develop a new multi-step model-free reinforcement learning algorithm, RLSPE(λ), for policy evaluation. Extending the classical least squares based policy evaluation algorithms, like LSPE(λ) and LSTD(λ) (Bertsekas & Ioffe, 1996; Nedić & Bertsekas, 2003; Bertsekas, 2012), to the robust case is very challenging due to the nonlinearity of the robust TD(λ) operator. We overcome this issue by a cleverly defined approximate robust TD(λ) operator that is amenable to online learning using least squares approaches. Also, as pointed out in (Bertsekas, 2011), convergence analysis of least squares style algorithms for RL is different from that of the standard temporal difference (TD) algorithm. Secondly, we develop a new robust policy iteration algorithm with provable guarantees on the performance of the policy at any iterate. In particular, we give a general weighted Euclidean norm bound on the error of the resulting policy. While similar results are available for the non-robust settings, this is the first work to provide such a characterization in the challenging setting of robust reinforcement learning.

2. Background and Problem Formulation

A Markov Decision Process is a tuple $M = (\mathcal{S}, \mathcal{A}, r, P, \alpha)$ where \mathcal{S} is the state space, \mathcal{A} is the action space, $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function, and $\alpha \in (0, 1)$ is the discount factor. The transition probability matrix $P_{s,a}(s')$ represents the probability of transitioning to state s' when action a is taken at state s . We consider a finite MDP setting where the cardinality of state and action spaces are finite (but very large). A (deterministic) policy π maps each state to an action. The value of a policy π evaluated at state s is given by

$$V_{\pi,P}(s) = \mathbb{E}_{\pi,P} \left[\sum_{t=0}^{\infty} \alpha^t r(s_t, a_t) \mid s_0 = s \right],$$

where $a_t \sim \pi(s_t)$ and $s_{t+1} \sim P_{s_t, a_t}(\cdot)$. The optimal value function and the optimal policy of an MDP with the transition probability P are defined as $V_P^* = \max_{\pi} V_{\pi, P}$ and $\pi_P^* = \arg \max_{\pi} V_{\pi, P}$.

The RMDP formulation considers a set of model parameters (uncertainty set) under the assumption that the actual parameters lie in this uncertainty set, and the algorithm computes a robust policy that performs best under the worst model. More precisely, instead of a fixed transition probability matrix P , we consider a set of transition probability matrices \mathcal{P} . We assume that the set \mathcal{P} satisfies the standard *rectangularity condition* (Iyengar, 2005). The objective is to find a policy that maximizes the worst-case performance. Formally, the *robust value function* V_{π} corresponding to a policy π and the *optimal robust value function* V^* are defined as (Iyengar, 2005; Nilim & El Ghaoui, 2005)

$$V_{\pi} = \inf_{P \in \mathcal{P}} V_{\pi, P}, \quad V^* = \sup_{\pi} \inf_{P \in \mathcal{P}} V_{\pi, P}. \quad (1)$$

The *optimal robust policy* π^* is such that the robust value function corresponding to it matches the optimal robust value function, that is, $V_{\pi^*} = V^*$.

A generic characterization of the set \mathcal{P} makes the RMDPs problems intractable to solve by model-free methods. In the standard model-free methods, the algorithm has access to a simulator that can simulate the next state given the current state and current action, according to a fixed transition probability matrix (that is unknown to the algorithm). However, generating samples according to each and every transition probability matrix from the set \mathcal{P} is clearly infeasible. To overcome this difficulty, we use the characterization of the uncertainty set used in (Roy et al., 2017).

Assumption 1 (Uncertainty Set). *Each $P \in \mathcal{P}$ can be represented as $P_{s,a}(\cdot) = P_{s,a}^o(\cdot) + U_{s,a}(\cdot)$ for some $U_{s,a} \in \mathcal{U}_{s,a}$, where $P_{s,a}^o(\cdot)$ is the unknown transition probability matrix corresponding to the nominal (simulator) model and $\mathcal{U}_{s,a}$ is a confidence region around it.*

Using the above characterization, we can write $\mathcal{P} = \{P^o + U : U \in \mathcal{U}\}$, where $\mathcal{U} = \cup_{s,a} \mathcal{U}_{s,a}$. So, \mathcal{U} is the set of all possible perturbations to the nominal model P^o .

An example of the uncertainty set \mathcal{U} can be the spherical uncertainty set with a radius parameter. Define $\mathcal{U}_{s,a} := \{x \mid \|x\|_2 \leq r, \sum_{s \in \mathcal{S}} x_s = 0, -P_{s,a}^o(s') \leq x_{s'} \leq 1 - P_{s,a}^o(s'), \forall s' \in \mathcal{S}\}$, for all $(s, a) \in (\mathcal{S}, \mathcal{A})$, for some $r > 0$. Notice that, this uncertainty set uses the knowledge of the nominal model P^o in its construction. In practice, we do not know P^o . So, in Section 3.2, we introduce an approximate uncertainty set without using this information.

We consider *robust Bellman operator for policy evaluation*,

defined as (Iyengar, 2005)

$$T_{\pi}(V)(s) = r(s, \pi(s)) + \alpha \inf_{P \in \mathcal{P}} \sum_{s'} P_{s, \pi(s)}(s') V(s'), \quad (2)$$

a popular approach to solve (1). Using our characterization of the uncertainty set, we can rewrite (2) as

$$T_{\pi}(V)(s) = r(s, \pi(s)) + \alpha \sum_{s'} P_{s, \pi(s)}^o(s') V(s') + \alpha \inf_{U \in \mathcal{U}_{s, \pi(s)}} \sum_{s'} U_{s, \pi(s)}(s') V(s'). \quad (3)$$

For any set \mathcal{B} and a vector v , define $\sigma_{\mathcal{B}}(v) = \inf\{u^{\top} v : u \in \mathcal{B}\}$. We denote $|\mathcal{S}|$ as the cardinality of the set \mathcal{S} . Let $\sigma_{\mathcal{U}_{s, \pi(s)}}(v)$ and r_{π} be the $|\mathcal{S}|$ dimensional column vectors defined as $(\sigma_{\mathcal{U}_{s, \pi(s)}}(v) : s \in \mathcal{S})^{\top}$ and $(r(s, \pi(s)) : s \in \mathcal{S})^{\top}$, respectively. Let P_{π}^o be the stochastic matrix corresponding to the policy π where for any $s, s' \in \mathcal{S}$, $P_{\pi}^o(s, s') = P_{s, \pi(s)}^o(s')$. Then, (3) can be written in the matrix form as

$$T_{\pi}(V) = r_{\pi} + \alpha P_{\pi}^o V + \alpha \sigma_{\mathcal{U}_{\pi}}(V). \quad (4)$$

It is known (Iyengar, 2005) that T_{π} is a contraction in sup norm and the robust value function V_{π} is the unique fixed point of T_{π} . The *robust Bellman operator* T can also be defined in the same way as in the non-robust setting,

$$T(V) = \max_{\pi} T_{\pi}(V). \quad (5)$$

It is also known (Iyengar, 2005) that T is a contraction in sup norm, and the optimal robust value function V^* is its unique fixed point.

The goal of the robust RL is to learn the optimal robust policy π^* without knowing the nominal model P^o or the uncertainty set \mathcal{P} .

3. Robust Least Squares Policy Evaluation

In this section, we develop the RLSPE(λ) algorithm for learning the robust value function.

3.1. Robust TD(λ) Operator and the Challenges

In RL, a very useful approach for analyzing the multi-step learning algorithms like TD(λ), LSTD(λ), and LSPE(λ) is to define a multi-step Bellman operator called TD(λ) operator (Tsitsiklis & Van Roy, 1997; Bertsekas, 2012). Following the same approach, we can define the robust TD(λ) operator as well. For a given policy π , and a parameter $\lambda \in [0, 1)$, the robust TD(λ) operator denoted by $T_{\pi}^{(\lambda)} : \mathbb{R}^{|\mathcal{S}|} \rightarrow \mathbb{R}^{|\mathcal{S}|}$ is defined as

$$T_{\pi}^{(\lambda)}(V) = (1 - \lambda) \sum_{m=0}^{\infty} \lambda^m T_{\pi}^{m+1}(V). \quad (6)$$

Note that for $\lambda = 0$, we recover T_{π} . The following result is straightforward.

Proposition 1 (informal). $T_\pi^{(\lambda)}$ is a contraction in sup norm and the robust value function V_π is its unique fixed point, for any $\alpha \in (0, 1)$, $\lambda \in [0, 1)$.

For RMDPs with very large state space, exact dynamic programming methods which involve the evaluation of (3) or (6) are intractable. A standard approach to overcome this issue is to approximate the value function using some function approximation architecture. Here we focus on linear function approximation architectures (Bertsekas, 2012). In linear function approximation architectures, the value function is represented as the weighted sum of features as, $\bar{V}(s) = \phi(s)^\top w$, $\forall s \in \mathcal{S}$, where $\phi(s) = (\phi_1(s), \phi_2(s), \dots, \phi_L(s))^\top$ is an L dimensional feature vector with $L < |\mathcal{S}|$, and $w = (w_1, \dots, w_L)^\top$ is a weight vector. In the matrix form, this can be written as $\bar{V} = \Phi w$ where Φ is an $|\mathcal{S}| \times L$ dimensional feature matrix whose s^{th} row is $\phi(s)^\top$. We assume linearly independent columns for Φ , i.e., $\text{rank}(\Phi) = L$.

The standard approach to find an approximate (robust) value function is to solve for a w_π , with $\bar{V}_\pi = \Phi w_\pi$, such that $\Phi w_\pi = \Pi T_\pi^{(\lambda)} \Phi w_\pi$, where Π is a projection onto the subspace spanned by the columns of Φ . The projection is with respect to a d -weighted Euclidean norm. This norm is defined as $\|V\|_d^2 = V^\top D V$, where D is a diagonal matrix with non-negative diagonal entries ($d(s)$, $s \in \mathcal{S}$), for any vector V . Under suitable assumptions, (Tamar et al., 2014) showed that ΠT_π is a contraction in a d -weighted Euclidean norm. We also use a similar assumption stated below.

Assumption 2. (i) For any given policy π , there exists an exploration policy $\pi_e = \pi_{\text{exp}}(\pi)$ and a $\beta \in (0, 1)$ such that $\alpha P_{s,\pi(s)}(s') \leq \beta P_{s,\pi_e(s)}^o(s')$, for all transition probability matrices $P \in \mathcal{P}$ and for all states $s, s' \in \mathcal{S}$.

(ii) There exists a steady state distribution $d_{\pi_e} = (d_{\pi_e}(s), s \in \mathcal{S})$ for the Markov chain with transition probability $P_{\pi_e}^o$ with $d_{\pi_e}(s) > 0, \forall s \in \mathcal{S}$.

In the following, we will simply use d instead of d_{π_e} .

Though the above assumption appears restrictive, it is necessary to show that ΠT_π is a contraction in the d -weighted Euclidean norm, as proved in (Tamar et al., 2014). Also, a similar assumption is used in proving the convergence of off-policy reinforcement learning algorithm (Bertsekas & Yu, 2009). In the robust case, we can expect a similar condition because we are learning a robust value function for a set of transition probability matrices instead of a single transition probability matrix. We can now show the following.

Proposition 2 (informal). Under Assumption 2, $\Pi T_\pi^{(\lambda)}$ is a contraction mapping in the d -weighted Euclidean norm for any $\lambda \in [0, 1)$.

The linear approximation based robust value function $\bar{V}_\pi = \Phi w_\pi$ can be computed using the iteration, $\Phi w_{k+1} =$

$\Pi T_\pi^{(\lambda)} \Phi w_k$. Since $\Pi T_\pi^{(\lambda)}$ is a contraction, w_k will converge to w^* . A closed form solution for w_{k+1} given w_k can be found by **least squares approach** as $w_{k+1} = \arg \min_w \|\Phi w - \Pi T_\pi^{(\lambda)} \Phi w_k\|_d^2$. It can be shown that (details are given in the supplementary material), we can get a closed form solution for w_{k+1} as

$$w_{k+1} = w_k + (\Phi^\top D \Phi)^{-1} \Phi^\top D (T_\pi^{(\lambda)} \Phi w_k - \Phi w_k). \quad (7)$$

This is similar to the projected equation approach (Bertsekas, 2012) in the non-robust setting. Even in the non-robust setting, iterations using the (7) is intractable for MDPs with large state space. Moreover, when the transition matrix is unknown, it is not feasible to use (7) exactly even for small RMDPs. Simulation-based model-free learning algorithms are developed for addressing this problem in the non-robust case. In particular, LSPE(λ) algorithm (Nedić & Bertsekas, 2003; Bertsekas, 2012) is used to solve the iterations of the above form.

However, compared to the non-robust setting, there are two significant challenges in learning the robust value function by using simulation-based model-free approaches.

(i) Non-linearity of the robust TD(λ) operator: The non-robust T_π operator and the TD(λ) operator do not involve any nonlinear operations. So, they can be estimated efficiently from simulation samples in a model-free way. However, the robust TD(λ) operator when expanded will have the following form (derivation is given in the supplementary material).

$$T_\pi^{(\lambda)}(V) = (1 - \lambda) \sum_{m=0}^{\infty} \lambda^m \left(\sum_{k=0}^m (\alpha P_\pi^o)^k r_\pi + (\alpha P_\pi^o)^{m+1} V + \alpha \sum_{k=0}^m (\alpha P_\pi^o)^k \sigma_{\mathcal{U}_\pi}(T_\pi^{(m-k)} V) \right). \quad (8)$$

The last term is very difficult to estimate using simulation-based model-free approaches due to the composition of operations $\sigma_{\mathcal{U}_\pi}$ and T_π . In addition, nonlinearity of the T_π operator by itself adds to the complexity.

(ii) Unknown uncertainty region \mathcal{U} : In our formulation, we assumed that the transition probability uncertainty set \mathcal{P} is given by $\mathcal{P} = P^o + \mathcal{U}$. So, for each $U \in \mathcal{U}$, $P^o + U$ should be a valid transition probability matrix. However, in the model-free setting, we do not know the nominal transition probability P^o . So, it is not possible to know \mathcal{U} exactly a priori. One can only use an approximation $\hat{\mathcal{U}}$ instead of \mathcal{U} . This can possibly affect the convergence of the learning algorithms.

3.2. Robust Least Squares Policy Evaluation (RLSPE(λ)) Algorithm

We overcome the challenges of learning the robust value function by defining an approximate robust TD(λ) operator,

and by developing a robust least squares policy evaluation algorithm based on that.

Let $\hat{\mathcal{U}}$ be the approximate uncertainty set we use instead of the actual uncertainty set. An example of the approximate uncertainty set $\hat{\mathcal{U}}$ can be the spherical uncertainty set defined *without* using the knowledge of the model P^o as $\hat{\mathcal{U}}_{s,a} := \{x \mid \|x\|_2 \leq r, \sum_{s \in \mathcal{S}} x_s = 0\}$ for all $(s, a) \in (\mathcal{S}, \mathcal{A})$. Note that $P^o + U$ for $U \in \hat{\mathcal{U}}$ need not be a valid transition probability matrix and this poses challenges both for the algorithm and analysis.

For a given policy π and a parameter $\lambda \in [0, 1)$, approximate robust TD(λ) operator denoted by $\tilde{T}_\pi^{(\lambda)} : \mathbb{R}^{|\mathcal{S}|} \rightarrow \mathbb{R}^{|\mathcal{S}|}$, is defined as

$$\begin{aligned} \tilde{T}_\pi^{(\lambda)}(V) = (1 - \lambda) \sum_{m=0}^{\infty} \lambda^m \left[\sum_{t=0}^m (\alpha P_\pi^o)^t r_\pi \right. \\ \left. + \alpha \sum_{t=0}^m (\alpha P_\pi^o)^t \sigma_{\hat{\mathcal{U}}_\pi}(V) + (\alpha P_\pi^o)^{m+1} V \right]. \end{aligned} \quad (9)$$

Note that even with $\hat{\mathcal{U}}_\pi = \mathcal{U}_\pi$, (9) is different from (8). We will show that this clever approximation helps to overcome the challenges due to the nonliterary associated with (8).

