

# 000 001 002 003 004 005 006 007 008 009 010 011 012 013 014 015 016 017 018 019 020 021 022 023 024 025 026 027 028 029 030 031 032 033 034 035 036 037 038 039 040 041 042 043 044 045 046 047 048 049 050 051 052 053 FINITE-TIME BOUNDS FOR DISTRIBUTIONALLY RO- BUST TD LEARNING WITH LINEAR FUNCTION APPROX- IMATION

006 **Anonymous authors**

007 Paper under double-blind review

## ABSTRACT

013 Distributionally robust reinforcement learning (DRRL) focuses on designing poli-  
014 cies that achieve good performance under model uncertainties. In particular, we  
015 are interested in maximizing the worst-case long-term discounted reward, where  
016 the data for RL comes from a nominal model while the deployed environment  
017 can deviate from the nominal model within a prescribed uncertainty set. Ex-  
018 isting convergence guarantees for robust temporal-difference (TD) learning for  
019 policy evaluation are limited to tabular MDPs or are dependent on restrictive  
020 discount-factor assumptions when function approximation is used. We present the  
021 first robust TD learning with linear function approximation, where robustness is  
022 measured with respect to the total-variation distance uncertainty set. Additionally,  
023 our algorithm is both model-free and does not require generative access to the MDP.  
024 Our algorithm combines a two-time-scale stochastic-approximation update with an  
025 outer-loop target-network update. We establish an  $\tilde{O}(1/\epsilon^2)$  sample complexity to  
026 obtain an  $\epsilon$ -accurate value estimate. Our results close a key gap between the empir-  
027 ical success of robust RL algorithms and the non-asymptotic guarantees enjoyed by  
028 their non-robust counterparts. The key ideas in the paper also extend in a relatively  
029 straightforward fashion to robust Q-learning with function approximation.

## 1 INTRODUCTION

030 Reinforcement learning (RL) aims to learn policies that maximize long-term reward. Standard RL  
031 methods learn the optimal strategy from trajectories generated by a simulator or the real environment,  
032 implicitly assuming that training and deployment environments share the same dynamics. Many  
033 applications face two issues: simulation-reality gaps and distribution shift between training and de-  
034 ployment. These call for policies that are robust to perturbations in the environment. Distributionally  
035 robust RL (DRRL) tackles this by assuming the true environment lies in an uncertainty set around a  
036 nominal model. It then learns a policy that maximizes the worst-case cumulative reward over that  
037 set, using data from trajectories corresponding to the nominal model. In this work, we focus on  
038 model-free DRRL with linear function approximation for the value function to deal with large state  
039 spaces.

040 In contrast to our model-free approach, model-based DRRL often proceeds by fitting an empirical  
041 transition model, defining an uncertainty set from it, and then optimizing for a robust policy (Shi &  
042 Chi, 2024; Wang & Zou, 2021; Xu et al., 2023; Panaganti & Kalathil, 2022; Yang et al., 2022; Zhou  
043 et al., 2021). In some model-based papers, access to a generative-model is assumed, which is not  
044 realistic in many cases (Wang & Zou, 2021; Xu et al., 2023). Whether one assumes generative access  
045 or not, the number of parameters that need to be estimated in a model-based approach grows with the  
046 cardinality of the state and action spaces, unless one makes additional structural assumptions on the  
047 model.

048 Another line of work focuses on model-free learning of robust policies, that is, learning without  
049 constructing an empirical transition matrix. In the tabular setting, Liang et al. (2023) analyzes  
050 Cressie–Read  $f$ -divergence–based uncertainty sets and establishes asymptotic convergence guarantees  
051 for robust temporal-difference (TD) learning. A complementary tabular result, Li et al. (2022), studies  
052 the  $R$ -contamination uncertainty set and exploits a distinctive property: the robust Bellman operator

054 in this model admits an unbiased stochastic estimator. The techniques developed there extend to any  
 055 uncertainty set that likewise permits an unbiased estimator of the robust Bellman operator, enabling  
 056 unbiased policy evaluation and, consequently, policy improvement in a model-free manner. However,  
 057 these papers do not consider function approximation, which is essential to deal with large state spaces.  
 058

059 When function approximation is introduced to represent the robust value function, the literature  
 060 typically proceeds along two directions with different limitations. One line of research constructs the  
 061 uncertainty set expressly so that the robust Bellman operator admits an unbiased estimator (Zhou  
 062 et al., 2023), allowing standard stochastic approximation arguments to go through or restrict to  
 063  $R$ -contamination uncertainty set (Wang & Zou, 2021). For  $R$ -contamination uncertainty set, Wang &  
 064 Zou (2021) investigates the TD-C algorithm under function approximation and provides finite-time  
 065 bounds for convergence to a stationary point of the associated objective, offering non-asymptotic  
 066 guarantees in a setting where the objective is nonconvex and only stationarity is generally attainable.  
 067 The other direction assumes extremely small discount factors to induce a contraction mapping for the  
 068 robust Bellman operator, which restores fixed-point uniqueness and enables convergence proofs Zhou  
 069 et al. (2023); Badrinath & Kalathil (2021); Tamar et al. (2014). Both approaches trade generality  
 070 for tractability: the first restricts attention to uncertainty sets with unbiased estimators and focuses  
 071 only on local optimality, while the second relies on unrealistically small discounting to guarantee  
 072 contraction.

073 Another line of work (Tang et al., 2024; Ma et al., 2022) for model-free DRRL considers linear  
 074 Markov decision process (MDP) for DRRL where the transition matrix of the underlying MDP has a  
 075 lower-dimensional structure. This reduces the complexity associated with large state spaces. In this  
 076 paper, we do not make such a modeling assumption.

077 In summary, most existing results on model-free robust RL are limited in at least one crucial way: they  
 078 prove only local or asymptotic convergence; focus on narrow uncertainty models (e.g., Liang et al.  
 079 (2023) observe on FrozenLake that  $R$ -contamination-based methods can mirror non-robust baselines  
 080 and even underperform due to over-conservatism); restricted to tabular settings; assume generative  
 081 access; or require extremely small discount factors. In particular, there are no finite-time guarantees  
 082 for robust TD with function approximation from a single trajectory under broad, practically motivated  
 083 uncertainty classes—such as those induced by total variation or Wasserstein- $\ell$  distances. At the same  
 084 time, practice-oriented deep-RL pipelines often use ad-hoc “robust TD” heuristics, leaving a sizable  
 085 gap between theory and deployment. This work closes a portion of that gap by establishing finite-time  
 086 guarantees for robust TD learning with function approximation under commonly used uncertainty sets,  
 087 without relying on generative sampling, vanishing discount factors, or purely asymptotic arguments.  
 088

**Contributions.** Our main contributions are summarized below.

1. **Finite-time guarantees for Robust TD Learning** For total variation and Wasserstein- $\ell$  uncertainty sets, we establish that the distributionally robust policy evaluation considered in the paper with linear function approximation admits non-asymptotic guarantees from a single trajectory. The robust TD method achieves an  $\epsilon$ -accurate value estimate with sample complexity  $\tilde{O}(1/\epsilon^2)$ .
2. **Overcoming projection mismatch via target networks.** While the robust Bellman operator is a contraction in  $\ell_\infty$  (Iyengar, 2005), function approximation induces a projected fixed-point equation that breaks direct contraction arguments. Prior approaches either remain tabular or require unrealistically small discount factors. We resolve this by incorporating a target-network mechanism—conceptually related to Munos & Szepesvári (2008) and, in the non-robust setting, Chen et al. (2023)—and prove stable, finite-time convergence of the resulting projected robust TD updates without restrictive discount-factor assumptions.
3. **Function approximation in the dual space.** Standard DRRL solvers compute the worst-case distribution at each step of an RL algorithm by using a dual formulation Iyengar (2005). However, this requires estimating a dual variable for each (state, action) pair, which is infeasible for large state spaces. To overcome this problem, we provide the first analysis of function approximation in the dual space.
4. **Robust Q-Learning.** The main technical contributions of the paper are in the proof of convergence and sample complexity bounds for robust TD learning with function approximation. It is straightforward to use these ideas to obtain finite-time bounds for robust Q-learning

108 with function approximation, which, to the best of our knowledge, has not been studied in  
 109 the literature. We refer the reader to the short argument in the Appendix (Section E).  
 110

111 Since our paper focuses on discounted-reward robust RL, we have not made an exhaustive comparison  
 112 of our work with work on average-reward robust RL; see, for example, Xu et al. (2025); Roch et al.  
 113 (2025); Chen et al. (2025). However, to the best of our knowledge, it is worth noting that there are no  
 114 performance guarantees even in the average-reward literature when function approximation is used.  
 115

## 116 2 MODEL AND PRELIMINARIES

117 **Model** We consider finite-state, finite-action, infinite-horizon discounted MDPs  $\mathcal{M} :=$   
 118  $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$ , where  $\mathcal{S}$  is the (finite) state space and  $\mathcal{A}$  is the (finite) action space. For any fi-  
 119 nite set  $\mathcal{X}$ , we denote by  $\Delta_{\mathcal{X}} := \left\{ \mu \in \mathbb{R}_+^{|\mathcal{X}|} : \sum_{x \in \mathcal{X}} \mu(x) = 1 \right\}$  the probability simplex over  $\mathcal{X}$ ; in  
 120 particular,  $\Delta_{\mathcal{S}}$  and  $\Delta_{\mathcal{A}}$  are the simplices over states and actions, respectively.  
 121

122 Throughout, we use lowercase letters  $s \in \mathcal{S}, a \in \mathcal{A}$  to denote deterministic (non-random) states and  
 123 actions, and uppercase letters  $S, S', A$  to denote random states and actions taking values in  $\mathcal{S}$  and  $\mathcal{A}$ . Given a state-action pair  $(s, a)$ , the transition kernel  $P(\cdot | s, a) \in \Delta_{\mathcal{S}}$  specifies the distribution of  
 124 the next-state random variable  $S' \sim P(\cdot | s, a)$ . The reward function is  $R : \mathcal{S} \times \mathcal{A} \rightarrow [-1, 1]$ , and  
 125  $\gamma \in (0, 1)$  is the discount factor. A (stochastic) policy  $\pi$  maps states to distributions over actions, that  
 126 is,  $\pi(\cdot | s) \in \Delta_{\mathcal{A}}$  for each  $s \in \mathcal{S}$ , and we write  $\pi(a | s)$  for the probability of choosing action  $a$  in  
 127 state  $s$ .  
 128

129 Let  $\{(S_t, A_t)\}_{t \geq 0}$  denote the state-action process for a policy  $\pi$ . Then for policy  $\pi$  and transition  
 130 model  $P$ , the (policy-dependent) state-action value is defined as  
 131

$$132 Q_P^\pi(s, a) := \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t R(S_t, A_t) \mid S_0 = s, A_0 = a, A_t \sim \pi(\cdot | S_t), S_{t+1} \sim P(\cdot | S_t, A_t) \right].$$

133 **Robust MDPs (RMDPs) and uncertainty sets.** Distributionally robust RL (DRRL) models trans-  
 134 ition uncertainty via an *uncertainty set* around a nominal kernel  $P_0$ . We adopt the standard  $(s, a)$ -  
 135 rectangular model (Iyengar, 2005; Nili & El Ghaoui, 2005):  
 136

$$137 \mathcal{P}_s^a = \left\{ q \in \Delta_{\mathcal{S}} : D(q, P_0(\cdot | s, a)) \leq \delta \right\}, \quad \mathcal{P} = \bigotimes_{(s, a) \in \mathcal{S} \times \mathcal{A}} \mathcal{P}_s^a, \quad (1)$$

138 where  $D(\cdot, \cdot)$  is a probability distance or divergence (e.g., total variation or Wasserstein- $\ell$ ), and  $\delta > 0$   
 139 is the radius. An RMDP is then the tuple (superscript ‘rob’ stands for ‘robust’ throughout the rest of  
 140 the paper)  
 141

$$142 \mathcal{M}^{\text{rob}} = (\mathcal{S}, \mathcal{A}, \mathcal{P}, R, \gamma).$$

143 **Robust value functions (fixed policy).** Given a fixed policy  $\pi$ , the *robust* state-action value function  
 144 is the worst-case value over  $\mathcal{P}$  :

$$145 Q^{\text{rob}, \pi}(s, a) := \min_{P \in \mathcal{P}} Q_P^\pi(s, a), \quad V^{\text{rob}, \pi}(s) := \sum_a \pi(a | s) Q^{\text{rob}, \pi}(s, a). \quad (2)$$

146 It satisfies the *robust Bellman equation*:

$$147 Q^{\text{rob}, \pi}(s, a) = R(s, a) + \gamma \min_{q \in \mathcal{P}_s^a} \sum_{s'} q(s' | s, a) \underbrace{\sum_{a'} \pi(a' | s') Q^{\text{rob}, \pi}(s', a')}_{=: V^{\text{rob}, \pi}(s')}.$$

148 Equivalently, defining the robust Bellman operator  $(\mathcal{T}^{\text{rob}, \pi} Q)(s, a) := R(s, a) + \gamma \sigma_{\mathcal{P}_s^a}(V^{\text{rob}, \pi}(s'))$   
 149 with

$$150 \sigma_{\mathcal{P}_s^a}(V) := \min_{q \in \mathcal{P}_s^a} \sum_{s'} q(s') V(s'), \quad V^{\text{rob}, \pi}(s') := \sum_{a'} \pi(a' | s') Q^{\text{rob}, \pi}(s', a'), \quad (4)$$

151 the fixed point relation is  $Q^{\text{rob}, \pi} = \mathcal{T}^{\text{rob}, \pi} Q^{\text{rob}, \pi}$ . We can write from the definitions,

$$152 0 \leq V^{\text{rob}, \pi}(s) \leq \frac{1}{1 - \gamma}, \forall s \in \mathcal{S}; \quad 0 \leq Q^{\text{rob}, \pi}(s, a) \leq \frac{1}{1 - \gamma}, \forall (s, a) \in \mathcal{S} \times \mathcal{A}.$$

162 For a fixed  $\pi$ , evaluating  $Q^{\text{rob}, \pi}$  reduces to solving Equation (3), which at each  $(s, a)$  requires solving  
 163 the inner problem Equation (4).

## 165 2.1 ROBUST TEMPORAL-DIFFERENCE LEARNING: CHALLENGES

167 **Function approximation.** Fix a policy  $\pi$ . We approximate the robust state-action value function  
 168 by a linear function class with the learnable parameter vector  $\theta \in \mathbb{R}^{n_\theta}$

$$170 \quad Q_\theta^{\text{rob}, \pi}(s, a) \approx \phi(s, a)^\top \theta, \quad \|\phi(s, a)\|_2 \leq 1, \forall (s, a) \in \mathcal{S} \times \mathcal{A}$$

172 with feature matrix  $\Phi \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times n_\theta}$ . Let  $d^\pi(s, a)$  be the stationary state-action distribution of  $(S_t, A_t)$   
 173 under  $\pi$ , and define  $D^\pi := \text{diag}(\{d^\pi(s, a)\}_{(s, a) \in \mathcal{S} \times \mathcal{A}})$ . Assume the weighted feature covariance is  
 174 well-conditioned:

$$175 \quad \Phi^\top D^\pi \Phi \succeq \mu I_{n_\theta} \quad \text{for some } \mu > 0.$$

177 Let  $\mathcal{W} := \{\Phi\theta : \theta \in \mathbb{R}^{n_\theta}\}$  and denote by  $\Pi : \mathbb{R}^{|\mathcal{S}||\mathcal{A}|} \rightarrow \mathcal{W}$  the  $D^\pi$ -orthogonal projection,

$$178 \quad \Pi f = \Phi(\Phi^\top D^\pi \Phi)^{-1} \Phi^\top D^\pi f.$$

180 For any scalar  $x \in \mathbb{R}$ , we define the clipping operator

$$182 \quad \text{Clip}(x) := \min \left\{ \max \left\{ x, -\frac{1}{1-\gamma} \right\}, \frac{1}{1-\gamma} \right\}.$$

185 When applied to a vector  $v \in \mathbb{R}^{n_\theta}$ ,  $\text{Clip}(v)$  denotes component-wise application of this operation.

186 We define the function approximation error for approximating the robust Q-function as:

$$188 \quad \epsilon_{\text{approx}} := \sup_{Q=\text{Clip}(\Phi\theta), \theta \in \mathbb{R}^{n_\theta}} \|\text{Clip}(\Pi \mathcal{T}^{\text{rob}, \pi}(Q)) - \mathcal{T}^{\text{rob}, \pi}(Q)\|_\infty. \quad (5)$$

191 **Key challenges in robust policy evaluation and our approach.** Model-free robust policy eval-  
 192 uation on a single trajectory typically hinges on a data-driven unbiased estimate  $\hat{\sigma}_{\mathcal{P}_s^a}(V)$  of the  
 193 inner-optimization objective defined in Equation (4). Except for special uncertainty sets (e.g.,  $R$ -  
 194 contamination), there is no direct plug-in *unbiased* single-sample estimator of this inner minimum,  
 195 which creates a bias in standard TD updates. To overcome this challenge, we use a *two-time-scale*  
 196 stochastic-approximation scheme in the inner loop of the algorithm: a fast time-scale solves for the  
 197 inner-optimization problem defined in Equation (4) in its equivalent dual form, while the slow loop  
 198 performs TD learning updates on  $\theta$  using the estimate of the inner-optimization objective of the fast  
 199 time-scale. Our two-time-scale algorithm is motivated by the algorithm in Liang et al. (2023), but the  
 200 key difference here is the use of function approximation which necessitates a different analysis.

201 While  $\mathcal{T}^{\text{rob}, \pi}$  is a  $\gamma$ -contraction in  $\ell_\infty$ -norm (Iyengar, 2005), function approximation introduces the  
 202 *projected* operator  $\Pi \mathcal{T}^{\text{rob}, \pi}$ , which is *not* known to be a contraction in any norm for typical  $\gamma \in (0, 1)$ .  
 203 Prior work by Zhou et al. (2023) circumvents this by imposing restrictive assumptions on  $\gamma$  which we  
 204 do not adopt. We address the non-contraction of  $\Pi \mathcal{T}^{\text{rob}, \pi}$  via a *target-network* mechanism prevalent  
 205 in deep RL, analyzed by Munos & Szepesvári (2008) and later used in the non-robust setting by Chen  
 206 et al. (2023), for Q-learning to overcome the contraction issue with the projected robust Bellman  
 207 operator. At outer iteration  $t$ , we freeze a target parameter  $\hat{\theta}_t$  and solve

$$208 \quad \Phi\theta = \Pi \mathcal{T}^{\text{rob}, \pi}(\Phi\hat{\theta}_t)$$

210 in the inner loop, then update the target in the outer loop. This decoupling stabilizes the projected  
 211 robust updates and enables our finite-time analysis under linear function approximation.

212 Standard DRRL literature solves the inner-optimization problem in Equation (4) in a corresponding  
 213 dual space. However, solving it for each state-action pair is impractical for problems with large state  
 214 and action spaces. We consider linear function approximation in the dual space of the optimization  
 215 problem in Equation 4 and provide the first finite-sample analysis under this function approximation  
 setup.

216 

### 3 ROBUST TD LEARNING WITH LINEAR FUNCTION APPROXIMATION

217 

#### 3.1 UNCERTAINTY SETS

220 Before presenting the robust policy evaluation algorithm, we discuss the uncertainty sets considered  
 221 in the paper: Total Variation (TV) uncertainty set and Wasserstein- $\ell$  uncertainty set.

