

# DISTRIBUTIONALLY ROBUST MODEL-BASED REINFORCEMENT LEARNING WITH LARGE STATE SPACES

**Shyam Sundhar Ramesh**

*Department of Electrical Engineering  
University College London*

SHYAM.RAMESH.22@UCL.AC.UK

**Pier Giuseppe Sessa**

*Department of Computer Science  
ETH Zurich*

SESSAP@ETHZ.CH

**Yifan Hu**

*College of Management of Technology  
EPFL*

YIFAN.HU@EPFL.CH

**Andreas Krause**

*Department of Computer Science  
ETH Zurich*

KRAUSEA@ETHZ.CH

**Ilija Bogunovic**

*Department of Electrical Engineering  
University College London*

I.BOGUNOVIC@UCL.AC.UK

## Abstract

Three major challenges in reinforcement learning are the complex dynamical systems with large state spaces, the costly data acquisition processes, and the deviation of real-world dynamics from the training environment deployment. To overcome these issues, we study distributionally robust Markov decision processes with continuous state spaces under the widely used Kullback–Leibler, chi-square, and total variation uncertainty sets. We propose a model-based approach that utilizes Gaussian Processes and the maximum variance reduction algorithm to efficiently learn multi-output nominal transition dynamics, leveraging access to a generative model (i.e., simulator). We further demonstrate the statistical sample complexity of the proposed method for different uncertainty sets. These complexity bounds are independent of the number of states and extend beyond linear dynamics, ensuring the effectiveness of our approach in identifying near-optimal distributionally-robust policies. The proposed method can be further combined with other model-free distributionally robust reinforcement learning methods to obtain a near-optimal robust policy. Experimental results demonstrate the robustness of our algorithm to distributional shifts and its superior performance in terms of the number of samples needed.

**Keywords:** Reinforcement Learning, Robustness, Generative Model, Sample Complexity

## 1. Introduction

The use of reinforcement learning (RL) algorithms is gaining momentum in various complex domains, including robotics, nuclear fusion, and molecular discovery. Data acquisition in such environments can be a challenging and resource-intensive process. Safety considerations may also limit the amount of data that can be collected through interactions with the environment. To address this issue, a commonly adopted approach is to train RL policies using a simulator (generative model) enabling RL agents to learn from a simulated environment.

Dealing with complex applications that involve large state spaces requires data-efficient learning, even when a simulator is available. However, achieving optimal policies using existing approaches often requires a significant amount of training data, making data-efficient learning an ongoing challenge. Additionally, when deploying a policy to a real-world system, it is crucial to ensure its performance remains reliable despite mismatches between the simulator and the real-world system. Such mismatches can arise from approximation errors, time-varying system parameters, or even due to adversarial influence. The resulting mismatch, known as the 'sim-to-real gap', can diminish the performance or impact the reliability of RL algorithms trained on a simulator model.

In this work, we examine the use of a generative model in *distributionally-robust model-based reinforcement learning*. Our aim is to find a distributionally-robust policy that is near-optimal by actively querying the simulator with a state-action pair selected by the learning algorithm. To achieve this, we introduce the kernelized Maximum Variance Reduction (MVR) algorithm, which identifies a state-action pair with the highest uncertainty according to the model to learn the nominal model dynamics. The algorithm produces a nominal dynamics estimate that is utilized within the robust Markov Decision Process (MDP) framework, where an uncertainty set that includes all models close to the learned one is considered. We provide a thorough characterization of statistical sample complexity rates by utilizing the learned model to generate a near-optimal robust policy.

## 2. Problem Setting

A discounted Markov Decision Process (MDP) is a tuple  $(\mathcal{S}, \mathcal{A}, P, r, \gamma)$ , with  $\mathcal{S}$  denoting the state space, the action space  $\mathcal{A}$ , and the probabilistic transition dynamics  $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ . Here,  $\Delta(\mathcal{S})$  denotes the set of all probability distributions over  $\mathcal{S}$ . The reward function  $r : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  characterizes the reward  $r(s, a)$  the learner receives upon playing  $a \in \mathcal{A}$  in  $s \in \mathcal{S}$ , and  $\gamma \in [0, 1]$  denotes the discount factor. The learner uses a policy  $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$  to select  $a \in \mathcal{A}$  upon observing the state  $s \in \mathcal{S}$ . We define the cumulative discounted reward as  $\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t)$  for known initial state  $s_0$  and  $s_t \sim P(s_{t-1}, a_{t-1})$  for  $t > 0$  and  $a_t \sim \pi(s_t)$ . The value function  $V_\pi$  and the state-action value function  $Q_\pi$  are given as follows:

$$V_\pi(s) = \mathbb{E}_{P,\pi} \left[ \sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \middle| s_0 = s \right], \quad Q_\pi(s, a) = r(s, a) + \mathbb{E}_{P,\pi} \left[ \sum_{t=1}^{\infty} \gamma^t r(s_t, a_t) \right],$$

where  $a_t \sim \pi(s_t)$  and  $s_{t+1} \sim P(s_t, a_t)$ . Finally, we define the optimal policy  $\pi^*$  corresponding to dynamics  $P$  which yields the optimal value function, i.e.,  $V_{\pi^*}(s) = \max_\pi V_\pi(s)$  for all  $s \in \mathcal{S}$ . We assume the standard generative (or random) access model, in which the learner can query transition data arbitrarily from a simulator, i.e., each query to the simulator  $(s_t, a_t)$  outputs a sample  $s_{t+1} \in \mathbb{R}^d$  where  $s_{t+1} \sim P(s_t, a_t)$ . In particular, we consider the following frequently used transition dynamics model:

$$s_{t+1} = f(s_t, a_t) + \omega_t, \tag{1}$$

where  $\omega_t \in \mathbb{R}^d$  represents independent additive transition noise and follows a Gaussian distribution with zero mean and covariance  $\sigma^2 I$ .

**Regularity assumptions:** We assume that  $f$  is *unknown* and continuous for tractability reasons which is a common assumption when dealing with continuous state spaces (e.g., (11; 17; 28)). Considering the multi-output definition of  $f$  and in line with the previous

work (e.g., (11; 17)), we define the modified state-action space  $\bar{\mathcal{X}}$  as  $\bar{\mathcal{X}} := \mathcal{S} \times \mathcal{A} \times [d]$ , where the last dimension  $i \in \{1, 2, \dots, d\}$  incorporates the index of the output state vector, i.e.,  $f(\cdot, \cdot) = (\tilde{f}(\cdot, \cdot, 1), \dots, \tilde{f}(\cdot, \cdot, d))$  where  $\tilde{f} : \bar{\mathcal{X}} \rightarrow \mathbb{R}$ . In particular, we assume that  $\tilde{f}$  belongs to a space of well-behaved functions (RKHS), denoted by  $\mathcal{H}$ . Further details regarding the assumption on  $f$  are detailed in the appendix. We refer to the simulator environment determined by  $f$  as the *nominal model*  $P_f$ , while the true environment encountered by the agent in the real world might not be the same (e.g., due to a sim-to-real gap). Consequently, we utilize the robust MDP framework to tackle this by considering an uncertainty set comprising of all models close to the nominal one.

**Robust Markov Decision Process (RMDP):** We consider the robust MDP setting that addresses the uncertainty in transition dynamics and considers a set of transition models called the *uncertainty set*. We use  $\mathcal{P}^f$  to denote the uncertainty set that satisfies the  $(s, a)$ -rectangularity condition (27) (as defined in Equation (2)), an assumption that is commonly used in the related literature (e.g., (41; 42; 62)). Similar to MDPs, we specify RMDP by a tuple  $(\mathcal{S}, \mathcal{A}, \mathcal{P}^f, r, \gamma)$  where the uncertainty set  $\mathcal{P}^f$  consists of all models close to a nominal model  $P_f$  in terms of a distance measure  $D$ :

$$\mathcal{P}_{s,a}^f = \{p \in \Delta(\mathcal{S}) : D(p || P_f(s, a)) \leq \rho\}, \quad \text{and} \quad \mathcal{P}^f = \bigotimes_{(s,a) \in \mathcal{S} \times \mathcal{A}} \mathcal{P}_{s,a}^f. \quad (2)$$

Here,  $D$  denotes some distance measure between probability distributions, and  $\rho > 0$  defines the radius of the uncertainty set. In the RMDP setting, the goal is to discover a policy that maximizes the cumulative discounted reward for the worst-case transition model within the *given* uncertainty set. Concretely, the robust value function  $V_{\pi, f}^R$  corresponding to a policy  $\pi$  and the optimal robust value are given as follows:

$$V_{\pi, f}^R(s) = \inf_{P \in \mathcal{P}^f} \mathbb{E}_{P, \pi} \left[ \sum_{t=1}^{\infty} \gamma^t r(s_t, a_t) \middle| s_0 = s \right], \quad V_{\pi^*, f}^R(s) = \max_{\pi} V_{\pi, f}^R(s) \quad \forall s \in \mathcal{S}. \quad (3)$$

The goal of the learner is to discover a near-optimal robust policy while minimizing the total number of samples  $N$ , i.e., queries to the nominal model (simulator). Concretely, for a fixed precision  $\epsilon > 0$ , the goal is to output a policy  $\hat{\pi}_N$  after collecting  $N$  samples, such that  $\|V_{\hat{\pi}_N, f}^R - V_{\pi^*, f}^R\|_{\infty} \leq \epsilon$ .

### 3. Sampling Algorithm

In this section, we outline our methodology for addressing the problem described in Section 2.

**Maximum variance reduction:** With certain assumptions on the loss function and noise distribution, the function estimation in RKHS is analogous to the Bayesian Gaussian process framework (47). When used with the same kernel function, this allows the construction of mean and variance estimates of  $\tilde{f} \in \mathcal{H}$  using Gaussian processes (eq. (8) and eq. (9)). Based on these, one can construct shrinking statistical confidence bounds that hold with probability at least  $1 - \delta$ , i.e., the following holds  $|\tilde{f}(x) - \mu_{n-1}(x)| \leq \beta_n(\delta)\sigma_{n-1}(x)$  for every  $n \geq 1$  and  $x \in \bar{\mathcal{X}}$ . Here  $\{\beta_i\}_{i=1}^n$  stands for the sequence of parameters that are suitably set (see Lemma 5) to ensure the validity of the confidence bounds. We use the maximum variance reduction (MVR) algorithm (Algorithm 1) to learn about the nominal model  $f$ . MVR works on the principle of reducing the maximum uncertainty measured by the posterior standard deviation of a GP model calculated by using previously collected data. At each iteration, MVR queries a state-action pair that has the highest uncertainty according to the model and

---

**Algorithm 1** Maximum Variance Reduction (MVR) for learning model dynamics

---

- 1: **Require:** Simulator  $f$ , kernel  $k$ , domain  $\mathcal{S} \times \mathcal{A}$
  - 2: Set  $\mu_0(s, a) = 0, \sigma_0(s, a) = 1$  for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$
  - 3: **for**  $i = 1, \dots, n$  **do**
  - 4:    $(s_i, a_i) = \arg \max_{(s,a) \in \mathcal{S} \times \mathcal{A}} \|\sigma_{i-1}(s, a)\|_2$
  - 5:   Observe  $s_{i+1} = f(s_i, a_i) + \omega_i$   
    (i.e., sample  $s_{i+1}$  from nominal  $P_f(s_i, a_i)$ )
  - 6:   Update to  $\mu_i$  and  $\sigma_i$  by using  $(s_i, a_i, s_{i+1})$   
    according to eq. (8) and eq. (9)
  - 7: **end for**
  - 8: **return** The dynamics estimate  $\hat{f}_n(\cdot, \cdot) = \mu_n(\cdot, \cdot)$
- 

uses the obtained observation to update the GP posterior. The algorithm outputs nominal dynamics estimate  $\hat{f}_n$  corresponding to the final GP posterior mean  $\mu_n$ . We defer further details of the Gaussian Process framework to the appendix.

#### 4. Sample Complexity

This section discusses the statistical sample complexity of the proposed MVR algorithm in distributionally robust MDPs rather than designing algorithms to find  $\hat{\pi}_N$  as done in (42; 21; 36). One can easily incorporate the MVR algorithm with the model-free algorithms from these previous works to find an optimal  $\hat{\pi}_N$  using  $\hat{f}_N$  allowing us to bypass the more costly simulator  $f$ . Given the optimal robust policies  $\hat{\pi}_N$  and  $\pi^*$  corresponding to the learned nominal dynamics  $\hat{f}_N$  by the MVR algorithm with  $N$  iterations and the true nominal dynamics  $f$ , respectively, we show the number of samples needed by the MVR algorithm to ensure that the following holds:

$$|V_{\hat{\pi}_N, f}^R(s) - V_{\pi^*, f}^R(s)| \leq \epsilon, \forall s \in \mathcal{S}. \quad (4)$$

**Theorem 1** (*Sample Complexity of MVR under KL uncertainty set*) Consider a robust MDP with nominal transition dynamics  $f$  satisfying the regularity assumptions from Section 2 and with uncertainty set defined as in Equation (2) w.r.t. KL divergence. For  $\pi^*$  denoting the robust optimal policy w.r.t. nominal transition dynamics  $f$  and  $\hat{\pi}_N$  denoting the robust optimal policy w.r.t. learned nominal transition dynamics  $\hat{f}_N$  via MVR (Algorithm 1), and  $\delta \in (0, 1), \epsilon \in (0, \frac{1}{1-\gamma})$ , it holds that  $\max_s |V_{\hat{\pi}_N, f}^R(s) - V_{\pi^*, f}^R(s)| \leq \epsilon$  with probability at least  $1 - \delta$  for any  $N$  such that

$$N = \mathcal{O}\left(\frac{e^{\frac{2-\gamma}{(1-\gamma)\alpha_{kl}}}}{\gamma^2 \beta_N^2(\delta) d^2 \Gamma_{Nd}}\right). \quad (5)$$

Theorem 1 shows the number of samples required from the nominal transition dynamics  $f$  (simulator) to construct a robust optimal policy reliably with high probability. The complexity bounds depend on the maximum information gain  $\Gamma_{Nd}$  (Equation (18) a kernel-dependent quantity that is frequently used in GP optimization), which is sublinear in  $N$  for many commonly used kernels ((52)). Furthermore, in our analysis, we use the confidence bounds from (54) with  $\beta_N^2(\delta)$  (see Lemma 6) which only exhibits a logarithmic dependence on  $N$ . An additional  $d$  factor that denotes the dimension of the state space  $\mathcal{S}$  in the obtained bound

comes from utilizing the multi-output (of dimension  $d$ ) GP framework to model the transition dynamics, which also appears in the regret bounds of similar works (11; 17; 18). Finally, the term  $\alpha_{kl} \in (0, \frac{1}{2(1-\gamma)\rho})$  is a problem-dependent parameter that is independent of  $N$ , which similarly appears in the guarantees of (41).