However, we emphasize that (9) is not an arbitrary definition. Note that, for $\hat{\mathcal{U}}_\pi = \mathcal{U}_\pi$, with $\lambda = 0$, we recover the operator T_π . Moreover, the robust value function V_π is a fixed point of $\tilde{T}_\pi^{(\lambda)}$ when $\hat{\mathcal{U}}_\pi = \mathcal{U}_\pi$ for any $\lambda \in [0, 1)$. We state this formally below.

Proposition 3. *Suppose $\hat{\mathcal{U}}_\pi = \mathcal{U}_\pi$. Then, for any $\alpha \in (0, 1)$ and $\lambda \in [0, 1)$, the robust value function V_π is a fixed point of $\tilde{T}_\pi^{(\lambda)}$, i.e., $\tilde{T}_\pi^{(\lambda)}(V_\pi) = V_\pi$.*

Intuitively, the convergence of any learning algorithm using the approximate robust TD(λ) operator will depend on the difference between the actual uncertainty set \mathcal{U} and its approximation $\hat{\mathcal{U}}$. To quantify this, we use the following metric. Let $\rho = \max_{s \in \mathcal{S}, a \in \mathcal{A}} \rho_{s,a}$ where

$$\rho_{s,a} = \max \left\{ \begin{aligned} &\max_{x \in \hat{\mathcal{U}}_{s,a}} \max_{y \in \mathcal{U}_{s,a} \setminus \hat{\mathcal{U}}_{s,a}} \|x - y\|_d / d_{\min}, \\ &\max_{x \in \mathcal{U}_{s,a}} \max_{y \in \hat{\mathcal{U}}_{s,a} \setminus \mathcal{U}_{s,a}} \|x - y\|_d / d_{\min} \end{aligned} \right\}$$

and $d_{\min} := \min_{s \in \mathcal{S}} d(s)$. By convention, we set $\rho_{s,a} = 0$ when $\hat{\mathcal{U}}_{s,a} = \mathcal{U}_{s,a}$ for all $(s, a) \in (\mathcal{S}, \mathcal{A})$. So, $\rho = 0$ if $\hat{\mathcal{U}} = \mathcal{U}$. Using this characterization and under some additional assumptions on the discount factor, we show that the approximate robust TD(λ) operator is a contraction in the d -weighted Euclidean norm.

Theorem 1. *Under Assumption 2, for any $V_1, V_2 \in \mathbb{R}^{|\mathcal{S}|}$ and $\lambda \in [0, 1)$,*

$$\|\Pi \tilde{T}_\pi^{(\lambda)} V_1 - \Pi \tilde{T}_\pi^{(\lambda)} V_2\|_d \leq c(\alpha, \beta, \rho, \lambda) \|V_1 - V_2\|_d, \quad (10)$$

where $c(\alpha, \beta, \rho, \lambda) = (\beta(2 - \lambda) + \rho\alpha)/(1 - \beta\lambda)$. So, if $c(\alpha, \beta, \rho, \lambda) < 1$, $\Pi \tilde{T}_\pi^{(\lambda)}$ is a contraction in the d -weighted Euclidean norm. Moreover, there exists a unique w_π such that $\Phi w_\pi = \Pi \tilde{T}_\pi^{(\lambda)}(\Phi w_\pi)$. Furthermore, for this w_π ,

$$\|V_\pi - \Phi w_\pi\|_d \leq \frac{1}{1 - c(\alpha, \beta, \rho, \lambda)} \left(\|V_\pi - \Pi V_\pi\|_d + \frac{\beta\rho \|V_\pi\|_d}{1 - \beta\lambda} \right). \quad (11)$$

We note that despite the assumption on the discount factor, we empirically show in Section 5 that our learning algorithm converges to a robust policy even if this assumption is violated. We also note that the upper bound in (11) quantifies the *error* of approximating the robust value function V_π with the approximate robust value function Φw_π . We will later use this error bound in characterizing performance of both RLSPE and RLSPI algorithms.

Using the contraction property of approximate robust TD(λ) operator, the linear approximation based robust value function $\bar{V}_\pi = \Phi w_\pi$ can be computed using the iteration, $\Phi w_{k+1} = \Pi \tilde{T}_\pi^{(\lambda)} \Phi w_k$. Similar to (7), we can get a closed form solution for w_{k+1} using **least squares approach** as

$$w_{k+1} = w_k + (\Phi^\top D \Phi)^{-1} \Phi^\top D (\tilde{T}_\pi^{(\lambda)} \Phi w_k - \Phi w_k). \quad (12)$$

This can be written in a more succinct matrix form as given below (derivation is given in the supplementary material).

$$w_{k+1} = w_k + B^{-1}(Aw_k + C(w_k) + b), \quad \text{where,} \quad (13)$$

$$A = \Phi^\top D (\alpha P_\pi^o - I) \sum_{m=0}^{\infty} (\alpha \lambda P_\pi^o)^m \Phi, \quad (14)$$

$$B = \Phi^\top D \Phi, \quad (15)$$

$$C(w) = \alpha \Phi^\top D \sum_{t=0}^{\infty} (\alpha \lambda P_\pi^o)^t \sigma_{\hat{\mathcal{U}}_\pi}(\Phi w), \quad (16)$$

$$b = \Phi^\top D \sum_{t=0}^{\infty} (\alpha \lambda P_\pi^o)^t r_\pi. \quad (17)$$

Iterations by evaluating (13) exactly is intractable for MDPs with large state space, and infeasible if we do not know the transition probability P_π^o . To address this issue, we propose a simulation-based model-free online reinforcement learning algorithm, which we call robust least squares policy evaluation (RLSPE(λ)) algorithm, for learning the robust value function.

RLSPE(λ) algorithm: Generate a sequence of states and rewards, $(s_t, r_t, t \geq 0)$, using the policy π . Update the

parameters as

$$w_{t+1} = w_t + \gamma_t B_t^{-1} (A_t w_t + b_t + C_t(w_t)), \text{ where,} \quad (18)$$

$$A_t = \frac{1}{t+1} \sum_{\tau=0}^t z_\tau (\alpha \phi^\top(s_{\tau+1}) - \phi^\top(s_\tau)), \quad (19)$$

$$B_t = \frac{1}{t+1} \sum_{\tau=0}^t \phi(s_\tau) \phi^\top(s_\tau), \quad (20)$$

$$C_t(w) = \frac{\alpha}{t+1} \sum_{\tau=0}^t z_\tau \sigma_{\hat{U}_{s_\tau, \pi(s_\tau)}}(\Phi w), \quad (21)$$

$$b_t = \frac{1}{t+1} \sum_{\tau=0}^t z_\tau r(s_\tau, \pi(s_\tau)), \quad (22)$$

$$z_\tau = \sum_{m=0}^{\tau} (\alpha \lambda)^{\tau-m} \phi(s_m), \quad (23)$$

where γ_t is a deterministic sequence of step sizes. We assume that the step size satisfies the the standard Robbins-Munro stochastic conditions for stochastic approximation, i.e., $\sum_{t=0}^{\infty} \gamma_t = \infty$, $\sum_{t=0}^{\infty} \gamma_t^2 < \infty$.

We use the on-policy version of the RLSPE(λ) algorithm in the above description. So, we implicitly assume that the given policy π is an exploration policy according to the Assumption 2. This is mainly for the clarity of the presentation and notational convenience. Also, this simplifies the presentation of the policy iteration algorithm introduced in the next section. An off-policy version of the above algorithm can be implemented using the techniques given in (Bertsekas & Yu, 2009). We now give the convergence result of the RLSPE(λ) algorithm.

Theorem 2. *Let Assumption 2 hold. Also, let $c(\alpha, \beta, \rho, \lambda) < 1$ so that $\Pi \tilde{T}_\pi^{(\lambda)}$ is a contraction according to Theorem 1. Let $\{w_t\}$ be the sequence generated by the RLSPE(λ) algorithm given in (18). Then, w_t converges to w_π with probability 1 where w_π satisfies the fixed point equation $\Phi w_\pi = \Pi \tilde{T}_\pi^{(\lambda)} \Phi w_\pi$.*

The key idea of the proof is to show that the RLSPE(λ) update (18) approximates the exact update equation (12) and both converge to the same value w_π . One particularly challenging task is in analyzing the behavior of the term $C_t(w)$ due to the non-linearity of the function $\sigma_{\hat{U}_\pi}(\cdot)$. We use the tools from stochastic approximation theory (Borkar, 2009; Nedić & Bertsekas, 2003) to show this rigorously after establishing the tractable properties of the function $\sigma_{\hat{U}_\pi}(\cdot)$.

Note that Theorem 2 and Theorem 1 together give an error bound for the converged solution of the RLSPE(λ) algorithm. More precisely, Theorem 2 shows the convergence of the RLSPE(λ) algorithm to w_π and Theorem 1 gives the bound on $\|V_\pi - \Phi w_\pi\|_d$, which is the error due to linear

function approximation. We will use this bound in the the convergence analysis of the RLSPI algorithm presented in the next section.

4. Robust Least Squares Policy Iteration

In this section, we introduce the robust least squares policy iteration (RLSPI) algorithm for finding the optimal robust policy. RLSPI algorithm can be thought as the robust version of the LSPI algorithm (Lagoudakis & Parr, 2003). RLSPI algorithm uses the RLSPE(λ) algorithm for policy evaluation. However, model-free policy improvement is difficult when working with value functions since the policy update step will require us to solve

$$\pi_{k+1} = \arg \max_{\pi} \tilde{T}_\pi^{(\lambda)}(\bar{V}_k), \quad (24)$$

where \bar{V}_k is the approximate robust value function corresponding to the policy π_k , in the $(k+1)^{\text{th}}$ policy iteration loop. To overcome this, we first introduce the robust state-action value function (Q-function).

For any given policy π and state-action pair (s, a) , we define the robust Q-value as,

$$Q_\pi(s, a) = \inf_{P \in \mathcal{P}} \mathbb{E}_P \left[\sum_{t=0}^{\infty} \alpha^t r(s_t, a_t) \mid s_0 = s, a_0 = a \right]. \quad (25)$$

Instead of learning the approximate robust value function \bar{V}_π , we can learn the approximate robust Q-value function \bar{Q}_π using RLSPE(λ). This can be done by defining the feature vector $\phi(s, a)$ where $\phi(s, a) = (\phi_1(s, a), \dots, \phi_L(s, a))^\top$ and the linear approximation of the form $\bar{Q}_\pi(s, a) = w^\top \phi(s, a)$ where w is a weight vector. The results from the previous section on the convergence of the RLSPE(λ) algorithm applies for the case of learning Q-value function as well.

RLSPI is a policy iteration algorithm that uses RLSPE(λ) for policy evaluation at each iteration. It starts with an arbitrary initial policy π_0 . At the k th iteration, RLSPE(λ) returns a weight vector that represents the approximate Q-value function $\bar{Q}_{\pi_k} = \Phi w_{\pi_k}$ corresponding to the policy π_k . The next policy π_{k+1} is the greedy policy corresponding to \bar{Q}_{π_k} , defined as $\pi_{k+1}(s) = \arg \max_{a \in \mathcal{A}} \bar{Q}_{\pi_k}(s, a)$. For empirical evaluation purposes, we terminate the policy iteration for some finite value K . RLSPI algorithm is summarized in Algorithm 1.

We make the following assumptions for the convergence analysis of the RLSPI algorithm. We note that we work with value functions instead of Q-value functions for notational convenience and consistency.

Assumption 3. (i) Each policy π_k is an exploration policy, i.e. $\pi_{exp}(\pi_k) = \pi_k$.
(ii) The Markov chain $P_{\pi_k}^o$ has a stationary distribution d_{π_k}

Algorithm 1 RLSPI Algorithm

- 1: Initialization: Policy evaluation weights error ϵ_0 , initial policy π_0 .
- 2: **for** $k = 0 \dots K$ **do**
- 3: Initialize the policy weight vector w_0 . Initialize time step $t \leftarrow 0$.
- 4: **repeat**
- 5: Observe the state s_t , take action $a_t = \pi_k(s_t)$, observe reward r_t and next state s_{t+1} .
- 6: Update the weight vector w_t according to RLSPE(λ) algorithm (c.f. (18)-(23))
- 7: $t \leftarrow t + 1$
- 8: **until** $\|w_t - w_{t-1}\|_2 < \epsilon_0$
- 9: $w_{\pi_k} \leftarrow w_t$
- 10: Update the policy

$$\pi_{k+1}(s) = \arg \max_{a \in \mathcal{A}} \phi(s, a)^\top w_{\pi_k}$$

11: **end for**

such that $d_{\pi_k}(s) > 0, \forall s \in \mathcal{S}$.

(iii) There exists a finite scalar δ such that $\|V_{\pi_k} - \Pi_{d_{\pi_k}} V_{\pi_k}\|_{d_{\pi_k}} < \delta$ for all k , where $\Pi_{d_{\pi_k}}$ is a projection onto the subspace spanned by the columns of Φ under the d_{π_k} -weighted Euclidean norm.

(iv) For any probability distribution μ , define another probability distribution $\mu_k = \mu H_k$ where H_k is a stochastic matrix defined with respect to π_k . Also assume that there exists a probability distribution $\bar{\mu}$ and finite positive scalars C_1, C_2 such that $\mu_k \leq C_1 \bar{\mu}$ and $d_{\pi_k} \geq \bar{\mu}/C_2$ for all k .

We note that these are the standard assumptions used in the RL literature to provide theoretical guarantees for approximate policy/value iteration algorithms with linear function approximation in the non-robust settings (Munos, 2003; Munos & Szepesvári, 2008; Lazaric et al., 2012). We make no additional assumptions even though we are addressing the more difficult robust RL problem. The specific form of the stochastic matrix H_k specified in Assumption 3.(iv) is deferred to the proof of Theorem 3 for brevity of the presentation.

We now give the asymptotic convergence result for the RLSPI algorithm. We assume that, similar to the non-robust setting (Munos, 2003), the policy evaluation step (inner loop) is run to the convergence. We only present the case where $\rho = 0$. The proof for the general case is straightforward, but involves much more detailed algebra. So, we omit those details for the clarity of presentation.

Theorem 3. *Let Assumption 2 and Assumption 3 hold. Let $\{\pi_k\}$ be the sequence of the policies generated by the RLSPI algorithm. Let V_{π_k} and $\bar{V}_k = \Phi w_{\pi_k}$ be true robust value function and the approximate robust value function*

corresponding to the policy π_k . Also, let V^* be the optimal robust value function. Then, with $c(\alpha, \beta, 0, \lambda) < 1$,

$$\begin{aligned} & \limsup_{k \rightarrow \infty} \|V^* - V_{\pi_k}\|_\mu \\ & \leq \frac{2\sqrt{C_1 C_2} c(\alpha, \beta, 0, \lambda)}{(1 - c(\alpha, \beta, 0, \lambda))^2} \limsup_{k \rightarrow \infty} \|V_{\pi_k} - \bar{V}_k\|_{d_{\pi_k}}. \end{aligned} \quad (26)$$

Moreover, from Theorem 1 and Assumption 3.(iii), we have

$$\limsup_{k \rightarrow \infty} \|V^* - V_{\pi_k}\|_\mu \leq \frac{2\sqrt{C_1 C_2} c(\alpha, \beta, 0, \lambda)}{(1 - c(\alpha, \beta, 0, \lambda))^3} \delta. \quad (27)$$

The above theorem, in particular (27), gives a (worst case) guarantee for the performance of the policy learned using the RLSPI algorithm. Note that the upper bound in (27) is a constant where δ represents the (unavoidable) error due to the linear function approximation. We also note that using ‘lim sup’ is necessary due to the policy chattering phenomenon in approximate policy iteration algorithms which exists even in the non-robust case (Bertsekas, 2012).

Instead of the asymptotic bound given in Theorem 3, we can actually get a bound for any K given in the RLSPI algorithm by modifying Assumption 3.(iv). We defer the precise statements of the assumption and theorem to Section C in the supplementary material due to page limitation.