222 **Total variation uncertainty set:** The total variation uncertainty set is defined as: for each  $(s, a)$ ,  
 223  $\mathcal{P}_s^{a, TV} = \{q \in \Delta_{\mathcal{S}} : \frac{1}{2} \|q - P_0(\cdot|s, a)\|_1 \leq \delta\}$ .  
 224

225 Simplifying (see: Appendix B) on the dual formulation originally given by Iyengar (2005) for the  
 226 Total Variation uncertainty set, we get the following equivalent dual optimization:

$$227 \sigma_{\mathcal{P}_s^a}(V) \equiv \max_{\lambda_s^a \in [\frac{-1}{1-\gamma}, \frac{1}{1-\gamma}]} \{\mathbb{E}_{S \sim P_0(\cdot|s, a)} [\min(V(S), \lambda_s^a)] - \delta \lambda_s^a\}.$$

230 **Wasserstein- $\ell$  uncertainty set:** The uncertainty set is defined as: for each  $(s, a)$ :  $\mathcal{P}_s^{a, W_\ell} = \{q \in$   
 231  $\Delta_{\mathcal{S}} : W_\ell(P_0(\cdot|s, a), q) \leq \delta\}$ , where  $\delta > 0$  is the uncertainty radius and  $W_\ell(P_0(\cdot|s, a), q)$  is the  
 232 Wasserstein- $\ell$  distance defined in detail in Appendix B.2.

233 The detailed analysis on TV and Wasserstein- $\ell$  uncertainty sets and the corresponding dual optimiza-  
 234 tion problem are given in the Appendix B.  
 235

236 

#### 3.2 ALGORITHM AND MAIN RESULTS

238 In this subsection, we present our robust policy evaluation algorithm and the main results of the paper.  
 239

240 

##### 3.2.1 ROBUST POLICY EVALUATION ALGORITHM

242 Our robust TD learning algorithm is presented in Algorithm 1. In the rest of this section, we  
 243 describe the algorithm and explain the notation used in the algorithm. In the outer loop (indexed by  
 244  $t = 0, \dots, T - 1$ ), we freeze a *target parameter*  $\hat{\theta}_t$ ; at the end of the inner loop we set  $\hat{\theta}_{t+1}$  to the  
 245 inner loop's final iterate. In the inner loop (indexed by  $k = 0, \dots, K - 1$ ) we approximately solve  
 246 for  $\theta$  satisfying:

$$247 \Phi\theta = \Pi \mathcal{T}^{\text{rob}, \pi}(\Phi\hat{\theta}_t),$$

248 using a two-time-scale stochastic approximation: a fast loop for the dual variables corresponding to  
 249 the inner-optimization problem 4, and a slow loop for the TD parameters. For a fixed outer loop  $t$ ,  
 250 the inner loop iterates are  $\theta_{t,k}$  for  $k \in [0, K - 1]$ .

251 At each inner loop iteration  $k$ , in a fast time-scale, we approximately solve the equivalent dual  
 252 optimization problem in (4) using a super-gradient ascent step. Instead of maintaining a separate dual  
 253 variable  $\lambda_s^a$  for each  $(s, a)$  (which would be tabular), we parameterize the dual variables  $\lambda_s^a$  with the  
 254 learnable parameter vector  $\nu \in \mathbb{R}^{n_\lambda}$  as

$$255 \lambda_s^a \approx \psi(s, a)^\top \nu, \quad \|\psi(s, a)\|_2 \leq 1, \forall (s, a) \in \mathcal{S} \times \mathcal{A},$$

257 with feature matrix  $\Psi \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times n_\lambda}$ .

259 Denote the robust value function estimate  $V_{\hat{\theta}_t}^{\text{rob}}$  evaluated at the target parameter  $\hat{\theta}_t$  as

$$261 V_{\hat{\theta}_t}^{\text{rob}}(s) = \sum_a \pi(a|s) \text{Clip} \left( \phi(s, a)^\top \hat{\theta}_t \right), \forall s \in \mathcal{S}. \quad (6)$$

263 The quantity  $V_{\hat{\theta}_t}^{\text{rob}}$  can be computed exactly for any fixed target parameter  $\hat{\theta}_t$ . In the case of  
 264 the TV distance uncertainty set, it suffices to compute  $V_{\hat{\theta}_t}^{\text{rob}}(s)$  *only for the state visited in each*  
 265 *inner-loop iteration*, rather than for all states. We update  $\nu_{t,k}$  with step-size  $\beta_k$  using a projected  
 266 *super-gradient ascent* on the dual objective with a super-gradient evaluated at the fresh data sample  
 267  $(S_{t,k}, A_{t,k}, S_{t,k+1})$ . Let  $B_\nu > 0$  be a fixed finite radius, and define

$$269 \mathcal{M}_\nu := \left\{ \nu \in \mathbb{R}^{n_\lambda} : \|\nu\|_2 \leq B_\nu \right\}.$$

270 The projection operator  $\text{Proj}_{\mathcal{M}_\nu}$  projects the dual parameter vector onto  $\mathcal{M}_\nu$ , ensuring that the  
 271 iterates remain bounded. Since  $\mathcal{M}_\nu$  is an  $\ell_2$  ball, this projection can be computed by simple norm  
 272 scaling.

273 In the algorithm,  $\bar{\nu}_{t,k}$  denotes the half-tail iterate-average of the dual parameter vector, i.e.,  
 274

$$275 \quad \bar{\nu}_{t,k} = \frac{1}{\lceil k/2 \rceil} \sum_{l=\lfloor k/2 \rfloor}^{k-1} \nu_{t,l} \quad (7)$$

278 which can be calculated easily by keeping track of the following two quantities:  $\sum_{l=0}^{k-1} \nu_{t,l}$  and  
 279  $\sum_{l=\lfloor k/2 \rfloor}^{k-1} \nu_{t,l}$ . While many elements of our algorithm have been used in implementations of robust  
 280 TD learning, to the best of our knowledge, such an averaging of the dual variables has not been used  
 281 previously. The averaging turns out to be crucial in obtaining finite-time bounds, since it allows us to  
 282 control the variance of the dual objective.  
 283

284 In the slow time-scale of the inner loop,  $\theta_{t,k}$  is updated using asynchronous stochastic approximation  
 285 with a step-size denoted by  $\alpha_k$  with a robust TD-target  $TD_{t,k+1}$ . The two-time-scale scheme ensures  
 286 that, at the slow scale, the dual variables appear near their sample-path equilibrium, yielding an  
 287 (asymptotically) unbiased robust TD target.  
 288

---

289 **Algorithm 1** Robust TD learning with Function Approximation

---

290 1: **Input:** Integers  $T, K$ . Initial  $\nu_0 \in \mathbb{R}^{n_\lambda}$ ,  $\theta_0 :=$  zero vector, fast time-scale step-sizes  $\beta_k =$   
 291  $\frac{\beta_0}{\sqrt{k+1}}$ , for some  $0 < \beta_0 < \infty$ , slow time-scale step-sizes  $\alpha_k = \frac{c}{(k+1)}$  for some  $0 < c < \infty$ ;  
 292  $\hat{\theta}_0 = \theta_0$ ,  $\theta_{0,0} = \theta_0$ , candidate policy  $\pi$ , Reward function  $R : (\mathcal{S} \times \mathcal{A}) \mapsto [-1, 1]$ , initial state  
 293  $S_{0,0}$ .  
 294 2: **for**  $t = 0, 1, \dots, T-1$  **do**  
 295 3:   **for**  $k = 0, 1, \dots, K-1$  **do**  
 296 4:     Take action  $A_{t,k}$  according to policy  $\pi$  and sample  $S_{t,k+1}$  ( $S_{t,k+1} \sim P_0(\cdot | S_{t,k}, A_{t,k})$ )  
 297 5:     **fast time-scale** ( $\beta_k$ )  
 298 6:     Compute  $\hat{G}(\psi(S_{t,k}, A_{t,k})^\top \nu_{t,k}; V_{\hat{\theta}_t}^{\text{rob}}, S_{t,k+1})$  from Equation (17) for TV uncertainty set  
 299     and Equation (20) for Wasserstein- $\ell$  uncertainty set  
 300 7:      $\nu_{t,k+1} = \text{Proj}_{\mathcal{M}_\nu}(\nu_{t,k} + \beta_k[\hat{G}(\psi(S_{t,k}, A_{t,k})^\top \nu_{t,k}; V_{\hat{\theta}_t}^{\text{rob}}, S_{t,k+1})\psi(S_{t,k}, A_{t,k})])$   
 301 8:     **Slow scale** ( $\alpha_k$ )  
 302 9:     Compute  $\bar{\nu}_{t,k}$  from Equation (7)  
 303 10:    Compute  $\hat{F}(\psi(S_{t,k}, A_{t,k})^\top \bar{\nu}_{t,k}; V_{\hat{\theta}_t}^{\text{rob}}, S_{t,k+1})$  from Equation (18) for TV uncertainty set  
 304     and Equation (21) for Wasserstein- $\ell$  uncertainty set  
 305 11:     $TD_{t,k+1} = R(S_{t,k}, A_{t,k}) + \gamma \hat{F}(\psi(S_{t,k}, A_{t,k})^\top \bar{\nu}_{t,k}; V_{\hat{\theta}_t}^{\text{rob}}, S_{t,k+1}) - \phi(S_{t,k}, A_{t,k})^\top \theta_{t,k}$   
 306 12:     $\theta_{t,k+1} = \theta_{t,k} + \alpha_k TD_{t,k+1} \phi(S_{t,k}, A_{t,k})$   
 307 13:   **end for**  
 308 14:    $\hat{\theta}_{t+1} = \theta_{t,K}$ ,  $S_{t+1,0} = S_{t,K}$ ,  $\theta_{t+1,0} = \theta_{t,K}$ ,  $\nu_{t+1,0} = \nu_{t,K}$ .  
 309 15: **end for**  
 310 16: **Output:**  $\hat{\theta}_T$

---

313 3.2.2 MAIN RESULT

315 We define the function approximation error for approximating the dual variables next. For compact-  
 316 ness of notation, denote for a value function  $V$ , for each  $(s, a)$ ,  $F_{s,a}^{*,V} := \sup_{\lambda_s^a} F(\lambda_s^a; V, P_0(\cdot | s, a))$   
 317 and  $F^V(\nu)_{s,a} := F(\psi(s, a)^\top \nu; V, P_0(\cdot | s, a))$ . Define

$$319 \quad \epsilon_{\text{approx}}^{\text{dual}} := \sup_{V: V(s) = \sum_a \pi(a|s) \text{Clip}(\phi(s, a)^\top \theta); \theta \in \mathbb{R}^{n_\theta}} \inf_{\nu \in \mathcal{M}_\nu} \|F^{*,V} - F^V(\nu)\|_\infty. \quad (8)$$

321 We make the following assumption on the policy  $\pi$ .  
 322

323 **Assumption 1.** *The policy  $\pi$  induces an irreducible and aperiodic Markov chain on  $\mathcal{S} \times \mathcal{A}$  under  
 324 the nominal transition kernel  $P_0$ .*

324 For  $\tau \leq k$ , define

$$325 \quad 326 \quad \eta_k^{t,\tau}(\cdot) := \mathbb{P}((S_{t,k}, A_{t,k}) \in \cdot \mid S_{t,0}, A_{t,0}, \dots, S_{t,k-\tau}, A_{t,k-\tau}).$$

327 Under Assumption 1, the Markov chain is geometrically mixing: there exist constants  $C_{\text{mix}} < \infty$   
328 and  $\rho \in (0, 1)$  such that

$$329 \quad 330 \quad \|\eta_k^{t,\tau} - d^\pi\|_{\text{TV}} \leq C_{\text{mix}} \rho^\tau, \quad \forall t, k, \tau.$$

331 Here  $C_{\text{mix}}$  and  $\rho$  depend only on the nominal model  $(P_0, \pi)$ .

332 Let  $\hat{Q}_t := \text{Clip}(\Phi\hat{\theta}_t)$  be the estimate of  $Q^{\text{rob},\pi}$  by Algorithm 1 at outer iteration  $t$ .  
333

334 In Theorem 1, we present our main result, which establishes the convergence of  $\hat{Q}_T$  to the robust  
335 value function  $Q^{\text{rob},\pi}$ , up to terms arising from function-approximation error.

336 **Theorem 1** (Finite-time bound: rates and dependencies (informal)). *Define*

$$337 \quad k_{\text{mix}} \\ 338 \quad := \min \left\{ m \in \mathbb{N} : \forall j \geq m, j \geq \max \left( \tau_\mu, 2 \left\lceil \frac{\log(\frac{C_{\text{mix}}}{\beta_0} \sqrt{j+1})}{\log(1/\rho)} \right\rceil, \left\lceil \frac{\log(C_{\text{mix}}(j+1)/c)}{\log(1/\rho)} \right\rceil \right) \right\},$$

342 where

$$343 \quad \tau_\mu := \left\lceil \frac{\log(C_{\text{mix}} \frac{1}{\mu})}{\log(1/\rho)} \right\rceil$$

345 Assume Assumption 1 holds, and we run  $K \geq k_{\text{mix}}$  inner iterations per outer iteration for either the  
346 TV uncertainty set or the Wasserstein- $\ell$  uncertainty set. Then, for any  $T \geq 1$ , we have

$$347 \quad \mathbb{E}[\|\hat{Q}_T - Q^{\text{rob},\pi}\|_\infty] \\ 348 \quad \leq \gamma^T \|\text{Clip}(\Phi\theta_0) - Q^{\text{rob},\pi}\|_\infty + \frac{\text{rate}_{\text{inner}}(K)}{(1-\gamma)^2} + \frac{\epsilon_{\text{approx}}}{1-\gamma} + \frac{2\sqrt{2} (1 + \frac{2}{K})^{\frac{\mu c}{4}} \epsilon_{\text{approx}}^{\text{dual}}}{\mu(1-\gamma)},$$

352 where the term  $\text{rate}_{\text{inner}}(K)$  is of the following order in terms of inner iteration number  $K$ :

$$353 \quad \text{rate}_{\text{inner}}(K) = \begin{cases} \mathcal{O}(K^{-\mu c/4}), & \alpha_k = \frac{c}{k+1}, \quad \mu c < 2, \\ \mathcal{O}((\log K)^{1/2} K^{-1/2}), & \alpha_k = \frac{c}{k+1}, \quad \mu c = 2, \\ \mathcal{O}(K^{-1/2}), & \alpha_k = \frac{c}{k+1}, \quad \mu c > 2, \end{cases}$$

359 where the notation  $\mathcal{O}$  captures the problem-dependent constants depending on  
360  $(\mu, \delta, C_{\text{mix}}, \rho, B_\nu, \beta_0, c)$ .

361 **Remark 1.** A fully constant-explicit version of Theorem 1 is provided in Theorem 2 in the Appendix.

363 Recall the slow time-scale step-size rule is  $\alpha_k = \frac{c}{k+1}, \forall k$ . The sample complexity to achieve an  
364  $\epsilon$ -approximate robust Q-function estimate can be derived in the following manner. Assume  $\mu c > 2$ . If  
365 we choose  $T = \mathcal{O}(\ln(\frac{1}{\epsilon(1-\gamma)}))$  and  $K = \mathcal{O}(\frac{1}{(\epsilon(1-\gamma)^2)^2})$ , we have  $\gamma^T \|\text{Clip}(\Phi\theta_0) - Q^{\text{rob},\pi}\|_\infty +$   
366  $\frac{\text{rate}_{\text{inner}}(K)}{(1-\gamma)^2} = \mathcal{O}(\epsilon)$ . This gives us the following sample complexity result.

368 **Corollary 1** (Sample Complexity). *Suppose the step-size rule  $\alpha_k = \frac{c}{1+k}$  is used with  $\mu c \geq 2$ . Then  
369 the sample complexity for Algorithm 1 achieves an element-wise  $\epsilon$ -accurate estimate of  $Q^{\text{rob},\pi}$  up to  
370 the function approximation error is*

$$371 \quad 372 \quad 373 \quad \mathcal{O} \left( \ln \left( \frac{1}{\epsilon(1-\gamma)} \right) \frac{1}{\epsilon^2(1-\gamma)^4} \right). \quad (9)$$

374 Similar sample complexity results can be obtained for other values of  $\mu c$ .

375 We note that the step-size rule  $\frac{c}{k+1}$  achieves the best sample complexity, but it requires  $c$  to be chosen  
376 sufficiently large. This is consistent with similar results in the non-robust RL literature; see, for  
377 example, Chen et al. (2023).

378 **4 KEY IDEAS AND PROOF OUTLINE**  
 379

380 While the detailed proof of Theorem 1 is presented in Appendix C, we provide the key ideas behind  
 381 the proof in this section.

382 We define the stacked reward vector  $r \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$  by

384  $r_{s,a} := R(s,a), \quad (s,a) \in \mathcal{S} \times \mathcal{A},$   
 385

386 using some fixed ordering of state-action pairs.

387 Fix an outer loop iteration  $t$ . Recall the definition  $F_{s,a}^{*,V} := \sup_{\lambda_s^a} F(\lambda_s^a; V, P_0(\cdot|s,a))$ . Define the  
 388 inner loop error for outer iteration index  $t$  as  $e_{t,k} := \theta_{t,k} - \theta_t^*$  with

389  $\theta_t^* := (\Phi^\top D^\pi \Phi)^{-1} \Phi^\top D^\pi [r + \gamma F^{*,V_{\theta_t}^{\text{rob}}}]$ . (10)  
 390

391 The next lemma bounds the expected estimation error at the final outer-loop iterate in terms of the  
 392 inner-loop error terms.

393 **Lemma 1.** *Under the setting in Theorem 1, Algorithm 1 guarantees*

394 
$$\mathbb{E} [\|\hat{Q}_T - Q^{\text{rob},\pi}\|_\infty] \leq \gamma^T \|\text{Clip}(\Phi\theta_0) - Q^{\text{rob},\pi}\|_\infty + \underbrace{\sum_{t=1}^T \gamma^{T-t} \mathbb{E} [\|e_{t,K}\|_\infty]}_{\text{Inner loop convergence error}} + \frac{\epsilon_{\text{approx}}}{1-\gamma}.$$

400 The proof of Lemma 1 is provided in Appendix C.1 and is inspired by the analysis in Chen et al.  
 401 (2023) for non-robust Q-learning.

402 In the analysis that follows, we establish that the inner loop error remains small (up to function  
 403 approximation error terms) in  $\ell_\infty$ -norm for sufficiently large  $k$ . We decompose the slow time-scale  
 404 update at inner loop  $k$  in Algorithm 1 into mean drift, noise and bias terms as

405  $\theta_{t,k+1} = \theta_{t,k} + \alpha_k [H(\theta_{t,k}) + b_{t,k}^\theta + n_{t,k+1}^\theta],$   
 406

407 where

408 
$$H(\theta_{t,k}) := \Phi^\top D^\pi [r + \gamma F^{*,V_{\theta_t}^{\text{rob}}} - \Phi\theta_{t,k}] \quad \underbrace{=}_{\text{from Equation (10)}} \Phi^\top D^\pi \Phi(\theta_t^* - \theta_{t,k}),$$
  
 409  
 410  $b_{t,k}^\theta := \gamma \Phi^\top D^\pi [F^{V_{\theta_t}^{\text{rob}}}(\bar{\nu}_{t,k}) - F^{*,V_{\theta_t}^{\text{rob}}}],$   
 411  
 412  $n_{t,k+1}^\theta := TD_{t,k+1} \phi(S_{t,k}, A_{t,k}) - H(\theta_{t,k}) - b_{t,k}^\theta.$   
 413

414 **Idealized recursion (without noise and bias).** The mean drift term corresponds to the deterministic  
 415 recursion:

416  $\theta_{t,k+1} = \theta_{t,k} + \alpha_k \Phi^\top D^\pi \Phi(\theta_t^* - \theta_{t,k}).$   
 417

418 This recursion admits  $\theta_t^*$  as its unique fixed point. Since the matrix  $\Phi^\top D^\pi \Phi$  is symmetric and  
 419 positive definite with minimum eigenvalue  $\mu > 0$ , in the absence of bias and noise terms, the iterates  
 420 satisfy

421  $\|\theta_{t,k+1} - \theta_t^*\|_2 \leq (1 - \alpha_k \mu) \|\theta_{t,k} - \theta_t^*\|_2,$  (11)  
 422

423 which implies geometric convergence of  $\theta_{t,k}$  to  $\theta_t^*$  at a rate governed by  $\mu$ .