We can compare our guarantees with the existing sample-complexity results in model-based distributionally robust RL which, however, only consider finite state-action spaces (41; 62; 59). In particular, when considering KL uncertainty sets, (41) obtain sample complexity of order  $\mathcal{O}\left(e^{\frac{\alpha_{kl}+2}{\alpha_{kl}(1-\gamma)}} \frac{\gamma^2 |\mathcal{S}|^2 |\mathcal{A}|}{(1-\gamma)^4 \rho^2 \epsilon^2}\right)$  up to logarithmic factors. Notably, the latter complexity bound explicitly depends on the cardinality of the state and action spaces  $|\mathcal{S}|$  and  $|\mathcal{A}|$ , thus scaling badly when  $\mathcal{S}$  and  $\mathcal{A}$  are large or continuous. Instead, the guarantee of Theorem 1 depends on the state-action space *only* through  $\Gamma_{Nd}$  which remains bounded even when these are continuous. This allows us to successfully extend the distributionally robust framework to continuous state spaces. Other terms in the bound of Theorem 1 such as  $\gamma$  (the discount factor),  $\rho$  (radius of the uncertainty set) have similar dependencies. Crucially, the dependence on the precision parameter  $\epsilon$  remains the same when compared to the guarantees provided for finite state-action setting. We relegate the proof of Theorem 1 and statistical sample complexities for  $\chi^2$  distance and TV distance uncertainty sets to the Appendix.

## 5. Experiments

The aim of our experiments is to show the effectiveness of the proposed distributionally-robust model-based approach. In particular, our goal is to evaluate the robustness of our policies against different perturbations of the environment’s parameters, and compare them with existing non-robust methods. Moreover, we compare our approach with model-free methods (robust and non-robust) which typically require a significantly larger number of interactions with the nominal environment. We consider the OpenAI’s gym (8) environments of swing-up Pendulum, Cartpole, and Reacher and test our approach against various perturbations.

**Module 1: Learning the model.** To learn the nominal environment, we utilize a setup similar to that of (38), but use the proposed Max Variance Reduction (MVR) method (Algorithm 1) instead. Similar to (38), we use a Gaussian process (GP) prior to model the transition dynamics  $f(s, a)$  (alternate models such as Neural Ensembles or Bayesian neural networks can be used to model the transition dynamics as done in, e.g., (17; 18)). As in continuous control problems the subsequent states are fairly close, we use our multi-output GP to model the difference  $f(s_t, a_t) - s_{t+1}$ .

**Module 2: Computing a robust policy.** Given a learned model  $\hat{f}_n$ , we compute the associated robust policy  $\hat{\pi}_n$  using the Robust Fitted Q-Iteration (RFQI) algorithm from (42). RFQI computes a robust policy from offline data by alternated maximization of a dual-variable function and a Q-function. We generate such offline data by using a  $\epsilon$ -greedy non-robust policy (using Soft Actor-Critic (25) or Model Predictive Control (9; 13)) which we train *on the learned model*  $\hat{f}_n$  from Module 1 and let interact with it for  $10^6/10^5$  steps. Note that this is crucially different from the vanilla RFQI (42) where the true nominal environment was used both for training such policy and for generating offline data. Indeed, this would require a significantly larger number of environment interactions.

We provide further implementation details regarding training, baselines, evaluation and hyperparameters in Appendix F and discuss the results below. In Figure 1 we plot the average performance (over 20 episodes) of the different baselines subject to different perturbation types

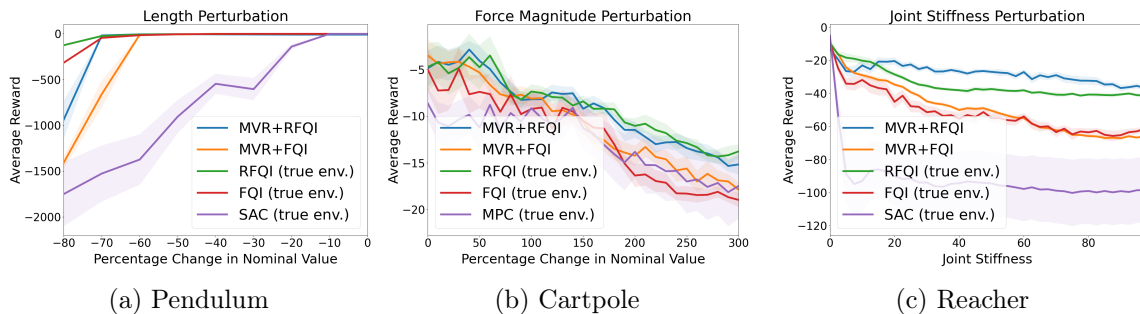


Figure 1: Average performance (over 20 episodes) on the considered environments, as a function of different perturbations: length perturbation for Pendulum, force magnitude perturbation for Cartpole, and perturbed joint stiffness for Reacher.

	MVR+RFQI (ours)	MVR+FQI	SAC	MPC	RFQI	FQI ((23))
Pendulum	60	60	$10^4$	-	$10^6 + 10^4$	$10^6 + 10^4$
Cartpole	150	150	-	2250/step	$10^5 \cdot 2250$	$10^5 \cdot 2250$
Reacher	2000	2000	$10^6$	-	$10^6 + 10^6$	$10^6 + 10^6$

Table 1: Number of interactions with the nominal environment to obtain the results of Figure 1. For MPC, a total of 2250 interactions are required at each step for planning multiple rollouts and selecting the best action. Both RFQI and FQI utilize  $10^6/10^5$  offline data generated by SAC or MPC.

and magnitudes for each environment. Results for other perturbations are relegated to Appendix F. In Table 1 we report the total number of interactions with the nominal environment required to compute the evaluated policies. We remark that MVR+RFQI and MVR+FQI interact with the environment only to learn a good Gaussian Process model via the MVR approach. Instead, the other model-free methods utilize the nominal environment throughout the whole training and, in case of RFQI and FQI, even to generate offline data. Notably, the policy computed by MVR+RFQI displays comparable performance to its model-free counterpart RFQI which, as shown in Table 1, requires a significantly larger number of environment interactions. This illustrates the sample-efficiency of the MVR approach in acquiring informative data and yielding good model estimates. Moreover, as the perturbation magnitude increases, MVR+RFQI generally achieves higher performance compared to MVR+FQI and the other non-robust methods, demonstrating the robustness of the computed policies.

## 6. Conclusions

We investigated distributionally robust reinforcement learning in the context of continuous state spaces and non-linear transition dynamics. Specifically, we proposed a model-based approach within the generative model setting, utilizing maximum variance reduction to learn nominal transition dynamics effectively. Our results include novel statistical sample complexity guarantees for commonly used uncertainty sets, required for identifying near-optimal distributionally robust policies in large state spaces. Through experiments conducted in popular RL-testing environments, we demonstrated the sample efficiency and robustness of our algorithm in the presence of distributional shifts. An important avenue for future research is the extension of our algorithm to the online and offline reinforcement learning settings.

## References

- [1] Y. Abbasi-Yadkori, P. Bartlett, K. Bhatia, N. Lazic, C. Szepesvari, and G. Weisz. Politex: Regret bounds for policy iteration using expert prediction. *International Conference on Machine Learning (ICML)*, 2019.
- [2] A. Agarwal, S. Kakade, and L. F. Yang. Model-based reinforcement learning with a generative model is minimax optimal. In *Conference on Learning Theory*, pages 67–83. PMLR, 2020.
- [3] S. M. Ali and S. D. Silvey. A general class of coefficients of divergence of one distribution from another. *Journal of the Royal Statistical Society: Series B (Methodological)*, 28(1):131–142, 1966.
- [4] O. M. Andrychowicz, B. Baker, M. Chociej, R. Jozefowicz, B. McGrew, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray, et al. Learning dexterous in-hand manipulation. *The International Journal of Robotics Research*, 39(1):3–20, 2020.
- [5] K. P. Badrinath and D. Kalathil. Robust reinforcement learning using least squares policy iteration with provable performance guarantees. *International Conference on Machine Learning (ICML)*, 2021.
- [6] A. Ben-Tal, D. Den Hertog, A. De Waegenaere, B. Melenberg, and G. Rennen. Robust solutions of optimization problems affected by uncertain probabilities. *Management Science*, 59(2):341–357, 2013.
- [7] I. Bogunovic, J. Scarlett, S. Jegelka, and V. Cevher. Adversarially robust optimization with Gaussian processes. *Conference on Neural Information Processing Systems (NeurIPS)*, 2018.
- [8] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. Openai gym, 2016.
- [9] E. F. Camacho and C. B. Alba. *Model predictive control*. Springer science & business media, 2013.
- [10] J. Chen and N. Jiang. Information-theoretic considerations in batch reinforcement learning. *International Conference on Machine Learning (ICML)*, 2019.
- [11] S. R. Chowdhury and A. Gopalan. Online learning in kernelized Markov decision processes. *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2019.
- [12] P. Christiano, Z. Shah, I. Mordatch, J. Schneider, T. Blackwell, J. Tobin, P. Abbeel, and W. Zaremba. Transfer from simulation to real world through learning deep inverse dynamics model. *arXiv preprint arXiv:1610.03518*, 2016.
- [13] K. Chua, R. Calandra, R. McAllister, and S. Levine. Deep reinforcement learning in a handful of trials using probabilistic dynamics models. *Conference on Neural Information Processing Systems (NeurIPS)*, 2018.

- [14] I. Clavera, J. Rothfuss, J. Schulman, Y. Fujita, T. Asfour, and P. Abbeel. Model-based reinforcement learning via meta-policy optimization. In *Conference on Robot Learning*, pages 617–629. PMLR, 2018.
- [15] N. Cressie and T. R. Read. Multinomial goodness-of-fit tests. *Journal of the Royal Statistical Society: Series B (Methodological)*, 46(3):440–464, 1984.
- [16] I. Csiszár. Information-type measures of difference of probability distributions and indirect observation. *studia scientiarum Mathematicarum Hungarica*, 2:229–318, 1967.
- [17] S. Curi, F. Berkenkamp, and A. Krause. Efficient model-based reinforcement learning through optimistic policy search and planning. *Conference on Neural Information Processing Systems (NeurIPS)*, 2020.
- [18] S. Curi, I. Bogunovic, and A. Krause. Combining pessimism with optimism for robust and efficient model-based deep reinforcement learning. *International Conference on Machine Learning (ICML)*, 2021.
- [19] M. P. Deisenroth and C. E. Rasmussen. Pilco: A model-based and data-efficient approach to policy search. In *ICML*, 2019.
- [20] E. Derman, D. Mankowitz, T. Mann, and S. Mannor. A bayesian approach to robust reinforcement learning. In *Uncertainty in Artificial Intelligence*, pages 648–658. PMLR, 2020.
- [21] E. Derman, D. J. Mankowitz, T. A. Mann, and S. Mannor. Soft-robust actor-critic policy-gradient. *arXiv preprint arXiv:1803.04848*, 2018.
- [22] J. Duchi and H. Namkoong. Learning models with uniform performance via distributionally robust optimization. *arXiv preprint arXiv:1810.08750*, 2018.
- [23] D. Ernst, P. Geurts, and L. Wehenkel. Tree-based batch mode reinforcement learning. *Journal of Machine Learning Research*, 6, 2005.
- [24] M. Gheshlaghi Azar, R. Munos, and H. J. Kappen. Minimax pac bounds on the sample complexity of reinforcement learning with a generative model. *Machine learning*, 91(3):325–349, 2013.
- [25] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. *International Conference on Machine Learning (ICML)*, 2018.
- [26] Z. Hu and L. J. Hong. Kullback-leibler divergence constrained distributionally robust optimization. *Available at Optimization Online*, pages 1695–1724, 2013.
- [27] G. N. Iyengar. Robust dynamic programming. *Mathematics of Operations Research*, 30(2):257–280, 2005.
- [28] S. Kakade, A. Krishnamurthy, K. Lowrey, M. Ohnishi, and W. Sun. Information theoretic regret bounds for online nonlinear control. *Conference on Neural Information Processing Systems (NeurIPS)*, 2020.