5. Experiments

We implemented our RLSPI algorithm using the MushroomRL library (D’Eramo et al., 2020), and evaluated its performance against Q-learning algorithm for an environment with discrete action space, deep deterministic policy gradient (DDPG) (Lillicrap et al., 2016) algorithm for continuous action space, and LSPI algorithm (Lagoudakis & Parr, 2003). For comparing with the performance of our RLSPI algorithm against another robust RL algorithm, we implemented the soft-robust algorithms proposed in (Derman et al., 2018) which use deep neural networks for function approximation.

We chose a spherical uncertainty set with a radius r . For such a set $\hat{\mathcal{U}}$, a closed form solution of $\sigma_{\hat{\mathcal{U}}}(\Phi w)$ can be computed for faster simulation. We note that in all the figures shown below, the quantity in the vertical axis is averaged over 100 runs, with the thick line showing the averaged value and the band around shows the ± 0.5 standard deviation. These figures act as the performance criteria for comparing results. We provide more details and additional experiment results in Section D of supplementary.

We used the CartPole, MountainCar, and Acrobot environments from OpenAI Gym (Brockman et al., 2016). We trained LSPI algorithm and our RLSPI algorithm on these environments with nominal parameters (default parameters in OpenAI Gym (Brockman et al., 2016)). We also

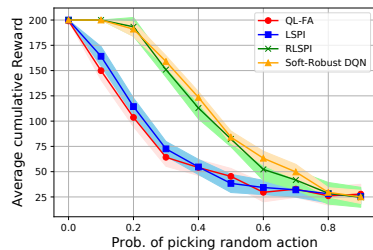


Figure 1: CartPole

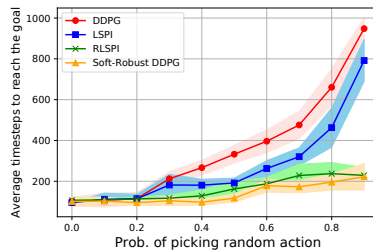


Figure 2: MountainCar

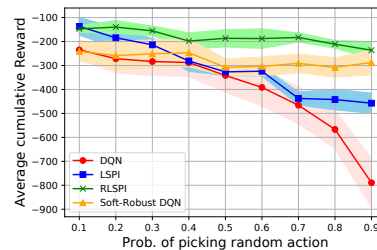


Figure 3: Acrobot

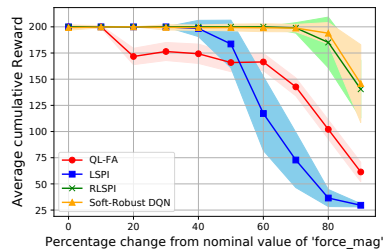


Figure 4: CartPole

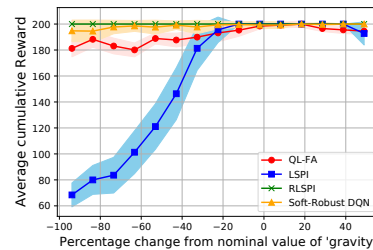


Figure 5: CartPole

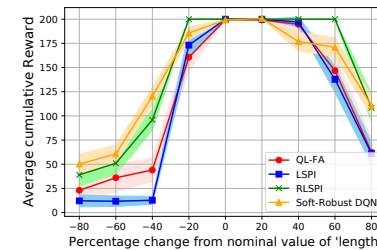


Figure 6: CartPole

trained Q-learning with linear function approximation and soft-robust deep Q-network (DQN) (Derman et al., 2018) algorithms on CartPole environment, DQN and soft-robust DQN (Derman et al., 2018) algorithms on Acrobot environment, and DDPG and soft-robust DDPG (Derman et al., 2018) algorithms on MountainCar environment. Then, to evaluate the robustness of the policies obtained, we changed the parameters of these environments and tested the performance of the learned policies on the perturbed environment.

In Figures 1-3, we show the robustness against action perturbations. In real-world settings, due to model mismatch or noise in the environments, the resulting action can be different from the intended action. We model this by picking a random action with some probability at each time step. Figure 1 shows the change in the average episodic reward against the probability of picking a random action for the CartPole environment. Figure 2 shows the average number of time steps to reach the goal in the MountainCar environment. Figure 3 shows the average episodic reward in the Acrobot environment. In all three cases, RLSPI algorithm shows robust performance against the perturbations.

Figures 4-6 show the test performance on CartPole, by changing the parameters *force_mag* (external force disturbance), *gravity*, *length* (length of pole on the cart). The nominal values of these parameters are 10, 9.8, and 0.5 respectively. RLSPI again exhibits robust performance.

The performance of our RLSPI algorithm is consistently superior to that of the non-robust algorithms. Moreover, the performance of RLSPI algorithm is comparable with that of the soft-robust algorithms (Derman et al., 2018), even

though the latter uses deep neural networks for function approximation while our algorithm uses only linear function approximation architecture. We also would like to emphasize that our work gives provable guarantees for the policy learned by the algorithm whereas (Derman et al., 2018) does not provide any such guarantees.

6. Conclusion and Future Work

We have presented an online model-free reinforcement learning algorithm to learn control policies that are robust to the parameter uncertainties of the model, for system with large state spaces. While there have been interesting empirical works on robust deep RL using neural network, they only provide convergence guarantees to a local optimum. Different from such empirical works, we proposed a learning based robust policy iteration algorithm called RLSPI algorithm with explicit theoretical guarantees on the performance of the learned policy. To the best of our knowledge, this is the first work that presents model-free reinforcement learning algorithm with function approximation for learning the optimal robust policy. We also empirically evaluated the performance of our RLSPI algorithm on standard benchmark RL problems.

In future, we plan to extend our theoretical results to nonlinear function approximation architectures. We also plan to characterize the sample complexity of robust reinforcement learning algorithms. Extending offline RL approaches to robust setting is another research area that we plan to pursue.

Acknowledgments

Dileep Kalathil gratefully acknowledges funding from the U.S. National Science Foundation (NSF) grants NSF-CRII-CPS-1850206 and NSF-CAREER-EPCN-2045783.

References

- Abbeel, P. and Ng, A. Y. Apprenticeship learning via inverse reinforcement learning. In *International Conference on Machine Learning*, pp. 1, 2004.
- Antos, A., Szepesvári, C., and Munos, R. Learning near-optimal policies with Bellman-residual minimization based fitted policy iteration and a single sample path. *Machine Learning*, 71(1):89–129, 2008.
- Arora, S., Du, S., Kakade, S., Luo, Y., and Saunshi, N. Provable representation learning for imitation learning via bi-level optimization. In *International Conference on Machine Learning*, pp. 367–376. PMLR, 2020.
- Bertsekas, D. P. Approximate policy iteration: A survey and some new methods. *Journal of Control Theory and Applications*, 9(3):310–335, 2011.
- Bertsekas, D. P. *Dynamic programming and optimal control, Vol - 2*. Athena scientific Belmont, MA, 2012.
- Bertsekas, D. P. and Ioffe, S. Temporal Differences-Based Policy Iteration and Applications in Neuro-Dynamic Programming. *Lab. for Info. and Decision Systems Report LIDS-P-2349, MIT, Cambridge, MA*, 14, 1996.
- Bertsekas, D. P. and Yu, H. Projected equation methods for approximate solution of large linear systems. *Journal of Computational and Applied Mathematics*, 227(1):27–50, 2009.
- Borkar, V. S. *Stochastic approximation: a dynamical systems viewpoint*, volume 48. Springer, 2009.
- Brockman, G., Cheung, V., Pettersson, L., Schneider, J., Schulman, J., Tang, J., and Zaremba, W. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- D’Eramo, C., Tateo, D., Bonarini, A., Restelli, M., and Peters, J. Mushroomrl: Simplifying reinforcement learning research. *arXiv preprint arXiv:2001.01102*, 2020. URL <https://github.com/MushroomRL/mushroom-rl>.
- Derman, E., Mankowitz, D. J., Mann, T. A., and Mannor, S. Soft-robust actor-critic policy-gradient. In *AUAI press for Association for Uncertainty in Artificial Intelligence*, pp. 208–218, 2018.
- Duan, Y., Jia, Z., and Wang, M. Minimax-optimal off-policy evaluation with linear function approximation. In *International Conference on Machine Learning*, pp. 2701–2709. PMLR, 2020.
- Iyengar, G. N. Robust dynamic programming. *Mathematics of Operations Research*, 30(2):257–280, 2005.
- Jin, C., Yang, Z., Wang, Z., and Jordan, M. I. Provably efficient reinforcement learning with linear function approximation. In *Conference on Learning Theory*, pp. 2137–2143, 2020.
- Kaufman, D. L. and Schaefer, A. J. Robust modified policy iteration. *INFORMS Journal on Computing*, 25(3):396–410, 2013.
- Kong, W., Somani, R., Song, Z., Kakade, S., and Oh, S. Meta-learning for mixed linear regression. In *International Conference on Machine Learning*, pp. 5394–5404. PMLR, 2020.
- Lagoudakis, M. G. and Parr, R. Least-squares policy iteration. *Journal of Machine Learning Research*, 4(Dec):1107–1149, 2003.
- Lazaric, A., Ghavamzadeh, M., and Munos, R. Finite-sample analysis of least-squares policy iteration. *Journal of Machine Learning Research*, 13(Oct):3041–3074, 2012.
- Lillicrap, T. P., Hunt, J. J., Pritzel, A., Heess, N., Erez, T., Tassa, Y., Silver, D., and Wierstra, D. Continuous control with deep reinforcement learning. In *ICLR (Poster)*, 2016.
- Lim, S. H. and Autef, A. Kernel-based reinforcement learning in robust Markov decision processes. In *International Conference on Machine Learning*, pp. 3973–3981, 2019.
- Lim, S. H., Xu, H., and Mannor, S. Reinforcement learning in robust Markov decision processes. In *Advances in Neural Information Processing Systems*, pp. 701–709, 2013.
- Mankowitz, D. J., Levine, N., Jeong, R., Abdolmaleki, A., Springenberg, J. T., Shi, Y., Kay, J., Hester, T., Mann, T., and Riedmiller, M. Robust reinforcement learning for continuous control with model misspecification. In *International Conference on Learning Representations*, 2020.
- Munos, R. Error bounds for approximate policy iteration. In *ICML*, volume 3, pp. 560–567, 2003.
- Munos, R. and Szepesvári, C. Finite-time bounds for fitted value iteration. *Journal of Machine Learning Research*, 9(27):815–857, 2008.

- Nedić, A. and Bertsekas, D. P. Least squares policy evaluation algorithms with linear function approximation. *Discrete Event Dynamic Systems*, 13(1-2):79–110, 2003.
- Nilim, A. and El Ghaoui, L. Robust control of Markov decision processes with uncertain transition matrices. *Operations Research*, 53(5):780–798, 2005.
- Pinto, L., Davidson, J., Sukthankar, R., and Gupta, A. Robust adversarial reinforcement learning. In *International Conference on Machine Learning*, pp. 2817–2826. PMLR, 2017.
- Puterman, M. L. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., NJ, 2005.
- Roy, A., Xu, H., and Pokutta, S. Reinforcement learning under model mismatch. In *Advances in Neural Information Processing Systems*, pp. 3043–3052, 2017.
- Tamar, A., Mannor, S., and Xu, H. Scaling up robust mdps using function approximation. In *International Conference on Machine Learning*, pp. 181–189, 2014.
- Tessler, C., Efroni, Y., and Mannor, S. Action robust reinforcement learning and applications in continuous control. In *International Conference on Machine Learning*, pp. 6215–6224, 2019.
- Tsitsiklis, J. N. and Van Roy, B. Analysis of temporal-difference learning with function approximation. In *Advances in Neural Information Processing Systems*, pp. 1075–1081, 1997.
- Vinitzky, E., Du, Y., Parvate, K., Jang, K., Abbeel, P., and Bayen, A. Robust reinforcement learning using adversarial populations. *arXiv preprint arXiv:2008.01825*, 2020.
- Wang, L., Cai, Q., Yang, Z., and Wang, Z. On the global optimality of model-agnostic meta-learning. In *International Conference on Machine Learning*, pp. 9837–9846. PMLR, 2020.
- Wang, R., Foster, D., and Kakade, S. M. What are the statistical limits of offline $\{rl\}$ with linear function approximation? In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=30EvkP2aQLD>.
- Wiesemann, W., Kuhn, D., and Rustem, B. Robust Markov decision processes. *Mathematics of Operations Research*, 38(1):153–183, 2013.
- Xu, H. and Mannor, S. Distributionally robust Markov decision processes. In *Advances in Neural Information Processing Systems*, pp. 2505–2513, 2010.
- Yang, L. and Wang, M. Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. In *International Conference on Machine Learning*, pp. 10746–10756. PMLR, 2020.
- Yu, P. and Xu, H. Distributionally robust counterpart in Markov decision processes. *IEEE Transactions on Automatic Control*, 61(9):2538–2543, 2015.
- Zhang, H., Chen, H., Xiao, C., Li, B., Boning, D., and Hsieh, C.-J. Robust deep reinforcement learning against adversarial perturbations on observations. In *Advances in Neural Information Processing Systems*, 2020.

Appendix

A. Proofs of the Results in Section 3.1

A.1. Proof of Proposition 1

We first restate Proposition 1 formally and then give the proof.

Proposition 4. (i) For any $V_1, V_2 \in \mathbb{R}^{|\mathcal{S}|}$ and $\lambda \in [0, 1)$, $\|T_\pi^{(\lambda)} V_1 - T_\pi^{(\lambda)} V_2\|_\infty \leq \frac{\alpha(1-\lambda)}{(1-\alpha\lambda)} \|V_1 - V_2\|_\infty$. So, $T_\pi^{(\lambda)}$ is a contraction in sup norm for any $\alpha \in (0, 1), \lambda \in [0, 1)$.

(ii) The robust value function V_π is the unique fixed point of $T_\pi^{(\lambda)}$, i.e., $T_\pi^{(\lambda)} V_\pi = V_\pi$, for all $\alpha \in (0, 1)$ and $\lambda \in [0, 1)$.

Proof. From (6) we have

$$\begin{aligned} \|T_\pi^{(\lambda)}(V_1) - T_\pi^{(\lambda)}(V_2)\|_\infty &= \|(1-\lambda) \sum_{m=0}^{\infty} \lambda^m (T_\pi^{m+1}(V_1) - T_\pi^{m+1}(V_2))\|_\infty \\ &\leq (1-\lambda) \sum_{m=0}^{\infty} \lambda^m \|T_\pi^{m+1}(V_1) - T_\pi^{m+1}(V_2)\|_\infty \\ &\stackrel{(a)}{\leq} (1-\lambda) \sum_{m=0}^{\infty} \lambda^m \alpha^{m+1} \|V_1 - V_2\|_\infty = \frac{\alpha(1-\lambda)}{(1-\alpha\lambda)} \|V_1 - V_2\|_\infty \end{aligned}$$

where (a) follows since T_π is a contraction operator with contraction modulus α . This proves (i). Since V_π is the unique fixed point of T_π , (ii) directly follows from (i) and the Banach Fixed Point Theorem (Puterman, 2005, Theorem 6.2.3). \square

A.2. Proof of Proposition 2

We first restate Proposition 2 formally and then give the proof.

Proposition 5. Under Assumption 2, for any $V_1, V_2 \in \mathbb{R}^{|\mathcal{S}|}$ and $\lambda \in [0, 1)$,

$\|\Pi T_\pi^{(\lambda)} V_1 - \Pi T_\pi^{(\lambda)} V_2\|_d \leq \frac{\beta(1-\lambda)}{(1-\beta\lambda)} \|V_1 - V_2\|_d$. So, $\Pi T_\pi^{(\lambda)}$ is a contraction mapping in d -weighted Euclidean norm for any $\beta \in (0, 1), \lambda \in [0, 1)$.