424 **Bias term analysis.** Recall that the bias term is given by

425  $b_{t,k}^\theta := \gamma \Phi^\top D^\pi [F^{V_{\theta_t}^{\text{rob}}}(\bar{\nu}_{t,k}) - F^{*,V_{\theta_t}^{\text{rob}}}]$ .  
 426

427 We show that this term becomes small for large  $k$ , up to a function approximation error  $\epsilon_{\text{approx}}^{\text{dual}}$ .

428 In the fast time-scale analysis, we prove that the stochastic update on  $\nu$  performs a super-gradient  
 429 ascent on the concave objective

431 
$$L_t(\nu) := \sum_{s,a} d^\pi(s,a) F(\psi(s,a)^\top \nu; V_{\theta_t}^{\text{rob}}, P_0(\cdot|s,a)),$$

432 which has bounded super-gradients. By a standard Lyapunov function argument for stochastic  
 433 approximation under a mixing Markov chain, we obtain the following guarantee on the iterates from  
 434 the fast time-scale for sufficiently large  $k$  (stated in detail in Lemma 4 in Appendix C):  
 435

$$436 \quad \mathbb{E} \left[ \max_{\nu \in \mathcal{M}_\nu} L_t(\nu) - L_t(\bar{\nu}_{t,k}) \right] \leq \frac{C_{\text{fast}}}{\sqrt{k}}, \quad (12)$$

438 where the constant  $C_{\text{fast}}$  is given in equation 25.  
 439

440 Using  $\|\phi(s, a)\|_2 \leq 1$  and  $F_{s,a}^{*,V_{\theta_t}^{\text{rob}}} \geq F^{V_{\theta_t}^{\text{rob}}}(\bar{\nu}_{t,k})_{s,a}$  for all  $(s, a)$ , we can write  
 441

$$442 \quad \|b_{t,k}^\theta\|_2 \leq \gamma \sum_{s,a} d^\pi(s, a) \left| F^{V_{\theta_t}^{\text{rob}}}(\bar{\nu}_{t,k})_{s,a} - F_{s,a}^{*,V_{\theta_t}^{\text{rob}}} \right| = \gamma \sum_{s,a} d^\pi(s, a) \left( F_{s,a}^{*,V_{\theta_t}^{\text{rob}}} - F^{V_{\theta_t}^{\text{rob}}}(\bar{\nu}_{t,k})_{s,a} \right)$$

$$443 \quad \leq \gamma \underbrace{\inf_{\nu \in \mathcal{M}_\nu} \sum_{s,a} d^\pi(s, a) \left( F_{s,a}^{*,V_{\theta_t}^{\text{rob}}} - F^{V_{\theta_t}^{\text{rob}}}(\nu)_{s,a} \right)}_{\leq \epsilon_{\text{dual}}^{\text{approx}}} + \gamma \underbrace{\left[ \sup_{\nu \in \mathcal{M}_\nu} L_t(\nu) - L_t(\bar{\nu}_{t,k}) \right]}_{\text{fast-scale objective gap}}.$$

444  
 445 **Handling the noise term.** Finally, to handle the noise terms  $n_{t,k+1}^\theta$ , we employ the approach in  
 446 Srikanth & Ying (2019), where a bound is obtained on the expectation of the error  $\|\theta_{t,k} - \theta_t^*\|_2^2$  con-  
 447 ditioned with respect to the filtration generated by the set  $(S_{t,0}, A_{t,0}, S_{t,1}, A_{t,1}, \dots, S_{t,k-\tau}, A_{t,k-\tau})$ .  
 448 By choosing a lag  $\tau$  such that the underlying Markov chain has mixed sufficiently, the effect of noise  
 449 can be controlled.  
 450

## 455 5 DISCUSSION

456 As mentioned in the introduction, we provide the first proof of convergence and finite-time bounds for  
 457 robust TD learning with function approximation, without making any assumptions on the underlying  
 458 model or imposing very restrictive assumptions on the discount factor. Some immediate extensions  
 459 and open problems are identified below:  
 460

- 461 1. The algorithm and the results can be extended to other families of distances between  
 462 probability distributions, such as the Cressie-Read family of  $f$ -divergences considered in  
 463 Liang et al. (2023), which admit duality representations that allow one to obtain unbiased  
 464 estimators of the quantities of interest. For the Cressie-Read family, this would require the  
 465 addition of one more time-scale, but the rest of the analysis would be similar. Our results  
 466 also apply to the  $R$ -contamination set, but the algorithm is even simpler in that case due to  
 467 the fact that the dual problem has a closed-form solution Xu et al. (2025).  
 468
- 469 2. Although the results in the main body of the paper have been presented for robust TD  
 470 learning, they can be easily extended to robust Q-learning with function approximation to  
 471 obtain optimal policies; see the Appendix.  
 472

## 473 REFERENCES

474 Kishan Panaganti Badrinath and Dileep Kalathil. Robust reinforcement learning using least squares  
 475 policy iteration with provable performance guarantees. In *International Conference on Machine  
 476 Learning*, pp. 511–520. PMLR, 2021.

477 Zaiwei Chen, John-Paul Clarke, and Siva Theja Maguluri. Target network and truncation overcome  
 478 the deadly triad in-learning. *SIAM Journal on Mathematics of Data Science*, 5(4):1078–1101,  
 479 2023.

480 Zijun Chen, Shengbo Wang, and Nian Si. Sample complexity of distributionally robust average-reward  
 481 reinforcement learning. *arXiv preprint arXiv:2505.10007*, 2025.

482 Rui Gao and Anton Kleywegt. Distributionally robust stochastic optimization with wasserstein  
 483 distance. *Mathematics of Operations Research*, 48(2):603–655, 2023.

486 Garud N Iyengar. Robust dynamic programming. *Mathematics of Operations Research*, 30(2):  
 487 257–280, 2005.

488

489 Yan Li, Guanghui Lan, and Tuo Zhao. First-order policy optimization for robust markov decision  
 490 process. *arXiv preprint arXiv:2209.10579*, 2022.

491 Zhipeng Liang, Xiaoteng Ma, Jose Blanchet, Jiheng Zhang, and Zhengyuan Zhou. Single-trajectory  
 492 distributionally robust reinforcement learning. *arXiv preprint arXiv:2301.11721*, 2023.

493

494 Xiaoteng Ma, Zhipeng Liang, Jose Blanchet, Mingwen Liu, Li Xia, Jiheng Zhang, Qianchuan Zhao,  
 495 and Zhengyuan Zhou. Distributionally robust offline reinforcement learning with linear function  
 496 approximation. *arXiv preprint arXiv:2209.06620*, 2022.

497 Rémi Munos and Csaba Szepesvári. Finite-time bounds for fitted value iteration. *Journal of Machine  
 498 Learning Research*, 9(5), 2008.

499

500 Arnab Nilim and Laurent El Ghaoui. Robust control of markov decision processes with uncertain  
 501 transition matrices. *Operations Research*, 53(5):780–798, 2005.

502 Kishan Panaganti and Dileep Kalathil. Sample complexity of robust reinforcement learning with  
 503 a generative model. In *International Conference on Artificial Intelligence and Statistics*, pp.  
 504 9582–9602. PMLR, 2022.

505

506 Zachary Roch, Chi Zhang, George Atia, and Yue Wang. A finite-sample analysis of distributionally  
 507 robust average-reward reinforcement learning. *arXiv preprint arXiv:2505.12462*, 2025.

508 Laixi Shi and Yuejie Chi. Distributionally robust model-based offline reinforcement learning with  
 509 near-optimal sample complexity. *Journal of Machine Learning Research*, 25(200):1–91, 2024.

510

511 Rayadurgam Srikant and Lei Ying. Finite-time error bounds for linear stochastic approximation and  
 512 learning. In *Conference on learning theory*, pp. 2803–2830. PMLR, 2019.

513 Aviv Tamar, Shie Mannor, and Huan Xu. Scaling up robust mdps using function approximation. In  
 514 *International conference on machine learning*, pp. 181–189. PMLR, 2014.

515

516 Cheng Tang, Zhishuai Liu, and Pan Xu. Robust offline reinforcement learning with linearly structured  
 517  $f$ -divergence regularization. *arXiv preprint arXiv:2411.18612*, 2024.

518

519 Yue Wang and Shaofeng Zou. Online robust reinforcement learning with model uncertainty. *Advances  
 520 in Neural Information Processing Systems*, 34:7193–7206, 2021.

521

522 Yang Xu, Washim Uddin Mondal, and Vaneet Aggarwal. Finite-sample analysis of policy evaluation  
 523 for robust average reward reinforcement learning. *arXiv preprint arXiv:2502.16816*, 2025.

523

524 Zaiyan Xu, Kishan Panaganti, and Dileep Kalathil. Improved sample complexity bounds for distribu-  
 525 tionally robust reinforcement learning. In *International Conference on Artificial Intelligence and  
 526 Statistics*, pp. 9728–9754. PMLR, 2023.

526

527 Wenhao Yang, Liangyu Zhang, and Zhihua Zhang. Toward theoretical understandings of robust  
 528 markov decision processes: Sample complexity and asymptotics. *The Annals of Statistics*, 50(6):  
 529 3223–3248, 2022.

530

531 Ruida Zhou, Tao Liu, Min Cheng, Dileep Kalathil, PR Kumar, and Chao Tian. Natural actor-critic  
 532 for robust reinforcement learning with function approximation. *Advances in neural information  
 533 processing systems*, 36:97–133, 2023.

533

534 Zhengqing Zhou, Zhengyuan Zhou, Qinxun Bai, Linhai Qiu, Jose Blanchet, and Peter Glynn.  
 535 Finite-sample regret bound for distributionally robust offline tabular reinforcement learning. In  
 536 *International Conference on Artificial Intelligence and Statistics*, pp. 3331–3339. PMLR, 2021.

537

538

539

540 **A CONTENTS**  
 541

542 The contents of the Appendix are as follows:  
 543

544 1. In Section B, we analyze the TV distance and Wasserstein- $\ell$  uncertainty sets in detail.  
 545 2. Section C proves the main result of the paper, that is, Theorem 1 in detail.  
 546 3. In Section E, we present the robust Q learning algorithm with linear function approximation  
 547 (Algorithm 2) and discuss how the theoretical analysis for robust TD learning can be  
 548 extended to the robust Q learning straightforwardly.  
 549

550 **B CONDITIONS FOLLOWED BY TV AND WASSERSTEIN- $\ell$  UNCERTAINTY SETS**  
 551

553 To aid the analysis, in the this section we outline a few properties of the uncertainty sets considered  
 554 in this paper: TV and Wasserstein- $\ell$  uncertainty sets. In Section 5, we discuss how our algorithm can  
 555 be trivially modified to satisfy a similar convergence guarantee for the R-contamination uncertainty  
 556 set and Cressie-Read family of f-divergences considered in Liang et al. (2023).  
 557

558 **Lemma 2.** *The TV and Wasserstein uncertainty sets considered in this paper satisfies the following  
 559 conditions. The optimization problem  $\sigma_{\mathcal{P}_s^a}(V)$  for a generic value function  $V$  as defined in Equation  
 560 (4) has an equivalent dual optimization problem corresponding to a dual variable  $\lambda_s^a$  :*

561 
$$\sigma_{\mathcal{P}_s^a}(V) \equiv \sup_{\lambda_s^a \geq 0} (F(\lambda_s^a; V, P_0(\cdot|s, a)))$$
  
 562

563 where  $F(\lambda_s^a; V, P_0(\cdot|s, a))$  is a  $\lambda_s^a$ -concave function with the following properties:  
 564

565 1. Let  $G(\lambda_s^a; V, P_0(\cdot|s, a))$  be a super-gradient of the concave function  $F(\lambda_s^a; V, P_0(\cdot|s, a))$ .  
 566 There exists an unbiased estimator  $\hat{G}(\lambda_s^a; V, S')$  of  $G(\lambda_s^a; V, P_0(\cdot|s, a))$  based on a sample  
 567 of the next state  $S' \sim P_0(\cdot|s, a)$ , that is,  
 568

569 
$$\mathbb{E}_{S' \sim P_0(\cdot|s, a)}[\hat{G}(\lambda_s^a; V, S')] = G(\lambda_s^a; V, P_0(\cdot|s, a)),$$
  
 570

571 and it satisfies  $|\hat{G}(\lambda_s^a; V, S')| \leq C_G < \infty$  for all  $\lambda_s^a \in \mathbb{R}$  for some constant  $C_G \geq 0$ .  
 572

573 2. There exists an unbiased estimator  $\hat{F}(\lambda_s^a; V, S')$  of the dual objective  $F(\lambda_s^a; V, P_0(\cdot|s, a))$   
 574 based on a sample of next state  $S' \sim P_0(\cdot|s, a)$ , that is,  
 575

576 
$$\mathbb{E}_{S' \sim P_0(\cdot|s, a)}[\hat{F}(\lambda_s^a; V, S')] = F(\lambda_s^a; V, P_0(\cdot|s, a)).$$
  
 577

578 Moreover, the estimator is uniformly bounded on bounded sets of  $\lambda_s^a$  : for every  $M > 0$   
 579 there exists a constant  $C_{F,M} < \infty$  such that, for all  $|\lambda_s^a| \leq M$  and all  $s' \in \mathcal{S}$ ,  
 580

581 
$$|\hat{F}(\lambda_s^a; V, s')| \leq C_{F,M}.$$
  
 582

583 Next, we discuss in detail the uncertainty sets considered in this paper, namely, TV distance uncer-  
 584 tainty and Wasserstein- $\ell$  uncertainty sets and prove the Lemma 2. For each uncertainty set,  
 585

586 1. We define the uncertainty set first. Then, we discuss and analyze the equivalent dual  
 587 optimization that corresponds to the inner-optimization problem defined in Equation 4.  
 588 2. We show the uncertainty set satisfies the conditions described in Lemma 2 and hence prove  
 589 Lemma 2.  
 590

591 **B.1 TOTAL VARIATION DISTANCE UNCERTAINTY SET**  
 592

593 The total variation uncertainty set is defined for each  $(s, a)$  pair as,  
 594

595 
$$\mathcal{P}_s^{aTV} = \{q \in \Delta_{\mathcal{S}} : \frac{1}{2} \|q - P_0(\cdot|s, a)\|_1 \leq \delta\}.$$

594 Next, we show that the optimization problem given in Equation 4 in the main body of the paper  
 595 with  $\mathcal{P}_s^{a,TV}$  as the uncertainty set satisfies the conditions described in Lemma 2. Let us rewrite the  
 596 optimization problem here for the TV distance uncertainty set.  
 597

$$598 \sigma_{\mathcal{P}_s^{a,TV}}(V) = \min_{q \in \mathcal{P}_s^{a,TV}} q^\top V.$$

600 From Lemma 4.3 in Iyengar (2005), we know that the above optimization problem can be solved  
 601 under the dual formulation :  
 602

$$603 \sigma_{\mathcal{P}_s^{a,TV}}(V) = \max_{f \in \mathbb{R}_+^{|\mathcal{S}|}} (\mathbb{E}_{S \sim P_0(\cdot|s,a)}[V(S) - f(S)] - \delta \text{span}(V - f)). \quad (13)$$

605 Next, we prove that the above dual optimization problem is equivalent to a scalar optimization  
 606 problem.  
 607

608 **Lemma 3.** *The optimization problem given in Equation (13) is equivalent to the following optimiza-  
 609 tion problem:*

$$610 \sigma_{\mathcal{P}_s^{a,TV}}(V) \equiv \delta \min_{s'} V(s') + \max_{\lambda_s^a \in [\min_{s'} V(s'), \max_{s'} V(s')]} \{\mathbb{E}_{S \sim P_0(\cdot|s,a)}[\min\{V(S), \lambda_s^a\}] - \delta \lambda_s^a\}. \quad (14)$$

614 **Proof. From the  $\mu$ -vector dual to a 1-D cut off problem:** The optimization problem in Equation  
 615 (13) can be written as

$$616 \max_{f \in \mathbb{R}_+^{|\mathcal{S}|}} \left\{ \mathbb{E}_{S \sim P_0(\cdot|s,a)}[V(S) - f(S)] - \delta [\max_{s'} (V(s') - f(s')) - \min_{s'} (V(s') - f(s'))] \right\}. \quad (15)$$

619 **Step 1 – restrict to “cut-off” vectors:** For any scalar  $z \in [\min_{s'} V(s'), \max_{s'} V(s')]$ , define

$$620 f_z(s) := [V(s) - z]_+ = \max\{0, V(s) - z\}.$$

623 Replacing an arbitrary feasible  $f$  by the corresponding  $f_{z:=\max_{s'} (V(s') - f(s'))}$  cannot decrease  
 624 the objective in equation 15, so an optimizer always has the form  $f_{z^*}$  for some  $z^* \in$   
 625  $[\min_{s'} V(s'), \max_{s'} V(s')]$ .