- [29] S. M. Kakade. *On the sample complexity of reinforcement learning*. University of London, University College London (United Kingdom), 2003.
- [30] M. Kearns, Y. Mansour, and A. Y. Ng. A sparse sampling algorithm for near-optimal planning in large markov decision processes. *Machine learning*, 49(2):193–208, 2002.
- [31] J. Kirschner, I. Bogunovic, S. Jegelka, and A. Krause. Distributionally robust Bayesian optimization. *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2020.
- [32] T. Lattimore, C. Szepesvari, and G. Weisz. Learning with good feature representations in bandits and in rl with a generative model. *International Conference on Machine Learning (ICML)*, 2020.
- [33] G. Li, Y. Wei, Y. Chi, Y. Gu, and Y. Chen. Breaking the sample size barrier in model-based reinforcement learning with a generative model. *Conference on Neural Information Processing Systems (NeurIPS)*, 2020.
- [34] X. Li, V. Mehta, J. Kirschner, I. Char, W. Neiswanger, J. Schneider, A. Krause, and I. Bogunovic. Near-optimal policy identification in active reinforcement learning. *International Conference on Learning Representations*, 2023.
- [35] S. H. Lim, H. Xu, and S. Mannor. Reinforcement learning in robust Markov decision processes. *Conference on Neural Information Processing Systems (NeurIPS)*, 2013.
- [36] D. J. Mankowitz, N. Levine, R. Jeong, Y. Shi, J. Kay, A. Abdolmaleki, J. T. Springenberg, T. Mann, T. Hester, and M. Riedmiller. Robust reinforcement learning for continuous control with model misspecification. *arXiv preprint arXiv:1906.07516*, 2019.
- [37] S. Mannor, O. Mebel, and H. Xu. Robust mdps with k-rectangular uncertainty. *Mathematics of Operations Research*, 41(4):1484–1509, 2016.
- [38] V. Mehta, B. Paria, J. Schneider, S. Ermon, and W. Neiswanger. An experimental design perspective on model-based reinforcement learning. *arXiv preprint arXiv:2112.05244*, 2021.
- [39] T. Nguyen, S. Gupta, H. Ha, S. Rana, and S. Venkatesh. Distributionally robust Bayesian quadrature optimization. *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2020.
- [40] A. Nilim and L. El Ghaoui. Robust control of markov decision processes with uncertain transition matrices. *Operations Research*, 53(5):780–798, 2005.
- [41] K. Panaganti and D. Kalathil. Sample complexity of robust reinforcement learning with a generative model. *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2022.
- [42] K. Panaganti, Z. Xu, D. Kalathil, and M. Ghavamzadeh. Robust reinforcement learning using offline data. *arXiv preprint arXiv:2208.05129*, 2022.

- [43] X. B. Peng, M. Andrychowicz, W. Zaremba, and P. Abbeel. Sim-to-real transfer of robotic control with dynamics randomization. In *2018 IEEE international conference on robotics and automation (ICRA)*, pages 3803–3810. IEEE, 2018.
- [44] M. Petrik and R. H. Russel. Beyond confidence regions: Tight Bayesian ambiguity sets for robust MDPs. *Conference on Neural Information Processing Systems (NeurIPS)*, 2019.
- [45] C. Pinneri, S. Sawant, S. Blaes, J. Achterhold, J. Stueckler, M. Rolinek, and G. Martius. Sample-efficient cross-entropy method for real-time planning. *arXiv preprint arXiv:2008.06389*, 2020.
- [46] L. Pinto, J. Davidson, R. Sukthankar, and A. Gupta. Robust adversarial reinforcement learning. *International Conference on Machine Learning (ICML)*, 2017.
- [47] C. E. Rasmussen and C. Williams. Gaussian processes for machine learning, vol. 1, 2006.
- [48] D. Rastogi, I. Koryakovskiy, and J. Kober. Sample-efficient reinforcement learning via difference models. In *Machine Learning in Planning and Control of Robot Motion Workshop at ICRA*, 2018.
- [49] A. Roy, H. Xu, and S. Pokutta. Reinforcement learning under model mismatch. *Conference on Neural Information Processing Systems (NeurIPS)*, 2017.
- [50] A. Shapiro. Distributionally robust stochastic programming. *SIAM Journal on Optimization*, 27(4):2258–2275, 2017.
- [51] R. Shariff and C. Szepesvári. Efficient planning in large MDPs with weak linear function approximation. *Conference on Neural Information Processing Systems (NeurIPS)*, 2020.
- [52] N. Srinivas, A. Krause, S. M. Kakade, and M. Seeger. Gaussian process optimization in the bandit setting: No regret and experimental design. *arXiv preprint arXiv:0912.3995*, 2009.
- [53] A. Tamar, S. Mannor, and H. Xu. Scaling up robust mdps using function approximation. *International Conference on Machine Learning (ICML)*, 2014.
- [54] S. Vakili, N. Bouziani, S. Jalali, A. Bernacchia, and D.-s. Shiu. Optimal order simple regret for gaussian process bandits. *Conference on Neural Information Processing Systems (NeurIPS)*, 2021.
- [55] Y. Wang and S. Zou. Online robust reinforcement learning with model uncertainty. *Conference on Neural Information Processing Systems (NeurIPS)*, 2021.
- [56] W. Wiesemann, D. Kuhn, and B. Rustem. Robust Markov decision processes. *Mathematics of Operations Research*, 38(1):153–183, 2013.
- [57] M. Wulfmeier, I. Posner, and P. Abbeel. Mutual alignment transfer learning. In *Conference on Robot Learning*, pages 281–290. PMLR, 2017.

- [58] H. Xu and S. Mannor. Distributionally robust Markov decision processes. *Conference on Neural Information Processing Systems (NeurIPS)*, 2010.
- [59] W. Yang, L. Zhang, and Z. Zhang. Towards theoretical understandings of robust markov decision processes: Sample complexity and asymptotics. *arXiv preprint arXiv:2105.03863*, 2021.
- [60] P. Yu and H. Xu. Distributionally robust counterpart in markov decision processes. *IEEE Transactions on Automatic Control*, 61(9):2538–2543, 2015.
- [61] H. Zhang, H. Chen, C. Xiao, B. Li, M. Liu, D. Boning, and C.-J. Hsieh. Robust deep reinforcement learning against adversarial perturbations on state observations. *Conference on Neural Information Processing Systems (NeurIPS)*, 2020.
- [62] Z. Zhou, Z. Zhou, Q. Bai, L. Qiu, J. Blanchet, and P. Glynn. Finite-sample regret bound for distributionally robust offline tabular reinforcement learning. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*. PMLR, 2021.

## Appendix A. Related Work

Reinforcement learning with a generative model is introduced in (30) wherein one assumes access to a simulator that outputs the next state given any state-action pair. (29) elucidate various uses for this generative setting and analyze it in further detail. For the finite MDP case, such a generative setting has been subsequently studied in various works such as (29; 24; 33) and, recently, by (2) who provide minimax optimality guarantees for the naive plug-in estimator based algorithm. For large state spaces, generative RL is typically combined with function approximation as studied, e.g., by (1; 51; 32; 34). Recently, (38) consider generative RL in continuous state-action spaces from an experimental perspective and showcase the relevance of this setting to the nuclear fusion dynamics research. In addition, (34) present an active exploration strategy that utilizes the least-squares value iteration. Their approach aims to identify a near-optimal policy across the entire state space, providing polynomial sample complexity guarantees that remain unaffected by the number of states. In contrast to these works, we use generative RL to discover *distributionally robust* policies through the modeling of unknown transition dynamics.

In model-based reinforcement learning, the model learned from a simulator encounters two issues well discussed in the literature, namely, the model-bias (19; 14) and the simulation to reality (sim2real) gap (4; 43; 36; 12; 48; 57). To address this from the perspective of distributional robustness, previous works (62; 41; 59) have considered distributional robustness aspects in the context of finite Markov decision processes (MDPs) using the robust MDP framework from (27; 40). Various other works utilize this robust MDP framework such as (58; 56; 60; 37; 5; 44) for the planning problem, and provide asymptotic guarantees for tabular and linear function approximators (35; 53; 49; 55). Our work is closely related to the recent works on distributionally robust RL (62; 41; 59). However, unlike ours, the sample complexity bounds established in these works rely on the number of states and actions, making them impractical for large or infinite state spaces. In the model-free setting, distributionally robust RL with large state space (though, still assumed to be finite) was considered by (42) in a function approximation setup. They assume access to offline data from the nominal transition dynamics and provide computational sample complexity bounds in terms of the size of the hypothesis space that is used to represent the set of state-action value functions (Q-function). Other works such as (46; 20; 36; 61) consider robustness aspects in deep reinforcement learning, but these approaches lack theoretical guarantees. To the best of our knowledge, our work is the first one to address the distributionally robust RL problem in the generative model setting with a *model-based* approach and *large* state spaces. Moreover, we are the first to consider general *non-linear* transition dynamics and derive provable sample complexity guarantees for such a setting.

Similar to previous works, we utilize the kernelized MDP framework from (11) to model transition dynamics with continuous states and actions by assuming that the transition function belongs to an associated Reproducing Kernel Hilbert Space (RKHS). Such continuous MDP formulations also appear in (17; 18), however, these works consider finite horizon MDPs while in our work we consider infinite horizon discounted MDPs. In particular, (18) propose an adversarially robust upper-confidence algorithm to optimize performance in the worst case. However, their algorithm provides robustness guarantees against adversarial perturbations to the transition dynamics. Our work differs from this perspective as we consider robustness

w.r.t. distributional shifts of the transition dynamics. Finally, in the related kernelized *bandit* setting, model-based distributionally robust algorithms are proposed in (31; 7; 39).

## Appendix B. Theoretical Guarantees of Maximum Variance Reduction (MVR)

In this section, we detail the assumptions on  $f$  and formally introduce the Gaussian process model. We describe the confidence bound results from (54) and adapt them to the case of multi-output GP models. Finally, we provide sample complexity guarantees for the MVR algorithm.

We recall the introduced notation  $\mathcal{X} = \mathcal{S} \times \mathcal{A}$  and remark that we use both  $(s_i, a_i)$  and  $x_i$  interchangeably in this section.

### B.1 Regularity Assumptions

We assume that  $f$  is *unknown* and continuous for tractability reasons which is a common assumption when dealing with continuous state spaces (e.g., (11; 17; 28)). Further on, we assume that  $f$  resides in the Reproducing Kernel Hilbert Space (RKHS). Considering the multi-output definition of  $f$  and in line with the previous work (e.g., (11; 17)), we define the modified state-action space  $\bar{\mathcal{X}}$  (over which the RKHS is defined) as  $\bar{\mathcal{X}} := \mathcal{S} \times \mathcal{A} \times [d]$ , where the last dimension  $i \in \{1, 2, \dots, d\}$  incorporates the index of the output state vector, i.e.,  $f(\cdot, \cdot) = (\tilde{f}(\cdot, \cdot, 1), \dots, \tilde{f}(\cdot, \cdot, d))$  where  $\tilde{f} : \bar{\mathcal{X}} \rightarrow \mathbb{R}$ . In particular, we assume that  $\tilde{f}$  belongs to a space of well-behaved functions, denoted by  $\mathcal{H}$ , induced by some continuous, positive definite kernel function  $k : \bar{\mathcal{X}} \times \bar{\mathcal{X}} \rightarrow \mathbb{R}$  and equipped with an inner product  $\langle \cdot, \cdot \rangle_k$ . All functions belonging to an RKHS  $\mathcal{H}$  satisfy the reproducing property defined w.r.t. the inner product  $\langle \cdot, \cdot \rangle_k : \langle \tilde{f}, k(x, \cdot) \rangle = \tilde{f}(x)$  for  $\tilde{f} \in \mathcal{H}$ . We also make the following common assumptions: (i) the kernel function  $k$  is bounded  $k(x, x') \leq 1$  for all  $x, x' \in \bar{\mathcal{X}}$  and  $\bar{\mathcal{X}}$  is a compact set ( $\mathcal{X} \subset \mathbb{R}^p$ ), and (ii) every function  $\tilde{f} \in \mathcal{H}$  has a bounded RKHS norm (induced by the inner product)  $\|\tilde{f}\|_k \leq B$ .

### B.2 Gaussian Process Model

Gaussian process (GP) is a non-parametric model that is often used to express uncertainty over functions on any set (e.g., RKHS). They allow to tractably construct posterior distribution over functions in the set to estimate the unknown non-linear function  $\tilde{f} : \mathcal{X} \rightarrow \mathbb{R}$  given data containing samples from function  $\tilde{f}$ . It follows the Bayesian methodology of calculating posterior given the prior and assumes that the function values at any finite subset of the domain  $\mathcal{X}$  follow the multivariate Gaussian distribution. One specifies a GP by a prior mean function and a covariance function usually defined using a kernel  $k(x, x')$  where  $x, x' \in \mathcal{X}$ .

Assuming that the samples of  $\tilde{f} : \mathcal{X} \rightarrow \mathbb{R}$  are noisy measurements of the underlying true function  $\tilde{f}$  with i.i.d. Gaussian noise  $\mathcal{N}(0, \lambda)$ , the posterior mean and covariance function of the posterior distribution can be explicitly calculated. In essence, for  $\{x_1, \dots, x_N\} \in \mathcal{X}$  and  $y_n = \tilde{f}(x_n) + \omega_n$ , the posterior mean, covariance and variance are given by:

$$\mu_n(x) = k_n(x)(K_n + I_n \lambda)^{-1} y_n, \quad (6)$$

$$\begin{aligned} k_n(x, x') &= k(x, x') - k_n(x)(K_n + I_n \lambda)^{-1} k_n^T(x'), \\ \sigma_n^2(x) &= k_n(x, x). \end{aligned} \quad (7)$$

Here  $K_n$  denotes the covariance matrix whose entries are  $[K_n]_{i,j} = k(x_i, x_j)$  with  $x_i, x_j \in \{x_1, \dots, x_N\}$  and  $k_n(x) = [k(x, x_1), \dots, k(x, x_N)]$  denotes the covariance vector whose entries are the covariance between  $x$  and  $x_j$  for all  $x_j \in \{x_1, \dots, x_N\}$ . The  $n \times n$  identity matrix is denoted as  $I_n$ .