Proof. From (6) we have

$$\begin{aligned} \|T_\pi^{(\lambda)}(V_1) - T_\pi^{(\lambda)}(V_2)\|_d &= \|(1-\lambda) \sum_{m=0}^{\infty} \lambda^m (T_\pi^{m+1}(V_1) - T_\pi^{m+1}(V_2))\|_d \\ &\leq (1-\lambda) \sum_{m=0}^{\infty} \lambda^m \|T_\pi^{m+1}(V_1) - T_\pi^{m+1}(V_2)\|_d \\ &\stackrel{(a)}{\leq} (1-\lambda) \sum_{m=0}^{\infty} \lambda^m \beta^{m+1} \|V_1 - V_2\|_d = \frac{\beta(1-\lambda)}{(1-\beta\lambda)} \|V_1 - V_2\|_d \end{aligned}$$

where (a) follows since ΠT_π is a contraction in the d -weighted Euclidean norm with contraction modulus β (Tamar et al., 2014, Corollary 4). From (Tsitsiklis & Van Roy, 1997), Π is a nonexpansive mapping in the d -weighted Euclidean norm. So, $\Pi T_\pi^{(\lambda)}$ has the stated property. \square

A.3. Derivation of (7)

Given w_k, w_{k+1} which satisfies the equation $\Phi w_{k+1} = \Pi T_\pi^{(\lambda)} \Phi w_k$ can be written as the solution of the minimization problem $w_{k+1} = \arg \min_w \|\Phi w - T_\pi^{(\lambda)} \Phi w_k\|_d^2$. Taking gradient w.r.t. w and equating to zero, we get $\Phi^\top D(\Phi w - T_\pi^{(\lambda)} \Phi w_k) = 0$, which implies $w_{k+1} = (\Phi^\top D\Phi)^{-1} \Phi^\top D T_\pi^{(\lambda)} \Phi w_k$. This can be written as

$$\begin{aligned} w_{k+1} &= (\Phi^\top D\Phi)^{-1} \Phi^\top D T_\pi^{(\lambda)} \Phi w_k = w_k + (\Phi^\top D\Phi)^{-1} (\Phi^\top D T_\pi^{(\lambda)} \Phi w_k - \Phi^\top D\Phi w_k) \\ &= w_k + (\Phi^\top D\Phi)^{-1} \Phi^\top D (T_\pi^{(\lambda)} \Phi w_k - \Phi w_k). \end{aligned}$$

A.4. Derivation of (8)

We have

$$\begin{aligned}
 TV &= r_\pi + \alpha P_\pi^o V + \alpha \sigma_{\mathcal{U}_\pi}(V) \\
 T^2V &= r_\pi + \alpha P_\pi^o(TV) + \alpha \sigma_{\mathcal{U}_\pi}(TV) \\
 &= (1 + \alpha P) r_\pi + (\alpha P_\pi^o)^2 V + \alpha (\alpha P_\pi^o) \sigma_{\mathcal{U}_\pi}(V) + \alpha \sigma_{\mathcal{U}_\pi}(TV).
 \end{aligned} \tag{28}$$

Suppose we have an induction hypothesis for $T^m V$, i.e.,

$$T^m V = \sum_{k=0}^{m-1} (\alpha P_\pi^o)^k r_\pi + (\alpha P_\pi^o)^m V + \alpha \sum_{k=0}^{m-1} (\alpha P_\pi^o)^k \sigma_{\mathcal{U}_\pi}(T^{(m-1-k)} V).$$

Now, to verify an induction step, observe that,

$$\begin{aligned}
 T^{m+1}V &= \sum_{k=0}^{m-1} (\alpha P_\pi^o)^k r_\pi + (\alpha P_\pi^o)^m TV + \alpha \sum_{k=0}^{m-1} (\alpha P_\pi^o)^k \sigma_{\mathcal{U}_\pi}(T^{(m-1-k)} TV) \\
 &\stackrel{(a)}{=} \sum_{k=0}^{m-1} (\alpha P_\pi^o)^k r_\pi + (\alpha P_\pi^o)^m (r_\pi + \alpha P_\pi^o V + \alpha \sigma_{\mathcal{U}_\pi}(V)) + \alpha \sum_{k=0}^{m-1} (\alpha P_\pi^o)^k \sigma_{\mathcal{U}_\pi}(T^{(m-1-k)} TV) \\
 &= \sum_{k=0}^m (\alpha P_\pi^o)^k r_\pi + (\alpha P_\pi^o)^{m+1} V + \alpha \sum_{k=0}^m (\alpha P_\pi^o)^k \sigma_{\mathcal{U}_\pi}(T^{(m-k)} V)
 \end{aligned}$$

where (a) follows from (28). Now, (8) directly follows from (6).

B. Proofs of the Results in Section 3.2

B.1. Proof of Proposition 3

Proof. With $\hat{\mathcal{U}}_\pi = \mathcal{U}_\pi$, we have,

$$\begin{aligned}
 \tilde{T}_\pi^{(\lambda)}(V_\pi) &= (1 - \lambda) \sum_{m=0}^{\infty} \lambda^m \left[\sum_{t=0}^m (\alpha P_\pi^o)^t r_\pi + \alpha \sum_{t=0}^m (\alpha P_\pi^o)^t \sigma_{\mathcal{U}_\pi}(V_\pi) + (\alpha P_\pi^o)^{m+1} V_\pi \right] \\
 &\stackrel{(a)}{=} (1 - \lambda) \sum_{m=0}^{\infty} \lambda^m \left[\sum_{t=0}^m (\alpha P_\pi^o)^t (r_\pi + \alpha \sigma_{\mathcal{U}_\pi}(V_\pi)) + (\alpha P_\pi^o)^{m+1} V_\pi \right] \\
 &= (1 - \lambda) \sum_{m=0}^{\infty} \lambda^m \left[\sum_{t=0}^m (\alpha P_\pi^o)^t (I - \alpha P_\pi^o) V_\pi + (\alpha P_\pi^o)^{m+1} V_\pi \right] \\
 &\stackrel{(b)}{=} (1 - \lambda) \sum_{m=0}^{\infty} \lambda^m V_\pi = V_\pi.
 \end{aligned}$$

Here, to get (a), we wrote $r_\pi + \alpha \sigma_{\mathcal{U}_\pi}(V_\pi) = (I - \alpha P_\pi^o) V_\pi$ since $T_\pi V_\pi = V_\pi$. (b) is from the telescopic sum of the previous equation. \square

B.2. Proof of Theorem 1

We have the following lemma which is similar to Lemma 4.2 in (Roy et al., 2017).

Lemma 1. For any vector $V \in \mathbb{R}^{|\mathcal{S}|}$ and for all $(s, a) \in \mathcal{S} \times \mathcal{A}$,

$$|\sigma_{\hat{\mathcal{U}}_{s,a}}(V) - \sigma_{\mathcal{U}_{s,a}}(V)| \leq \rho \|V\|_d.$$

Proof. First note that, for any $x, y \in \mathbb{R}^{|\mathcal{S}|}$ we have

$$x^\top y \leq (x^\top D y) / d_{\min} \stackrel{(a)}{\leq} (\|x\|_d \|y\|_d) / d_{\min}, \tag{29}$$

where (a) follows from Cauchy-Schwarz inequality with respect to $\|\cdot\|_d$ norm.

Consider any $p \in \widehat{\mathcal{U}}_{s,a}$ and $q \in \mathcal{U}_{s,a} \setminus \widehat{\mathcal{U}}_{s,a}$. For any $V \in \mathbb{R}^{|S|}$, we have

$$\begin{aligned} p^\top V &= q^\top V + (p - q)^\top V \geq \sigma_{\mathcal{U}_{s,a}}(V) + (p - q)^\top V \geq \sigma_{\mathcal{U}_{s,a}}(V) + \min_{x \in \widehat{\mathcal{U}}_{s,a}} \min_{y \in \mathcal{U}_{s,a} \setminus \widehat{\mathcal{U}}_{s,a}} (x - y)^\top V \\ &= \sigma_{\mathcal{U}_{s,a}}(V) - \max_{x \in \widehat{\mathcal{U}}_{s,a}} \max_{y \in \mathcal{U}_{s,a} \setminus \widehat{\mathcal{U}}_{s,a}} (y - x)^\top V \stackrel{(b)}{\geq} \sigma_{\mathcal{U}_{s,a}}(V) - \max_{x \in \widehat{\mathcal{U}}_{s,a}} \max_{y \in \mathcal{U}_{s,a} \setminus \widehat{\mathcal{U}}_{s,a}} (\|y - x\|_d \|V\|_d) / d_{\min} \\ &\geq \sigma_{\mathcal{U}_{s,a}}(V) - \rho_{s,a} \|V\|_d, \end{aligned}$$

where (b) follows from (29). By taking infimum on both sides with respect to $p \in \widehat{\mathcal{U}}_{s,a}$, we get,

$$\sigma_{\widehat{\mathcal{U}}_{s,a}}(V) \geq \sigma_{\mathcal{U}_{s,a}}(V) - \rho_{s,a} \|V\|_d.$$

We can also get $\sigma_{\mathcal{U}_{s,a}}(V) \geq \sigma_{\widehat{\mathcal{U}}_{s,a}}(V) - \rho_{s,a} \|V\|_d$ by similar arguments. Combining, we get,

$$|\sigma_{\widehat{\mathcal{U}}_{s,a}}(V) - \sigma_{\mathcal{U}_{s,a}}(V)| \leq \rho_{s,a} \|V\|_d.$$

Since $\rho = \max_{s,a} \rho_{s,a}$, we get the desired result. \square

We will use the following result which follows directly from (Tsitsiklis & Van Roy, 1997, Lemma 1).

Lemma 2. Under Assumption 2, for any V , we have $\|P_{\pi_e}^o V\|_d \leq \|V\|_d$.

Remark 1. The inequality in the Assumption 2 can be written as, $\alpha P_{s,\pi(s)}^o(s') + \alpha U_{s,\pi(s)}(s') \leq \beta P_{s,\pi_e(s)}^o(s')$. From this, we can conclude that $\alpha U_{s,\pi(s)}(s') \leq \beta P_{s,\pi_e(s)}^o(s')$.

Next, we show the following.

Lemma 3. For any $V_1, V_2 \in \mathbb{R}^{|S|}$,

$$\|\sigma_{\widehat{\mathcal{U}}_\pi}(V_1) - \sigma_{\widehat{\mathcal{U}}_\pi}(V_2)\|_d \leq \left(\frac{\beta}{\alpha} + \rho\right) \|V_1 - V_2\|_d.$$

Proof. For any $s \in \mathcal{S}$ we have

$$\begin{aligned} \sigma_{\widehat{\mathcal{U}}_{s,\pi(s)}}(V_2) - \sigma_{\widehat{\mathcal{U}}_{s,\pi(s)}}(V_1) &= \inf_{q \in \widehat{\mathcal{U}}_{s,\pi(s)}} q^\top V_2 - \inf_{\tilde{q} \in \widehat{\mathcal{U}}_{s,\pi(s)}} \tilde{q}^\top V_1 = \inf_{q \in \widehat{\mathcal{U}}_{s,\pi(s)}} \sup_{\tilde{q} \in \widehat{\mathcal{U}}_{s,\pi(s)}} q^\top V_2 - \tilde{q}^\top V_1 \\ &\geq \inf_{q \in \widehat{\mathcal{U}}_{s,\pi(s)}} q^\top (V_2 - V_1) = \sigma_{\widehat{\mathcal{U}}_{s,\pi(s)}}(V_2 - V_1) \stackrel{(a)}{\geq} \sigma_{\mathcal{U}_{s,\pi(s)}}(V_2 - V_1) - \rho \|V_1 - V_2\|_d, \end{aligned} \quad (30)$$

where (a) follows from Lemma 1. By definition, for any arbitrary $\epsilon > 0$, there exists a $U_{s,\pi(s)} \in \mathcal{U}_{s,\pi(s)}$ such that

$$U_{s,\pi(s)}^\top (V_2 - V_1) - \epsilon \leq \sigma_{\mathcal{U}_{s,\pi(s)}}(V_2 - V_1). \quad (31)$$

Using (31) and (30),

$$\begin{aligned} \alpha(\sigma_{\widehat{\mathcal{U}}_{s,\pi(s)}}(V_1) - \sigma_{\widehat{\mathcal{U}}_{s,\pi(s)}}(V_2)) &\leq \alpha U_{s,\pi(s)}^\top (V_1 - V_2) + \alpha\epsilon + \rho\alpha \|V_1 - V_2\|_d \\ &\leq \alpha U_{s,\pi(s)}^\top (V_2 - V_1) + \alpha\epsilon + \rho\alpha \|V_1 - V_2\|_d \\ &\stackrel{(b)}{\leq} \beta(P_{s,\pi_e(s)}^o)^\top (V_1 - V_2) + \alpha\epsilon + \rho\alpha \|V_1 - V_2\|_d \end{aligned} \quad (32)$$

where (b) follows from Remark 1. Since ϵ is arbitrary, we get,

$$\alpha(\sigma_{\widehat{\mathcal{U}}_{s,\pi(s)}}(V_1) - \sigma_{\widehat{\mathcal{U}}_{s,\pi(s)}}(V_2)) \leq \beta(P_{s,\pi_e(s)}^o)^\top (V_1 - V_2) + \rho\alpha \|V_1 - V_2\|_d.$$

By exchanging the roles of V_1 and V_2 , we get $\alpha(\sigma_{\widehat{U}_{s,\pi(s)}}(V_2) - \sigma_{\widehat{U}_{s,\pi(s)}}(V_1)) \leq \beta(P_{s,\pi_e}^o)^\top |V_1 - V_2| + \rho\alpha\|V_1 - V_2\|_d$. Combining these and writing compactly in vector form, we get

$$\alpha|\sigma_{\widehat{U}_\pi}(V_1) - \sigma_{\widehat{U}_\pi}(V_2)| \leq \beta P_{\pi_e}^o |V_1 - V_2| + \rho\alpha\|V_1 - V_2\|_d \mathbf{1},$$

where $\mathbf{1} = (1, 1, \dots, 1)^\top$, an $|\mathcal{S}|$ -dimensional unit vector. Since $\alpha|\sigma_{\widehat{U}_\pi}(V_1) - \sigma_{\widehat{U}_\pi}(V_2)| \geq 0$, by the property of the norm, we get

$$\begin{aligned} \alpha\|\sigma_{\widehat{U}_\pi}(V_1) - \sigma_{\widehat{U}_\pi}(V_2)\|_d &\leq \|\beta P_{\pi_e}^o |V_1 - V_2| + \rho\alpha\|V_1 - V_2\|_d \mathbf{1}\|_d \\ &\stackrel{(c)}{\leq} \beta\|P_{\pi_e}^o |V_1 - V_2|\|_d + \rho\alpha\|V_1 - V_2\|_d \mathbf{1}\|_d \\ &\stackrel{(d)}{\leq} \beta\|V_1 - V_2\|_d + \rho\alpha\|V_1 - V_2\|_d \end{aligned}$$

where (c) follows from triangle inequality and (d) from Lemma 2 and from the fact that $\|\mathbf{1}\|_d = 1$. Dividing by α and rearranging, we get the desired result. \square

Proof of Theorem 1. We first observe

$$\|\alpha P_\pi^o |V|\|_d \stackrel{(a)}{\leq} \|\beta P_{\pi_e}^o |V|\|_d \stackrel{(b)}{\leq} \beta\|V\|_d, \quad (33)$$

where (a) follows from the Assumption 2 and (b) from the Lemma 2.