626 **Step 2 – plug  $f_z$  into the objective.** Because  $V(s) - f_z(s) = \min\{V(s), z\}$ ,

$$628 \max_s (V - f_z) = z, \quad \min_s (V - f_z) = \min_{s'} V(s'),$$

630 and

$$631 \mathbb{E}_{S \sim P_0}[V(S) - f_z(S)] = \mathbb{E}_{S \sim P_0}[\min\{V(S), z\}].$$

632 Substituting these identities into equation 15 yields the *scalar* optimization  
 633

$$634 \sigma_{\mathcal{P}_s^{a,TV}}(V) = \delta \min_{s'} V(s') + \max_{z \in [\min_{s'} V(s'), \max_{s'} V(s')]} \{\mathbb{E}_{S \sim P_0(\cdot|s,a)}[\min\{V(S), z\}] - \delta z\}. \quad (16)$$

637  $\square$

640 As we are dealing with  $V$  functions for which  $V(s) \in \{\frac{-1}{1-\gamma}, \frac{1}{1-\gamma}\}$ , the optimum dual variable lies  
 641 in:  $\lambda_s^a \in \{\frac{-1}{1-\gamma}, \frac{1}{1-\gamma}\}$  and we can equivalently write from Lemma 3,  
 642

$$643 \sigma_{\mathcal{P}_s^{a,TV}}(V) = \delta \min_{s'} V(s') + \max_{\lambda_s^a \in \{\frac{-1}{1-\gamma}, \frac{1}{1-\gamma}\}} \{\mathbb{E}_{S \sim P_0(\cdot|s,a)}[\min\{V(S), \lambda_s^a\}] - \delta \lambda_s^a\}.$$

645 It is easy to verify that the concave objective has a super-gradient:

$$646 G^{TV}(\lambda_s^a; V, P_0(\cdot|s,a)) := \mathbb{P}_{S \sim P_0(\cdot|s,a)}[V(S) \geq \lambda_s^a] - \delta.$$

648 An unbiased estimate of the super-gradient for a value of  $\lambda_s^a$  and the value function  $V$  from a next  
 649 state  $S' \sim P_0(\cdot|s, a)$  can be given as:  
 650

$$651 \hat{G}^{TV}(\lambda_s^a; V, S') := \mathbf{1}_{V(S') \geq \lambda_s^a} - \delta. \quad (17)$$

652 An unbiased estimate of the dual objective for a value of  $\lambda_s^a$  and the value function  $V$  from a next  
 653 state  $S' \sim P_0(\cdot|s, a)$  can be given as  
 654

$$655 \hat{F}^{TV}(\lambda_s^a; V, S') = \delta \min_{s'} V(s') + \min(V(S'), \lambda_s^a) - \delta \lambda_s^a. \quad (18)$$

657 As we have  $|V(s)| \leq \frac{1}{1-\gamma}$ ,  $\forall s \in \mathcal{S}$ , its easy to see that,  
 658

$$659 |\hat{G}^{TV}(\lambda_s^a; V, S')| \leq C_G^{TV} := \max(\delta, 1 - \delta), \forall \lambda_s^a \in \mathbb{R},$$

661 and, for any  $0 < M < \infty$ ,

$$663 |\hat{F}^{TV}(\lambda_s^a; V, S')| \leq C_{F,M}^{TV} := (1 + \delta) \left( M + \frac{1}{1 - \gamma} \right), \forall \lambda_s^a \in [-M, M].$$

## 665 B.2 WASSERSTEIN- $\ell$ UNCERTAINTY SET

667 We define the Wasserstein- $\ell$  uncertainty set for each  $(s, a)$  pair as:  
 668

$$669 \mathcal{P}_s^{aW_\ell} = \{q \in \Delta_{\mathcal{S}} : W_\ell(P_0(\cdot|s, a), q) \leq \delta\},$$

670 where  $\delta > 0$  is the uncertainty radius and  $W_\ell(P_0(\cdot|s, a), q)$  is the Wasserstein- $\ell$  distance defined next.  
 671 Consider the generic metric space  $(\mathcal{S}, d)$  by defining some distance metric  $d$ . For some parameter  
 672  $\ell \in [1, \infty)$ , and two distributions  $p, q \in \Delta_{\mathcal{S}}$ , define the Wasserstein- $\ell$  distance between them as  
 673  $W_\ell(p, q) = \inf_{N \in \Gamma(p, q)} \|d\|_{N, \ell}$ , where  $\Gamma(p, q)$  denotes the distribution over  $\mathcal{S} \times \mathcal{S}$  with marginal  
 674 distributions  $p, q$  and  $\|d\|_{N, \ell} = (\mathbb{E}_{(X, Y) \sim N}[d(X, Y)^\ell])^{1/\ell}$ . Let us use the distance matrix with  
 675 normalization, ensuring  $|d(s, s')| \leq 1, \forall (s, s')$ .  
 676

677 Next, we show that the following optimization problem with  $\mathcal{P}_s^{aW_\ell}$  as the uncertainty set satisfies the  
 678 conditions described in Lemma 2.

$$679 \sigma_{\mathcal{P}_s^{aW_\ell}}(V) = \min_{q \in \mathcal{P}_s^{aW_\ell}} q^\top V.$$

682 From Gao & Kleywegt (2023), we know that the above optimization problem can be solved under the  
 683 dual formulation :

$$685 \sigma_{\mathcal{P}_s^a}(V) = \sup_{\lambda_s^a \geq 0} \left( -\lambda_s^a \delta^\ell + \mathbb{E}_{P_0(\cdot|s, a)} [\inf_y (V(y) + \lambda_s^a d(S, y)^\ell)] \right).$$

687 As the state space  $\mathcal{S}$  is finite, we can replace the inner-optimization  $[\inf_y (V(y) + \lambda_s^a d(S, y)^\ell)]$  with  
 688  $[\min_y (V(y) + \lambda_s^a d(S, y)^\ell)]$ . Next, we show that the optimum dual variable of the above optimization  
 689 problem lies inside a compact set  $[0, \lambda_M^{W_\ell}]$  with  $\lambda_M^{W_\ell} := \frac{\text{span}(V)}{\delta^\ell}$ .  
 690

692 As point-wise minimum of affine functions is concave, the above optimization problem is a concave  
 693 optimization problem. It is easy to verify that the concave objective has a super-gradient:  
 694

$$G^{W_\ell}(\lambda_s^a; V, P_0(\cdot|s, a)) = -\delta^\ell + \mathbb{E}_{X \sim P_0(\cdot|s, a)} [d(X, y_{\lambda_s^a}^*(X))^\ell], \quad (19)$$

696 where,

$$697 y_{\lambda_s^a}^*(x) \in \arg \min_y [V(y) + \lambda_s^a d(x, y)^\ell].$$

699 Let us fix an  $S = s$  and its minimizer  $y_{\lambda_s^a}^*(x)$  for the inner-optimization  $[\inf_y (V(y) + \lambda_s^a d(s, y)^\ell)]$ .  
 700 Because the candidate  $y = s$  is always feasible,  
 701

$$V(y_{\lambda_s^a}^*(s)) + \lambda_s^a d(s, y_{\lambda_s^a}^*)^\ell \leq V(s).$$

702 Rearrange:

$$703 \quad 704 \quad 705 \quad d(s, y_{\lambda_s^a}^*(s))^\ell \leq \frac{V(s) - V(y_{\lambda_s^a}^*(s))}{\lambda_s^a} \leq \frac{\text{span}(V)}{\lambda_s^a}.$$

706 Taking expectation in Equation 19 and using the above equation gives

$$707 \quad 708 \quad 709 \quad G^{W_\ell}(\lambda_s^a; V, P_0(\cdot|s, a)) \leq -\delta^\ell + \frac{\text{span}(V)}{\lambda_s^a}.$$

710 Now, for any  $\lambda_s^a > \lambda_M^{W_\ell} = \frac{\text{span}(V)}{\delta^\ell}$ , we have,

$$711 \quad 712 \quad 713 \quad G^{W_\ell}(\lambda_s^a; V, P_0(\cdot|s, a)) \leq 0.$$

714 Due to the concavity of the objective, a non-positive super-gradient means the function is non-  
715 increasing for all  $\lambda_s^a > \lambda_M^{W_\ell}$ . Combining the observation with the boundedness of the objective for  
716 bounded  $\lambda_s^a$ , we conclude that the supremum is attained and lies in  $[0, \lambda_M^{W_\ell}]$ .

717 An unbiased estimate of the super-gradient for a value of  $\lambda_s^a$  and the value function  $V$  from a next  
718 state  $S' \sim P_0(\cdot|s, a)$  can be given as:

$$719 \quad 720 \quad \hat{G}^{W_\ell}(\lambda_s^a; V, S') = -\delta^\ell + d(S', y^{*'})^\ell, \quad (20)$$

721 where,

$$722 \quad 723 \quad y^{*'} = \arg \min_y [V(y) + \lambda_s^a d(S', y)^\ell].$$

724 An unbiased estimate of the dual objective for a value of  $\lambda_s^a$  and the value function  $V$  from a next  
725 state  $S' \sim P_0(\cdot|s, a)$  can be given as

$$726 \quad 727 \quad \hat{F}^{W_\ell}(\lambda_s^a; V, S') = -\lambda_s^a \delta^\ell + V(y^{*'}) + \lambda_s^a d(S', y^{*'})^\ell. \quad (21)$$

728 If we assume  $|V(s)| \leq \frac{1}{1-\gamma}$ ,  $\forall s \in \mathcal{S}$ , its easy to show that,

$$729 \quad 730 \quad |\hat{G}^{W_\ell}(\lambda_s^a; V, S')| \leq C_G^{W_\ell} := 1 + \delta^\ell, \forall \lambda_s^a \in \mathbb{R},$$

731 and, for any  $0 < M < \infty$ ,

$$732 \quad 733 \quad |\hat{F}^{W_\ell}(\lambda_s^a; V, S')| \leq C_{F,M}^{W_\ell} := (\delta^\ell + 1)M + \frac{1}{1-\gamma}, \forall \lambda_s^a \in [-M, M].$$

## 737 C CONVERGENCE ANALYSIS OF ALGORITHM 1 AND THE PROOF OF 738 THEOREM 1

741 In this section, we provide the proof of Theorem 1. The proof follows in a similar manner described  
742 in the proof sketch in the main body of the paper. We start with proving Lemma 1 which establishes  
743 the convergence of the outer loop iterates in terms of inner loop convergence error. Subsequently, we  
744 establish the convergence of the inner loop. Finally we combine them to prove Theorem 1.

### 745 C.1 OUTER LOOP CONVERGENCE ANALYSIS: PROOF OF LEMMA 1

746 In this subsection, we prove Lemma 1. The proof is inspired by the analysis in Chen et al. (2023) for  
747 non-robust Q-learning. The analysis of the outer loop follows from the paper (Chen et al., 2023). To  
748 write the bound for the outer loop, we have to start with a few notations as used in the mentioned  
749 paper. Recall, the function approximation error  $\epsilon_{\text{approx}}$  is defined as:

$$750 \quad 751 \quad \epsilon_{\text{approx}} := \sup_{Q=\text{Clip}(\Phi\theta), \theta \in \mathbb{R}^{n_\theta}} \|\text{Clip}(\Pi \mathcal{T}^{\text{rob}, \pi}(Q)) - \mathcal{T}^{\text{rob}, \pi}(Q)\|_\infty.$$

752 Also, recall the definition of  $\theta_t^*$  from Equation 10.

753 Recall the fact that  $Q^{\text{rob}, \pi} = \mathcal{T}^{\text{rob}, \pi}(Q^{\text{rob}, \pi})$ .

756 Then, for any  $t = 1, 2, \dots, T$ , we have,

$$\begin{aligned}
 758 \quad \|\hat{Q}_t - Q^{\text{rob}, \pi}\|_\infty &= \|\text{Clip}(\Phi\hat{\theta}_t) - \mathcal{T}^{\text{rob}, \pi}(Q^{\text{rob}, \pi})\|_\infty \\
 759 \quad &= \underbrace{\|(\mathcal{T}^{\text{rob}, \pi}(\hat{Q}_{t-1}) - \mathcal{T}^{\text{rob}, \pi}(Q^{\text{rob}, \pi}))\|_\infty}_I \\
 760 \quad &\quad + \underbrace{\|(\text{Clip}(\Phi\hat{\theta}_t) - \text{Clip}(\Pi\mathcal{T}^{\text{rob}, \pi}(\hat{Q}_{t-1})))\|_\infty}_{II} \\
 761 \quad &\quad + \underbrace{\|(\mathcal{T}^{\text{rob}, \pi}(\hat{Q}_{t-1}) - \text{Clip}(\Pi\mathcal{T}^{\text{rob}, \pi}(\hat{Q}_{t-1})))\|_\infty}_{\leq \epsilon_{\text{approx}}}.
 \end{aligned}$$

767 **First Term:**

$$I = \|(\mathcal{T}^{\text{rob}, \pi}(\hat{Q}_{t-1}) - \mathcal{T}^{\text{rob}, \pi}(Q^{\text{rob}, \pi}))\|_\infty \leq \gamma \|\hat{Q}_{t-1} - Q^{\text{rob}, \pi}\|_\infty,$$

772 as the robust bellman operator is a  $\gamma$ -contraction with respect to the  $\infty$ -norm (Iyengar, 2005).

773 **Second Term:**

$$\begin{aligned}
 776 \quad II &= \|(\text{Clip}(\Phi\hat{\theta}_t) - \text{Clip}(\Pi\mathcal{T}^{\text{rob}, \pi}(\hat{Q}_{t-1})))\|_\infty \\
 777 \quad &\leq \underbrace{\|\Phi\hat{\theta}_t - \Pi\mathcal{T}^{\text{rob}, \pi}(\hat{Q}_{t-1})\|_\infty}_{(a)} \\
 778 \quad &= \underbrace{\|\Phi(\theta_{t-1, K} - \theta_{t-1}^*)\|_\infty}_{(b)} \\
 779 \quad &\leq \max_{s, a} \|\phi(s, a)\|_2 \|\theta_{t-1, K} - \theta_{t-1}^*\|_2 \\
 780 \quad &\leq \underbrace{\|\theta_{t-1, K} - \theta_{t-1}^*\|_2}_{(c)}.
 \end{aligned}$$

786 where (a) is using the non-expansive property of the clipping operator with respect to  $\|\cdot\|_\infty$ ; for (b),  
787 recall the definition of  $\theta_t^*$  in the inner loop in Equation 10; for (c), assume  $\|\phi(s, a)\|_2 \leq 1, \forall (s, a) \in$   
788  $\mathcal{S} \times \mathcal{A}$ .

789 Hence, we get:

$$\|\hat{Q}_t - Q^{\text{rob}, \pi}\|_\infty \leq \gamma \|\hat{Q}_{t-1} - Q^{\text{rob}, \pi}\|_\infty + \|\theta_{t-1, K} - \theta_{t-1}^*\|_2 + \epsilon_{\text{approx}}.$$

793 Unroll the recursion and take the expectation:

$$\mathbb{E}\|\hat{Q}_T - Q^{\text{rob}, \pi}\|_\infty \leq \gamma^T \|\hat{Q}_0 - Q^{\text{rob}, \pi}\|_\infty + \sum_{t=1}^T \gamma^{T-t} \mathbb{E}[\|\theta_{t-1, K} - \theta_{t-1}^*\|_2] + \frac{\epsilon_{\text{approx}}}{1 - \gamma}.$$

## 799 C.2 INNER LOOP CONVERGENCE ANALYSIS

801 The purpose of this subsection is to bound the term  $\mathbb{E}[\|\theta_{t-1, K} - \theta_{t-1}^*\|_2]$  for any fixed outer loop  
802 iteration index  $t$ .

803 Using the notations stated in the Lemma 2, we instantiate different problem-dependent constants,  
804 namely  $C_G$  and  $C_{F, M}$  as follows for different uncertainty sets. In the fast time-scale, we have  
805  $\|\nu_{t, k}\|_2 \leq B_\nu$ . This implies,  $|\psi(s, a)^\top \nu_{t, k}| \leq B_\nu$  for all  $(s, a)$  pair. Hence, for the rest of the  
806 Appendix denote  $M := B_\nu$  and similarly  $C_{F, M} := C_{B_\nu}$ .

807 Recall from the Appendix B, we know that,

1. For total variation uncertainty sets,  $C_{B_\nu} = (1 + \delta) \left( M + \frac{1}{1 - \gamma} \right)$ .

810 2. For Wasserstein- $\ell$  uncertainty sets,  $C_{B_\nu} = (1 + \delta^\ell)B_\nu + \frac{1}{1-\gamma}$ .  
 811

812 Also, from the Appendix B, we know that,  
 813

814 1. For total variation uncertainty sets,  $C_G = \max(\delta, 1 - \delta)$ .  
 815 2. For Wasserstein- $\ell$  uncertainty sets,  $C_G = 1 + \delta^\ell$ .  
 816

817 In this subsection, we show that for each outer iteration  $t$ , the inner loop parameter  $\theta_{t,k}$  converges to  
 818  $\theta_t^*$  as defined in Equation (10). Recall the definition  $F_{s,a}^{*,V} := \sup_{\lambda_s^a} F(\lambda_s^a; V, P_0(\cdot|s, a))$ . We denote  
 819 the inner loop error for outer iteration index  $t$  as  $e_{t,k} := \theta_{t,k} - \theta_t^*$  with  
 820

$$821 \theta_t^* := (\Phi^\top D^\pi \Phi)^{-1} \Phi^\top D^\pi [r + \gamma F^{*,V_{\theta_t}^{\text{rob}}}] . \quad (22)$$

822 Using earlier notation, the dual objective corresponding to an  $(s, a)$ -pair for a target value function  
 823  $V_{\hat{\theta}_t}^{\text{rob}}$  is  
 824

$$825 \max_{\lambda_s^a} F(\lambda_s^a; V_{\hat{\theta}_t}^{\text{rob}}, P_0(\cdot|s, a)).$$

826 For the rest of the discussion in this subsection, let us fix an outer loop iteration  $t$  and the target  
 827 parameter  $\hat{\theta}_t$  is treated as a deterministic vector. For a given outer loop index  $t$ , for all inner loop  
 828 iterations  $k \geq 1$  let the filtration  $\mathcal{F}_{t,k}$  be the sigma algebra generated by the transitions sampled till  
 829 inner loop iteration index  $k - 1$  that is, on the set  $\{S_{t,j}, A_{t,j}, S_{t,j+1} : 0 \leq j \leq k - 1\}$ .  
 830

831 Observe that the pair process  $Z_{t,k} := (S_{t,k}, A_{t,k})$  is a Markov chain. We define another filtration  
 832  $\mathcal{G}_{t,k}$  as the sigma algebra over the set  $\{Z_{t,0}, Z_{t,1}, \dots, Z_{t,k}\}$ .  
 833

### 834 C.2.1 ANALYSIS ON THE FAST TIME-SCALE:

835 The fast time-scale update is given as  
 836

$$837 \nu_{t,k+1} = \text{Proj}_{\mathcal{M}_\nu} (\nu_{t,k} + \beta_k [\hat{G}(\psi(S_{t,k}, A_{t,k})^\top \nu_{t,k}; V_{\hat{\theta}_t}^{\text{rob}}, S_{t,k+1}) \psi(S_{t,k}, A_{t,k})]).$$

838 Define the diagonal matrix  $D_{t,k} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}||\mathcal{A}|}$  with each diagonal element as  
 839  $D_{t,k}((s, a), (s, a)) = \mathbf{1}_{(s,a)=(S_{t,k},A_{t,k})}$ .  
 840

841 For each outer iteration  $t$  and dual vector  $\nu_{t,k}$ , we define a vector  $\bar{g}_t(\nu_{t,k}) \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ , indexed by  
 842  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , via  
 843

$$844 \bar{g}_t(\nu_{t,k})_{s,a} := \mathbb{E}_{S' \sim P_0(\cdot|s,a)} [\hat{G}(\psi(s, a)^\top \nu_{t,k}; V_{\hat{\theta}_t}^{\text{rob}}, S')].$$

845 Here  $S'$  denotes the next-state random variable with distribution  $P_0(\cdot | s, a)$ .  
 846

847 Also, define the stochastic update vector  $X_{t,k} \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$  defined as  
 848

$$849 [X_{t,k}]_{s,a} := \mathbf{1}_{(s,a)=(S_{t,k},A_{t,k})} \cdot \hat{G}(\psi(s, a)^\top \nu_{t,k}; V_{\hat{\theta}_t}^{\text{rob}}, S_{t,k+1}) \psi(s, a).$$

850 We split the update into stationary drift and different noise terms as:  
 851

$$852 \nu_{t,k+1} = \text{Proj}_{\mathcal{M}_\nu} \left( \nu_{t,k} + \beta_k \left[ \Psi^\top D^\pi \bar{g}_t(\nu_{t,k}) + \underbrace{X_{t,k} - \mathbb{E}[X_{t,k} | \mathcal{G}_{t,k}]}_{m_{k+1}^\nu} + \underbrace{\mathbb{E}[X_{t,k} | \mathcal{G}_{t,k}] - \Psi^\top D^\pi \bar{g}_t(\nu_{t,k})}_{\zeta_{k+1}^\nu} \right] \right).$$

853 We see that,  
 854

$$855 \zeta_{k+1}^\nu = \mathbb{E}[X_{t,k} | \mathcal{G}_{t,k}] - \Psi^\top D^\pi \bar{g}_t(\nu_{t,k}).$$

856 So the update now becomes,  
 857

$$858 \nu_{t,k+1} = \text{Proj}_{\mathcal{M}_\nu} (\nu_{t,k} + \beta_k [\Psi^\top D^\pi \bar{g}_t(\nu_{t,k}) + m_{k+1}^\nu + \zeta_{k+1}^\nu]).$$

864 In the above equation,  $m_{k+1}^\nu$  denotes the state-innovation noise that is a martingale difference on the  
 865 filtration  $\mathcal{G}_{t,k}$ .

866 Hence,

868

$$869 \mathbb{E}[m_{k+1}^\nu | \mathcal{G}_{t,k}] = 0.$$

870

871 We analyze the finite time convergence of the fast time-scale first. We show that the fast time-scale  
 872 update is equivalent to performing a stochastic gradient super-gradient ascent on the following  
 873 objective function:

874

$$875 L_t(\nu) := \sum_{s,a} d^\pi(s,a) F(\psi(s,a)^\top \nu; V_{\hat{\theta}_t}, P_0(\cdot|s,a)).$$

876

877 It is easy to show that  $L_t(\nu)$  is concave on  $\nu$ . Let  $\nu_t^*$  be one maximizer of the above objective function  
 878  $L_t(\nu)$  in the domain  $\mathcal{M}_\nu$ . Notice that the fast time-scale update depends on the target parameter  
 879 vector  $\hat{\theta}_t$  and independent of the slow time-scale parameters for a fixed outer loop.