We consider multi-output GPs to model the unknown function  $f$  that outputs states in  $\mathbb{R}^d$ . Similar to Equation (6) and Equation (7), we get analogous expressions for the multi-output case in Equation (8) and Equation (9).

**Multi-output Gaussian process:** Under the assumptions of Section 2, modeling uncertainty and learning the transition model  $f$  can be performed via the Gaussian process framework. A Gaussian process  $GP(\mu(\cdot), k(\cdot, \cdot))$  over the input domain  $\bar{\mathcal{X}}$ , is a collection of random variables  $(\tilde{f}(x))_{x \in \bar{\mathcal{X}}}$  whose every finite subset  $(\tilde{f}(x_i))_{i=1}^n, n \in \mathbb{N}$ , follows multivariate Gaussian distribution with mean  $\mathbb{E}[\tilde{f}(x_i)] = \mu(x_i)$  and covariance  $\mathbb{E}[(\tilde{f}(x_i) - \mu(x_i))(\tilde{f}(x_j) - \mu(x_j))] = k(x_i, x_j)$  for every  $1 \leq i, j \leq n$ . Standard algorithms implicitly use a zero-mean  $GP(0, k(\cdot, \cdot))$  as the prior distribution over  $\tilde{f}$ , i.e.  $\tilde{f} \sim GP(0, k(\cdot, \cdot))$ , and assume that the noise variables are drawn independently across  $t$  from  $\mathcal{N}(0, \lambda)$  with  $\lambda > 0$ . Considering the multi-output definition of  $f(\cdot, \cdot) = (\tilde{f}(\cdot, \cdot, 1), \dots, \tilde{f}(\cdot, \cdot, d))$ , we build  $d$  copies of the dataset such that  $\mathcal{D}_{1:n,l} = \{(s_i, a_i, l), s_{i+1,l}\}_{i=1}^n$  each with  $n$  transitions from a particular state-action pair  $(s, a)$  to component  $l$  of next state. For  $x_i = (s_i, a_i)$  and  $y_{i,l} = s_{i+1,l}$ , the posterior mean, covariance and variance for  $\tilde{f}(x, l)$  are given by:

$$\mu_{nd}(x, l) = k_{nd}(x, l)(K_{nd} + I_{nd}\lambda)^{-1}y_{nd}, \quad (8)$$

$$\begin{aligned} k_{nd}((x, l), (x', l)) &= k((x, l), (x', l)) - \\ & k_{nd}(x, l)(K_{nd} + I_{nd}\lambda)^{-1}k_{nd}^T(x', l), \\ \sigma_{nd}^2(x, l) &= k_{nd}((x, l), (x, l)). \end{aligned} \quad (9)$$

Here  $K_{nd}$  denotes the covariance matrix of dimensions  $nd \times nd$  whose entries are  $k((x_i, l), (x_j, l'))$  with  $1 \leq i, j \leq n$  and  $1 \leq l, l' \leq d$ .  $k_{nd}(x, l) = [k((x, l), (x_i, l'))]_{1 \leq i \leq n, 1 \leq l' \leq d}$  denotes the covariance vector and  $y_{nd} = [y_{i,l}]_{1 \leq i \leq n, 1 \leq l \leq d}$  denotes the output vector.

Correspondingly, the posterior mean and variance for  $f$  would be

$$\mu_n(s, a) = (\mu_{nd}(s, a, 1), \dots, \mu_{nd}(s, a, d)), \quad (10)$$

$$\sigma_n(s, a) = (\sigma_{nd}(s, a, 1), \dots, \sigma_{nd}(s, a, d)). \quad (11)$$

### B.3 Non-adaptive Multi-output Confidence Bounds

Our Algorithm 1 uses the maximum variance reduction rule to learn about the transition dynamics. As seen in our analysis (see Theorem 10), we are interested in constructing confidence intervals for  $f$  only at the end of  $n$  iterations (i.e., after taking  $n$  samples), and hence, we do not require anytime confidence bounds (e.g., as in (52)). Moreover, in our algorithm, the current decision  $(s_i, a_i)$  does not depend on the previous noise realizations. By focusing on the single-output case first, the following confidence lemma from (54), can be used to construct confidence intervals with  $\beta(\delta)$  independent of  $n$  which holds w.h.p. for a fixed  $x \in \mathcal{X}$ :

**Lemma 2** *Given  $n$  noisy observations of  $\tilde{f} : \mathcal{X} \rightarrow \mathbb{R}$  with  $\|f\|_k \leq B$  where noise  $\{\omega_1, \dots, \omega_n\}$  is independent of inputs  $\{x_1, \dots, x_n\}$ , for  $\beta(\delta) = B + \frac{\sigma}{\lambda} \sqrt{2 \log(2/\delta)}$ , and  $\mu_n, \sigma_n$  as defined in*

Equation (6) and Equation (7), the following holds for a fixed  $x \in \mathcal{X}$  with probability at least  $1 - \delta$ ,

$$|f(x) - \mu_n(x)| \leq \beta(\delta)\sigma_n(x).$$

To extend this result over the entire input set  $x \in \mathcal{X}$ , the authors in (54) use a discretization assumption which ensures that there exists a discretization  $\mathcal{D}_n$  such that  $\tilde{f}(x) - \tilde{f}([x]_n) \leq \frac{1}{\sqrt{n}}$ , where  $[x]_n = \arg \min_{x' \in \mathcal{D}_n} \|x - x'\|_2$  and  $|\mathcal{D}_n| \leq CB^d n^{d/2}$  for  $C$  being independent of  $n$  and  $B$  (RKHS norm bound). Consequently, they obtain the following lemma providing uniform confidence bounds:

**Lemma 3** ((54, Theorem-3)) *Given  $n$  noisy observations of  $\tilde{f} : \mathcal{X} \rightarrow \mathbb{R}$ ,  $\mathcal{X} \subset \mathbb{R}$  satisfying  $\|\tilde{f}\|_k \leq B$  where noise  $\{\omega_1, \dots, \omega_n\}$  is independent of inputs  $\{x_1, \dots, x_n\} \subset \mathcal{X}$  and when there exists discretization  $\mathcal{D}_n$  of  $\mathcal{X}$  with  $|\mathcal{D}_n| \leq CB^d n^{d/2}$ , for  $\beta(\delta) = B + \frac{\sigma}{\lambda} \sqrt{2 \log(2/\delta)}$  and  $\beta_n(\delta) = 2B + \beta(\frac{\delta}{3C(B + \sqrt{n}\beta(2\delta/3n))^{d_n d/2}})$ ,  $\mu_n, \sigma_n$  as defined in Equation (6) and Equation (7), the following holds for all  $x \in \mathcal{D}_n$  with probability at least  $1 - \delta$ ,*

$$|\tilde{f}(x) - \mu_n(x)| \leq \beta_n(\delta)\sigma_n(x).$$

To extend this result to multiple dimensions as required in our work, we take the same discretization assumption as in (54). But considering the multi-output definition of  $f$ , we define the modified state-action space  $\bar{\mathcal{X}}$ . This is in line with (11), which also has a similar multi-output setting. We define the modified state-action space as  $\bar{\mathcal{X}} := \mathcal{S} \times \mathcal{A} \times \{1, 2, \dots, d\}$  where the last dimension  $i \in \{1, 2, \dots, d\}$  incorporates the index of the output vector, in the sense that  $f(\cdot, \cdot) = (\tilde{f}(\cdot, \cdot, 1), \dots, \tilde{f}(\cdot, \cdot, d))$  where  $\tilde{f} : \bar{\mathcal{X}} \rightarrow \mathbb{R}$ . We then detail the discretization assumption as in (54) w.r.t.  $\tilde{f}$  (see also Section 2 for more details).

**Assumption 4** *For every  $n \in \mathbb{N}$  and  $\tilde{f} \in \mathcal{H}_k(\mathcal{S} \times \mathcal{A} \times \mathcal{I})$  there exists a discretization  $\mathcal{D}_n(\mathcal{S} \times \mathcal{A})$  of  $\mathcal{S} \times \mathcal{A}$  such that  $\tilde{f}(s, a, i) - \tilde{f}([s, a]_n, i) \leq \frac{1}{\sqrt{n}}$ , where  $[s, a]_n = \arg \min_{(s', a') \in \mathcal{D}_n(\mathcal{S} \times \mathcal{A})} \|(s, a) - (s', a')\|_2$ ,  $i \in \mathcal{I}$ , and  $|\mathcal{D}_n(\mathcal{S} \times \mathcal{A})| \leq CB^p n^{p/2}$  ( $|\mathcal{D}_n(\mathcal{S} \times \mathcal{A} \times \mathcal{I})| \leq CB^p n^{p/2} d$ ) for  $C$  being independent of  $n$  and  $B$ , and  $\mathcal{S} \times \mathcal{A} \subset \mathbb{R}^p$ .*

Assumption 4 allows us to provide bounds for  $\|f(s, a) - \mu_n(s, a)\|_2$  for all  $(s, a) \in \mathcal{S}$  using Lemma 3. Note that Assumption 4 does not discretize the modified state-action space ( $\bar{\mathcal{X}} = \mathcal{S} \times \mathcal{A} \times \{1, 2, \dots, d\}$ ) but instead discretizes  $\mathcal{S} \times \mathcal{A}$  for each  $i \in \mathcal{I}$ . Hence,  $|\mathcal{D}_n(\mathcal{S} \times \mathcal{A} \times \mathcal{I})| \leq CB^p n^{p/2} d$ , and  $\beta_n(\delta)$  will change accordingly. We describe the following lemma detailing the same.

**Lemma 5** *Under Assumption 4 with  $\beta_n(\delta)$  as in Lemma 3 and training a Gaussian process model on observations up to iteration  $n$  ( $\{s_1, \dots, s_n\}$ ) and their corresponding inputs ( $\{(s_0, a_0), \dots, (s_{n-1}, a_{n-1})\}$ ), it holds with probability at least  $1 - \delta$ ,*

$$\|f(s, a) - \mu_n(s, a)\|_2 \leq \beta_n(\delta)\sqrt{d}\|\sigma_n([s, a]_n)\|_2 + \frac{2d}{\sqrt{n}},$$

*uniformly for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$  and  $[s, a]_n = \arg \min_{(s', a') \in \mathcal{D}_n(\mathcal{S} \times \mathcal{A})} \|(s, a) - (s', a')\|_2$ .*

**Proof** For any  $(s, a) \in \mathcal{S} \times \mathcal{A}$ ,

$$\begin{aligned} & \|f(s, a) - \mu_n(s, a)\|_2 \\ &= \sqrt{\sum_{i=1}^d (\tilde{f}(s, a, i) - \mu_n(s, a, i))^2} \end{aligned} \quad (12)$$

$$\begin{aligned} &= \sqrt{\sum_{i=1}^d |\tilde{f}(s, a, i) - \tilde{f}([s, a]_n, i) + \tilde{f}([s, a]_n, i) - \mu_n([s, a]_n, i) + \mu_n([s, a]_n, i) - \mu_n(s, a, i)|^2} \\ &\leq \sum_{i=1}^d \left( |\tilde{f}(s, a, i) - \tilde{f}([s, a]_n, i)| + |\tilde{f}([s, a]_n, i) - \mu_n([s, a]_n, i)| + |\mu_n([s, a]_n, i) - \mu_n(s, a, i)| \right) \end{aligned} \quad (13)$$

$$\leq \left( \sum_{i=1}^d (|\tilde{f}([s, a]_n, i) - \mu_n([s, a]_n, i)|) \right) + \frac{2d}{\sqrt{n}} \quad (14)$$

$$\leq \beta_n(\delta) \left( \sum_{i=1}^d (\sigma_n([s, a]_n, i)) \right) + \frac{2d}{\sqrt{n}} \quad (15)$$

$$\leq \beta_n(\delta) \sqrt{d} \sqrt{\sum_{i=1}^d (\sigma_n([s, a]_n, i))^2} + \frac{2d}{\sqrt{n}} \quad (16)$$

$$\leq \beta_n(\delta) \sqrt{d} \|\sigma_n([s, a]_n)\|_2 + \frac{2d}{\sqrt{n}}. \quad (17)$$

In Equation (13), Equation (16) we use  $\|x\|_2 \leq \|x\|_1 \leq \sqrt{d}\|x\|_2$ . And Equation (14) and Equation (15) follow from Assumption 4 (since  $\tilde{f}, \mu_n \in \mathcal{H}_k(\mathcal{S} \times \mathcal{A} \times \mathcal{I})$ ) and Lemma 3, respectively. ■

## B.4 Sample Complexity Guarantees

Our objective is to obtain a uniform upper bound on the model precision  $\|\mu_n(s, a) - f(s, a)\|_2$  for all state-action pairs  $(s, a)$  while accounting for the errors induced by discretization. Here,  $\mu_n(\cdot, \cdot)$  is obtained from Algorithm 1. We achieve this by using Lemma 5 to obtain a bound in terms of maximum information gain (Equation (18)).

To characterize the precision of the learned model, we use the maximum information gain (52)

$$\Gamma_n(\bar{\mathcal{X}}) = \max_{x_1, \dots, x_n \in \bar{\mathcal{X}}} 0.5 \log \det(I_n + \lambda^{-1} K_n), \quad (18)$$

a kernel-dependent quantity that is frequently used in GP optimization. For many commonly used kernels,  $\Gamma_n$  is sublinear in  $n$ , which implies that the predictive uncertainties are shrinking sufficiently fast, and thus  $\hat{f}_n$  is capable of generalizing well across the entire domain. This is formalized in the following lemma.