Notice that for any finite t we have,

$$\|(\alpha P_\pi^o)^t |V|\|_d = \|\alpha P_\pi^o (\alpha P_\pi^o)^{t-1} |V|\|_d \stackrel{(c)}{\leq} \beta\|(\alpha P_\pi^o)^{t-1} |V|\|_d \quad (34)$$

where (c) follows from (33). Using this repeatedly, we get,

$$\|(\alpha P_\pi^o)^t |V|\|_d \leq \beta^t \|V\|_d. \quad (35)$$

Now,

$$\begin{aligned} &\|\widetilde{T}_\pi^{(\lambda)}(V_1) - \widetilde{T}_\pi^{(\lambda)}(V_2)\|_d \\ &= \|(1 - \lambda) \sum_{m=0}^{\infty} \lambda^m [\alpha \sum_{t=0}^m (\alpha P_\pi^o)^t (\sigma_{\widehat{U}_\pi}(V_1) - \sigma_{\widehat{U}_\pi}(V_2)) + (\alpha P_\pi^o)^{m+1} (V_1 - V_2)]\|_d \\ &\leq (1 - \lambda) \sum_{m=0}^{\infty} \lambda^m [\alpha \sum_{t=0}^m \|(\alpha P_\pi^o)^t |\sigma_{\widehat{U}_\pi}(V_1) - \sigma_{\widehat{U}_\pi}(V_2)|\|_d + \|(\alpha P_\pi^o)^{m+1} |V_1 - V_2|\|_d] \\ &\stackrel{(d)}{\leq} (1 - \lambda) \sum_{m=0}^{\infty} \lambda^m [\alpha \sum_{t=0}^m (\beta)^t \|\sigma_{\widehat{U}_\pi}(V_1) - \sigma_{\widehat{U}_\pi}(V_2)\|_d + \beta^{m+1} \|V_1 - V_2\|_d] \\ &\stackrel{(e)}{\leq} (1 - \lambda) \sum_{m=0}^{\infty} \lambda^m [(\beta + \rho\alpha) \sum_{t=0}^m (\beta)^t \|V_1 - V_2\|_d + \beta^{m+1} \|V_1 - V_2\|_d] \end{aligned} \quad (36)$$

$$\begin{aligned} &= \left[\frac{(\beta + \rho\alpha)}{(1 - \beta)} \left(1 - \frac{\beta(1 - \lambda)}{(1 - \beta\lambda)} \right) + \frac{\beta(1 - \lambda)}{(1 - \beta\lambda)} \right] \|V_1 - V_2\|_d \\ &= \frac{\beta(2 - \lambda) + \rho\alpha}{(1 - \beta\lambda)} \|V_1 - V_2\|_d, \end{aligned} \quad (37)$$

where (d) follows from (35) and (e) follows from Lemma 3.

From (Bertsekas, 2012), Π is a non-expansive operator in $\|\cdot\|_d$. Thus,

$$\|\Pi \widetilde{T}_\pi^{(\lambda)} V_1 - \Pi \widetilde{T}_\pi^{(\lambda)} V_2\|_d \leq \|\widetilde{T}_\pi^{(\lambda)}(V_1) - \widetilde{T}_\pi^{(\lambda)}(V_2)\|_d \leq c(\alpha, \beta, \rho, \lambda) \|V_1 - V_2\|_d. \quad (38)$$

where $c(\alpha, \beta, \rho, \lambda) = (\beta(2 - \lambda) + \rho\alpha)/(1 - \beta\lambda)$. This concludes the proof of getting (10).

For proving (11), first denote the operator $\tilde{T}_\pi^{(\lambda)}$ as $\bar{T}_\pi^{(\lambda)}$ when $\hat{U}_\pi = U_\pi$. Now observe that

$$\begin{aligned}
 \|\tilde{T}_\pi^{(\lambda)}(V) - \bar{T}_\pi^{(\lambda)}(V)\|_d &= \|(1 - \lambda) \sum_{m=0}^{\infty} \lambda^m [\alpha \sum_{t=0}^m (\alpha P_\pi^o)^t (\sigma_{U_\pi}(V) - \sigma_{\hat{U}_\pi}(V))]\|_d \\
 &\leq (1 - \lambda) \sum_{m=0}^{\infty} \lambda^m [\alpha \sum_{t=0}^m \|(\alpha P_\pi^o)^t (\sigma_{U_\pi}(V) - \sigma_{\hat{U}_\pi}(V))\|_d] \\
 &\stackrel{(f)}{\leq} (1 - \lambda) \sum_{m=0}^{\infty} \lambda^m [\alpha \sum_{t=0}^m \beta^t \|\sigma_{U_\pi}(V) - \sigma_{\hat{U}_\pi}(V)\|_d] \\
 &\stackrel{(g)}{\leq} (1 - \lambda) \sum_{m=0}^{\infty} \lambda^m \left[\frac{(1 - \beta^{m+1})}{1 - \beta} \alpha \rho \|V\|_d \right] = \frac{\alpha \rho \|V\|_d}{1 - \beta\lambda}
 \end{aligned} \tag{39}$$

where (f) follows (35) and (g) from Lemma 1.

Now,

$$\begin{aligned}
 \|V_\pi - \Phi w_\pi\|_d &\leq \|V_\pi - \Pi V_\pi\|_d + \|\Pi V_\pi - \Phi w_\pi\|_d \\
 &\stackrel{(h)}{=} \|V_\pi - \Pi V_\pi\|_d + \|\Pi \tilde{T}_\pi^{(\lambda)} V_\pi - \Pi \tilde{T}_\pi^{(\lambda)} \Phi w_\pi\|_d \\
 &\stackrel{(i)}{\leq} \|V_\pi - \Pi V_\pi\|_d + \|\Pi \tilde{T}_\pi^{(\lambda)} V_\pi - \Pi \tilde{T}_\pi^{(\lambda)} V_\pi\|_d + \|\Pi \tilde{T}_\pi^{(\lambda)} V_\pi - \Pi \tilde{T}_\pi^{(\lambda)} \Phi w_\pi\|_d \\
 &\stackrel{(j)}{\leq} \|V_\pi - \Pi V_\pi\|_d + \frac{\beta \rho \|V_\pi\|_d}{1 - \beta\lambda} + c(\alpha, \beta, \rho, \lambda) \|V_\pi - \Phi w_\pi\|_d
 \end{aligned}$$

We get (h) because $\tilde{T}_\pi^{(\lambda)} V_\pi = V_\pi$ from Proposition 3 and $\Pi \tilde{T}_\pi^{(\lambda)} \Phi w_\pi = \Phi w_\pi$ by the premise of the proposition, (i) by triangle inequality, (j) from (39) and (10). Rearranging, we get,

$$\|V_\pi - \Phi w_\pi\|_d \leq \frac{1}{1 - c(\alpha, \beta, \rho, \lambda)} \left(\|V_\pi - \Pi V_\pi\|_d + \frac{\beta \rho \|V_\pi\|_d}{1 - \beta\lambda} \right). \tag{40}$$

This completes the proof of Theorem 1. \square

B.3. Derivation of (13)

For any bounded mapping W , observe that

$$\sum_{t=0}^{\infty} (\alpha \lambda P_\pi^o)^t W = (1 - \lambda) \sum_{m=0}^{\infty} \lambda^m \sum_{t=0}^m (\alpha P_\pi^o)^t W,$$

yielded by exchanging summations. Using this observation with equations (9) and (12) we get (13).

B.4. Proof of Theorem 2

To prove this theorem, we will use the following result from (Nedić & Bertsekas, 2003). Let $\|\cdot\|$ denote the standard Euclidean norm.

Proposition 6. (Nedić & Bertsekas, 2003, Proposition 4.1) Consider a sequence $\{x_t\}$ generated by the update equation

$$x_{t+1} = x_t + \gamma_t (h_t(x_t) + e_t),$$

where $h_t : \mathbb{R}^n \rightarrow \mathbb{R}^n$, γ_t is a positive deterministic stepsize, and e_t is a random noise vector. Let $f : \mathbb{R}^n \rightarrow \mathbb{R}$ be a continuously differentiable function. Assume the following:

(i) Function f is positive, i.e., $f(x) \geq 0$. Also, f has Lipschitz continuous gradient, i.e., there exists some scalar $L > 0$ such that

$$\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|, \quad \forall x, y \in \mathbb{R}^n. \tag{41}$$

(ii) Let $\mathcal{F}_t = \{x_0, x_1, \dots, x_t\}$. There exists positive scalars c_1, c_2 , and c_3 such that

$$(\nabla f(x_t))^\top \mathbb{E}[h_t(x_t)|\mathcal{F}_t] \leq -c_1 \|\nabla f(x_t)\|^2, \quad \forall t, \quad (42)$$

$$\|\mathbb{E}[e_t|\mathcal{F}_t]\| \leq c_2 \epsilon_t (1 + \|\nabla f(x_t)\|), \quad \forall t, \quad (43)$$

$$\mathbb{E}[\|h_t(x_t) + e_t\|^2|\mathcal{F}_t] \leq c_3 (1 + \|\nabla f(x_t)\|^2), \quad \forall t, \quad (44)$$

where ϵ_t is a positive deterministic scalar.

(ii) The deterministic sequences $\{\gamma_t\}$ and $\{\epsilon_t\}$ satisfy

$$\sum_{t=0}^{\infty} \gamma_t = \infty, \quad \sum_{t=0}^{\infty} \gamma_t^2 < \infty, \quad \sum_{t=0}^{\infty} \gamma_t \epsilon_t^2 < \infty, \quad \lim_{t \rightarrow \infty} \epsilon_t = 0. \quad (45)$$

Then, with probability 1:

(i) The sequence $\{f(x_t)\}$ converges.

(ii) The sequence $\{\nabla f(x_t)\}$ converges to zero.

(iii) Every limit point of $\{x_t\}$ is a stationary point of f .

Using the above result, we now prove Theorem 2.

Proof of Theorem 2. We rewrite (18) as $w_{t+1} = w_t + \gamma_t(h_t(w_t) + e_t)$, where,

$$\begin{aligned} h_t(w_t) &= B^{-1}(Aw_t + b + C(w_t)), \\ e_t &= B_t^{-1}(A_t w_t + b_t + C_t(w_t)) - B^{-1}(Aw_t + b + C(w_t)). \end{aligned}$$

Define

$$f(w) = \frac{1}{2}(w - w_\pi)^\top \Phi^\top D \Phi (w - w_\pi).$$

We now verify the conditions given in Proposition 6.

Verifying (41): By definition $f(w) \geq 0$. Also, $\nabla f(w) = \Phi^\top D \Phi (w - w_\pi)$. Hence,

$$\|\nabla f(w) - \nabla f(w')\| = \|\Phi^\top D \Phi (w - w')\| \leq \|\Phi^\top D \Phi\|_{\text{op}} \|w - w'\|,$$

where $\|\cdot\|_{\text{op}}$ is the operator norm corresponding to the Euclidean space. Set $L = \|\Phi^\top D \Phi\|_{\text{op}}$. From (ii) in Assumption 2 and Φ being full rank, we know that $\Phi^\top D \Phi$ is a positive definite matrix and hence L is a positive scalar. This verifies (41).

Verifying (42): It is straightforward to verify that $Aw_t + b + C(w_t) = \Phi^\top D(\tilde{T}_\pi^{(\lambda)} \Phi w_t - \Phi w_t)$ by comparing (12) - (16). Then,

$$\begin{aligned} (\nabla f(w_t))^\top \mathbb{E}[h_t(w_t)|\mathcal{F}_t] &= (w_t - w_\pi)^\top (Aw_t + b + C(w_t)) \\ &= (w_t - w_\pi)^\top \Phi^\top D(\tilde{T}_\pi^{(\lambda)} \Phi w_t - \Phi w_t) < 0, \end{aligned}$$

where the last inequality is by using Lemma 9 from (Tsitsiklis & Van Roy, 1997) with the fact that, from Theorem 1, $\Pi \tilde{T}_\pi^{(\lambda)}$ is a contraction. Since $(\nabla f(w_t))^\top \mathbb{E}[h_t|\mathcal{F}_t]$ is strictly negative, we can find a positive scalar c_2 that satisfies (42).

Verifying (43): We can write $e_t = \Delta_t^1 w_t + \Delta_t^2 + \Delta_t^3(w_t)$, where,

$$\Delta_t^1 = B_t^{-1} A_t - B^{-1} A, \quad \Delta_t^2 = B_t^{-1} b_t - B^{-1} b, \quad \Delta_t^3(w_t) = B_t^{-1} C_t(w_t) - B^{-1} C(w_t).$$

From Proposition 2.1 of (Nedić & Bertsekas, 2003), we have

$$\|\mathbb{E}[\Delta_t^1]\| \leq \frac{\bar{c}_1}{\sqrt{t+1}}, \quad \|\mathbb{E}[\Delta_t^2]\| \leq \frac{\bar{c}_2}{\sqrt{t+1}} \quad \forall t. \quad (46)$$

So, we will now bound $\mathbb{E}[\Delta_t^3(w_t)|\mathcal{F}_t]$.

$$\begin{aligned} \|\mathbb{E}[\Delta_t^3(w_t)|\mathcal{F}_t]\| &= \|\mathbb{E}[B_t^{-1}C_t(w_t) - B^{-1}C(w_t)|\mathcal{F}_t]\| \\ &= \|\mathbb{E}[B_t^{-1}C_t(w_t) - B^{-1}C_t(w_t) + B^{-1}C_t(w_t) - B^{-1}C(w_t)|\mathcal{F}_t]\| \\ &\leq \mathbb{E}[\|(B_t^{-1} - B^{-1})\| \|C_t(w_t)\| | \mathcal{F}_t] + \|B^{-1}\| \mathbb{E}[\|C_t(w_t) - C(w_t)\| | \mathcal{F}_t]. \end{aligned} \quad (47)$$

First consider $\|C_t(w_t)\|$. We can bound

$$\|z_\tau\| \leq \sum_{m=0}^{\tau} (\alpha\lambda)^{\tau-m} \|\phi(s_m)\| \stackrel{(a)}{\leq} \bar{c}_{31} \quad (48)$$

$$\begin{aligned} |\sigma_{\hat{U}_{s_\tau, \pi(s_\tau)}}(\Phi w_t)| &= \left| \inf_{u \in \hat{U}_{s_\tau, \pi(s_\tau)}} u^\top (\Phi w_t) \right| \leq \sup_{u \in \hat{U}_{s_\tau, \pi(s_\tau)}} |u^\top (\Phi w_t)| \\ &\leq \sup_{u \in \hat{U}_{s_\tau, \pi(s_\tau)}} \|u\| \|\Phi w_t\| \stackrel{(b)}{\leq} \bar{c}_{32} \|w_t\| \end{aligned} \quad (49)$$

where \bar{c}_{31} and \bar{c}_{32} are finite positive scalars. For (a), we used the fact that $\sup_s \|\phi(s)\| < \infty$. For (b), we used the fact that $\sup_{u \in \hat{U}_{s, \pi(s)}} \|u\| < \infty$ since $\hat{U}_{s, \pi(s)}$ is a finite set for any $s \in \mathcal{S}$ and $\|\Phi w_t\| \leq c' \|w_t\|$ for some finite positive scalar c' . Using (48) and (49),

$$\|C_t(w_t)\| \leq \frac{\alpha}{(t+1)} \sum_{\tau=0}^t \|z_\tau\| \sup_{u \in \hat{U}_{s_\tau, \pi(s_\tau)}} |\sigma_{\hat{U}_{s_\tau, \pi(s_\tau)}}(\Phi w_t)| \leq \frac{\alpha}{(t+1)} \sum_{\tau=0}^t \|z_\tau\| \bar{c}_{32} \|w_t\|. \quad (50)$$

From Lemma 4.3 of (Nedić & Bertsekas, 2003), we have $\mathbb{E}[\|(B_t^{-1} - B^{-1})\|] \leq \bar{c}_{34}/\sqrt{(t+1)}$. Using this and (50), we get,

$$\mathbb{E}[\|(B_t^{-1} - B^{-1})\| \|C_t(w_t)\|] \leq \frac{\bar{c}_{35}}{\sqrt{t+1}} \|w_t\|. \quad (51)$$

For bounding $\mathbb{E}[\|C_t(w_t) - C(w_t)\|]$ in (47), we define V_t and n_t as,

$$V_t = \frac{\alpha}{(t+1)} \sum_{m=0}^t \phi(s_m) \sum_{\tau=t+1}^{\infty} \left[(\alpha\lambda P_\pi^\circ)^{\tau-m} \sigma_{\hat{U}_\pi}(\Phi w_t) \right] (s_m), \quad n_t(s) = \sum_{\tau=0}^t \mathbb{I}\{s_\tau = s\}.$$