880

881 Our goal is to bound the sub-optimality gap of the dual objective (the weighted objective  $L_t(\nu)$ ) for  
 882 each iteration in the inner loop. We will be able to use the error in estimating the dual objective from  
 883 the fast time-scale as a bias in the slow time-scale to get a sample complexity bound for the inner  
 884 loop of the algorithm 1.

885

Let us define the dual objective sub-optimality at  $\nu = \nu_{t,k}$  for the fast time-scale as :

886

$$887 L_{t,k} := \sum_{s,a} d^\pi(s,a) [F(\psi(s,a)^\top \nu_t^*; V_{\hat{\theta}_t}^{\text{rob}}, P_0(\cdot|s,a)) - F(\psi(s,a)^\top \nu_{t,k}; V_{\hat{\theta}_t}^{\text{rob}}, P_0(\cdot|s,a))].$$

888

Define the Lyapunov function for the fast time-scale as

889

$$890 e_{t,k}^\nu = \|\nu_{t,k} - \nu_t^*\|_2^2.$$

891

Using the non-expansiveness of projection,

892

$$893 e_{t,k+1}^\nu \leq \|\nu_{t,k} + \beta_k \Psi^\top D^\pi \bar{g}_t(\nu_{t,k}) + \beta_k (m_{k+1}^\nu + \zeta_{k+1}^\nu) - \nu_t^*\|_2^2.$$

894

We can write:

895

$$896 \begin{aligned} 897 e_{t,k+1}^\nu &\leq \|(\nu_{t,k} - \nu_t^*)\|_2^2 + 2\beta_k (\nu_{t,k} - \nu_t^*)^\top \Psi^\top D^\pi \bar{g}_t(\nu_{t,k}) \\ 898 &\quad + 2\beta_k (\nu_{t,k} - \nu_t^*)^\top (m_{k+1}^\nu + \zeta_{k+1}^\nu) \\ 899 &\quad + \beta_k^2 \|\Psi^\top D^\pi \bar{g}_t(\nu_{t,k}) + m_{k+1}^\nu + \zeta_{k+1}^\nu\|_2^2. \end{aligned}$$

900

We simplify the term  $(\nu_{t,k} - \nu_t^*)^\top \Psi^\top D^\pi \bar{g}_t(\nu_{t,k})$  first.

901

$$902 \begin{aligned} 903 (\nu_{t,k} - \nu_t^*)^\top \Psi^\top D^\pi \bar{g}_t(\nu_{t,k}) &= \sum_{s,a \in \mathcal{S} \times \mathcal{A}} d^\pi(s,a) \bar{g}_t(\nu_{t,k})(s,a) (\nu_{t,k} - \nu_t^*)^\top \psi(s,a) \\ 904 &= (\nu_{t,k} - \nu_t^*)^\top \nabla_{\nu=\nu_{t,k}} L_t(\nu). \end{aligned}$$

905

Hence, from the first order optimality condition on a concave objective, we can write:

906

$$907 (\nu_{t,k} - \nu_t^*)^\top \Psi^\top (D^\pi) \bar{g}_t(\nu_{t,k}) \leq -L_{t,k}$$

908

Hence, we get,

909

$$910 \begin{aligned} 911 e_{t,k+1}^\nu &\leq \|(\nu_{t,k} - \nu_t^*)\|_2^2 - 2\beta_k L_{t,k} \\ 912 &\quad + 2\beta_k (\nu_{t,k} - \nu_t^*)^\top (m_{k+1}^\nu) + 2\beta_k (\nu_{t,k} - \nu_t^*)^\top (\zeta_{k+1}^\nu) \\ 913 &\quad + \beta_k^2 \|\Psi^\top D^\pi \bar{g}_t(\nu_{t,k}) + m_{k+1}^\nu + \zeta_{k+1}^\nu\|_2^2. \end{aligned}$$

914

918 Now we condition on a lagged filtration  $\mathcal{G}_{t,k-\tau}$  where  $\tau \leq k$  would be chosen later.  
919  
920

$$\begin{aligned} 921 \mathbb{E}[e_{t,k+1}^\nu | \mathcal{G}_{t,k-\tau}] &\leq \mathbb{E}[e_{t,k}^\nu | \mathcal{G}_{t,k-\tau}] - 2\beta_k \mathbb{E}[L_{t,k} | \mathcal{G}_{t,k-\tau}] \\ 922 &\quad + 2\beta_k \mathbb{E}[(\nu_{t,k} - \nu_t^*)^\top (m_{k+1}^\nu) | \mathcal{G}_{t,k-\tau}] + 2\beta_k \mathbb{E}[(\nu_{t,k} - \nu_t^*)^\top (\zeta_{k+1}^\nu) | \mathcal{G}_{t,k-\tau}] \\ 923 &\quad + \beta_k^2 \mathbb{E}[\|\Psi^\top D^\pi \bar{g}_t(\nu_{t,k}) + m_{k+1}^\nu + \zeta_{k+1}^\nu\|_2^2 | \mathcal{G}_{t,k-\tau}]. \\ 924 \end{aligned}$$

925 Let us first bound the  $\beta_k^2$  terms. Recall from the conditions described in Lemma 2,  $\|\bar{g}_t(\nu_{t,k})\|_\infty \leq C_G$ .  
926 As  $\|\psi(s, a)\|_2 \leq 1, \forall (s, a) \in \mathcal{S} \times \mathcal{A}$ , one can easily show that  $\|\Psi^\top D^\pi \bar{g}_t(\nu_{t,k})\|_2 \leq C_G$  and  
927  $\|m_{k+1}^\nu\|_2 \leq 2C_G$  and  $\|\zeta_{k+1}^\nu\|_2 \leq 2C_G$ . Hence,

$$928 \beta_k^2 \mathbb{E}[\|\Psi^\top D^\pi \bar{g}_t(\nu_{t,k}) + m_{k+1}^\nu + \zeta_{k+1}^\nu\|_2^2 | \mathcal{G}_{t,k-\tau}] \leq 25\beta_k^2 C_G^2.$$

929 Now we work on the cross terms. Let us start with  $2\beta_k \mathbb{E}[(\nu_{t,k} - \nu_t^*)^\top D^\pi(m_{k+1}^\nu) | \mathcal{G}_{t,k-\tau}]$ .  
930

931 We write:

$$\begin{aligned} 934 2\beta_k \mathbb{E}[(\nu_{t,k} - \nu_t^*)^\top D^\pi(m_{k+1}^\nu) | \mathcal{G}_{t,k-\tau}] &= 2\beta_k \mathbb{E}[\mathbb{E}[(\nu_{t,k} - \nu_t^*)^\top D^\pi(m_{k+1}^\nu) | \mathcal{G}_{t,k}] | \mathcal{G}_{t,k-\tau}] \\ 935 &= 2\beta_k \mathbb{E}[(\nu_{t,k} - \nu_t^*)^\top D^\pi[\mathbb{E}[(m_{k+1}^\nu) | \mathcal{G}_{t,k}] | \mathcal{G}_{t,k-\tau}] \\ 936 &= 0. \\ 937 \end{aligned}$$

938 Now we focus on the term  $2\beta_k \mathbb{E}[(\nu_{t,k} - \nu_t^*)^\top (\zeta_{k+1}^\nu) | \mathcal{G}_{t,k-\tau}]$ .  
939

940 We use the shorthand

$$941 Z_{t,k} := (S_{t,k}, A_{t,k}) \in \mathcal{S} \times \mathcal{A}$$

942 for the state-action pair at outer iteration  $t$  and inner iteration  $k$ .  
943

944 Define the vector  $e_{Z_{t,k}}$   
945

$$946 e_{Z_{t,k}}(s, a) := \mathbf{1}_{S_{t,k}, A_{t,k} = s, a}.$$

947 Recall the definition  
948

$$949 \eta_k^{t,\tau}(\cdot) := \mathbb{P}((S_{t,k}, A_{t,k}) \in \cdot | S_{t,0}, A_{t,0}, \dots, S_{t,k-\tau}, A_{t,k-\tau})$$

950 From Assumption 1, we know the following holds:  
951

$$953 \|\eta_k^{t,\tau} - d^\pi\|_{\text{TV}} \leq C_{\text{mix}} \rho^\tau \quad (0 < \rho < 1).$$

955 Thus, we can write:  
956

$$\begin{aligned} 957 2\beta_k \mathbb{E}[(\nu_{t,k} - \nu_t^*)^\top (\Psi^\top (e_{Z_{t,k}} - d^\pi) \odot \bar{g}_t(\nu_{t,k})) | \mathcal{G}_{t,k-\tau}] \\ 958 &= 2\beta_k \mathbb{E}[(\Psi(\nu_{t,k} - \nu_t^*))^\top ((e_{Z_{t,k}} - d^\pi) \odot \bar{g}_t(\nu_{t,k})) | \mathcal{G}_{t,k-\tau}] \\ 959 &\leq 2\beta_k \mathbb{E}[\|(\Psi(\nu_{t,k} - \nu_t^*))\|_\infty \|((e_{Z_{t,k}} - d^\pi) \odot \bar{g}_t(\nu_{t,k}))\|_1 | \mathcal{G}_{t,k-\tau}] \\ 960 &\leq 2\beta_k \mathbb{E}[2B_\nu C_G \|((e_{Z_{t,k}} - d^\pi)\|_1 | \mathcal{G}_{t,k-\tau}] \\ 961 &\leq 8\beta_k B_\nu C_G C_{\text{mix}} \rho^\tau. \\ 962 \end{aligned}$$

963 Putting it together, we have:  
964

$$966 \mathbb{E}[e_{t,k+1}^\nu | \mathcal{G}_{t,k-\tau}] \leq \mathbb{E}[e_{t,k}^\nu | \mathcal{G}_{t,k-\tau}] - 2\beta_k \mathbb{E}[L_{t,k} | \mathcal{G}_{t,k-\tau}] + 8\beta_k B_\nu C_G C_{\text{mix}} \rho^\tau + 25\beta_k^2 g_m^2.$$

968 Now we choose  $\forall l \in (\lfloor k/2 \rfloor, k-1)$ ,  $\tau = \tau_{\beta_l} := \lceil \frac{\log(\frac{C_{\text{mix}}}{\beta_0} \sqrt{l+1})}{\log(\frac{1}{\rho})} \rceil$ , and assume  $k \geq 2\tau_{\beta_k}$  then,  
969 taking unconditional expectation gives us: (as  $C_{\text{mix}} \rho^{\tau_{\beta_l}} \leq \beta_l$ )  $\forall l \in (\lfloor k/2 \rfloor, k-1)$ :  
970

$$971 2\beta_k \mathbb{E}[L_l] \leq \mathbb{E}[e_{t,l}^\nu] - \mathbb{E}[e_{t,l+1}^\nu] + \beta_l^2 (8B_\nu C_G + 25C_G^2).$$

972 Next, we use telescoping for iterates over the index  $l$  from  $\lfloor \frac{k}{2} \rfloor$  to  $k-1$ .  
 973

974 
$$2 \sum_{l=\lfloor k/2 \rfloor}^{k-1} \beta_l \mathbb{E}[L_l] \leq 4B_\nu^2 + (8B_\nu C_G + 25C_G^2) \sum_{l=\lfloor k/2 \rfloor}^{k-1} \beta_l^2, \quad (23)$$
  
 975  
 976  
 977

978 where, we used that  $\|e_{t,\lfloor k/2 \rfloor}\|_2 \leq 4B_\nu^2$  due to the projection step  $\text{Proj}_{\mathcal{M}_\nu}$ . Recall, the fast time-scale  
 979 passes the following suffix-average of the dual parameter vector iterates to the slow time-scale at  
 980 each iterate  $k$ :

981  
 982 
$$\bar{\nu}_{t,k} = \frac{1}{\lceil k/2 \rceil} \sum_{l=\lfloor k/2 \rfloor}^{k-1} \nu_{t,k}.$$
  
 983  
 984  
 985

986 We use the step-size rule of  $\beta_k = \frac{\beta_0}{\sqrt{k+1}}$ .  
 987

988 Similar to the definition of  $L_{t,k}$ , let us define

989 
$$\bar{L}_{t,k} := \sum_{s,a} d^\pi(s,a) [F(\psi(s,a)^\top \nu_t^*; V_{\theta_t}^{\text{rob}}, P_0(\cdot|s,a)) - F(\psi(s,a)^\top \bar{\nu}_{t,k}; V_{\theta_t}^{\text{rob}}, P_0(\cdot|s,a))] \quad (24)$$
  
 990  
 991

992 Hence, using Jensen's inequality, we write:

993  
 994 
$$\mathbb{E}[\bar{L}_{t,k}] \leq \frac{1}{\lceil k/2 \rceil} \sum_{l=\lfloor k/2 \rfloor}^{k-1} \mathbb{E}[L_{t,l}]$$
  
 995  
 996 
$$\stackrel{(a)}{\leq} \frac{1}{\lceil k/2 \rceil} \frac{\sqrt{k}}{\beta_0} \sum_{l=\lfloor k/2 \rfloor}^{k-1} \beta_l \mathbb{E}[L_{t,k}] \leq \frac{2}{k} \frac{\sqrt{k}}{\beta_0} \sum_{l=\lfloor k/2 \rfloor}^{k-1} \beta_l \mathbb{E}[L_{t,l}]$$
  
 997  
 998 
$$\stackrel{(b)}{\leq} \frac{4B_\nu^2}{\beta_0 \sqrt{k}} + \frac{(8B_\nu C_G + 25C_G^2)}{\beta_0 \sqrt{k}} \sum_{l=\lfloor k/2 \rfloor}^{k-1} \beta_l^2$$
  
 999  
 1000  
 1001  
 1002  
 1003  
 1004 
$$= \frac{4B_\nu^2}{\beta_0 \sqrt{k}} + \frac{(8B_\nu C_G + 25C_G^2)}{\beta_0 \sqrt{k}} \sum_{l=\lfloor k/2 \rfloor}^{k-1} \frac{\beta_0^2}{l+1}$$
  
 1005  
 1006  
 1007 
$$\stackrel{(c)}{\leq} \frac{4B_\nu^2}{\beta_0 \sqrt{k}} + \frac{(8B_\nu C_G + 25C_G^2) (\beta_0^2 (1 + \ln(k) - \ln(k/2)))}{\beta_0 \sqrt{k}}$$
  
 1008  
 1009  
 1010  
 1011

1012 where

1013 
$$C_{\text{fast}} = \frac{(4B_\nu^2 + \beta_0^2 (8B_\nu C_G + 25C_G^2) \ln(2e))}{\beta_0}. \quad (25)$$
  
 1014

1015 In (a), we used  $\beta_k \geq \frac{\beta_0}{\sqrt{k}}$ ,  $\forall k \leq k-1$ . In (b), we used Equation 23. In (c), we used the following  
 1016 identity:  $\ln(k) \leq 1 + \frac{1}{2} + \dots + \frac{1}{k} \leq 1 + \ln(k)$ .  
 1017

1018 In summary, we have the following guarantee from the fast time-scale:

1019 **Lemma 4.** *Fix an outer loop  $t \geq 0$ . The following holds for the fast time-scale iterates of the*  
 1020 *Algorithm 1: If  $k \geq 2 \left\lceil \frac{\log(\frac{C_{\text{mix}}}{\beta_0} \sqrt{k+1})}{\log(\frac{1}{\rho})} \right\rceil$ ,*  
 1021

1022  
 1023 
$$\mathbb{E}[\bar{L}_{t,k}] \leq \frac{C_{\text{fast}}}{\sqrt{k}}, \quad (26)$$
  
 1024  
 1025

where,  $\bar{L}_{t,k}$  is defined in Equation (24) and  $C_{\text{fast}}$  is given in Equation (25).

1026 C.2.2 SLOW TIME-SCALE ANALYSIS  
1027

1028 Next, we will prove the convergence of the slow time-scale iterates  $\theta_{t,k}$  to  $\theta_t^*$  for sufficiently large  $k$ .  
1029 We denote the inner loop error for outer iteration index  $t$  as  $e_{t,k} := \theta_{t,k} - \theta_t^*$  with

1030 1031 
$$\theta_t^* := (\Phi^\top D^\pi \Phi)^{-1} \Phi^\top D^\pi [r + \gamma F^{*,V_{\theta_t}^{\text{rob}}}] \quad (27)$$
  
1032

1033 Recall the notation  $Z_{t,k} = (S_{t,k}, A_{t,k})$ . The slow update is

1034 1035 
$$\theta_{t,k+1} = \theta_{t,k} + \alpha_k T D_{t,k+1} \phi(Z_{t,k}),$$

1036 1037 
$$T D_{t,k+1} = R(Z_{t,k}) + \gamma \hat{F}(\psi(Z_{t,k})^\top \bar{\nu}_{t,k}; V_{\hat{\theta}}^{\text{rob}}, S_{t,k+1}) - \phi(Z_{t,k})^\top \theta_{t,k},$$

1038 with  $\bar{\nu}_{t,k}$  the suffix average produced by the fast time-scale updates.

1039 For each outer iteration  $t$  and dual vector  $\nu_{t,k}$ , we define a vector  $\bar{f}_t(\nu_{t,k}) \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}|}$ , indexed by  
1040  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , via

1041 1042 
$$[\bar{f}_t(\nu_{t,k})]_{s,a} := \mathbb{E}_{S' \sim P_0(\cdot | s, a)} [\hat{F}(\psi(s, a)^\top \nu_{t,k}; V_{\hat{\theta}}^{\text{rob}}, S')].$$

1043 Here  $S'$  denotes the next-state random variable with distribution  $P_0(\cdot | s, a)$ .

1044 We decompose the term  $T D_{t,k+1} \phi(Z_{t,k})$  as

1045 1046 
$$H_{t,k}^\theta + b_{t,k}^\theta + \xi_{t,k+1}^\theta + m_{t,k+1}^\theta,$$

1047 where

1048 1049 
$$H_{t,k}^\theta := \Phi^\top D^\pi [r + \gamma F^{*,V_{\theta_t}^{\text{rob}}} - \Phi \theta_{t,k}] \underbrace{=}_{\text{from Equation (10)}} \Phi^\top D^\pi \Phi (\theta_t^* - \theta_{t,k}),$$

1050 1051 
$$b_{t,k}^\theta := \gamma \Phi^\top D^\pi [F^{V_{\theta_t}^{\text{rob}}}(\bar{\nu}_{t,k}) - F^{*,V_{\theta_t}^{\text{rob}}}],$$

1052 1053 
$$\xi_{t,k+1}^\theta := \Phi^\top (D_{t,k} - D^\pi) [r + \gamma F^{V_{\theta_t}^{\text{rob}}}(\bar{\nu}_{t,k}) - \Phi \theta_{t,k}]$$

1054 1055 
$$= \Phi^\top (D_{t,k} - D^\pi) [r + \gamma \bar{f}_t(\bar{\nu}_{t,k}) - \Phi \theta_{t,k}],$$

1056 1057 
$$m_{t,k+1}^\theta := \gamma \Phi^\top D_{t,k} (\hat{F}(\psi(Z_{t,k})^\top \bar{\nu}_{t,k}; S_{k+1}, V_{\hat{\theta}}^{\text{rob}}) - \bar{f}_t(\bar{\nu}_{t,k})).$$

1058 Note  $\mathbb{E}[m_{t,k+1}^\theta | \mathcal{G}_{t,k}] = 0$  and, by tower property of conditional expectation,  $\mathbb{E}[e_{t,k}^\top m_{t,k+1}^\theta | \mathcal{G}_{t,k-\tau}] = 0$  for a lag  $\tau \leq k$ .