**Lemma 6** For  $\beta_n(\delta)$  set as in Lemma 3 and  $\mathcal{I}_d$  denoting  $\{1, 2, \dots, d\}$ , the MVR algorithm (Algorithm 1) outputs the dynamics estimate  $\hat{f}_n(\cdot, \cdot) = \mu_n(\cdot, \cdot)$  such that the following holds uniformly for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$  with probability at least  $1 - \delta$ ,

$$\|\mu_n(s, a) - f(s, a)\|_2 \leq \mathcal{O}\left(\frac{\beta_n(\delta)2ed}{\sqrt{n}}\sqrt{\Gamma_{nd}(\mathcal{S} \times \mathcal{A} \times \mathcal{I}_d)}\right).$$

The preceding lemma asserts that we can effectively estimate the unknown dynamics by utilizing the pure exploration procedure and that the error in the model reduces as we increase the number of samples. In the subsequent section, we leverage this finding to establish the minimum number of samples needed to obtain a distributionally robust policy that is close to optimal.

**Proof** From Lemma 5, it holds that with probability at least  $1 - \delta$  uniformly for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ :

$$\begin{aligned} \|\mu_n(s, a) - f(s, a)\|_2 &\leq \beta_n(\delta)\sqrt{d}\|\sigma_n([s, a]_n)\|_2 + \frac{2d}{\sqrt{n}} \\ &\leq \beta_n(\delta)\sqrt{d}\max_{(s,a) \in \mathcal{S} \times \mathcal{A}} \|\sigma_n(s, a)\|_2 + \frac{2d}{\sqrt{n}} \\ &\leq \beta_n(\delta)\sqrt{d}\|\sigma_n(s_n, a_n)\|_2 + \frac{2d}{\sqrt{n}} \\ &\leq \frac{2d}{\sqrt{n}} + \frac{\beta_n(\delta)}{n}\sqrt{d}\sum_{j=1}^n \|\sigma_j(s_n, a_n)\|_2 \\ &\leq \frac{2d}{\sqrt{n}} + \frac{\beta_n(\delta)}{n}\sqrt{d}\sum_{j=1}^n \|\sigma_j(s_j, a_j)\|_2 \end{aligned} \tag{19}$$

$$\begin{aligned} &\leq \frac{\beta_n(\delta)}{\sqrt{n}}\sqrt{d}\sqrt{\sum_{j=1}^n \|\sigma_j(s_j, a_j)\|_2^2} + \frac{2d}{\sqrt{n}} \\ &\leq \frac{\beta_n(\delta)2ed}{\sqrt{n}}\sqrt{\Gamma_{nd}(\mathcal{S} \times \mathcal{A} \times \mathcal{I}_d)} + \frac{2d}{\sqrt{n}} \end{aligned} \tag{20}$$

$$= \mathcal{O}\left(\frac{\beta_n(\delta)2ed}{\sqrt{n}}\sqrt{\Gamma_{nd}(\mathcal{S} \times \mathcal{A} \times \mathcal{I}_d)}\right). \tag{21}$$

Here, Equation (19) follows from the decision rule in line-4 of Algorithm 1 and Equation (20) is obtained using standard bound for the sum of variances in the case of multi-output GPs from (18, Lemma-7) and (11, Lemma-11).  $\blacksquare$

## Appendix C. Proof Outline

We begin by defining the robust Bellman operator (27) in terms of the robust state-action value function  $Q_{\pi, f}^R$  as follows:

$$Q_{\pi, f}^R(s, a) = r(s, a) + \gamma \inf_{D(p||P_f(s, a)) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\pi, f}^R(s') \right]. \tag{22}$$

**Step (i):** The first step is to bound the approximation error of policy  $\hat{\pi}_n$  (i.e., the left-hand side of Equation (4)) by the sum of two error terms:  $|V_{\hat{\pi}_N, f}^R(s) - V_{\hat{\pi}_N, \hat{f}_N}^R(s)|$  and  $|V_{\hat{\pi}_N, \hat{f}_N}^R(s) - V_{\pi^*, f}^R(s)|$ . Utilizing the robust Bellman Equation (22), bounding such errors boils down to bounding differences of the form:

$$\max_s \left| \inf_{\text{KL}(p||P_f(s, \hat{\pi}_N(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_N, f}^R(s') \right] - \inf_{\text{KL}(p||P_{\hat{f}_N}(s, \hat{\pi}_N(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_N, f}^R(s') \right] \right|. \quad (23)$$

where  $P_f(s, a)$  denotes the Gaussian transition distribution with mean  $f(s, a)$  and covariance  $\sigma^2 I$ .

**Step (ii):** The major challenge of bounding Equation (23) lies in the inner infinite-dimensional minimization problems over distributions. To overcome this, we can reformulate such problems into single-dimensional ones using duality (26; 62; 41) according to the following lemma.

**Lemma 7** (Variant of (26)) *For random variable  $X$  and function  $V$  satisfying that  $V(X)$  has a finite Moment Generating function, it holds for all  $\rho > 0$ :*

$$\inf_{P: \text{KL}(P||P_0) \leq \rho} \mathbb{E}_{X \sim P} [V(X)] = \sup_{\alpha \geq 0} \{-\alpha \log(\mathbb{E}_{X \sim P_0} [e^{-\frac{V(X)}{\alpha}}]) - \alpha \rho\}. \quad (24)$$

Let  $H(V, P) := \sup_{\alpha \geq 0} \{-\alpha \log(\mathbb{E}_{X \sim P} [e^{-\frac{V(X)}{\alpha}}]) - \alpha \rho\}$ . Thus, applying Lemma 7, we rewrite Equation (23) as the difference of two single-dimensional convex optimization problems with expectations over  $P_f$  and  $P_{\hat{f}_N}$ , respectively:

$$\begin{aligned} & \max_s \left| H(V_{\hat{\pi}_N, f}^R, P_f(s, \hat{\pi}_N(s))) - H(V_{\hat{\pi}_N, f}^R, P_{\hat{f}_N}(s, \hat{\pi}_N(s))) \right| \\ & \leq \max_{V(\cdot) \in \mathcal{V}} \max_{s, a} \left| H(V, P_f(s, a)) - H(V, P_{\hat{f}_N}(s, a)) \right| \\ & \leq \max_{V(\cdot) \in \mathcal{V}} \max_{s, a} \max_{\alpha \in [\underline{\alpha}, \bar{\alpha}]} c \left| \mathbb{E}_{s' \sim P_f(s, a)} \left[ e^{-\frac{V(s')}{\alpha}} \right] - \mathbb{E}_{s' \sim P_{\hat{f}_N}(s, a)} \left[ e^{-\frac{V(s')}{\alpha}} \right] \right|, \end{aligned} \quad (25)$$

where  $c, \underline{\alpha}, \bar{\alpha} > 0$  are constants,  $\mathcal{V}$  denotes the value functional space, and the last inequality holds due to certain structural properties of the single-dimensional optimization problem in the RHS of Equation (24).

**Step (iii):** Finally, we bound Equation (25) using the difference between the estimated model  $\hat{f}_N$  and the true  $f$ , which is characterized by Lemma 6, in Appendix D. Moreover, to address the outer maximum over all value functions, states, and actions, we incorporate a covering number argument.

**Other uncertainty sets:** We further obtain the statistical sample complexities for  $\chi^2$  distance and TV distance uncertainty sets. We note that the analysis follows similar steps as the ones of Theorem 1. The major difference lies in incorporating and handling the dual forms of  $\chi^2$ /TV uncertainty sets in our analysis which differ from the one of Lemma 7. For  $\chi^2$  uncertainty set, we utilize the dual formulation that appears in (22), while for TV uncertainty sets we follow the approach of (59). As before, we can upper bound Lemma 23 via covering number arguments and the distance between the nominal transition dynamics  $f$  and the learned transition dynamics  $\hat{f}_N$  by using Theorem 6. Below, we outline the statistical sample complexity in the case of  $\chi^2$  and TV uncertainty sets in Theorems 8 and 9, respectively.

**Proposition 8** (*Sample Complexity of MVR under  $\chi^2$  uncertainty set*) Under the setup of Theorem 1 with uncertainty set defined w.r.t.  $\chi^2$  distance, it holds that  $\max_s |V_{\hat{\pi}_N, f}^R(s) - V_{\pi^*, f}^R(s)| \leq \epsilon$  with probability at least  $1 - \delta$  for any  $N$  such that

$$N = \mathcal{O}\left(\left(\frac{1 + 2\rho}{\sqrt{1 + 2\rho} - 1}\right)^4 \frac{\gamma^4 \beta_N^2(\delta) d^2 \Gamma_{Nd}}{(1 - \gamma)^8 \epsilon^4}\right). \quad (26)$$

**Proposition 9** (*Sample Complexity of MVR under TV uncertainty set*)

Under the setup of Theorem 1 with uncertainty set defined w.r.t. TV distance, it holds that  $\max_s |V_{\hat{\pi}_N, f}^R(s) - V_{\pi^*, f}^R(s)| \leq \epsilon$  with probability at least  $1 - \delta$  for any  $N$  such that

$$N = \mathcal{O}\left(\frac{(2 + \rho)^2 \gamma^2 \beta_N^2(\delta) d^2 \Gamma_{Nd}}{\rho^2 (1 - \gamma)^4 \epsilon^2}\right). \quad (27)$$

We relegate the proofs of Theorems 8 and 9 to Appendices E.1 and E.2. In comparison to the exponential dependence on  $\frac{1}{1-\gamma}$  for KL uncertainty set in Theorem 1, we note that for both  $\chi^2$ /TV uncertainty sets, we obtain *polynomial* dependence on  $\frac{1}{1-\gamma}$ . In the context of the TV uncertainty set, the dependency on  $\epsilon$  in Theorem 9 remains consistent with the finite state case ((41)). However, in the  $\chi^2$  case, the bound presented in Theorem 8 exhibits a worse dependence on  $\epsilon$  compared to the result derived in (41). This difference arises because we refrain from utilizing the same dual reformulation lemmas from (27), as they are applicable exclusively to finite state-action settings. Improving these rates is an interesting direction for future work.

## Appendix D. Sample Complexity Bounds for KL Uncertainty Sets

**Theorem 10** (*Sample Complexity of MVR under KL uncertainty set*) Consider a robust MDP with nominal transition dynamics  $f$  satisfying the regularity assumptions from Section 2 and with uncertainty set defined as in Equation (2) w.r.t. KL divergence. For  $\pi^*$  denoting the robust optimal policy w.r.t. nominal transition dynamics  $f$  and  $\hat{\pi}_N$  denoting the robust optimal policy w.r.t. learned nominal transition dynamics  $\hat{f}_N$  via MVR (Algorithm 1), and  $\delta \in (0, 1)$ ,  $\epsilon \in (0, \frac{1}{1-\gamma})$ , it holds that  $\max_s |V_{\hat{\pi}_N, f}^R(s) - V_{\pi^*, f}^R(s)| \leq \epsilon$  with probability at least  $1 - \delta$  for any  $N$  such that

$$N = \mathcal{O}\left(e^{\frac{2-\gamma}{(1-\gamma)\alpha_{\text{kl}}}} \frac{\gamma^2 \beta_N^2(\delta) d^2 \Gamma_N d}{(1-\gamma)^4 \rho^2 \epsilon^2}\right). \quad (28)$$

**Proof Step (i):** As detailed in the proof outline of Section 4, in order to bound  $|V_{\hat{\pi}_N, f}^R(s) - V_{\pi^*, f}^R(s)|$ , we begin by adding and subtracting  $V_{\hat{\pi}_n, \hat{f}_n}^R(s)$  which is the robust value function w.r.t. the nominal transition dynamics  $\hat{f}_n$  and its corresponding optimal policy  $\hat{\pi}_n$ . Then, we split the difference into two terms as follows:

$$|V_{\hat{\pi}_N, f}^R(s) - V_{\pi^*, f}^R(s)| = \underbrace{|V_{\hat{\pi}_N, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s)|}_{(i)} + \underbrace{|V_{\hat{\pi}_n, \hat{f}_n}^R(s) - V_{\pi^*, f}^R(s)|}_{(ii)}. \quad (29)$$

In order to not disturb the flow of the proof we bound (i) and (ii) separately Lemma 11 and Lemma 12 respectively. From Lemma 11, we obtain that

$$\begin{aligned} (i) &\leq \max_s \left| V_{\hat{\pi}_N, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s) \right| \\ &\leq \frac{\gamma}{1-\gamma} \max_s \left| \inf_{\text{KL}(p||P_{f(s, \hat{\pi}_n(s))}) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_N, f}^R(s') \right] - \inf_{\text{KL}(p||P_{\hat{f}_n(s, \hat{\pi}_n(s))}) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, \hat{f}_n}^R(s') \right] \right|. \end{aligned} \quad (30)$$

And from Lemma 12, we obtain that

$$\begin{aligned} (ii) &\leq \max_s \left| V_{\hat{\pi}_n, \hat{f}_n}^R(s) - V_{\pi^*, f}^R(s) \right| \\ &\leq \frac{\gamma}{1-\gamma} \max_s \left| \inf_{\text{KL}(p||P_{\hat{f}_n(s, \hat{\pi}_n(s))}) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\pi^*, f}^R(s') \right] - \inf_{\text{KL}(p||P_{f(s, \hat{\pi}_n(s))}) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\pi^*, f}^R(s') \right] \right|. \end{aligned} \quad (31)$$

Note that both these terms in Equations (30) and (31) are of the form mentioned in the **Step (i)** of Section 4.