Note that $n_t(s)$ is the number of visits to state s until time t . Then,

$$\begin{aligned} \mathbb{E}[C_t(w_t)|\mathcal{F}_t] &= \mathbb{E} \left[\frac{\alpha}{(t+1)} \sum_{\tau=0}^t \sum_{m=0}^{\tau} (\alpha\lambda)^{\tau-m} \phi(s_m) \sigma_{\hat{U}_{s_\tau, \pi(s_\tau)}}(\Phi w_t) | \mathcal{F}_t \right] \\ &= \mathbb{E} \left[\frac{\alpha}{(t+1)} \sum_{\tau=0}^t \sum_{m=0}^{\tau} (\alpha\lambda)^{\tau-m} \phi(s_m) \mathbb{E} \left[\sigma_{\hat{U}_{s_\tau, \pi(s_\tau)}}(\Phi w_t) | \mathcal{F}_m \right] | \mathcal{F}_t \right] \\ &\stackrel{(a)}{=} \mathbb{E} \left[\frac{\alpha}{(t+1)} \sum_{\tau=0}^t \sum_{m=0}^{\tau} (\alpha\lambda)^{\tau-m} \phi(s_m) ((P_\pi^\circ)^{\tau-m} \sigma_{\hat{U}_\pi}(\Phi w_t))(s_m) | \mathcal{F}_t \right] \\ &\stackrel{(b)}{=} \mathbb{E} \left[\frac{\alpha}{(t+1)} \sum_{m=0}^t \phi(s_m) \sum_{\tau=m}^t ((\alpha\lambda P_\pi^\circ)^{\tau-m} \sigma_{\hat{U}_\pi}(\Phi w_t))(s_m) | \mathcal{F}_t \right] \\ &= \mathbb{E} \left[\frac{\alpha}{(t+1)} \sum_{m=0}^t \phi(s_m) \sum_{\tau=m}^{\infty} ((\alpha\lambda P_\pi^\circ)^{\tau-m} \sigma_{\hat{U}_\pi}(\Phi w_t))(s_m) | \mathcal{F}_t \right] - \mathbb{E}[V_t | \mathcal{F}_t] \\ &\stackrel{(c)}{=} \mathbb{E} \left[\frac{\alpha}{(t+1)} \sum_{s \in \mathcal{S}} n_t(s) \phi(s) \sum_{i=0}^{\infty} ((\alpha\lambda P_\pi^\circ)^i \sigma_{\hat{U}_\pi}(\Phi w_t))(s) | \mathcal{F}_t \right] - \mathbb{E}[V_t | \mathcal{F}_t], \end{aligned} \quad (52)$$

where (a) follows since the transition probability matrix governing along policy π is P_π^o , (b) by exchanging the order of summation, and (c) by using the definition of $n_t(s)$. Note that,

$$\begin{aligned}\mathbb{E}[C(w_t)|\mathcal{F}_t] &= \mathbb{E}[\alpha\Phi^\top D \sum_{i=0}^{\infty} (\alpha\lambda P_\pi^o)^i \sigma_{\hat{U}_\pi}(\Phi w_t)|\mathcal{F}_t] \\ &= \mathbb{E}[\alpha \sum_{s \in \mathcal{S}} d_s \phi(s) \sum_{i=0}^{\infty} ((\alpha\lambda P_\pi^o)^i \sigma_{\hat{U}_\pi}(\Phi w_t))(s)|\mathcal{F}_t].\end{aligned}\quad (53)$$

So, using (53) and (52),

$$\begin{aligned}& \|\mathbb{E}[(C_t(w_t) - C(w_t)) | \mathcal{F}_t]\| \\ &= \left\| \alpha \sum_{s \in \mathcal{S}} \left(\frac{\mathbb{E}[n_t(s)]}{t+1} - d_s \right) \phi(s) \sum_{i=0}^{\infty} ((\alpha\lambda P_\pi^o)^i \sigma_{\hat{U}_\pi}(\Phi w_t))(s) - \mathbb{E}[V_t | \mathcal{F}_t] \right\| \\ &\stackrel{(d)}{\leq} \bar{c}_{36} \sum_{s \in \mathcal{S}} \left| \left(\frac{\mathbb{E}[n_t(s)]}{t+1} - d_s \right) \right| \|w_t\| + \mathbb{E}[\|V_t\| | \mathcal{F}_t] \stackrel{(e)}{\leq} \frac{\bar{c}_{37}}{t+1} \|w_t\| + \frac{\bar{c}_{38}}{t+1} \|w_t\|\end{aligned}\quad (54)$$

For getting (d), note that $\|\phi(s) \sum_{i=0}^{\infty} ((\alpha\lambda P_\pi^o)^i \sigma_{\hat{U}_\pi}(\Phi w_t))(s)\| \leq c' \|w_t\|$ for some positive constant c' , using (49) and the fact that the summation is bounded due to the discounted factor $(\alpha\lambda)$. For getting (e), we use the result from Lemma 4.2 (Nedić & Bertsekas, 2003) that $\left| \frac{\mathbb{E}[n_t(s)]}{t+1} - d_s \right| \leq \frac{c'}{t+1}$ for some positive number c' . Also, it is straightforward to show that $\mathbb{E}[\|V_t\| | \mathcal{F}_t] \leq c' \frac{\|w_t\|}{t+1}$ for some positive number c' using (49) and the fact that the summation is bounded due to the discounted factor $(\alpha\lambda)$.

Using (51) and (54) in (47), we get

$$\|\mathbb{E}[\Delta_t^3(w_t)]\| \leq \frac{\bar{c}_3}{\sqrt{t+1}} \|w_t\|. \quad (55)$$

Notice that

$$\begin{aligned}\|w_t\| &\leq \|w_t - w_\pi\| + \|w_\pi\| \leq \|(\Phi^\top D\Phi)^{-1}\| \|(\Phi^\top D\Phi)(w_t - w_\pi)\| + \|w_\pi\| \\ &\leq \bar{c}_{39}(1 + \|\nabla f(w_t)\|).\end{aligned}\quad (56)$$

Now,

$$\begin{aligned}\|\mathbb{E}[e_t | \mathcal{F}_t]\| &= \|\mathbb{E}[\Delta_t^1]\| \|w_t\| + \|\mathbb{E}[\Delta_t^2]\| + \|\mathbb{E}[\Delta_t^3(w_t) | \mathcal{F}_t]\| \\ &\stackrel{(f)}{\leq} \frac{\bar{c}_1}{\sqrt{t+1}} \|w_t\| + \frac{\bar{c}_2}{\sqrt{t+1}} + \frac{\bar{c}_3}{\sqrt{t+1}} \|w_t\| \stackrel{(g)}{\leq} \frac{c_2}{\sqrt{t+1}} (1 + \|\nabla f(w_t)\|)\end{aligned}\quad (57)$$

where (f) is by using (46) and (55) and (g) is by using (56). This completes the verification of the condition (43).

Verifying (44): From definition, we have $h_t(w_t) + e_t = B_t^{-1}(A_t w_t + b_t + C_t(w_t))$, for all t . Using similar steps as before, it is straightforward to show that $\|B_t\| \leq c_{41}$, $|A_t| \leq c_{42}$, $|b_t| \leq c_{43}$ for some positive scalars c_{41}, c_{42}, c_{43} . From (50) we have $\|c_t(w_t)\| \leq \bar{c}_{33} \|w_t\|$. Combining these, we get

$$\|h_t(w_t) + e_t\| \leq c_{44}(1 + \|w_t\|). \quad (58)$$

Now, from (56) and (41), $\|h_t(w_t) + e_t\|^2 \leq c_{45}(1 + \|\nabla f(w_t)\|)^2$ for all t . So, finally we have,

$$\mathbb{E}[\|h_t(w_t) + e_t\|^2 | \mathcal{F}_t] = \mathbb{E}[\|h_t(w_t) + e_t\|^2 | w_t] \leq c_3(1 + \|\nabla f(w_t)\|)^2 \quad \forall t,$$

showing that (44) is satisfied.

Verifying (45): From (57), $\epsilon_t = 1/\sqrt{t+1}$. So, this condition is satisfied. ‘ So, all the assumption of Proposition 6 are satisfied. Hence the result of that Proposition is true. In particular, $\nabla f(w_t) = \Phi^\top D\Phi(w_t - w_\pi)$ converges to 0. Since $\Phi^\top D\Phi$ is positive definite, this implies that $w_t \rightarrow w_\pi$. \square

C. Proof of the Results in Section 4

Let $\bar{V}_k = \Phi w_{\pi_k}$ and V^* be the optimal robust value function. Define

$$e_k = V_{\pi_k} - \bar{V}_k, \quad l_k = V^* - V_{\pi_k}, \quad g_k = V_{\pi_{k+1}} - V_{\pi_k}. \quad (59)$$

Interpretations of these expressions: Since the robust value function in the k^{th} iteration \bar{V}_k is used as a surrogate for the robust value function V_{π_k} , e_k quantifies the *approximation error*. g_k signifies the *gain* of value functions between iterations k and $k+1$. Finally, l_k encapsulates the *loss* in the value function because of using policy π_k instead of the optimal policy.

Let $|x|$ denote element-wise absolute values of vector $x \in \mathbb{R}^{|\mathcal{S}|}$. We first prove the following result. This parallels to the result for the nonrobust setting in (Munos, 2003).

Lemma 4. *We have*

$$|l_{k+1}| \leq c(\alpha, \beta, 0, \lambda) \bar{H}_*(|l_k| + |e_k|) + c(\alpha, \beta, 0, \lambda) \bar{H}_k(|g_k| + |e_k|) \quad \text{and} \quad (60)$$

$$|g_k| \leq c(\alpha, \beta, 0, \lambda) (I - c(\alpha, \beta, 0, \lambda) \bar{H}_{k+1})^{-1} (\bar{H}_{k+1} + \bar{H}_k) |e_k| \quad (61)$$

where the stochastic matrices \bar{H}_* , \bar{H}_k , and \bar{H}_{k+1} are defined in (64)-(65).

Proof. As before, denote the operator $\tilde{T}_\pi^{(\lambda)}$ as $\bar{T}_\pi^{(\lambda)}$ when $\hat{U}_\pi = U_\pi$. Now, similar to (36), for any policy π and V_1, V_2 , we get that

$$\begin{aligned} & |\bar{T}_\pi^{(\lambda)}(V_1) - \bar{T}_\pi^{(\lambda)}(V_2)| \\ &= |(1-\lambda) \sum_{m=0}^{\infty} \lambda^m [\alpha \sum_{t=0}^m (\alpha P_\pi^o)^t (\sigma_{U_\pi}(V_1) - \sigma_{U_\pi}(V_2)) + (\alpha P_\pi^o)^{m+1} (V_1 - V_2)]| \\ &\leq (1-\lambda) \sum_{m=0}^{\infty} \lambda^m [\alpha \sum_{t=0}^m |(\alpha P_\pi^o)^t| |\sigma_{U_\pi}(V_1) - \sigma_{U_\pi}(V_2)| + |(\alpha P_\pi^o)^{m+1}| |V_1 - V_2|] \\ &\stackrel{(a)}{\leq} (1-\lambda) \sum_{m=0}^{\infty} \lambda^m [\alpha \sum_{t=0}^m (\beta P_{\pi_e}^o)^t |\sigma_{U_\pi}(V_1) - \sigma_{U_\pi}(V_2)| + |(\beta P_{\pi_e}^o)^{m+1}| |V_1 - V_2|] \\ &\stackrel{(b)}{\leq} (1-\lambda) \sum_{m=0}^{\infty} \lambda^m [\beta \sum_{t=0}^m (\beta P_{\pi_e}^o)^t |V_1 - V_2| + |(\beta P_{\pi_e}^o)^{m+1}| |V_1 - V_2|] \end{aligned} \quad (62)$$

where (a) follows from Assumption 2 and π_e (dependent on π) being the exploration policy. (b) follows from (32) in Lemma 3.

Recall that the optimal robust value function V^* and the optimal robust policy π^* satisfy the equation $\bar{T}_{\pi^*}^{(\lambda)} V^* = V^*$. Using this,

$$\begin{aligned} l_{k+1} &= V^* - V_{\pi_{k+1}} = \bar{T}_{\pi^*}^{(\lambda)} V^* - \bar{T}_{\pi_{k+1}}^{(\lambda)} V_{\pi_{k+1}} \\ &= (\bar{T}_{\pi^*}^{(\lambda)} V^* - \bar{T}_{\pi^*}^{(\lambda)} V_{\pi_k}) + (\bar{T}_{\pi^*}^{(\lambda)} V_{\pi_k} - \bar{T}_{\pi^*}^{(\lambda)} \bar{V}_k) + (\bar{T}_{\pi^*}^{(\lambda)} \bar{V}_k - \bar{T}_{\pi_{k+1}}^{(\lambda)} \bar{V}_k) \\ &\quad + (\bar{T}_{\pi_{k+1}}^{(\lambda)} \bar{V}_k - \bar{T}_{\pi_{k+1}}^{(\lambda)} V_{\pi_k}) + (\bar{T}_{\pi_{k+1}}^{(\lambda)} V_{\pi_k} - \bar{T}_{\pi_{k+1}}^{(\lambda)} V_{\pi_{k+1}}) \\ &\stackrel{(c)}{\leq} (\bar{T}_{\pi^*}^{(\lambda)} V^* - \bar{T}_{\pi^*}^{(\lambda)} V_{\pi_k}) + (\bar{T}_{\pi^*}^{(\lambda)} V_{\pi_k} - \bar{T}_{\pi^*}^{(\lambda)} \bar{V}_k) \\ &\quad + (\bar{T}_{\pi_{k+1}}^{(\lambda)} \bar{V}_k - \bar{T}_{\pi_{k+1}}^{(\lambda)} V_{\pi_k}) + (\bar{T}_{\pi_{k+1}}^{(\lambda)} V_{\pi_k} - \bar{T}_{\pi_{k+1}}^{(\lambda)} V_{\pi_{k+1}}) \\ &\stackrel{(d)}{\leq} (1-\lambda) \sum_{m=0}^{\infty} \lambda^m [\beta \sum_{t=0}^m (\beta P_{\pi^*}^o)^t (|l_k| + |e_k|) + (\beta P_{\pi^*}^o)^{m+1} (|l_k| + |e_k|)] \\ &\quad + (1-\lambda) \sum_{m=0}^{\infty} \lambda^m [\beta \sum_{t=0}^m (\beta P_{\pi_{k+1}}^o)^t (|g_k| + |e_k|) + (\beta P_{\pi_{k+1}}^o)^{m+1} (|g_k| + |e_k|)] \\ &\stackrel{(e)}{=} c(\alpha, \beta, 0, \lambda) \bar{H}_*(|l_k| + |e_k|) + c(\alpha, \beta, 0, \lambda) \bar{H}_{k+1} (|g_k| + |e_k|) \end{aligned} \quad (63)$$

Here (c) follows because π_{k+1} is the greedy policy w.r.t. \bar{V}_k and hence $\bar{T}_{\pi^*}^{(\lambda)} \bar{V}_k - \bar{T}_{\pi_{k+1}}^{(\lambda)} \bar{V}_k \leq 0$. (d) follows from (62) noting (i) in Assumption 3. Finally, (e) follows by taking

$$\bar{H}_* = \frac{1}{c(\alpha, \beta, 0, \lambda)} (1 - \lambda) \sum_{m=0}^{\infty} \lambda^m [\beta \sum_{t=0}^m (\beta P_{\pi^*}^o)^t + (\beta P_{\pi^*}^o)^{m+1}], \text{ and} \quad (64)$$

$$\bar{H}_j = \frac{1}{c(\alpha, \beta, 0, \lambda)} (1 - \lambda) \sum_{m=0}^{\infty} \lambda^m [\beta \sum_{t=0}^m (\beta P_{\pi_j}^o)^t + (\beta P_{\pi_j}^o)^{m+1}], \text{ for } j \geq 1. \quad (65)$$

Note that, matrices \bar{H}_* and \bar{H}_j are stochastic matrices. This follows easily by verifying $\bar{H}_* \mathbf{1} = \mathbf{1}$, $\bar{H}_j \mathbf{1} = \mathbf{1}$ using algebra analysis as in (37).