1059 1060 Recall the definition

1061 1062 
$$\eta_k^{t,\tau}(\cdot) := \mathbb{P}((S_{t,k}, A_{t,k}) \in \cdot | S_{t,0}, A_{t,0}, \dots, S_{t,k-\tau}, A_{t,k-\tau})$$

1063 From Assumption 1, we know the following holds:

1064 1065 
$$\|\eta_k^{t,\tau} - d^\pi\|_{\text{TV}} \leq C_{\text{mix}} \rho^\tau \quad (0 < \rho < 1).$$

1066 1067 For the fixed lag (to be chosen later)  $\tau \geq 1$  and define  $\mathcal{H}_{t,k}$  as the sigma algebra over the set  
1068  $\{\mathcal{G}_{t,k-\tau}, \theta_{t,k}, \bar{\nu}_{t,k}\}.$

1069 1070 Conditioning on  $\mathcal{H}_{t,k}$  “freezes”  $e_{t,k} := \theta_{t,k} - \theta_t^*$  and

1071 1072 
$$y_{t,k} := r + \gamma \bar{f}_t(\bar{\nu}_{t,k}) - \Phi \theta_{t,k}.$$

1073 1074 We use: for any signed vector  $w$  on  $\mathcal{S} \times \mathcal{A}$ ,

1075 1076 
$$\|\Phi^\top w\|_2 = \left\| \sum_z w_{s,a} \phi(s, a) \right\|_2 \leq \sum_{s,a} |w_{s,a}| = \|w\|_1, \quad (28)$$

1080 because each row vector satisfies  $\|\phi(s, a)\|_2 \leq 1$ . We also write the Markov noise term in terms of  
 1081  $y_{t,k}$  as

$$1082 \quad \xi_{t,k+1}^\theta = \Phi^\top (D_{t,k} - D^\pi) y_{t,k}.$$

1083 Finally set  $Y_0 := 1 + \gamma C_{F,M} + \frac{1}{(1-\gamma)\sqrt{\mu}}$  so that

$$1085 \quad \|y_{t,k}\|_\infty \leq Y_0 + \|e_{t,k}\|_2, \quad (29)$$

1087 using Lemma 7 and,  $r(\cdot) \in [0, 1]$ ,  $\|\bar{f}_t(\cdot)\|_\infty \leq C_{F,M}$ , and  $\|\phi(z)\|_2 \leq 1$ .

1088 Recall the definition  $e_{t,k} = \theta_{t,k} - \theta_t^*$  with

$$1090 \quad \theta_t^* := (\Phi^\top D^\pi \Phi)^{-1} \Phi^\top D^\pi [r + \gamma F^{*, V_{\theta_t}^{\text{rob}}}].$$

1092 (I) ONE-STEP LYAPUNOV EXPANSION UNDER A CONDITIONAL EXPECTATION WITH A  
 1093 FILTRATION UNDER A GENERIC LAG  $\tau$

1094 With  $x_{t,k} := \|e_{t,k}\|^2$ , for  $k \geq \tau$ , as  $\mathbb{E}[e_{t,k}^\top m_{t,k+1}^\theta \mid \mathcal{G}_{t,k-\tau}] = 0$ , we can write, if  $k \geq \tau$

$$1096 \quad \begin{aligned} \mathbb{E}[x_{k+1} \mid \mathcal{G}_{t,k-\tau}] &= \mathbb{E}\left[\|e_{t,k} + \alpha_k(H_{t,k}^\theta + b_{t,k}^\theta + \xi_{t,k+1}^\theta + m_{t,k+1}^\theta)\|^2 \mid \mathcal{G}_{t,k-\tau}\right] \\ 1097 &= x_{t,k} + 2\alpha_k \mathbb{E}[e_{t,k}^\top H_{t,k}^\theta \mid \mathcal{G}_{t,k-\tau}] + 2\alpha_k \mathbb{E}[e_{t,k}^\top b_{t,k}^\theta \mid \mathcal{G}_{t,k-\tau}] \\ 1098 &\quad + 2\alpha_k \mathbb{E}[e_{t,k}^\top \xi_{t,k+1}^\theta \mid \mathcal{G}_{t,k-\tau}] \\ 1099 &\quad + \alpha_k^2 \mathbb{E}[\|H_{t,k}^\theta + b_{t,k}^\theta + \xi_{t,k+1}^\theta + m_{t,k+1}^\theta\|^2 \mid \mathcal{G}_{t,k-\tau}]. \end{aligned} \quad (30)$$

1103 (II) MAIN DRIFT

1104 Recall, we denote  $\mu$  as the minimum eigenvalue of the matrix  $\Phi^\top D^\pi \Phi$  and from Assumption 1,  
 1105  $\mu > 0$ .

1107 Since  $e_{t,k}^\top H_{t,k}^\theta = -e_{t,k}^\top (\Phi^\top D^\pi \Phi) e_{t,k} \leq -\mu \|e_{t,k}\|^2$ ,

$$1109 \quad 2\alpha_k \mathbb{E}[e_{t,k}^\top H_{t,k}^\theta \mid \mathcal{G}_{t,k-\tau}] \leq -2\mu \alpha_k \mathbb{E}[x_{t,k} \mid \mathcal{G}_{t,k-\tau}]. \quad (31)$$

1110 (III) CROSS TERM CORRESPONDING TO BIAS

1112 By conditional Cauchy–Schwarz and Young inequality,

$$1114 \quad 2\alpha_k \mathbb{E}[e_{t,k}^\top b_{t,k}^\theta \mid \mathcal{G}_{t,k-\tau}] \leq \frac{\mu}{2} \alpha_k \mathbb{E}[x_{t,k} \mid \mathcal{G}_{t,k-\tau}] + \frac{2\alpha_k}{\mu} \mathbb{E}[\|b_{t,k}^\theta\|^2 \mid \mathcal{G}_{t,k-\tau}]. \quad (32)$$

1116 (IV) CROSS TERM CORRESPONDING TO MARKOV NOISE

1118 **Lemma 5** (Cross with Markov noise). *For any  $\tau \geq 1$ , if  $k \geq \tau$*

$$1121 \quad \begin{aligned} 2\alpha_k \mathbb{E}[e_{t,k}^\top \xi_{t,k+1}^\theta \mid \mathcal{G}_{t,k-\tau}] &\leq \left(\frac{\mu}{2} + 4 \min(1, C_{\text{mix}} \rho^\tau)\right) \alpha_k \mathbb{E}[\|e_{t,k}\|_2^2 \mid \mathcal{G}_{t,k-\tau}] \\ 1122 &\quad + \frac{8 Y_0^2}{\mu} \alpha_k \min(1, C_{\text{mix}} \rho^\tau)^2. \end{aligned}$$

1125 *Proof.* By the tower property,

$$1127 \quad \mathbb{E}[e_{t,k}^\top \xi_{t,k+1}^\theta \mid \mathcal{G}_{t,k-\tau}] = \mathbb{E}[\mathbb{E}[e_{t,k}^\top \Phi^\top (D_{t,k} - D^\pi) y_{t,k} \mid \mathcal{H}_{t,k}] \mid \mathcal{G}_{t,k-\tau}].$$

1129 Given  $\mathcal{H}_{t,k}$ , the only randomness is  $Z_{t,k} \sim \eta_k^{t,\tau}$ . Hence

$$1130 \quad \mathbb{E}[\Phi^\top (D_{t,k} - D^\pi) y_{t,k} \mid \mathcal{H}_{t,k}] = \Phi^\top (\eta_k^{t,\tau} - d^\pi) \odot y_{t,k}.$$

1132 Therefore,

$$1133 \quad \left| \mathbb{E}[e_{t,k}^\top \xi_{t,k+1}^\theta \mid \mathcal{G}_{t,k-\tau}] \right| \leq \mathbb{E}[\|e_{t,k}\|_2 \|\Phi^\top (\eta_k^{t,\tau} - d^\pi) \odot y_{t,k}\|_2 \mid \mathcal{G}_{t,k-\tau}].$$

1134 Apply Equation (28) and  $\|(\eta_k^{t,\tau} - d^\pi) \odot y_{t,k}\|_1 \leq \|\eta_k^{t,\tau} - d^\pi\|_1 \|y_{t,k}\|_\infty \leq$   
 1135  $2 \min(1, C_{\text{mix}} \rho^\tau) \|y_{t,k}\|_\infty$ :

1136

$$1137 \left| \mathbb{E}[e_{t,k}^\top \xi_{t,k+1}^\theta \mid \mathcal{G}_{t,k-\tau}] \right| \leq 2 \min(1, C_{\text{mix}} \rho^\tau) \mathbb{E}[\|e_{t,k}\|_2 \|y_{t,k}\|_\infty \mid \mathcal{G}_{t,k-\tau}].$$

1138

1139 Multiply by  $2\alpha_k$  and split  $\|y_{t,k}\|_\infty$  using Equation (29):

1140

$$1141 2\alpha_k \mathbb{E}[\cdot] \leq 4\alpha_k \min(1, C_{\text{mix}} \rho^\tau) \mathbb{E}[\|e_{t,k}\|_2 Y_0 \mid \mathcal{G}_{t,k-\tau}] + 4\alpha_k \min(1, C_{\text{mix}} \rho^\tau) \mathbb{E}[\|e_{t,k}\|_2 \|e_{t,k}\|_2 \mid \mathcal{G}_{t,k-\tau}].$$

1142

1143 For the first term in the right hand side, use Young's inequality:

1144

$$1145 4\alpha_k \min(1, C_{\text{mix}} \rho^\tau) Y_0 \|e_{t,k}\|_2 \leq \alpha_k \left( \frac{\mu}{2} \|e_{t,k}\|_2^2 + \frac{8 \min(1, C_{\text{mix}} \rho^\tau)^2 Y_0^2}{\mu} \right).$$

1146

1147 Combining,

1148

$$1149 2\alpha_k \mathbb{E}[e_{t,k}^\top \xi_{t,k+1}^\theta \mid \mathcal{G}_{t,k-\tau}] \leq \left( \frac{\mu}{2} + 4 \min(1, C_{\text{mix}} \rho^\tau) \right) \alpha_k \mathbb{E}[\|e_{t,k}\|_2^2 \mid \mathcal{G}_{t,k-\tau}] + \frac{8 Y_0^2}{\mu} \alpha_k \min(1, C_{\text{mix}} \rho^\tau)^2.$$

1150

□

1151

## 1152 (v) REMAINING SECOND ORDER TERMS

1153

1154 Now from the fact that  $\|\phi(s, a)\|_2^2 \leq 1$ , and from the conditions described in Lemma 2, we can write

1155

1156

$$1157 \|H_{t,k}^\theta\|^2 \leq x_{t,k},$$

1158

$$\mathbb{E}[\|m_{t,k+1}^\theta\|^2 \mid \mathcal{G}_{t,k-\tau}] \leq 4\gamma^2 C_{F,M}^2.$$

1159

1160

$$\mathbb{E}[\|\xi_{t,k+1}^\theta\|^2 \mid \mathcal{G}_{t,k-\tau}] \leq 8\mathbb{E}[x_{t,k} \mid \mathcal{G}_{t,k-\tau}] + 8Y_0^2$$

1161

1162 Therefore,

1163

$$1164 \alpha_k^2 \mathbb{E}[\|H_{t,k}^\theta + b_{t,k}^\theta + \xi_{t,k+1}^\theta + m_{t,k+1}^\theta\|^2 \mid \mathcal{G}_{t,k-\tau}]$$

1165

$$\leq \alpha_k^2 \left( 68 \mathbb{E}[x_{t,k} \mid \mathcal{G}_{t,k-\tau}] + 32\gamma^2 C_{F,M}^2 + 32Y_0^2 + 2\mathbb{E}[\|b_{t,k}^\theta\|^2 \mid \mathcal{G}_{t,k-\tau}] \right). \quad (33)$$

1166

1167 **Lemma 6** (Bias second order at  $1/k$ ). *For  $k \geq \max(\tau, 2\tau_{\beta_k})$ ,*

1168

1169

$$1170 \mathbb{E}[\|b_{t,k}^\theta\|_2^2 \mid \mathcal{G}_{t,k-\tau}] \leq \frac{2C_{\text{bias}}}{k} + 2(\epsilon_{\text{approx}}^{\text{dual}})^2,$$

1171

1172 with the explicit constant

1173

$$1174 C_{\text{bias}} := C_{\text{fast}}^2 + 2C_{F,M}^2 + 64C_{F,M}^2 \frac{C_{\text{mix}} \rho}{1 - \rho},$$

1175

1176 and  $C_{\text{fast}}$  as in Lemma 4.

1177 *Proof.* Using  $\|\phi(s, a)\|_2 \leq 1$  for all  $(s, a)$ , we can write using Equation (28),

1178

1179

$$1180 \|b_{t,k}^\theta\|_2 \leq \gamma \sum_{s,a} d^\pi(s, a) \left| F_{s,a}^{V_{\theta_t}^{\text{rob}}} (\bar{\nu}_{t,k})_{s,a} - F_{s,a}^{*,V_{\theta_t}^{\text{rob}}} \right|$$

1181

$$1182 = \gamma \sum_{s,a} d^\pi(s, a) \left( F_{s,a}^{*,V_{\theta_t}^{\text{rob}}} - F_{s,a}^{V_{\theta_t}^{\text{rob}}} (\bar{\nu}_{t,k})_{s,a} \right)$$

1183

$$\leq \gamma \inf_{\nu \in \mathcal{M}_\nu} \underbrace{\sum_{s,a} d^\pi(s, a) \left( F_{s,a}^{*,V_{\theta_t}^{\text{rob}}} - F_{s,a}^{V_{\theta_t}^{\text{rob}}} (\nu)_{s,a} \right)}_{\leq \epsilon_{\text{approx}}^{\text{dual}}} + \gamma \underbrace{\overline{L}_{t,k}}_{\text{fast-scale objective gap}}.$$

1184

1185

1186

1187

1188

$$\mathbb{E}[\|b_{t,k}^\theta\|_2^2 \mid \mathcal{G}_{t,k-\tau}] \leq 2\mathbb{E}[(\bar{L}_{t,k})^2 \mid \mathcal{G}_{t,k-\tau}] + 2(\epsilon_{approx}^{dual})^2.$$

1190

1191 To bound the RHS at the  $1/k$  scale, use the suffix average  $\bar{L}_{t,k} \leq \frac{1}{m} \sum_{j=k-m}^{k-1} L_{t,j}$  with  $m = \lfloor k/2 \rfloor$ .  
1192 We also use from Section B that and  $0 \leq L_{t,j} \leq 2C_{F,M}$ . Write

1193

$$\mathbb{E}[(\bar{L}_{t,k})^2 \mid \mathcal{G}_{t,k-\tau}] \leq \frac{1}{m^2} \left\{ \underbrace{\text{Var} \left( \sum_{j=k-m}^{k-1} L_{t,j} \mid \mathcal{G}_{t,k-\tau} \right)}_{(I)} + \underbrace{\left( \mathbb{E} \sum_{j=k-m}^{k-1} L_{t,j} \mid \mathcal{G}_{t,k-\tau} \right)^2}_{(II)} \right\}.$$

1194

1195

1196

1197

1198 *Term (II) (mean square).* By Lemma 4, for  $k \geq \tau_{\beta_k}$ ,  $\frac{1}{m} \mathbb{E} \sum_{j=k-m}^{k-1} L_{t,j} \leq \mathbb{E}[\bar{L}_{t,k}] \leq C_{\text{fast}}/\sqrt{k}$ , so  
1199  $(II)/m^2 \leq C_{\text{fast}}^2/k$ .

1200

1201

1202

1203

$$\begin{aligned} \text{Var} \left( \sum_{j=k-m}^{k-1} L_{t,j} \mid \mathcal{G}_{t,k-\tau} \right) &= \sum_j \text{Var}(L_{t,j} \mid \mathcal{G}_{t,k-\tau}) + 2 \sum_{h=1}^{m-1} (m-h) \text{Cov}(L_{t,j}, L_{t,j+h} \mid \mathcal{G}_{t,k-\tau}) \\ &\leq m C_{F,M}^2 + 32C_{F,M}^2 m \sum_{h=1}^{\infty} C_{\text{mix}} \rho^h \leq m \left( C_{F,M}^2 + 32C_{F,M}^2 \frac{C_{\text{mix}} \rho}{1-\rho} \right). \end{aligned}$$

1204

1205

1206

1207

1208

1209

1210

1211

1212

1213

1214

Dividing by  $m^2$  and using  $m \geq k/2$  gives

$$\frac{(I)}{m^2} \leq \frac{2}{k} \left( C_{F,M}^2 + 32C_{F,M}^2 \frac{C_{\text{mix}} \rho}{1-\rho} \right).$$

Combining the two terms yields

$$\mathbb{E}[(\bar{L}_{t,k})^2 \mid \mathcal{G}_{t,k-\tau}] \leq \frac{C_{\text{fast}}^2}{k} + \frac{2}{k} \left( C_{F,M}^2 + 32C_{F,M}^2 \frac{C_{\text{mix}} \rho}{1-\rho} \right),$$

1215

1216

which is the claimed bound with the displayed  $C_{\text{bias}}$ .  $\square$

1217

1218

1219

## FINAL RECURSION FOR THE SLOW TIME-SCALE

1220

1221

1222

From Equation (30) and the statements of the Lemma 5, Lemma 6 and from the analysis above, we can write the following recursion for  $\|e_{t,k}\|_2$ .

1223

1224

1225

We first characterize a suitable choice of the lag  $\tau$  to use in the recursion for  $\|e_{t,k}\|_2$ . We make a few observations. If  $\tau \geq \lceil \frac{\log(C_{\text{mix}}(k+1)/(c))}{\log(1/\rho)} \rceil$  then  $C_{\text{mix}} \rho^\tau \leq \alpha_k$ .

1226

1227

Also, if  $\tau \geq \tau_\mu := \lceil \frac{\log(C_{\text{mix}} \frac{1}{\mu})}{\log(1/\rho)} \rceil$ , then  $C_{\text{mix}} \rho^\tau \leq \mu$ .

1228

1229

1230

1231

1232

From the above two observations, a suitable choice for  $\tau$  in each inner iteration  $k$  is  $\tau_{\alpha_k} := \max\{\tau_{\alpha_k}, \tau_\mu\}$ . As we are conditioning on filtration  $\mathcal{G}_{t,k-\tau}$ , we need  $k \geq \tau_k$ . Moreover, to use Lemma 4, we need  $k \geq 2\tau_{\beta_k} = \lceil \frac{\log(\frac{C_{\text{mix}}}{\beta_0} \sqrt{k+1})}{\log(\frac{1}{\rho})} \rceil$ .