**Step (ii):** Next, corresponding to **Step (ii)** of the proof outline in Section 4, we use Lemma 7 to bound Equations (30) and (31). Denote  $M := \frac{1}{1-\gamma} \geq \max_s V_{\pi^*}^R(s)$  for convenience. Using Equation (30) and Lemma 14 (internally using Lemma 7), conditioned on the event of Lemma 14 holding true, it holds that

$$\begin{aligned} (i) &\leq \max_s \left| V_{\hat{\pi}_N, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s) \right| \\ &\leq \frac{1}{1-\gamma} \max_s \left| \gamma \inf_{\text{KL}(p||P_{f(s, \hat{\pi}_n(s))}) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_N, f}^R(s') \right] - \gamma \inf_{\text{KL}(p||P_{\hat{f}_n(s, \hat{\pi}_n(s))}) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, \hat{f}_n}^R(s') \right] \right| \end{aligned}$$

$$\leq \max_{s,a} \left( 2\gamma \frac{M^2}{\rho} e^{\frac{M}{\underline{\alpha}}} \max_{\alpha \in [\underline{\alpha}, \frac{M}{\rho}]} \left| \mathbb{E}_{s' \sim P_{\hat{f}_n}(s,a)} \left[ e^{-\frac{V_{\hat{\pi}_n, f}^R(s')}{\alpha}} \right] - \mathbb{E}_{s' \sim P_f(s,a)} \left[ e^{-\frac{V_{\pi^*, f}^R(s')}{\alpha}} \right] \right| \right). \quad (32)$$

$$\leq \max_{V(\cdot) \in \mathcal{V}} \max_{s,a} \left( 2\gamma \frac{M^2}{\rho} e^{\frac{M}{\underline{\alpha}}} \max_{\alpha \in [\underline{\alpha}, \frac{M}{\rho}]} \left| \mathbb{E}_{s' \sim P_{\hat{f}_n}(s,a)} \left[ e^{-\frac{V(s')}{\alpha}} \right] - \mathbb{E}_{s' \sim P_f(s,a)} \left[ e^{-\frac{V(s')}{\alpha}} \right] \right| \right). \quad (33)$$

We can bound (ii) similarly.

$$(ii) \leq \max_s \left| V_{\hat{\pi}_n, \hat{f}_n}^R(s) - V_{\pi^*, f}^R(s) \right| \quad (34)$$

$$\leq \max_{V(\cdot) \in \mathcal{V}} \max_{s,a} \left( 2\gamma \frac{M^2}{\rho} e^{\frac{M}{\underline{\alpha}}} \max_{\alpha \in [\underline{\alpha}, \frac{M}{\rho}]} \left| \mathbb{E}_{s' \sim P_{\hat{f}_n}(s,a)} \left[ e^{-\frac{V(s')}{\alpha}} \right] - \mathbb{E}_{s' \sim P_f(s,a)} \left[ e^{-\frac{V(s')}{\alpha}} \right] \right| \right). \quad (35)$$

**Step (iii):** Next, we want to utilize the learning error bound (Equation (21)) that bounds the difference between the means of true nominal transition dynamics  $P_f$  and learned nominal transition dynamics  $P_{\hat{f}_n}$  to bound Equations (33) and (35).

We begin by bounding the difference  $\left| \mathbb{E}_{s' \sim P_{\hat{f}_n}(s,a)} \left[ e^{-\frac{V(s')}{\alpha}} \right] - \mathbb{E}_{s' \sim P_f(s,a)} \left[ e^{-\frac{V(s')}{\alpha}} \right] \right|$ , by the difference in means of  $P_f$  and  $P_{\hat{f}_n}$  in Lemma 15. Since Equation (33) has a max over all value functions, we introduce a covering number argument in Lemma 17 to reform it to a max over the functions in the  $\zeta$ -covering set. We then use Lemma 15 to obtain bounds in terms of maximum information gain  $\Gamma_{Nd}$  (Equation (18)) and  $\zeta$ . Further details regarding the covering number argument are deferred to Lemma 17. Then, we apply the result of Lemma 17 with  $\zeta = 1$  (defined in Lemma 17) on Equation (33). Then, it holds that

$$(i) \leq \max_s \left| V_{\hat{\pi}_n, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s) \right| = \mathcal{O} \left( 2 \frac{M^2}{\rho} e^{\frac{M}{\alpha_{kl}}} e^{\frac{1}{\alpha_{kl}}} \frac{\beta_n(\delta) \sqrt{2ed^2 \Gamma_{nd}}}{\sigma \sqrt{n}} \right), \quad (36)$$

where  $\alpha_{kl}$  is a problem-dependent constant denoting the minimum value of  $\underline{\alpha}$  defined in Lemma 14. A similar constant also appears in the sample complexity bounds provided in (41; 62). Note that  $\beta_n$ , which appears in Lemma 2, has a logarithmic dependence on  $n$ . Similarly, from Equation (35) and Lemmas 15,17, we obtain

$$(ii) \leq \max_s \left| V_{\hat{\pi}_n, \hat{f}_n}^R(s) - V_{\pi^*, f}^R(s) \right| = \mathcal{O} \left( 2\gamma \frac{M^2}{\rho} e^{\frac{M}{\alpha_{kl}}} e^{\frac{1}{\alpha_{kl}}} \frac{\beta_n(\delta) \sqrt{2ed^2 \Gamma_{nd}}}{\sigma \sqrt{n}} \right). \quad (37)$$

Note that we want to bound  $V_{\hat{\pi}_n, f}^R(s) - V_{\pi^*, f}^R(s) = (i) + (ii)$  over all  $s \in \mathcal{S}$ . Using  $\max_s \left| V_{\hat{\pi}_n, f}^R(s) - V_{\pi^*, f}^R(s) \right| \leq \max_s \left| V_{\hat{\pi}_n, \hat{f}_n}^R(s) - V_{\pi^*, f}^R(s) \right| + \max_s \left| V_{\hat{\pi}_n, \hat{f}_n}^R(s) - V_{\hat{\pi}_n, f}^R(s) \right|$  and substituting  $M$  by  $1/(1-\gamma)$ , we obtain from Equation (36) and Equation (37)

$$\max_s \left| V_{\hat{\pi}_n, f}^R(s) - V_{\pi^*, f}^R(s) \right| = \mathcal{O} \left( \gamma e^{\frac{1}{(1-\gamma)\alpha_{kl}}} e^{\frac{1}{\alpha_{kl}}} \frac{\beta_n(\delta) d \sqrt{2e \Gamma_{nd}}}{(1-\gamma)^2 \rho \sigma \sqrt{n}} \right).$$

Finally, to ensure that  $\max_s \left| V_{\hat{\pi}_n, f}^R(s) - V_{\pi^*, f}^R(s) \right| \leq \epsilon$ , it suffices to have

$$\max_s \left| V_{\hat{\pi}_n, f}^R(s) - V_{\pi^*, f}^R(s) \right| = \mathcal{O} \left( \gamma e^{\frac{1}{(1-\gamma)\alpha_{kl}}} e^{\frac{1}{\alpha_{kl}}} \frac{\beta_n(\delta) d \sqrt{2e \Gamma_{nd}}}{(1-\gamma)^2 \rho \sigma \sqrt{n}} \right) = \epsilon.$$

By inverting the previously obtained result, we arrive at

$$n = \mathcal{O}\left(e^{\frac{2}{(1-\gamma)\alpha_{kl}}} e^{\frac{2}{\alpha_{kl}}} \frac{\gamma^2 \beta_n^2(\delta) d^2 \Gamma_{nd}}{(1-\gamma)^4 \rho^2 \epsilon^2}\right).$$

■

**Lemma 11** (*Simplification using robust Bellman equation*) Denote  $(i) := \left| V_{\hat{\pi}_n, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s) \right|$  for  $V_{\hat{\pi}_n, f}^R$  being the robust value function of policy  $\hat{\pi}_n$  w.r.t. true nominal transition dynamics  $f$  and  $V_{\hat{\pi}_n, \hat{f}_n}^R$  being the robust value function of policy  $\hat{\pi}_n$  w.r.t. learned nominal transition dynamics  $\hat{f}$ . Then the following holds,

$$\begin{aligned} (i) &= \left| V_{\hat{\pi}_n, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s) \right| \\ &\leq \max_s \left| V_{\hat{\pi}_n, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s) \right| \\ &\leq \frac{\gamma}{1-\gamma} \max_s \left| \inf_{D(p||P_f(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] - \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] \right|. \end{aligned} \quad (38)$$

**Proof** Since both the quantities are w.r.t. the same policy, using the definition of the robust  $Q$ -function and the robust Bellman equation (see Equation (22)), we obtain:

$$\begin{aligned} (i) &= |V_{\hat{\pi}_n, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s)| \quad (39) \\ &= |Q_{\hat{\pi}_n, f}^R(s, \hat{\pi}_n(s)) - Q_{\hat{\pi}_n, \hat{f}_n}^R(s, \hat{\pi}_n(s))| \\ &= |r(s, \hat{\pi}_n(s)) - r(s, \hat{\pi}_n(s))| \\ &\quad + \gamma \inf_{D(p||P_f(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] - \gamma \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, \hat{f}_n}^R(s') \right] \\ &= |\gamma \inf_{D(p||P_f(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] - \gamma \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, \hat{f}_n}^R(s') \right]| \quad (40) \end{aligned}$$

Adding and subtracting  $\gamma \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right]$  to Equation (40), we obtain the following two terms:

$$\begin{aligned} (i_a) &= \left| \gamma \inf_{D(p||P_f(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] - \gamma \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] \right|, \\ (i_b) &= \left| \gamma \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] - \gamma \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, \hat{f}_n}^R(s') \right] \right|. \end{aligned}$$

Now, we use Lemma 13 to bound  $(i_b)$ . We have:

$$(i_b) = \left| \gamma \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] - \gamma \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, \hat{f}_n}^R(s') \right] \right|$$

$$\stackrel{\text{Lemma13}}{\leq} \gamma \max_s \left| V_{\hat{\pi}_n, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s) \right| \quad (\text{Lemma13}). \quad (41)$$

Plugging Equation (41) into Equation (39) and using the fact that  $(i) = (i_a) + (i_b)$ , we have

$$(i) = |V_{\hat{\pi}_n, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s)| \quad (42)$$

$$\begin{aligned} &\leq (i_a) + \gamma \max_s \left| V_{\hat{\pi}_n, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s) \right| \\ &= \left| \gamma \inf_{D(p||P_f(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] - \gamma \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] \right| \\ &\quad + \gamma \max_s \left| V_{\hat{\pi}_n, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s) \right|. \end{aligned} \quad (43)$$

Taking maximum over states in Equation (42) and Equation (43) we have

$$\begin{aligned} &\max_s \left| V_{\hat{\pi}_n, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s) \right| \\ &\leq \max_s \left| \gamma \inf_{D(p||P_f(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] - \gamma \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] \right| \\ &\quad + \gamma \max_s \left| V_{\hat{\pi}_n, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s) \right|. \end{aligned}$$

Moving  $\gamma \max_s \left| V_{\hat{\pi}_n, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s) \right|$  to the LHS and dividing  $(1 - \gamma)$  on both sides, it holds that

$$\begin{aligned} (i) &\leq \max_s \left| V_{\hat{\pi}_n, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s) \right| \\ &\leq \frac{\gamma}{1 - \gamma} \max_s \left| \inf_{D(p||P_f(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] - \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] \right|. \end{aligned} \quad (44)$$

■

**Lemma 12** (*Simplification using robust Bellman equation*) Denote  $(ii) := \left| V_{\hat{\pi}_n, \hat{f}_n}^R(s) - V_{\pi^*, f}^R(s) \right|$  for  $V_{\hat{\pi}_n, \hat{f}_n}^R$  being the robust value function of policy  $\hat{\pi}_n$  w.r.t. learned nominal transition dynamics  $\hat{f}_n$  and  $V_{\pi^*, f}^R$  being the robust value function of policy  $\pi^*$  w.r.t. true nominal transition dynamics  $f$ . Then the following holds,

$$\begin{aligned} (ii) &= \left| V_{\hat{\pi}_n, \hat{f}_n}^R(s) - V_{\pi^*, f}^R(s) \right| \\ &\leq \max_s \left| V_{\hat{\pi}_n, \hat{f}_n}^R(s) - V_{\pi^*, f}^R(s) \right| \\ &\leq \frac{\gamma}{1 - \gamma} \max_s \left| \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\pi^*, f}^R(s') \right] - \inf_{D(p||P_f(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\pi^*, f}^R(s') \right] \right|. \end{aligned} \quad (45)$$

**Proof** We first note that  $Q_{\pi^*,f}^R(s, \hat{\pi}_n(s)) \leq Q_{\pi^*,f}^R(s, \pi^*(s))$  as  $\pi^*$  is the robust optimal policy for the nominal transition dynamics  $f$  (see Equation (3)). As a result, we have

$$\begin{aligned}
(ii) &= |V_{\hat{\pi}_n, \hat{f}_n}^R(s) - V_{\pi^*, f}^R(s)| \tag{46} \\
&= |Q_{\hat{\pi}_n, \hat{f}_n}^R(s, \hat{\pi}_n(s)) - Q_{\pi^*, f}^R(s, \pi^*(s))| \\
&\leq |Q_{\hat{\pi}_n, \hat{f}_n}^R(s, \hat{\pi}_n(s)) - Q_{\pi^*, f}^R(s, \hat{\pi}_n(s))| \\
&= |r(s, \hat{\pi}_n(s)) - r(s, \pi^*(s))| \\
&\quad + \gamma \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} [V_{\hat{\pi}_n, \hat{f}_n}^R(s')] - \gamma \inf_{D(p||P_f(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} [V_{\pi^*, f}^R(s')] | \\
&= |\gamma \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} [V_{\hat{\pi}_n, \hat{f}_n}^R(s')] - \gamma \inf_{D(p||P_f(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} [V_{\pi^*, f}^R(s')] | \tag{47}
\end{aligned}$$