The same argument can be repeated to get

$$-l_{k+1} \leq c(\alpha, \beta, 0, \lambda) \bar{H}_*(|l_k| + |e_k|) + c(\alpha, \beta, 0, \lambda) \bar{H}_{k+1}(|g_k| + |e_k|).$$

Combining, we get

$$|l_{k+1}| \leq c(\alpha, \beta, 0, \lambda) \bar{H}_*(|l_k| + |e_k|) + c(\alpha, \beta, 0, \lambda) \bar{H}_{k+1}(|g_k| + |e_k|).$$

Now,

$$\begin{aligned} g_k &= \bar{T}_{\pi_{k+1}}^{(\lambda)} V_{\pi_{k+1}} - \bar{T}_{\pi_k}^{(\lambda)} V_{\pi_k} \\ &= \bar{T}_{\pi_{k+1}}^{(\lambda)} V_{\pi_{k+1}} - \bar{T}_{\pi_{k+1}}^{(\lambda)} V_{\pi_k} + \bar{T}_{\pi_{k+1}}^{(\lambda)} V_{\pi_k} - \bar{T}_{\pi_{k+1}}^{(\lambda)} \bar{V}_k + (\bar{T}_{\pi_{k+1}}^{(\lambda)} \bar{V}_k - \bar{T}_{\pi_k}^{(\lambda)} \bar{V}_k) + \bar{T}_{\pi_k}^{(\lambda)} \bar{V}_k - \bar{T}_{\pi_k}^{(\lambda)} V_{\pi_k} \\ &\geq \bar{T}_{\pi_{k+1}}^{(\lambda)} V_{\pi_{k+1}} - \bar{T}_{\pi_{k+1}}^{(\lambda)} V_{\pi_k} + \bar{T}_{\pi_{k+1}}^{(\lambda)} V_{\pi_k} - \bar{T}_{\pi_{k+1}}^{(\lambda)} \bar{V}_k + \bar{T}_{\pi_k}^{(\lambda)} \bar{V}_k - \bar{T}_{\pi_k}^{(\lambda)} V_{\pi_k} \end{aligned}$$

where the last inequality follows because π_{k+1} is the greedy policy w.r.t. \bar{V}_k and hence $(\bar{T}_{\pi_{k+1}}^{(\lambda)} \bar{V}_k - \bar{T}_{\pi_k}^{(\lambda)} \bar{V}_k) \geq 0$. From (62), noting (i) in Assumption 3, we have

$$\begin{aligned} -g_k &\leq (1 - \lambda) \sum_{m=0}^{\infty} \lambda^m [\beta \sum_{t=0}^m (\beta P_{\pi_{k+1}}^o)^t (|g_k| + |e_k|) + (\beta P_{\pi_{k+1}}^o)^{m+1} (|g_k| + |e_k|)] \\ &\quad + (1 - \lambda) \sum_{m=0}^{\infty} \lambda^m [\beta \sum_{t=0}^m (\beta P_{\pi_k}^o)^t (|e_k|) + (\beta P_{\pi_k}^o)^{m+1} (|e_k|)] \\ &\leq c(\alpha, \beta, 0, \lambda) \bar{H}_{k+1}(|g_k| + |e_k|) + c(\alpha, \beta, 0, \lambda) \bar{H}_k(|e_k|). \end{aligned}$$

Repeating the above argument for $-g_k = \bar{T}_{\pi_k}^{(\lambda)} V_{\pi_k} - \bar{T}_{\pi_{k+1}}^{(\lambda)} V_{\pi_{k+1}}$, we get

$$|g_k| \leq c(\alpha, \beta, 0, \lambda) \bar{H}_{k+1}(|g_k| + |e_k|) + c(\alpha, \beta, 0, \lambda) \bar{H}_k(|e_k|).$$

So, $|g_k| \leq c(\alpha, \beta, 0, \lambda) (I - c(\alpha, \beta, 0, \lambda) \bar{H}_{k+1})^{-1} (\bar{H}_{k+1} + \bar{H}_k) |e_k|$. Thus, proving the lemma. \square

Proof of Theorem 3. From the Lemma 4, taking lim sup on both sides of (60) we have

$$\begin{aligned} \limsup_{k \rightarrow \infty} |l_k| &\leq \limsup_{k \rightarrow \infty} c(\alpha, \beta, 0, \lambda) (I - c(\alpha, \beta, 0, \lambda) \bar{H}_*)^{-1} (\bar{H}_* + \bar{H}_k) |e_k| \\ &\quad + c(\alpha, \beta, 0, \lambda) (I - c(\alpha, \beta, 0, \lambda) \bar{H}_*)^{-1} \bar{H}_k |g_k| \\ &\stackrel{(a)}{\leq} \limsup_{k \rightarrow \infty} c(\alpha, \beta, 0, \lambda) (I - c(\alpha, \beta, 0, \lambda) \bar{H}_*)^{-1} (\bar{H}_* + \bar{H}_k) |e_k| \\ &\quad + c^2(\alpha, \beta, 0, \lambda) (I - c(\alpha, \beta, 0, \lambda) \bar{H}_*)^{-1} \bar{H}_k (I - c(\alpha, \beta, 0, \lambda) \bar{H}_{k+1})^{-1} (\bar{H}_{k+1} + \bar{H}_k) |e_k| \\ &\stackrel{(b)}{=} \frac{2c(\alpha, \beta, 0, \lambda)}{(1 - c(\alpha, \beta, 0, \lambda))^2} \limsup_{k \rightarrow \infty} H_k |e_k| \end{aligned} \quad (66)$$

where (a) follows from (61) in Lemma 4. (b) follows by taking

$$H_k = (1 - c(\alpha, \beta, 0, \lambda))^2 (I - c(\alpha, \beta, 0, \lambda) \bar{H}_*)^{-1} \left(\frac{(\bar{H}_* + \bar{H}_k)}{2} + c(\alpha, \beta, 0, \lambda) \bar{H}_k (I - c(\alpha, \beta, 0, \lambda) \bar{H}_{k+1})^{-1} \frac{(\bar{H}_{k+1} + \bar{H}_k)}{2} \right). \quad (67)$$

Notice that H_k is a stochastic matrix. To see this, we know that, if P_i , $i = \{1, 2, 3, 4\}$ are stochastic matrices and $c < 1$, then $P_1 P_2 P_3 P_4$, $(P_1 + P_2)/2$, and $(1 - c)(I - cP_1)^{-1}$ are valid stochastic matrices as well. Now, it is easy to verify that $H_k \mathbf{1} = \mathbf{1}$. Then, $\mu_k = \mu H_k$ is a valid probability distribution.

Let x^2 denote element-wise squares of vector $x \in \mathbb{R}^{|\mathcal{S}|}$ and also let $\|x\|_\mu^2 = \int_{s \in \mathcal{S}} |x(s)|^2 d\mu(s) = \mu|x|^2$. Now, from (66) we have

$$\begin{aligned} \limsup_{k \rightarrow \infty} \|l_k\|_\mu^2 &= \limsup_{k \rightarrow \infty} \mu |l_k|^2 \leq \frac{4c^2(\alpha, \beta, 0, \lambda)}{(1 - c(\alpha, \beta, 0, \lambda))^4} \limsup_{k \rightarrow \infty} \mu [H_k |e_k|]^2 \\ &\stackrel{(c)}{\leq} \frac{4c^2(\alpha, \beta, 0, \lambda)}{(1 - c(\alpha, \beta, 0, \lambda))^4} \limsup_{k \rightarrow \infty} \mu H_k |e_k|^2 \\ &= \frac{4c^2(\alpha, \beta, 0, \lambda)}{(1 - c(\alpha, \beta, 0, \lambda))^4} \limsup_{k \rightarrow \infty} \mu_k |e_k|^2 = \frac{4c^2(\alpha, \beta, 0, \lambda)}{(1 - c(\alpha, \beta, 0, \lambda))^4} \limsup_{k \rightarrow \infty} \|e_k\|_{\mu_k}^2 \\ &\stackrel{(d)}{\leq} \frac{4C_1 C_2 c^2(\alpha, \beta, 0, \lambda)}{(1 - c(\alpha, \beta, 0, \lambda))^4} \limsup_{k \rightarrow \infty} \|e_k\|_{d_{\pi_k}}^2. \end{aligned}$$

Here (c) follows from Jensen's inequality. To see this, let

$$H_k = \begin{pmatrix} -q_1 - \\ -q_2 - \\ \vdots \\ -q_{|\mathcal{S}|} - \end{pmatrix} \text{ where } q_i \in \mathbb{R}^{|\mathcal{S}|}, \text{ for all } i, \text{ are probability vectors.}$$

For any $x \in \mathbb{R}^{|\mathcal{S}|}$, let $x(j)$ denote the j^{th} coordinate value in x . Now, for each $i \in \{1, 2, \dots, |\mathcal{S}|\}$, define $|\mathcal{S}|$ -discrete valued random variable X_i such that it takes value $|e_k|(j)$ with probability $q_i(j)$ for all $j \in \{1, 2, \dots, |\mathcal{S}|\}$. Thus, from Jensen's inequality, we have

$$[H_k |e_k|]^2 = ((\mathbb{E}[X_1])^2, (\mathbb{E}[X_2])^2, \dots, (\mathbb{E}[X_{|\mathcal{S}|}])^2)^\top \leq (\mathbb{E}[X_1^2], \mathbb{E}[X_2^2], \dots, \mathbb{E}[X_{|\mathcal{S}|}^2])^\top = H_k |e_k|^2. \quad (68)$$

(d) follows by noting that for any $x \in \mathbb{R}^{|\mathcal{S}|}$, from (iv) in Assumption 3, we have $\|x\|_{\mu_k}^2 \leq C_1 \|x\|_\mu^2 \leq C_1 C_2 \|x\|_{d_{\pi_k}}^2$ for all k . Thus proving (26) of this theorem.

Now, since $\rho = 0$, (27) of this theorem directly follows from (iii) in Assumption 3 and (11) in Theorem 1. This completes the proof of this theorem. \square

Now we make an alternative assumption to Assumption 3.(iv) to get a guarantee for any K^{th} iteration of the RLSPI algorithm.

Assumption 4. For an arbitrary sequence of stationary policies $\{\pi_i\}_{i \geq 1}$, let some probability distributions μ and $\bar{\mu}$ satisfy

$$C_{\mu, \bar{\mu}} = (1 - c(\alpha, \beta, 0, \lambda)) \sum_{m \geq 1} (c(\alpha, \beta, 0, \lambda))^m c_{\mu, \bar{\mu}}(m) < \infty, \quad (69)$$

where for any given $m \geq 1$ the coefficients are defined as

$$c_{\mu, \bar{\mu}}(m) = \sup_{\pi_1, \dots, \pi_m} \left\| \frac{d(\mu \bar{H}_1 \bar{H}_2 \dots \bar{H}_m)}{d\bar{\mu}} \right\|_\infty, \quad (70)$$

and \bar{H}_k is a stochastic matrix that depends on $P_{\pi_k}^o$ for any $1 \leq k \leq m$. Also assume that $d_{\pi_k} \geq \bar{\mu}/C$ for all k .

We note that $c_{\mu, \bar{\mu}}(m)$ can potentially diverge to ∞ , but $C_{\mu, \bar{\mu}}$ is finite as long as $(c(\alpha, \beta, 0, \lambda))^m$ converges to 0 at a faster rate. Assumption 4 being similar as in the non-robust setting, we refer the reader to (Antos et al., 2008; Lazaric et al., 2012) for its detailed interpretation.

Here is a result that provides a guarantee for the performance of the policy learned in K^{th} iteration of the RLSPI algorithm.

Theorem 4. *Let Assumption 2, Assumption 3.(i)-(iii), and Assumption 4 hold. Let the range of reward function r be $(0, R_{\max}]$. Let $\{\pi_k\}_{k=1}^K$ be the sequence of policies generated by the RLSPI algorithm for some $K \geq 1$. Let V_{π_k} and $\bar{V}_k = \Phi w_{\pi_k}$ be true robust value function and the approximate robust value function corresponding to the policy π_k . Also, let V^* be the optimal robust value function. Then, with $c(\alpha, \beta, 0, \lambda) < 1$,*

$$\|V^* - V_{\pi_K}\|_{\mu} \leq \frac{2\sqrt{2}(c(\alpha, \beta, 0, \lambda))^{(K+1)/2}}{(1 - c(\alpha, \beta, 0, \lambda))^{3/2}} R_{\max} + \frac{2\sqrt{2}c(\alpha, \beta, 0, \lambda)\sqrt{CC_{\mu, \bar{\mu}}}}{(1 - c(\alpha, \beta, 0, \lambda))^2} \max_{0 \leq k < K} \|V_{\pi_k} - \bar{V}_k\|_{d_{\pi_k}}. \quad (71)$$

Moreover, from Theorem 1 and Assumption 3.(iii), as $K \rightarrow \infty$ we have

$$\|V^* - V_{\pi_K}\|_{\mu} \leq \frac{2\sqrt{2}c(\alpha, \beta, 0, \lambda)\sqrt{CC_{\mu, \bar{\mu}}}}{(1 - c(\alpha, \beta, 0, \lambda))^3} \delta. \quad (72)$$

Proof of Theorem 4. Let $\zeta = c(\alpha, \beta, 0, \lambda)$. From the Lemma 4, using (61) in (60) we have

$$\begin{aligned} |l_{k+1}| &\leq \zeta \bar{H}_* |l_k| + \zeta (\bar{H}_* + \bar{H}_k) |e_k| + \zeta \bar{H}_k |g_k| \\ &\leq \zeta \bar{H}_* |l_k| + \zeta (\bar{H}_* + \bar{H}_k) |e_k| + \zeta^2 \bar{H}_k (I - \zeta \bar{H}_{k+1})^{-1} (\bar{H}_{k+1} + \bar{H}_k) |e_k| \\ &\stackrel{(a)}{=} \zeta \bar{H}_* |l_k| + \frac{2\zeta}{1 - \zeta} \bar{H}_k |e_k|, \end{aligned} \quad (73)$$

where (a) follows since

$$H_k = (1 - \zeta) \left(\frac{(\bar{H}_* + \bar{H}_k)}{2} + \zeta \bar{H}_k (I - \zeta \bar{H}_{k+1})^{-1} \frac{(\bar{H}_{k+1} + \bar{H}_k)}{2} \right). \quad (74)$$

Taking $(K - 1)$ -recursions of (73) we get

$$|l_K| \leq (\zeta \bar{H}_*)^K |l_0| + \frac{2\zeta}{1 - \zeta} \sum_{k=0}^{K-1} (\zeta \bar{H}_*)^{K-k-1} \bar{H}_k |e_k|. \quad (75)$$

Note that since $\alpha \leq \beta$, we have that $\alpha \leq \zeta$. Now, since the rewards are in $(0, R_{\max}]$, we also have that $|l_0| = |V^* - V_{\pi_0}| \leq \frac{2R_{\max}}{(1-\alpha)} \mathbb{1} \leq \frac{2R_{\max}}{(1-\zeta)} \mathbb{1}$. Thus, bounding (75) further we get

$$|l_K| \leq \frac{2\zeta K \zeta}{(1 - \zeta)^2} \left[\alpha_K G_K R_{\max} \mathbb{1} + \sum_{k=0}^{K-1} \alpha_k G_k |e_k| \right], \quad (76)$$

where $\zeta_K = \zeta^{K-1}(2 - \zeta - \zeta^K) \leq 2$, the positive coefficients α_k are

$$\alpha_k = \frac{\zeta^{K-k-1}(1 - \zeta)}{\zeta_K}, \text{ for } 0 \leq k \leq K - 1, \text{ and } \alpha_K = \frac{\zeta^{K-1}(1 - \zeta)}{\zeta_K}, \quad (77)$$

and the operators G_k are

$$G_k = (\bar{H}_*)^{K-k-1} H_k, \text{ for } 0 \leq k \leq K - 1, \text{ and } G_K = (\bar{H}_*)^K.$$

Note that $\sum_{k=0}^K \alpha_k = 1$ and G_k for $0 \leq k \leq K$ are stochastic matrices.