1233

Hence, we use the following definitions:

1234

1235

1236

1237

1238

1239

1240

1241

$$\tau_{\alpha_k} := \lceil \frac{\log(C_{\text{mix}}(k+1)/c)}{\log(1/\rho)} \rceil, \quad (34)$$

$$\tau_{\beta_k} := \lceil \frac{\log(\frac{C_{\text{mix}}}{\beta_0} \sqrt{k+1})}{\log(1/\rho)} \rceil, \quad (35)$$

$$\tau_\mu := \lceil \frac{\log(C_{\text{mix}}/\mu)}{\log(1/\rho)} \rceil \quad (36)$$

1242  $\forall k \geq \max(\tau_{\alpha_k}, 2\tau_{\beta_k}, \tau_\mu)$ :

$$\begin{aligned} \mathbb{E}[\|e_{t,k+1}\|_2^2 \mid \mathcal{G}_{t,k-\tau}] &\leq (1 - \mu \alpha_k + C_e \alpha_k^2) \mathbb{E}[\|e_{t,k}\|_2^2 \mathcal{G}_{t,k-\tau}] \\ &+ \left(\frac{2\alpha_k}{\mu} + 2\alpha_k^2\right) \left(\frac{2C_{bias}}{k} + 2(\epsilon_{approx}^{dual})^2\right) + C_{cross} \alpha_k^2 + C_0 \alpha_k^2, \end{aligned}$$

1248

1249 and hence, after taking total expectation,  $\forall k \geq \max(\tau, 2\tau_{\beta_k}, \tau_\mu)$ 

$$\begin{aligned} \mathbb{E}\|e_{t,k+1}\|_2^2 &\leq (1 - \mu \alpha_k + C_e \alpha_k^2) \mathbb{E}\|e_{t,k}\|_2^2 + \\ &+ \left(\frac{2\alpha_k}{\mu} + 2\alpha_k^2\right) \left(\frac{2C_{bias}}{k} + 2(\epsilon_{approx}^{dual})^2\right) + C_{cross} \alpha_k^2 + C_0 \alpha_k^2. \end{aligned} \quad (37)$$

1255

1256 The constants in the above recursion are given as:

1257  $C_e := 72,$

1259  $C_0 := 32\gamma^2 C_{F,M}^2 + 32 \left(1 + \gamma C_{F,M} + \frac{1}{(1-\gamma)\sqrt{\mu}}\right)^2,$

1262  $C_{cross} := 8 \left(1 + \gamma C_{F,M} + \frac{1}{(1-\gamma)\sqrt{\mu}}\right)^2,$

1264  $C_{bias} := C_{fast}^2 + 2C_{F,M}^2 + 64C_{F,M}^2 \frac{C_{mix}\rho}{1-\rho},$

1266  $C_{fast} := \frac{(4B_\nu^2 + \beta_0^2(8B_\nu C_G + 25C_G^2) \ln(2e))}{\beta_0}.$

1268

1269 Now we derive the last iterate convergence of the error  $\|e_{t,k+1}^\theta\|_2$ .1270 We also know from the definition of  $\varepsilon_{approx}^{dual}$ ,

1273  $|\varepsilon_{approx}^{dual}| \leq \frac{1}{1-\gamma} + C_{F,M}.$

1275

1276 In compact notations,  $\forall k \geq \max(\tau_{\alpha_k}, 2\tau_{\beta_k}, \tau_\mu)$ 

1278

$$\begin{aligned} \mathbb{E}\|e_{t,k+1}\|_2^2 &\leq (1 - \mu \alpha_k + C_e \alpha_k^2) \mathbb{E}\|e_{t,k}\|_2^2 + \\ &\alpha_k^2 \left( \left(\frac{4}{c\mu} + 2\right) 2C_{bias} + C_{cross} + C_0 + 4 \left(\frac{1}{1-\gamma} + C_{F,M}\right)^2 \right) \\ &+ \frac{4\alpha_k(\epsilon_{approx}^{dual})^2}{\mu} \end{aligned}$$

1286

1287 Or,

1289  $\mathbb{E}\|e_{t,k+1}\|_2^2 \leq (1 - \mu \alpha_k + C_e \alpha_k^2) \mathbb{E}\|e_{t,k}\|_2^2 + \alpha_k^2 C_1 + \frac{4\alpha_k(\epsilon_{approx}^{dual})^2}{\mu} \quad (38)$

1292 with  $C_1 = \left( \left(\frac{8}{c\mu} + 4\right) C_{bias} + C_{cross} + C_0 + 4 \left(\frac{1}{1-\gamma} + C_{F,M}\right)^2 \right).$

1293 Next, we derive the final iterate convergence from the above recursion for the step-size rule  $\alpha_k = \frac{c}{1+k}$ ,  $\forall k \in [0, K-1]$ :

1296 C.3 LAST-ITERATE CONVERGENCE FOR  $\alpha_k = \frac{c}{k+1}$   
1297  
1298 **Recursion on  $\mathbb{E}[\|e_{t,k}\|_2^2]$ :** From Equation (38), we have that, whenever  
1299  
1300 
$$k \geq k_{\text{mix}} := \max(\tau_{\alpha_k}, 2\tau_{\beta_k}, \tau_\mu),$$
  
1301  
1302  
1303 
$$\mathbb{E}[\|e_{t,k+1}\|_2^2] \leq \left(1 - \mu\alpha_k + C_e \alpha_k^2\right) \mathbb{E}[\|e_{t,k}\|_2^2] + \alpha_k^2 C_1 + \frac{4}{\mu} \alpha_k (\varepsilon_{\text{approx}}^{\text{dual}})^2, \quad \alpha_k = \frac{c}{k+1}. \quad (39)$$
  
1304  
1305

1306 Let

$$o_k := 1 - \mu\alpha_k + C_e \alpha_k^2, \quad k_0 := \max\left\{\left\lceil \frac{2C_e c}{\mu} \right\rceil, \left\lceil \mu c \right\rceil\right\},$$

1307 Then for all  $j \geq k_0$ ,  $0 < o_j \leq 1 - \frac{\mu c}{2}/(j+1) \leq 1$ . When  $k_{\text{mix}} < k_0$ , define the finite pre-burn  
1308 product before the contraction kicks in as follows  
1309  
1310

$$H_{\text{pre}} := \prod_{j=k_{\text{mix}}}^{k_0-1} o_j, \quad \text{and set } H_{\text{pre}} = 1 \text{ if } k_{\text{mix}} \geq k_0.$$

1311 Assume the final iterate index  $K > k_{\text{mix}}$   
1312  
1313

1314 **Tail-product bounds.** For  $m \leq u$ , set  $G_m^u := \prod_{j=m}^u o_j$ . Then  
1315  
1316

$$1317 G_{k_0}^{K-1} \leq \prod_{j=k_0}^{K-1} \left(1 - \frac{\mu c}{j+1}\right) \leq \left(\frac{k_0+1}{K}\right)^{\frac{\mu c}{2}}, \quad \text{and if } k_{\text{mix}} \geq k_0, \quad G_{k_{\text{mix}}}^{K-1} \leq \left(\frac{k_{\text{mix}}+1}{K}\right)^{\frac{\mu c}{2}}. \quad (40)$$

1318  
1319 **Unrolling from  $k_{\text{mix}}$ .** Fix  $K > k_{\text{mix}}$ . Unrolling equation 39 from  $t = k_{\text{mix}}$  to  $K-1$  yields  
1320  
1321

$$1322 \mathbb{E}[\|e_{t,K}\|_2^2] \leq \underbrace{\mathbb{E}[\|e_{t,k_{\text{mix}}}\|_2^2] G_{k_{\text{mix}}}^{K-1}}_{\text{initial term}} + \sum_{t=k_{\text{mix}}}^{K-1} \alpha_t^2 C_1 G_{t+1}^{K-1} + \frac{4}{\mu} (\varepsilon_{\text{approx}}^{\text{dual}})^2 \sum_{t=k_{\text{mix}}}^{K-1} \alpha_t G_{t+1}^{K-1}. \quad (41)$$

1323  
1324 **Splitting at  $k_0$ .** Split each sum at  $k_0$ :

$$1325 \sum_{t=k_{\text{mix}}}^{K-1} (\cdot) = \underbrace{\sum_{t=k_{\text{mix}}}^{k_0-1} (\cdot)}_{\text{finite "pre-window"}} + \underbrace{\sum_{t=k_0}^{K-1} (\cdot)}_{\text{tail}}.$$

1326  
1327 **Initial term.** If  $k_{\text{mix}} \geq k_0$ , then by equation 40  
1328

$$1329 \mathbb{E}[\|e_{t,k_{\text{mix}}}\|_2^2] G_{k_{\text{mix}}}^{K-1} \leq \mathbb{E}[\|e_{t,k_{\text{mix}}}\|_2^2] \left(\frac{k_{\text{mix}}+1}{K}\right)^{\frac{\mu c}{2}}.$$

1330 If  $k_{\text{mix}} < k_0$ , then  
1331

$$1332 \mathbb{E}[\|e_{t,k_{\text{mix}}}\|_2^2] G_{k_{\text{mix}}}^{K-1} \leq \mathbb{E}[\|e_{t,k_{\text{mix}}}\|_2^2] H_{\text{pre}} \left(\frac{k_0+1}{K}\right)^{\frac{\mu c}{2}}.$$

1333  
1334 **Variance sum.** Define the finite pre-window constant  
1335

$$1336 U_{\text{pre}}^{(\text{mix})}(k_{\text{mix}}, k_0) := \sum_{t=k_{\text{mix}}}^{k_0-1} \frac{c^2 C_1}{(t+1)^2} \prod_{j=t+1}^{k_0-1} o_j \quad (\text{define it as 0 if } k_{\text{mix}} \geq k_0).$$

1350

Then

$$1352 \quad \sum_{t=k_{\text{mix}}}^{K-1} \frac{c^2 C_1}{(t+1)^2} G_{t+1}^{K-1} \leq U_{\text{pre}}^{(\text{mix})}(k_{\text{mix}}, k_0) \left( \frac{k_0+1}{K} \right)^{\frac{\mu c}{2}} + \sum_{t=k_0}^{K-1} \frac{c^2 C_1}{(t+1)^2} \left( \frac{t+2}{K} \right)^{\frac{\mu c}{2}}.$$

1354

Using integral approximation, we can prove that the above term has the following upper bound:

$$1356 \quad \sum_{t=k_0}^{K-1} \frac{c^2 C_1}{(t+1)^2} \left( \frac{t+2}{K} \right)^{\frac{\mu c}{2}} \leq \begin{cases} \frac{c^2 C_1}{\frac{\mu c}{2} - 1} \frac{1}{K} + \frac{c^2 C_1}{\frac{\mu c}{2} - 1} \frac{(k_0+2)^{\frac{\mu c}{2}-1}}{K^{\frac{\mu c}{2}}}, & \frac{\mu c}{2} > 1, \\ \frac{c^2 C_1}{K} \left( 1 + \ln \frac{K}{k_0+2} \right), & \frac{\mu c}{2} = 1, \\ \frac{c^2 C_1}{1 - \frac{\mu c}{2}} \frac{(k_0+2)^{\frac{\mu c}{2}-1}}{K^{\frac{\mu c}{2}}}, & 0 < \frac{\mu c}{2} < 1. \end{cases}$$

1363

1364 **Dual-approximation sum.** Define the finite pre-window constant

$$1366 \quad B_{\text{pre}}^{(\text{mix})}(k_{\text{mix}}, k_0) := \sum_{t=k_{\text{mix}}}^{k_0-1} \alpha_t \prod_{j=t+1}^{k_0-1} o_j \quad (\text{again 0 if } k_{\text{mix}} \geq k_0).$$

1368

Then

$$1370 \quad \sum_{t=k_{\text{mix}}}^{K-1} \alpha_t G_{t+1}^{K-1} \leq B_{\text{pre}}^{(\text{mix})}(k_{\text{mix}}, k_0) \left( \frac{k_0+1}{K} \right)^{\frac{\mu c}{2}} + \sum_{t=k_0}^{K-1} \frac{c}{t+1} \left( \frac{t+2}{K} \right)^{\frac{\mu c}{2}}.$$

1373

For the tail, an integral comparison yields

$$1374 \quad \sum_{t=k_0}^{K-1} \frac{c}{t+1} \left( \frac{t+2}{K} \right)^{\frac{\mu c}{2}} \leq \frac{c}{K^{\frac{\mu c}{2}}} \int_{k_0}^K (x+2)^{\frac{\mu c}{2}-1} dx = \frac{c}{\frac{\mu c}{2} K^{\mu c/2}} \left[ (K+2)^{\mu c/2} - (k_0+2)^{\mu c/2} \right].$$

1377

We also know that,

$$1379 \quad |\varepsilon_{\text{approx}}^{\text{dual}}| \leq \frac{1}{1-\gamma} + C_{F,M}.$$

1380

We can also write,

$$1382 \quad \frac{c}{\frac{\mu c}{2} K^{\mu c/2}} \left[ (K+2)^{\mu c/2} - (k_0+2)^{\mu c/2} \right] \\ 1383 \quad \leq \frac{2}{\mu} \left( 1 + \frac{2}{K} \right)^{\frac{\mu c}{2}}$$

1387

Therefore,

$$1389 \quad \frac{4}{\mu} (\varepsilon_{\text{approx}}^{\text{dual}})^2 \sum_{t=k_{\text{mix}}}^{K-1} \alpha_t G_{t+1}^{K-1} \\ 1390 \quad \leq \frac{4}{\mu} \left( \frac{1}{1-\gamma} + C_{F,M} \right)^2 B_{\text{pre}}^{(\text{mix})}(k_{\text{mix}}, k_0) \left( \frac{k_0+1}{K} \right)^{\frac{\mu c}{2}} + \frac{8}{\mu^2} \left( 1 + \frac{2}{K} \right)^{\frac{\mu c}{2}} (\varepsilon_{\text{approx}}^{\text{dual}})^2.$$

1394

**Bounding the initial iterate at  $k_{\text{mix}}$ .** Recall that the slow update of Algorithm 1 satisfies

$$1396 \quad \theta_{t,k+1} = \theta_{t,k} + \alpha_k TD_{k+1} \phi(Z_{t,k}), \quad \|TD_{k+1} \phi(Z_{t,k})\|_2 \leq 1 + C_{F,M} + \|\theta_{t,k}\|_2,$$

1397

with  $\theta_{t,0} = \mathbf{0}$ . Define  $u_k := \|\theta_{t,k}\|_2$ . Then

$$1399 \quad u_{k+1} \leq u_k + \alpha_k (1 + C_{F,M} + u_k) = (1 + \alpha_k) u_k + \alpha_k (1 + C_{F,M}).$$

1400

Introduce the shifted sequence  $v_k := u_k + (1 + C_{F,M})$ . We then have

$$1402 \quad v_{k+1} = u_{k+1} + (1 + C_{F,M}) \leq (1 + \alpha_k) u_k + \alpha_k (1 + C_{F,M}) + (1 + C_{F,M}) \\ 1403 \quad = (1 + \alpha_k) (u_k + 1 + C_{F,M}) = (1 + \alpha_k) v_k.$$

1404 Iterating from  $k = 0$  to  $k_{\text{mix}} - 1$  and using  $\theta_{t,0} = \mathbf{0}$  (so  $u_0 = 0$  and  $v_0 = 1 + C_{F,M}$ ) yields  
1405

$$1406 \quad 1407 \quad v_{k_{\text{mix}}} \leq (1 + C_{F,M}) \prod_{i=0}^{k_{\text{mix}}-1} (1 + \alpha_i) = (1 + C_{F,M}) \prod_{i=0}^{k_{\text{mix}}-1} \left(1 + \frac{c}{i+1}\right).$$

1409 Using  $1 + x \leq e^x$  and the harmonic-sum bound  $\sum_{i=0}^{k_{\text{mix}}-1} \frac{1}{i+1} \leq 1 + \ln k_{\text{mix}}$ , we obtain  
1410

$$1411 \quad 1412 \quad \prod_{i=0}^{k_{\text{mix}}-1} \left(1 + \frac{c}{i+1}\right) \leq \exp\left(\sum_{i=0}^{k_{\text{mix}}-1} \frac{c}{i+1}\right) \leq \exp(c(1 + \ln k_{\text{mix}})) = e^c k_{\text{mix}}^c.$$

1414 Therefore  
1415

$$1416 \quad u_{k_{\text{mix}}} = \|\theta_{t,k_{\text{mix}}}\|_2 \leq v_{k_{\text{mix}}} \leq (1 + C_{F,M}) e^c k_{\text{mix}}^c.$$

1417 Finally, recall that  $e_{t,k} := \theta_{t,k} - \theta^{*,t}$ , where  $\theta^{*,t}$  is the (time- $t$ ) fixed point of the mean ODE. Using  
1418  $\|e_{t,k_{\text{mix}}}\|_2^2 \leq 2\|\theta_{t,k_{\text{mix}}}\|_2^2 + 2\|\theta^{*,t}\|_2^2$  and the standard bound  $\|\theta^{*,t}\|_2 \leq (1 - \gamma)^{-1}/\sqrt{\mu}$ , we obtain  
1419

$$1420 \quad 1421 \quad \mathbb{E}[\|e_{t,k_{\text{mix}}}\|_2^2] \leq 2(1 + C_{F,M})^2 e^{2c} k_{\text{mix}}^{2c} + \frac{2}{\mu(1 - \gamma)^2} \\ 1422 \quad 1423 \quad =: C_{\text{init}}^{(\text{mix},1)}.$$

1424 This constant will be used as the initial-error term in the last-iterate bound for the  $\alpha_k = c/(k+1)$   
1425 schedule.  
1426

1427 **Final bound.** Combining the pieces, for any  $K > k_{\text{mix}}$ ,  
1428

$$1429 \quad \mathbb{E}[\|e_{t,K}\|_2^2] \\ 1430 \quad \leq \begin{cases} C_{\text{init}}^{(\text{mix},1)} \left(\frac{k_{\text{mix}}+1}{K}\right)^{\frac{\mu c}{2}}, & \text{if } k_{\text{mix}} \geq k_0, \\ C_{\text{init}}^{(\text{mix},1)} H_{\text{pre}} \left(\frac{k_0+1}{K}\right)^{\frac{\mu c}{2}}, & \text{if } k_{\text{mix}} < k_0, \end{cases} \\ 1431 \quad + U_{\text{pre}}^{(\text{mix})}(k_{\text{mix}}, k_0) \left(\frac{k_0+1}{K}\right)^{\frac{\mu c}{2}} + \begin{cases} \frac{c^2 C_1}{\frac{\mu c}{2} - 1} \frac{1}{K} + \frac{c^2 C_1}{\frac{\mu c}{2} - 1} \frac{(k_0+2)^{\frac{\mu c}{2}-1}}{K^{\frac{\mu c}{2}}}, & \frac{\mu c}{2} > 1, \\ \frac{c^2 C_1}{K} \left(1 + \ln \frac{K}{k_0+2}\right), & \frac{\mu c}{2} = 1, \\ \frac{c^2 C_1}{1 - \frac{\mu c}{2}} \frac{(k_0+2)^{\frac{\mu c}{2}-1}}{K^{\frac{\mu c}{2}}}, & 0 < \frac{\mu c}{2} < 1, \end{cases} \\ 1432 \quad 1433 \quad 1434 \quad 1435 \quad 1436 \quad 1437 \quad 1438 \quad 1439 \quad 1440 \quad 1441 \quad 1442 \quad 1443 \quad 1444 \quad 1445 \quad 1446 \quad 1447 \quad 1448 \quad 1449 \quad 1450 \quad 1451 \quad 1452 \quad 1453 \quad 1454 \quad 1455 \quad 1456 \quad 1457 \quad \frac{4}{\mu} \left(\frac{1}{1 - \gamma} + C_{F,M}\right)^2 B_{\text{pre}}^{(\text{mix})}(k_{\text{mix}}, k_0) \left(\frac{k_0+1}{K}\right)^{\frac{\mu c}{2}} + \frac{8}{\mu^2} \left(1 + \frac{2}{K}\right)^{\frac{\mu c}{2}} (\varepsilon_{\text{approx}}^{\text{dual}})^2.$$

#### C.4 FORMAL VERSION OF THEOREM 1

1448 **Theorem 2** (Main finite-time guarantee under function approximation). *Let  $\hat{Q}_t := \Phi\hat{\theta}_t$  be the  
1449 estimate of  $Q^{\text{rob},\pi}$  returned by Algorithm 1 at outer iteration  $t$ .*  
1450

$$1451 \quad k_{\text{mix}} \\ 1452 \quad := \min \left\{ m \in \mathbb{N} : \forall j \geq m, j \geq \max \left( \tau_\mu, 2 \left\lceil \frac{\log(\frac{C_{\text{mix}}}{\beta_0} \sqrt{j+1})}{\log(1/\rho)} \right\rceil, \left\lceil \frac{\log(C_{\text{mix}}(j+1)/(c))}{\log(1/\rho)} \right\rceil \right) \right\},$$