Adding and subtracting  $\gamma \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} [V_{\pi^*, f}^R(s')]$  to Equation (47), we obtain the following two terms:

$$\begin{aligned}
(ii_a) &= |\gamma \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} [V_{\hat{\pi}_n, \hat{f}_n}^R(s')] - \gamma \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} [V_{\pi^*, f}^R(s')] |, \\
(ii_b) &= |\gamma \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} [V_{\pi^*, f}^R(s')] - \gamma \inf_{D(p||P_f(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} [V_{\pi^*, f}^R(s')] |.
\end{aligned}$$

Now, we use Lemma 13 to bound  $(ii_a)$ . We have:

$$\begin{aligned}
(ii_a) &= |\gamma \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} [V_{\hat{\pi}_n, \hat{f}_n}^R(s')] - \gamma \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} [V_{\pi^*, f}^R(s')] | \\
&\leq \gamma \max_s |V_{\pi^*, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s)|. \tag{48}
\end{aligned}$$

Plugging Equation (48) into Equation (46) and using the fact that  $(ii) = (ii_a) + (ii_b)$ , we have

$$\begin{aligned}
(ii) &= |V_{\pi^*, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s)| \tag{49} \\
&\leq (ii_b) + \max_s |V_{\pi^*, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s)| \\
&= |\gamma \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} [V_{\pi^*, f}^R(s')] - \gamma \inf_{D(p||P_f(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} [V_{\pi^*, f}^R(s')] | \\
&\quad + \gamma \max_s |V_{\pi^*, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s)|. \tag{50}
\end{aligned}$$

Taking maximum over states in Equation (49) and Equation (50) and following similar steps as in Equation (44), we have

$$(ii) \leq \max_s |V_{\pi^*, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s)|$$



$$\begin{aligned}
&\leq \max_s \left| \gamma \inf_{D(p||P_f(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] - \gamma \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] \right| \\
&\quad + \gamma \max_s \left| V_{\hat{\pi}_n, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s) \right| \\
&\leq \frac{\gamma}{1 - \gamma} \max_s \left| \inf_{D(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] - \inf_{D(p||P_f(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] \right|.
\end{aligned} \tag{51}$$

■

**Lemma 13** (from (41, Lemma 1)) *Let  $V_1$  and  $V_2$  be two value functions from  $\mathcal{S} \rightarrow [0, 1/(1 - \gamma)]$ . Let  $D$  be any distance measure between probability distributions (e.g., KL-divergence,  $\chi^2$ -divergence, or variation distance defined in Equation (2)). The following inequality (1-Lipschitz w.r.t.  $V$ ) holds true*

$$\left| \inf_{D(p||P_{\hat{f}}(s, a)) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_1(s') \right] - \inf_{D(p||P_{\hat{f}}(s, a)) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_2(s') \right] \right| \leq \max_{s'} |V_2(s') - V_1(s')|.$$

**Proof** We want to bound

$$\left| \inf_{D(p||P_{\hat{f}}(s, a)) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_1(s') \right] - \inf_{D(p||P_{\hat{f}}(s, a)) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_2(s') \right] \right|.$$

Notice that

$$\begin{aligned}
&\inf_{D(p||P_{\hat{f}}(s, a)) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_1(s') \right] - \inf_{D(p||P_{\hat{f}}(s, a)) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_2(s') \right] \\
&= \inf_{D(p||P_{\hat{f}}(s, a)) \leq \rho} \sup_{D(p'||P_{\hat{f}}(s, a)) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_1(s') \right] - \mathbb{E}_{s' \sim p'} \left[ V_2(s') \right] \\
&\geq \inf_{D(p||P_{\hat{f}}(s, a)) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_1(s') \right] - \mathbb{E}_{s' \sim p} \left[ V_2(s') \right] \\
&= \inf_{D(p||P_{\hat{f}}(s, a)) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_1(s') - V_2(s') \right],
\end{aligned}$$

where the inequality follows from the property of supremum. By the definition of inf, for any  $\epsilon > 0$ , there exists some distribution  $q$  s.t.  $D(q||P_{\hat{f}}(s, a)) \leq \rho$  satisfying

$$\mathbb{E}_{s' \sim q} \left[ V_1(s') - V_2(s') \right] - \epsilon \leq \inf_{D(p||P_{\hat{f}}(s, a)) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_1(s') - V_2(s') \right].$$

Then, we have

$$\begin{aligned}
&\inf_{D(p||P_{\hat{f}}(s, a)) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_2(s') \right] - \inf_{D(p||P_{\hat{f}}(s, a)) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_1(s') \right] \\
&\leq - \inf_{D(p||P_{\hat{f}}(s, a)) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_1(s') - V_2(s') \right] \\
&\leq - \mathbb{E}_{s' \sim q} \left[ V_1(s') - V_2(s') \right] + \epsilon
\end{aligned}$$

$$\begin{aligned}
&\leq \mathbb{E}_{s' \sim q} [V_2(s') - V_1(s')] + \epsilon \\
&\leq \max_{s'} |V_2(s') - V_1(s')| + \epsilon.
\end{aligned} \tag{52}$$

Let  $\epsilon \rightarrow 0$ , we obtain one side of the desired bound.

One can similarly bound  $\inf_{D(p||P_{\hat{f}}(s,a)) \leq \rho} \mathbb{E}_{s' \sim p} [V_1(s')] - \inf_{D(p||P_{\hat{f}}(s,a)) \leq \rho} \mathbb{E}_{s' \sim p} [V_2(s')]$  by just interchanging  $V_1$  and  $V_2$  everywhere. Combining this argument with Equation (52), we obtain

$$\left| \inf_{D(p||P_{\hat{f}}(s,a)) \leq \rho} \mathbb{E}_{s' \sim p} [V_1(s')] - \inf_{D(p||P_{\hat{f}}(s,a)) \leq \rho} \mathbb{E}_{s' \sim p} [V_2(s')] \right| \leq \max_{s'} |V_2(s') - V_1(s')|.$$

■

**Lemma 14** (*Simplification using Lemma 7 reformulation*) For any value function  $V(\cdot) : \mathcal{S} \rightarrow [0, 1/(1 - \gamma)]$ , define the event  $\mathbf{E}$  as follows:

$$\begin{aligned}
\max_s \left| \inf_{KL(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] - \inf_{KL(p||P_f(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] \right| \leq \\
\max_{s,a} 2 \frac{M}{\rho} e^{\frac{M}{\underline{\alpha}}} \max_{\alpha \in [\underline{\alpha}, \bar{\alpha}]} \left| \mathbb{E}_{s' \sim P_{\hat{f}_n}(s,a)} [e^{\frac{-V(s')}{\alpha}}] - \mathbb{E}_{s' \sim P_f(s,a)} [e^{\frac{-V(s')}{\alpha}}] \right|.
\end{aligned}$$

Then, for any  $n > \{\max_{s,a} N'(\rho, P_f(s, a)), \max_{s,a} N''(\rho, P_f(s, a))\}$  where  $N'(\rho, P_f(s, a)) = \mathcal{O}\left(\frac{\beta_n^2(\delta) 2ed^2 \Gamma_{nd}}{(\kappa - e^{-\rho})^2}\right)$  and  $N''(\rho, P_f(s, a)) = \mathcal{O}\left(\frac{4M^2 e^{\frac{2M}{\underline{\alpha}}} \beta_n^2(\delta) 2ed^2 \Gamma_{nd}}{(\rho\tau)^2}\right)$  with  $\bar{\alpha} = \frac{M}{\rho}$ ,  $M = \frac{1}{1-\gamma}$ ,  $\kappa$  defined in Equation (68),  $\tau$  defined in Equation (71), and  $\underline{\alpha} = \alpha^*/2$  defined in Equation (57), the event  $\mathbf{E}$  holds true with probability at least  $1 - \delta$ .

**Proof** (A similar proof as in (62, Lemma-4)). First note that,

$$\begin{aligned}
\max_s \left| \inf_{KL(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] - \inf_{KL(p||P_f(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] \right| \leq \\
\max_{s,a} \left| \inf_{KL(p||P_{\hat{f}_n}(s,a)) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] - \inf_{KL(p||P_f(s,a)) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] \right|. \tag{53}
\end{aligned}$$

Recall (26, Theorem-1) for distributionally robust optimization with a random variable  $X$  and a random function  $H$ . One can rewrite an infinite-dimensional optimization problem as a scalar optimization problem:

$$\sup_{P: KL(p||P_0) \leq \rho} \mathbb{E}_{X \sim P} [H(X)] = \inf_{\alpha \geq 0} \{\alpha \log(\mathbb{E}_{X \sim P_0} [e^{\frac{H(X)}{\alpha}}]) + \alpha \rho\}. \tag{54}$$

For now, we focus on bounding  $\left| \inf_{KL(p||P_{\hat{f}_n}(s,a)) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] - \inf_{KL(p||P_f(s,a)) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] \right|$  for one particular  $(s, a)$ . For brevity, we write  $P_f(s, a)$  and  $P_{\hat{f}_n}(s, a)$  as  $P_f$  and  $P_{\hat{f}_n}$ , respectively. By Equation (54), we have

$$\inf_{P: KL(p||P_f) \leq \rho} \mathbb{E}_{s' \sim P} [V(s')] = \max_{\alpha \geq 0} \{-\alpha \log(\mathbb{E}_{s' \sim P_f} [e^{\frac{-V(s')}{\alpha}}]) - \alpha \rho\}, \tag{55}$$

$$\inf_{P:KL(p||\hat{P}_{\hat{f}_n})\leq\rho} \mathbb{E}_{s'\sim P}[V(s')] = \max_{\alpha\geq 0}\{-\alpha\log(\mathbb{E}_{s'\sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}]) - \alpha\rho\}. \quad (56)$$

For the finite state-action space setting, (62, Lemma-4) characterizes the property of the optimal  $\alpha^*$ . Following a similar proof strategy, we denote

$$\alpha^* = \arg \max_{\alpha\geq 0}\{-\alpha\log(\mathbb{E}_{s'\sim P_f}[e^{-\frac{V(s')}{\alpha}}]) - \alpha\rho\}, \quad (57)$$

and

$$\hat{\alpha}_n^* = \arg \max_{\alpha\geq 0}\{-\alpha\log(\mathbb{E}_{s'\sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}]) - \alpha\rho\}. \quad (58)$$

To ensure that  $\max_{\alpha\geq 0}\{-\alpha\log(\mathbb{E}_{s'\sim P_f}[e^{-\frac{V(s')}{\alpha}}]) - \alpha\rho\} - \max_{\alpha\geq 0}\{-\alpha\log(\mathbb{E}_{s'\sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}]) - \alpha\rho\}$  is small enough, we need to show that  $\alpha^*$  and  $\hat{\alpha}_n^*$  are close enough. For this, one considers two different cases,  $\alpha^* = 0$  and  $\alpha^* > 0$ .

**Case-1:** In Case-1, we investigate the conditions for  $\hat{\alpha}_n^* = 0$  given that  $\alpha^* = 0$ . According to (26, Proposition-2), for  $\alpha^* = 0$  to occur, the random variable  $Y := V(s')$  where  $s' \sim \mathcal{N}(f(s, a), \sigma^2 I)$  must satisfy three conditions namely, (i)  $Y$  must be bounded, (ii)  $Y$  must have finite mass at its essential infimum, and (iii) the finite mass at essential infimum should be greater than  $e^{-\rho}$ . So we want to verify whether these conditions hold true for  $\hat{Y}_n := V(s')$  where  $s' \sim \mathcal{N}(\hat{f}_n(s, a), \sigma^2 I)$  when  $Y$  satisfies these conditions.

We restate definition of the essential infimum for a real-valued random variable  $Y$ , denoted as  $\text{ESI}(Y)$ .

$$\text{ESI}(Y) = \sup\{t \in \mathbb{R} : \mathbb{P}\{Y < t\} = 0\}. \quad (59)$$

We first show that  $Y = V(s')$  where  $s' \sim \mathcal{N}(f(s, a), \sigma^2 I)$  and  $\hat{Y}_n = V(s')$  where  $s' \sim \mathcal{N}(\hat{f}_n(s, a), \sigma^2 I)$  have the same essential infimum. By the definition of  $\text{ESI}(Y)$ , for any  $\epsilon > 0$ , it holds that

$$\mathbb{P}\{\text{ESI}(Y) \leq Y < \text{ESI}(Y) + \epsilon\} > 0, \quad \mathbb{P}\{Y < \text{ESI}(Y)\} = 0. \quad (60)$$

It implies for  $Y = V(s')$  and  $s' \sim \mathcal{N}(f(s, a), \sigma^2 I)$  that

$$\mathbb{P}_{s' \sim \mathcal{N}(f(s, a), \sigma^2 I)}\{s' \in \mathbb{R}^d : \text{ESI}(Y) \leq Y = V(s') < \text{ESI}(Y) + \epsilon\} > 0, \quad (61)$$

$$\mathbb{P}_{s' \sim \mathcal{N}(f(s, a), \sigma^2 I)}\{s' \in \mathbb{R}^d : Y = V(s') < \text{ESI}(Y)\} = 0. \quad (62)$$

It further implies that, the set  $\{s' \in \mathbb{R}^d : \text{ESI}(Y) \leq V(s') < \text{ESI}(Y) + \epsilon\}$  must have a Lebesgue measure greater than 0 and  $\{s' \in \mathbb{R}^d : V(s') < \text{ESI}(Y)\}$  must have a Lebesgue measure equal to 0 since  $s' \sim \mathcal{N}(f(s, a), \sigma^2 I)$  is a continuous distribution.