Now, from (76) we have

$$\begin{aligned}
 \|l_K\|_\mu^2 &\stackrel{(b)}{\leq} \frac{4\zeta_K^2\zeta^2}{(1-\zeta)^4} \mu \left[\alpha_K (G_K R_{\max} \mathbb{1})^2 + \sum_{k=0}^{K-1} \alpha_k (G_k |e_k|)^2 \right] \\
 &\stackrel{(c)}{\leq} \frac{4\zeta_K^2\zeta^2}{(1-\zeta)^4} \mu \left[\alpha_K G_K R_{\max}^2 \mathbb{1} + \sum_{k=0}^{K-1} \alpha_k G_k |e_k|^2 \right] \\
 &= \frac{4\zeta_K^2\zeta^2}{(1-\zeta)^4} \left[\alpha_K R_{\max}^2 + \mu \sum_{k=0}^{K-1} \alpha_k G_k |e_k|^2 \right] \\
 &\stackrel{(d)}{\leq} \frac{4\zeta_K^2\zeta^2}{(1-\zeta)^4} \left[\alpha_K R_{\max}^2 + (1-\zeta) \sum_{k=0}^{K-1} \alpha_k \sum_{m \geq 0} \zeta^m c_{\mu, \bar{\mu}}(m+K-k) \|e_k\|_{\bar{\mu}}^2 \right] \\
 &\stackrel{(e)}{\leq} \frac{4\zeta_K^2\zeta^2}{(1-\zeta)^4} \left[\alpha_K R_{\max}^2 + (1-\zeta) \|e\|_{\bar{\mu}}^2 \sum_{k=0}^{K-1} \alpha_k \sum_{m \geq 0} \zeta^m c_{\mu, \bar{\mu}}(m+K-k) \right] \\
 &\stackrel{(f)}{\leq} \frac{4\zeta_K^2\zeta^2}{(1-\zeta)^4} \left[\zeta^{K-1} (1-\zeta) R_{\max}^2 + (1-\zeta) \|e\|_{\bar{\mu}}^2 \sum_{k=0}^{K-1} \zeta^{K-k-1} (1-\zeta) \sum_{m \geq 0} \zeta^m c_{\mu, \bar{\mu}}(m+K-k) \right] \\
 &\stackrel{(g)}{\leq} \frac{8\zeta^2}{(1-\zeta)^4} [\zeta^{K-1} (1-\zeta) R_{\max}^2 + \|e\|_{\bar{\mu}}^2 C_{\mu, \bar{\mu}}] \\
 &= \frac{8\zeta^{K+1}}{(1-\zeta)^3} R_{\max}^2 + \frac{8\zeta^2}{(1-\zeta)^4} \|e\|_{\bar{\mu}}^2 C_{\mu, \bar{\mu}} \stackrel{(h)}{\leq} \frac{8\zeta^{K+1}}{(1-\zeta)^3} R_{\max}^2 + \frac{8\zeta^2 C_{\mu, \bar{\mu}}}{(1-\zeta)^4} \|e\|_{d_{\pi_k}}^2,
 \end{aligned}$$

where (b) follows from Jensen's inequality, i.e., $f(\sum_{k=0}^K \alpha_k x_k) \leq \sum_{k=0}^K \alpha_k f(x_k)$ for any convex function f . Also, (c) follows from Jensen's inequality, similar to (68), with stochastic matrix G_k . (d) follows from the definition of the coefficients $c_{\mu, \bar{\mu}}(m)$ (70), (e) follows by taking e such that $\|e\|_{\bar{\mu}}^2 = \max_{0 \leq k \leq K-1} \|e_k\|_{\bar{\mu}}^2$, (f) follows from (77), (g) follows since $\zeta_K \leq 2$ and the definition of $C_{\mu, \bar{\mu}}$ (69), and (h) follows by noting that for any $x \in \mathbb{R}^{|\mathcal{S}|}$, from Assumption 4, we have $\|x\|_{\bar{\mu}}^2 \leq C \|x\|_{d_{\pi_k}}^2$ for all k . Using the fact that $a^2 + b^2 \leq (a+b)^2$ for $a \geq 0, b \geq 0$ completes the proof of (71) in this theorem.

Now, since $\rho = 0$, (72) of this theorem directly follows from (iii) in Assumption 3 and (11) in Theorem 1. This completes the proof of this theorem. \square

D. Experiments

In all the experiments reported, we use a spherical uncertainty set $\{x : \|x\|_2 \leq r\}$ where r is the radius parameter. For such a set, we can compute a closed form solution for $\sigma_{\hat{U}_{s, \pi(s)}}(\Phi w)$ as $\sqrt{r w^\top \Phi^\top \Phi w}$ (Roy et al., 2017). Note that, we can precompute $\Phi^\top \Phi$ once and reuse it in every iteration of the RLSPI Algorithm, thus saving the computational overhead.

Chain MDP: We first consider a tabular MDP problem represented in the Figure 7 for verifying the convergence of RLSPI algorithm. This MDP consists of 10 states depicted by circles here. We have two actions, that is, move left or right. The actions fail to remain in a given direction with probability 0.1, depicted by the red arrows. Thus, with probability 0.9, actions succeed to be in a given direction, depicted by the blue (action left being unchanged) and green (action right being unchanged) arrows. Finally, visiting states of yellow color, that is 0 and 9, are rewarded 1, and visiting other states are rewarded 0.

(Lagoudakis & Parr, 2003) observes that learning algorithms often attain sub-optimal policies under such MDPs due to the randomization of actions (as depicted by red arrows in the Figure 7). It is also straightforward that the optimal policy of this MDP is moving left for states 0 through 4 and moving right for states 5 through 9. We train RLSPI algorithm on this MDP with $\alpha = 0.9, \lambda = 0$. We use the space spanned by polynomials, degree up to 2, as the feature space and set $\delta = 0.1$ (error of weights as mentioned in Step 8 of the RLSPI algorithm). We select r as 0.01 times the constant $\|\Phi^\top \Phi\|_F^{-1}$ where $\|\cdot\|_F$ is the Frobenius norm.

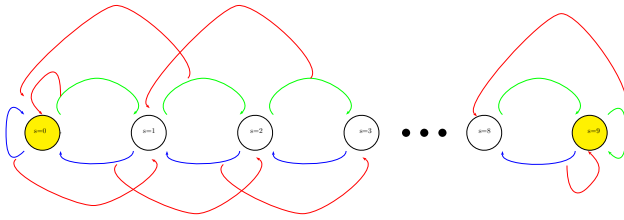


Figure 7

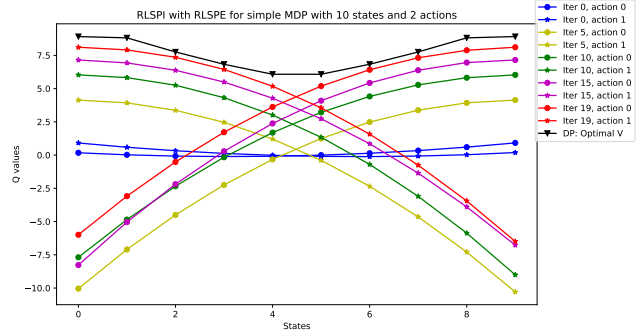


Figure 8

Figure 8 shows how the Q-value functions in RLSPI algorithm training evolve as the iteration progress. From this, we note that RLSPI algorithm is able to find the optimal policy with relatively less number of iterations. From this figure, we also note that the Q-value functions corresponding to the optimal policy in RLSPI algorithm converges to the optimal robust value function.

Examples from OpenAI Gym (Brockman et al., 2016): We now provide more details for the OpenAI Gym experiments demonstrated in Section 5. We use the radial basis functions (RBFs) for the purpose of feature spaces in our experiments. The general expression for RBFs is $\psi(x) = \exp(-\|x - \mu\|^2/\sigma)$ where the RBF parameters μ and σ are chosen before running the experiment. Here x is a concatenation of states and actions when both \mathcal{S}, \mathcal{A} are continuous spaces. In this case, the feature map is simply defined as $\phi(s, a) = \psi((s, a))$ where (s, a) represents the concatenation operation. While working with experiments whose action space is discrete, we naturally choose $\phi(s, a)$ to be the vector $(\mathbb{1}(a = 1)\psi(s), \dots, \mathbb{1}(a = |\mathcal{A}|\psi(s))^\top$ where $\mathbb{1}(E)$ is the indicator function which produces value 1 if the event E is true, and 0 otherwise. After a few trials, we observed that using few (typically 3-4) uniformly spaced RBFs in each dimension of x , here x is as described before, with approximate overlap percentage of 2.5 works suitably for getting the desired results shown here and in Section 5. Figure 9 illustrates this for the case of using two(= n) uniformly spaced RBFs with one dimensional x variable. For this illustration, we have the low(l) and high(h) values of x to be -0.5 and 0.5 respectively. Thus, the centers (i.e., parameter μ) of the two uniformly spaced RBFs in Figure 9 are -0.5 and 0.5 . We select the parameter σ as $(h - l)^2/n^3 = 0.125$. We execute this idea on the OpenAI Gym environments. We also experiment on FrozenLake8x8 OpenAI Gym environment, for which we use the tabular feature space since the state space is discrete.

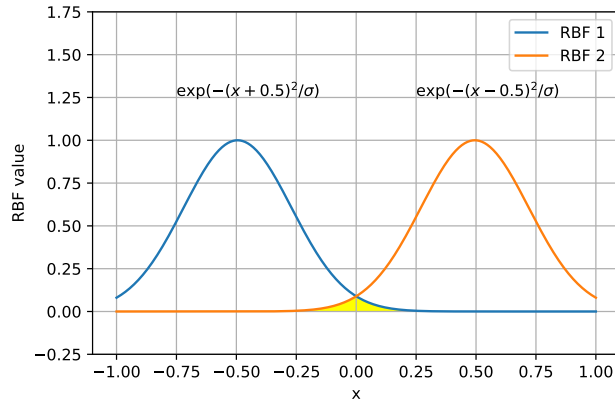


Figure 9: We have discretized x with 10^4 points here. The overlap percentage (ratio of the area in the yellow region to the area of either RBF 1 or RBF 2) is 2.88.

A short description of the OpenAI Gym tasks **CartPole**, **MountainCar**, **Acrobot**, and **FrozenLake8x8** we used are as follows.

CartPole: By a hinge, a pole is attached to a cart, which moves along a one-dimensional path. The motion of the cart is controllable, which is either to move it left or right. The pole starts upright, and the goal is to prevent it from falling over. A reward of +1 is provided for every time-step that the pole remains upright. CartPole consists of a 4-dimension continuous state space with 2 discrete actions.

MountainCar: A car is placed in the valley and there exists a flag on top of the hill. The goal is to reach the flag. The control signal is the acceleration and deceleration in continuous domain. A reward of 0 is provided if the car reaches goal,

otherwise it is provided -1 . MountainCar consists of a 3-dimension $\mathcal{S} \times \mathcal{A}$ continuous space.

Acrobot: Two poles attached to each other by a free moving joint and one of the poles is attached to a hinge on a wall. Initially, the poles are hanging downwards. An action, positive and negative torques can be applied to the movable joint. A reward of -1 is provided every time-step until the end of the lower pole reaches a given height, at termination 0 reward is provided. The goal is to maximize the reward gathered. Acrobot consists of a 6-dimension continuous state space with 3 discrete actions.

FrozenLake8x8: A grid of size 8×8 consists of some tiles which lead to the agent falling into the water. The agent is rewarded 1 after reaching a goal tile without falling and rewarded 0 in every other timestep.

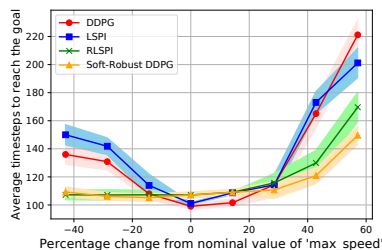


Figure 10

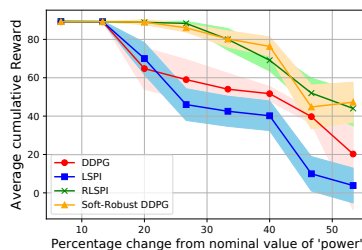


Figure 11

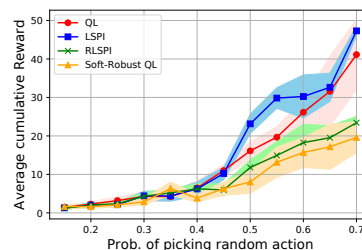


Figure 12

In Section 5, we provided performance evaluation curves in Figures 1-6. Here, we provide more results.

Figure 10 shows the average time steps to reach the goal in MountainCar environment as we change the parameter max_speed . The nominal value of this parameter is 0.07. As the parameter deviates from the nominal value, the performance of the policy obtained by the LSPI algorithm degrades quickly whereas the performance of the policy obtained by the RLSPI algorithm is fairly robust. Figure 11 shows the average cumulative reward on the MountainCar environment as we change the parameter $power$. The nominal value of this parameter is 15×10^{-4} . We again note that the RLSPI algorithm showcases robust performance. Figure 12 shows the ratio of average time to reach the goal and the number of trajectories which actually reach the goal on the FrozenLake8x8 environment against probability of picking a random action. Note that for large values of this probability all algorithms take more time to reach the goal or often fall into the water. Here again, RLSPI shows robust performance. Intuitively, perturbation in the parameters (like the action space, CartPole’s $force_mag$, $gravity$, $length$, MountainCar’s max_speed , $power$) of the environment is captured by the uncertainty set in the RMDP framework. Thus we see good performances of the robust algorithms like our RLSPI algorithm, Soft-Robust algorithms (Derman et al., 2018), and Robust Q-learning algorithm (Roy et al., 2017) compared to the non-robust algorithms.

In each policy iteration loop, in both LSPI and RLSPI algorithms, we generate t trajectories of horizon length h using the last updated policy (the initial policy π_0 is random.) We generally stop the simulation after 10-20 policy iteration loops. The details of the hyper-parameters are shown in Table 1 in addition to λ being set to zero.

OpenAI Gym Environment	Discount α	Weights error ϵ_0	t	h
CartPole	0.95	0.01	150	200
MountainCar	0.95	0.05	1000	20
Acrobot	0.98	0.1	100	200
FrozenLake8x8	0.99	0.01	3000	200

Table 1: Details of hyper-parameters in LSPI and RLSPI algorithms experiments.

Here are the details on the Q-learning based algorithms. The Q-learning algorithm with linear function approximation on CartPole uses 64 parameterized (centers and variances) RBFs chosen by the Adam optimizer. Both Q-learning and Soft-Robust Q-learning algorithms for FrozenLake8x8 uses the tabular method, instead of deep neural or linear function architectures. We use the usual decaying-epsilon-greedy for the exploration policies, such that it exponentially decays to 0.01 at half-way through the total number of training episodes (Epochs in Table 2) starting from 0.99. We provide the hyper-parameters in Table 2 like the discount factor, size of hidden layers h starting from the first hidden layer in the given

array, and size of the batch of tuples (state, action, next state, reward) chosen uniformly from the experience buffer of size 5000 to update the neural network (Batch in Table 2). For all neural networks, we used the *relu* activation functions. Note that the DDPG algorithm uses two same sized neural networks for actor and critic.

OpenAI Gym Environment	Discount α	Hidden layers h	Batch	Epochs
CartPole	0.9	[200, 100]	200	3000
MountainCar	0.99	[300, 200]	40	1000
Acrobot	0.995	[128, 64, 64]	500	10000
FrozenLake8x8	0.95	-	-	80000

Table 2: Details of hyper-parameters in Q-learning based algorithms experiments.

For completeness, we also point out some weaknesses of the experiments we have done. Firstly, we are not optimizing over the parameter r which is the radius of the spherical set associated with the uncertainty. We believe that performing a hyper-parameter search for the best r will make the policy obtained by the RLSPI further robust. Secondly, since we are focusing on the linear approximation architecture for developing the theoretical understanding of model-free robust RL, the experiments may not be immediately scalable to very high dimensional OpenAI Gym environments which typically require nonlinear approximation architecture.

To end this section, we mention the software configurations used to generate these results: *Python3.7 with OpenAI Gym (Brockman et al., 2016) and few basic libraries (non-exhaustive) like numpy, scipy, matplotlib*. Also, the hardware configurations used was *macOS High Sierra Version 10.13.6, 16 GB LPDDR3, Intel Core i7*.