1455 where,  
1456

$$1457 \quad \tau_\mu := \left\lceil \frac{\log(C_{\text{mix}} \frac{1}{\mu})}{\log(1/\rho)} \right\rceil$$

1458 Assume 1 holds, and we run  $K \geq k_{\text{mix}}$  inner iterations per outer iteration for either the TV or the  
 1459 Wasserstein- $\ell$  uncertainty sets.. Then, for any horizon  $T \geq 1$ ,  
 1460

$$\begin{aligned} 1461 \mathbb{E}[\|\hat{Q}_T - Q^{\text{rob},\pi}\|_\infty] \\ 1462 &\leq \gamma^T \|\Phi\theta_0 - Q^{\text{rob},\pi}\|_\infty + \frac{A_{\text{sched}}(K)}{1-\gamma} + \frac{\epsilon_{\text{approx}}}{1-\gamma} + \frac{2\sqrt{2}(1+\frac{2}{K})^{\frac{\mu c}{4}}\epsilon_{\text{approx}}^{\text{dual}}}{\mu(1-\gamma)}. \end{aligned} \quad (42)$$

1465 Here  $A_{\text{sched}}(K) \geq 0$  is the schedule-dependent residual, which takes one of the following explicit  
 1466 forms depending on the range of  $c$ .  
 1467

1468 Recall the definition  $k_0 = \max\left\{\left\lceil\frac{144c}{\mu}\right\rceil, \left\lceil\mu c\right\rceil\right\}$  and  
 1469

1470 Set

$$\begin{aligned} 1471 o_j &:= \left(1 - \frac{c\mu}{j+1}\right) + \frac{72c^2}{(j+1)^2}, \quad H_{\text{pre}} := \prod_{j=k_{\text{mix}}}^{k_0-1} o_j, \\ 1472 U_{\text{pre}}^{(\text{mix})}(k_{\text{mix}}, k_0) &:= \sum_{t=k_{\text{mix}}}^{k_0-1} \frac{c^2 C_1}{(t+1)^2} \prod_{j=t+1}^{k_0-1} o_j, \quad B_{\text{pre}}^{(\text{mix})}(k_{\text{mix}}, k_0) := \sum_{t=k_{\text{mix}}}^{k_0-1} \alpha_t \prod_{j=t+1}^{k_0-1} o_j, \\ 1473 \end{aligned}$$

1474 with the convention that  $H_{\text{pre}} = 1$ ,  $U_{\text{pre}}^{(\text{mix})} = 0$ ,  $B_{\text{pre}}^{(\text{mix})} = 0$  when  $k_{\text{mix}} \geq k_0$ . Also set  
 1475

$$C_{\text{init}}^{(\text{mix},1)} := 2(1+C_{F,M})^2 e^{2c} k_{\text{mix}}^{\frac{\mu c}{2}} + \frac{2}{\mu(1-\gamma)^2}.$$

1476 Then  $A_{\text{sched}}(K) = \sqrt{\Xi_1(K)}$  with  
 1477

$$\begin{aligned} 1478 \Xi_1(K) &:= \begin{cases} C_{\text{init}}^{(\text{mix},1)} \left(\frac{k_{\text{mix}}+1}{K}\right)^{\frac{\mu c}{2}}, & k_{\text{mix}} \geq k_0, \\ C_{\text{init}}^{(\text{mix},1)} H_{\text{pre}} \left(\frac{k_0+1}{K}\right)^{\frac{\mu c}{2}}, & k_{\text{mix}} < k_0, \\ + U_{\text{pre}}^{(\text{mix})}(k_{\text{mix}}, k_0) \left(\frac{k_0+1}{K}\right)^{\frac{\mu c}{2}} \end{cases} \\ 1479 &+ \begin{cases} \frac{c^2}{\frac{\mu c}{2}-1} \frac{1}{K} C_1 + \frac{c^2}{\frac{\mu c}{2}-1} \frac{(k_0+2)^{\frac{\mu c}{2}-1}}{K^{\frac{\mu c}{2}}} C_1, & \frac{\mu c}{2} > 1, \\ \frac{c^2}{K} \left(1 + \ln \frac{K}{k_0+2}\right) C_1, & \frac{\mu c}{2} = 1, \\ \frac{c^2}{1-\frac{\mu c}{2}} \frac{(k_0+2)^{\frac{\mu c}{2}-1}}{K^{\frac{\mu c}{2}}} C_1, & 0 < \frac{\mu c}{2} < 1, \\ + \frac{4}{\mu} \left(\frac{1}{1-\gamma} + C_{F,M}\right)^2 B_{\text{pre}}^{(\text{mix})}(k_{\text{mix}}, k_0) \left(\frac{k_0+1}{K}\right)^{\frac{\mu c}{2}}. \end{cases} \end{aligned}$$

## 1499 Restating the Constants

$$\begin{aligned} 1500 C_0 &= 32\gamma^2 C_{F,M}^2 + 32 \left(1 + \gamma C_{F,M} + \frac{1}{(1-\gamma)\sqrt{\mu}}\right)^2, \\ 1501 C_{\text{cross}} &= 8 \left(1 + \gamma C_{F,M} + \frac{1}{(1-\gamma)\sqrt{\mu}}\right)^2, \\ 1502 C_{\text{fast}} &= \frac{4B_\nu^2 + \beta_0^2(8B_\nu C_G + 25C_G^2) \ln(2e)}{\beta_0}, \\ 1503 C_{\text{bias}} &= C_{\text{fast}}^2 + 2C_{F,M}^2 + 64C_{F,M}^2 \frac{C_{\text{mix}}\rho}{1-\rho}, \\ 1504 C_1 &= \left( \left(\frac{8}{c\mu} + 4\right) C_{\text{bias}} + C_{\text{cross}} + C_0 + 4 \left(\frac{1}{1-\gamma} + C_{F,M}\right)^2 \right). \end{aligned}$$

1512 **D REMAINING PROOFS**  
 1513

1514 **D.1 BOUND ON  $\theta_t^*$**   
 1515

1516 We drop the superscript  $t$  from  $\theta_t^*$  as  $t$  is fixed through the discussion of this subsection. Recall,

1517 
$$1518 \theta_t^* = (\Phi^\top D^\pi \Phi)^{-1} \Phi^\top D^\pi [r + \gamma \bar{f}_t(\lambda^*)]$$

1519 **Lemma 7** (Bound on the optimal weight vector). *Let*

1520 
$$1521 \theta^* = (\Phi^\top D^\pi \Phi)^{-1} \Phi^\top D^\pi [r + \gamma \bar{f}_t(\lambda^*)],$$

1522 where

1523 

- 1524  $\Phi \in \mathbb{R}^{|\mathcal{S}||\mathcal{A}| \times d}$  has full column rank and row vectors  $\phi(s, a)$  satisfying  $\|\phi(s, a)\|_2 \leq 1$ ;
- 1525  $\Phi^\top D^\pi \Phi$  is positive definite with min-eigenvalue  $\nu$ . The diagonal entries of the matrix  $D^\pi$  satisfy  $d_i \geq d_{\min} > 0$  and  $\sum_i d_i = 1$ ;
- 1526  $\cdot$  each entry of  $r$  obeys  $|r_i| \leq 1$ ;
- 1527  $\cdot$  each entry of  $\bar{f}_t(\lambda^*)$  obeys  $|\bar{f}_t(\lambda^*)_i| \leq 1/(1 - \gamma)$ .

1528 Then

1529 
$$1530 \|\theta^*\|_2 \leq \frac{1}{1 - \gamma} \frac{1}{\sqrt{\mu}}.$$

1531 *Proof.* Set

1532 
$$1533 C := \Phi^\top D^\pi \Phi, \quad v := r + \gamma \bar{f}_t(\lambda^*).$$

1534 By definition of  $\theta^*$ ,

1535 
$$C\theta^* = \Phi^\top D^\pi v.$$

1536 Multiply by  $\theta^{*\top}$  on the left:

1537 
$$\theta^{*\top} C\theta^* = \theta^{*\top} \Phi^\top D^\pi v = (\Phi\theta^*)^\top D^\pi v.$$

1538 Let  $y := \Phi\theta^*$ . Then

1539 
$$\theta^{*\top} C\theta^* = y^\top D^\pi v.$$

1540 On the other hand,

1541 
$$\theta^{*\top} C\theta^* = \theta^{*\top} \Phi^\top D^\pi \Phi\theta^* = (\Phi\theta^*)^\top D^\pi (\Phi\theta^*) = y^\top D^\pi y.$$

1542 Thus

1543 
$$y^\top D^\pi y = y^\top D^\pi v.$$

1544 Apply Cauchy–Schwarz in the  $D^\pi$ –weighted inner product:

1545 
$$1546 y^\top D^\pi v = (D^{\pi 1/2} y)^\top (D^{\pi 1/2} v) \leq \|D^{\pi 1/2} y\|_2 \|D^{\pi 1/2} v\|_2 = (y^\top D^\pi y)^{1/2} (v^\top D^\pi v)^{1/2}.$$

1547 If  $y^\top D^\pi y = 0$ , then  $\theta^* = 0$  and the desired bound is trivial, so assume  $y^\top D^\pi y > 0$  and divide both sides by  $(y^\top D^\pi y)^{1/2}$ :

1548 
$$1549 (y^\top D^\pi y)^{1/2} \leq (v^\top D^\pi v)^{1/2} \implies y^\top D^\pi y \leq v^\top D^\pi v.$$

1550 Recalling  $y^\top D^\pi y = \theta^{*\top} C\theta^*$ , we obtain

1551 
$$\theta^{*\top} C\theta^* \leq v^\top D^\pi v.$$

1552 **Upper bound on  $v^\top D^\pi v$ .** For each component,

1553 
$$1554 |v_i| \leq |r_i| + \gamma |\bar{f}_t(\lambda^*)_i| \leq 1 + \frac{\gamma}{1 - \gamma} = \frac{1}{1 - \gamma}.$$

1555 Hence

1556 
$$1557 v^\top D^\pi v = \sum_i d_i v_i^2 \leq \sum_i d_i \left( \frac{1}{1 - \gamma} \right)^2 = \left( \frac{1}{1 - \gamma} \right)^2.$$

1566 **Lower bound via the minimum eigenvalue.** Since  $C = \Phi^\top D^\pi \Phi \succeq \mu I$ ,

$$1568 \quad \theta^{*\top} C \theta^* \geq \mu \|\theta^*\|_2^2.$$

1570 Combining the upper and lower bounds,

$$1572 \quad \mu \|\theta^*\|_2^2 \leq \theta^{*\top} C \theta^* \leq \left( \frac{1}{1-\gamma} \right)^2,$$

1574 so

$$1575 \quad \|\theta^*\|_2 \leq \frac{1}{1-\gamma} \frac{1}{\sqrt{\mu}}.$$

1577  $\square$

## 1579 E ROBUST Q-LEARNING

1581 In this section, we discuss a robust Q-learning algorithm with function approximation that finds the  
 1582 optimal policy for the worst-case transition kernel in the uncertainty set considered in this paper. We  
 1583 first define the optimal state-action value function  $Q^{\text{rob},*}$  as the state-action value function of the best  
 1584 admissible policy to maximize  $Q^{\text{rob},\pi}$  for each  $(s, a)$ -pair.

$$1585 \quad Q^{\text{rob},*}(s, a) = \max_{\pi} Q^{\text{rob},\pi}(s, a), \forall (s, a) \in \mathcal{S} \times \mathcal{A}.$$

1588 It is shown in prior literature (Iyengar, 2005) that  $Q^{\text{rob},*}$  satisfies the following equation, which is  
 1589 called the robust Bellman optimality equation

$$1591 \quad Q^{\text{rob},*}(s, a) = R(s, a) + \gamma \min_{q \in \mathcal{P}_s^a} \sum_{s'} q(s') \underbrace{\max_{a'} Q^{\text{rob},*}(s', a')}_{=: V^{\text{rob},*}(s')}.$$

1595 Equivalently, define the robust Bellman optimality operator  $(\mathcal{T}^{\text{rob},*}Q)(s, a) := R(s, a) +$   
 1596  $\gamma \sigma_{\mathcal{P}_s^a}(V^{\text{rob},*})$  with

$$1597 \quad V^{\text{rob},*}(s') := \max_{a'} Q(s', a'), \quad (44)$$

1598 and  $\sigma_{\mathcal{P}_s^a}(V)$  is given in Equation (4). Iyengar (2005) proved that the robust Bellman optimality  
 1599 operator is  $\gamma$ -contraction in  $\ell_\infty$  norm.

1601 Now, we discuss how the TD learning algorithm presented in Algorithm 1 in the main body of the  
 1602 paper can be extended to estimate  $Q^{\text{rob},*}$  in a relatively straightforward manner. Similar to the TD  
 1603 learning setup, assume that we can sample data corresponding to a behavioral policy  $\pi_b$  from the  
 1604 nominal model  $P_0$ . Also, assume that the policy  $\pi_b$  satisfies Assumption 1.

1605 The goal here is to approximate  $Q^{\text{rob},*}$  by  $\Phi \theta^*$  for an appropriately chosen  $\theta^*$ . Our Q-learning  
 1606 algorithm is presented in Algorithm 2. The algorithm computes an estimate  $\hat{\theta}_t$  of this parameter at  
 1607 each iteration  $t$  of the outer loop. The quantity  $V_{\hat{\theta}_t}^{\text{rob},*}$  in the description of the algorithm is given by

$$1609 \quad V_{\hat{\theta}_t}^{\text{rob},*}(s) = \max_a \text{Clip} \left( \phi(s, a)^\top \hat{\theta}_t \right), \forall s \in \mathcal{S}. \quad (45)$$

1611 **Difference between Algorithm 2 and Algorithm 1:** The only difference between the robust Q-  
 1612 learning algorithm in Algorithm 2 and the robust TD learning algorithm in Algorithm 1 is that, we  
 1613 use  $V_{\hat{\theta}_t}^{\text{rob},*}$  instead of  $V_{\hat{\theta}_t}^{\text{rob}}$  in the calculation of the dual super-gradient in line 6 and the calculation  
 1614 of the dual objective in line 10 in Algorithm 2.

1616 **Finite-Time Performance Bound for the Robust Q-Learning (Algorithm 2):** Recall that we  
 1617 established a finite-time performance bound for the robust TD learning in Theorem 1. By following  
 1618 the steps of the proof of that theorem, it is easy to see that an analogous guarantee holds for the  
 1619 estimate of  $Q^{\text{rob},*}$  produced by Algorithm 2. The reason that the proof is identical is that the robust  
 Bellman optimality operator is a  $\gamma$ -contraction in the  $\ell_\infty$  norm as was the robust Bellman operator

1620 for a fixed policy. The only difference is that the function approximation error for approximating  
 1621 the Q-function should now be defined as the error in approximating  $Q^{\text{rob},*}$  by the class of functions  
 1622  $\{\Phi\theta : \theta \in \mathbb{R}^{n_\theta}\}$  :

1623

$$1624 \epsilon_{\text{approx}}^* := \sup_{Q=\text{Clip}(\Phi\theta), \theta \in \mathbb{R}^{n_\theta}} \|\text{Clip}(\Pi \mathcal{T}^{\text{rob},*}(Q)) - \mathcal{T}^{\text{rob},*}(Q)\|_\infty. \quad (46)$$

1625

1626 The above definition is completely analogous to the TD learning setting in the main body of the paper,  
 1627 but with  $\mathcal{T}^{\text{rob},*}$  instead of  $\mathcal{T}^{\text{rob},\pi}$  for a policy  $\pi$ .

1628 Thus, the sample complexity of robust Q-learning is of the same order as that of robust TD-learning  
 1629 up to a function approximation error.

1630

---

1631 **Algorithm 2** Robust Q-learning with Function Approximation

---

1632 1: **Input:** Integers  $T, K$ . Initial  $\nu_0 \in \mathbb{R}^{n_\lambda}$ ,  $\theta_0 :=$  zero vector, fast time-scale step-sizes  $\beta_k =$   
 1633  $\frac{\beta_0}{\sqrt{k+1}}$ , slow time-scale step-sizes  $\alpha_k = \frac{c}{(k+1)}$  for some  $c : 0 < c < \infty$ ;  $\hat{\theta}_0 = \theta_0$ ,  $\theta_{0,0} = \theta_0$ ,  
 1634 behavioral policy  $\pi_b$ , Reward function  $R : (\mathcal{S} \times \mathcal{A}) \mapsto [-1, 1]$ , initial state  $S_{0,0}$ .

1635 2: **for**  $t = 0, 1, \dots, T-1$  **do**

1636 3:   **for**  $k = 0, 1, \dots, K-1$  **do**

1637 4:     Take action  $A_{t,k}$  according to the behavioral policy  $\pi_b$  and sample  $S_{t,k+1}$  ( $S_{t,k+1} \sim$   
 1638  $P_0(\cdot | S_{t,k}, A_{t,k})$ )

1639 5:     **fast time-scale** ( $\beta_k$ )

1640 6:     Compute  $\hat{G}(\psi(S_{t,k}, A_{t,k})^\top \nu_{t,k}; V_{\hat{\theta}_t}^{\text{rob},*}, S_{t,k+1})$  from Equation (17) for TV uncertainty set  
 1641 and Equation (20) for Wasserstein- $\ell$  uncertainty set

1642 7:      $\nu_{t,k+1} = \text{Proj}_{\mathcal{M}_\nu}(\nu_{t,k} + \beta_k [\hat{G}(\psi(S_{t,k}, A_{t,k})^\top \nu_{t,k}; V_{\hat{\theta}_t}^{\text{rob},*}, S_{t,k+1}) \psi(S_{t,k}, A_{t,k})])$

1643 8:     **Slow time-scale** ( $\alpha_k$ )

1644 9:     Compute  $\bar{\nu}_{t,k}$  from Equation (7)

1645 10:    Compute  $\hat{F}(\psi(S_{t,k}, A_{t,k})^\top \bar{\nu}_{t,k}; V_{\hat{\theta}_t}^{\text{rob},*}, S_{t,k+1})$  from Equation (18) for TV uncertainty set  
 1646 and Equation (21) for Wasserstein- $\ell$  uncertainty set

1647 11:     $TD_{t,k+1} = R(S_{t,k}, A_{t,k}) + \gamma \hat{F}(\psi(S_{t,k}, A_{t,k})^\top \bar{\nu}_{t,k}; V_{\hat{\theta}_t}^{\text{rob},*}, S_{t,k+1}) - \phi(S_{t,k}, A_{t,k})^\top \theta_{t,k}$

1648 12:     $\theta_{t,k+1} = \theta_{t,k} + \alpha_k TD_{t,k+1} \phi(S_{t,k}, A_{t,k})$

1649 13:   **end for**

1650 14:    $\hat{\theta}_{t+1} = \theta_{t,K}$ ,  $S_{t+1,0} = S_{t,K}$ ,  $\theta_{t+1,0} = \theta_{t,K}$ ,  $\nu_{t+1,0} = \nu_{t,K}$ .

1651 15: **end for**

1652 16: **Output:**  $\hat{\theta}_T$

---

1653

1654

1655

1656

1657

1658

1659

1660

1661

1662

1663

1664

1665

1666

1667

1668

1669

1670

1671

1672

1673

1674 **USE OF LARGE LANGUAGE MODEL**

1675

1676 The authors used large language models (e.g., ChatGPT) to polish the language in certain parts of the  
1677 paper. All technical content, proofs, and conclusions are the sole work of the authors.

1678

1679

1680

1681

1682

1683

1684

1685

1686

1687

1688

1689

1690

1691

1692

1693

1694

1695

1696

1697

1698

1699

1700

1701

1702

1703

1704

1705

1706

1707

1708

1709

1710

1711

1712

1713

1714

1715

1716

1717

1718

1719

1720

1721

1722

1723

1724

1725

1726

1727