Due to this fact that the set  $\{s' \in \mathbb{R}^d : \text{ESI}(Y) \leq V(s') < \text{ESI}(Y) + \epsilon\}$  has a Lebesgue measure greater than zero and noting that  $\mathcal{N}(\hat{f}_n(s, a), \sigma^2 I)$  is also a continuous distribution with the same support as of  $\mathcal{N}(f(s, a), \sigma^2 I)$  (i.e., the probability density function of  $\mathcal{N}(\hat{f}_n(s, a), \sigma^2 I)$  is positive whenever probability density function of  $\mathcal{N}(f(s, a), \sigma^2 I)$  is positive), it holds that

$$\mathbb{P}_{s' \sim \mathcal{N}(\hat{f}_n(s, a), \sigma^2 I)}\{s' \in \mathbb{R}^d : \text{ESI}(Y) \leq \hat{Y}_n = V(s') < \text{ESI}(Y) + \epsilon\} > 0. \quad (63)$$

A similar argument follows for

$$\mathbb{P}_{s' \sim \mathcal{N}(\hat{f}_n(s,a), \sigma^2 I)} \{s' \in \mathbb{R}^d : \hat{Y}_n = V(s') < \text{ESI}(Y)\} = 0. \quad (64)$$

In essence, Equations (63) and (64) imply,

$$\mathbb{P}\{\text{ESI}(Y) \leq \hat{Y}_n < \text{ESI}(Y) + \epsilon\} = 0, \quad \mathbb{P}\{\hat{Y}_n < \text{ESI}(Y)\} > 0.$$

Hence, from the definition of  $\text{ESI}(\cdot)$  in Equations (59) and (60), we have  $\text{ESI}(Y) = \text{ESI}(\hat{Y}_n)$ .

As a result, for  $\alpha^* = 0$  to occur and for  $Y = V(s')(s' \sim \mathcal{N}(f(s,a), \sigma^2 I))$  to have finite mass at the essential infimum (condition-(ii)), i.e.,  $\mathbb{P}\{Y = \text{ESI}(Y)\} > 0$ , it requires that

$$\mathbb{P}_{s' \sim \mathcal{N}(f(s,a), \sigma^2 I)} \{s' \in \mathbb{R}^d : Y = V(s') = \text{ESI}(Y)\} > 0.$$

This will further require that the set  $\{s' \in \mathbb{R}^d : Y = V(s') = \text{ESI}(Y)\}$  must have a Lebesgue measure greater than 0. Following a similar argument as to have obtained Equation (63) (the probability density function of  $\mathcal{N}(\hat{f}_n(s,a), \sigma^2 I)$  is positive whenever probability density function of  $\mathcal{N}(f(s,a), \sigma^2 I)$  is positive), the set  $\{s' \in \mathbb{R}^d : Y = V(s') = \text{ESI}(Y)\}$  having Lebesgue measure greater than 0, will imply

$$\mathbb{P}_{s' \sim \mathcal{N}(\hat{f}_n(s,a), \sigma^2 I)} \{s' \in \mathbb{R}^d : \hat{Y}_n = V(s') = \text{ESI}(Y)\} > 0, \quad (65)$$

and

$$\mathbb{P}\{\hat{Y}_n = \text{ESI}(Y)\} > 0 \quad (66)$$

Since  $\text{ESI}(Y) = \text{ESI}(\hat{Y}_n)$ , Equations (65) and (66) imply

$$\mathbb{P}\{\hat{Y}_n = \text{ESI}(\hat{Y}_n)\} > 0, \quad (67)$$

Hence, if  $\mathbb{P}\{Y = \text{ESI}(Y)\} > 0$  holds true, it also holds that  $\mathbb{P}\{\hat{Y}_n = \text{ESI}(\hat{Y}_n)\} > 0$ . This implies that whenever  $Y$  has a finite mass at its essential infimum,  $\hat{Y}_n$  also has finite mass at its essential infimum (condition-(ii) satisfied).

But, recall that according to (26, Proposition-2) for  $\alpha^* = 0$  to occur, the finite mass which  $Y$  has at its essential infimum should also be greater than  $e^{-\rho}$  (condition-(iii)). Hence, one has to check if  $Y$  satisfies

$$\mathbb{P}_{s' \sim \mathcal{N}(f(s,a), \sigma^2 I)} \{s' \in \mathbb{R}^d : Y = V(s') = \text{ESI}(Y)\} > e^{-\rho}, \quad (68)$$

what is the condition that  $Y_n$  satisfies

$$\mathbb{P}_{s' \sim \mathcal{N}(\hat{f}_n(s,a), \sigma^2 I)} \{s' \in \mathbb{R}^d : \hat{Y}_n = V(s') = \text{ESI}(\hat{Y}_n)\} > e^{-\rho},$$

so that  $\hat{\alpha}_n^* = 0$  whenever  $\alpha^* = 0$ . Denote  $\kappa := \mathbb{P}_{s' \sim \mathcal{N}(f(s,a), \sigma^2 I)} \{s' \in \mathbb{R}^d : Y = V(s') = \text{ESI}(Y)\}$ ,  $\kappa_n := \mathbb{P}_{s' \sim \mathcal{N}(\hat{f}_n(s,a), \sigma^2 I)} \{s' \in \mathbb{R}^d : \hat{Y}_n = V(s') = \text{ESI}(\hat{Y}_n)\}$ , and  $S_{min} := \{s' \in \mathbb{R}^d : V(s') = \text{ESI}(Y) = \text{ESI}(\hat{Y}_n)\}$ . If  $\kappa > e^{-\rho}$  and  $\kappa - \kappa_n \leq \frac{\kappa - e^{-\rho}}{2}$ , then it will hold that  $\kappa_n > e^{-\rho}$ .

$$|\kappa - \kappa_n| = \left| \int_{S_{min}} \frac{1}{\sqrt{(2\pi\sigma^2)^d}} (e^{-\frac{\|s' - f(s,a)\|^2}{\sigma^2}} - e^{-\frac{\|s' - \hat{f}_n(s,a)\|^2}{\sigma^2}}) dx \right|$$

$$\begin{aligned}
&\leq \int_{S_{min}} \frac{1}{\sqrt{(2\pi\sigma^2)^d}} \left| e^{-\frac{\|s'-f(s,a)\|^2}{\sigma^2}} - e^{-\frac{\|s'-\hat{f}_n(s,a)\|^2}{\sigma^2}} \right| dx \\
&\leq \int_{\mathbb{R}^d} \frac{1}{\sqrt{(2\pi\sigma^2)^d}} \left| e^{-\frac{\|s'-f(s,a)\|^2}{\sigma^2}} - e^{-\frac{\|s'-\hat{f}_n(s,a)\|^2}{\sigma^2}} \right| dx \\
&\leq \|f(s,a) - \hat{f}_n(s,a)\|_2 \quad (\text{Lemma15}) \\
&\leq \mathcal{O}\left(\frac{\beta_n(\delta)\sqrt{2ed^2\Gamma_{nd}}}{\sqrt{n}}\right),
\end{aligned}$$

We need  $\mathcal{O}\left(\frac{\beta_n(\delta)\sqrt{2ed^2\Gamma_{nd}}}{\sqrt{n}}\right) \leq \frac{\kappa - e^{-\rho}}{2}$ , which in turn requires  $n = \mathcal{O}\left(\frac{\beta_n^2(\delta)2ed^2\Gamma_{nd}}{(\frac{\kappa - e^{-\rho}}{2})^2}\right) = N'(\rho, P_f(s, a))$ . Hence, for  $n > \max_{s,a} N'(\rho, P_f(s, a))$  with probability at least  $1 - \delta$ , it holds that

$$\kappa_n > e^{-\rho},$$

for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$  whenever  $\kappa > e^{-\rho}$ , implying  $\alpha_n^* = 0$  whenever  $\alpha^* = 0$ .

**Case-2:** Consider the case of  $\alpha^* > 0$ . The idea is to bound both  $\alpha^*$  and  $\hat{\alpha}_n^*$  by a set  $[\underline{\alpha}, \bar{\alpha}]$  and bound  $\max_{\alpha \geq 0} \{(-\alpha \log(\mathbb{E}_{s' \sim P_f(s, \pi(s))}[e^{-\frac{V(s')}{\alpha}}]) - \alpha\rho) - (-\alpha \log(\mathbb{E}_{s' \sim P_f(s, \pi'(s))}[e^{-\frac{V(s')}{\alpha}}]) - \alpha\rho)\}$  for  $\alpha$  taking values within set  $[\underline{\alpha}, \bar{\alpha}]$ . We first provide an upper bound for  $\alpha^*$  as  $\frac{M}{\rho}$  where  $M = \frac{1}{1-\gamma}$  denoting the maximum value of  $V(s')$ .

$$\begin{aligned}
\max_{\alpha \geq 0} \{-\alpha \log(\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]) - \alpha\rho\} &\geq \lim_{\alpha \rightarrow 0} [-\alpha \log(\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]) - \alpha\rho] \\
&= \text{ESI}(V(s')|_{s' \sim P_f}) \quad (\text{Lemma16}) \\
&\geq 0.
\end{aligned} \tag{69}$$

Since  $\max_s V(s) \leq M$ , we have

$$-\alpha \log(\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]) - \alpha\rho \leq -\alpha \log(e^{-\frac{M}{\alpha}}) - \alpha\rho = M - \alpha\rho.$$

It implies for  $\alpha > \frac{M}{\rho}$  that

$$-\alpha \log(\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]) - \alpha\rho < 0. \tag{70}$$

By Equation (69), since  $\max_{\alpha \geq 0} \{-\alpha \log(\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]) - \alpha\rho\} \geq 0$ ,  $\arg \max_{\alpha \geq 0} \{-\alpha \log(\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]) - \alpha\rho\}$  cannot be greater than  $\frac{M}{\rho}$  due to Equation (70) holding for all  $\alpha > \frac{M}{\rho}$ . Hence, we have  $\alpha^* \leq \frac{M}{\rho}$ . A similar argument holds for  $\hat{\alpha}_n^*$  and it holds that  $\hat{\alpha}_n^* \leq \frac{M}{\rho}$ .

Denote  $\underline{\alpha} := \alpha^*/2$ ,  $\bar{\alpha} := \frac{M}{\rho}$ , and

$$\tau := \min \left\{ \underline{\alpha} \log(\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\underline{\alpha}}}] + \alpha\rho, \bar{\alpha} \log(\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\bar{\alpha}}}] + \bar{\alpha}\rho) \right\} - \alpha^* \log(\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha^*}}]) - \alpha^*\rho.$$

We first show that,

$$\left| \log\left(\frac{\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}]}{\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]}\right) \right| \leq e^{\frac{M}{\alpha}} |\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}] - \mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]|. \quad (71)$$

Consider 2 cases:  $\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}] \geq \mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]$  and  $\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}] > \mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}]$

**Case-1:**  $\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}] \geq \mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]$ :

$$\begin{aligned} \left| \log\left(\frac{\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}]}{\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]}\right) \right| &= \log\left(\frac{\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}]}{\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]}\right) \\ &= \log\left(1 + \frac{\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}] - \mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]}{\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]}\right) \\ &\leq \frac{\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}] - \mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]}{\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]} \\ &\leq e^{\frac{M}{\alpha}} (\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}] - \mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]). \end{aligned}$$

**Case-2:**  $\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}] < \mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]$ :

$$\begin{aligned} \left| \log\left(\frac{\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}]}{\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]}\right) \right| &= \log\left(\frac{\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]}{\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}]}\right) \\ &= \log\left(1 + \frac{\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}] - \mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}]}{\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}]}\right) \\ &\leq \frac{\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}] - \mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}]}{\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}]} \\ &\leq e^{\frac{M}{\alpha}} (\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}] - \mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}]). \end{aligned}$$

Hence, Equation (71) holds. Then, for  $\alpha \in [\underline{\alpha}, \bar{\alpha}]$ , we have

$$\begin{aligned} &|(\alpha \log(\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}]) + \alpha\rho) - (\alpha \log(\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]) + \alpha\rho)| \\ &= \alpha \left| \log\left(1 + \frac{\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}] - \mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]}{\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]}\right) \right| \\ &\stackrel{(i)}{\leq} \alpha e^{\frac{M}{\alpha}} |\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}] - \mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]| \end{aligned} \quad (72)$$

$$\begin{aligned}
&\stackrel{(ii)}{\leq} \alpha e^{\frac{M}{\alpha}} \|f(s, a) - \hat{f}_n(s, a)\| \quad (\text{Lemma 15}) \\
&\stackrel{(iii)}{\leq} \mathcal{O}\left(\alpha e^{\frac{M}{\alpha}} \beta_n(\delta) \sqrt{\frac{2ed^2\Gamma_{nd}}{n}}\right) \quad (\text{from Equation (21)}). \tag{73}
\end{aligned}$$

Here (i) holds from Equation (71), (ii) from Lemma 15 and (iii) from Equation (21).

We further show that  $\hat{\alpha}_n^* \in [\underline{\alpha}, \bar{\alpha}]$ . The first step in achieving that is to restrict  $n > N''(\rho, P_f(s, a)) = \mathcal{O}\left(4 \frac{M^2 e^{\frac{2M}{\alpha}} \beta_n^2(\delta) 2ed^2\Gamma_{nd}}{(\rho\tau)^2}\right)$ . It implies that if  $\mathcal{O}\left(\alpha e^{\frac{M}{\alpha}} \beta_n(\delta) \sqrt{\frac{2ed^2\Gamma_{nd}}{n}}\right) < \tau/2$  and for  $n > \max_{s,a} N''(\rho, P_f(s, a))$  from Equation (73) with probability at least  $1 - \delta$ , for all  $(s, a) \in \mathcal{S} \times \mathcal{A}$ , we have

$$\max_{\underline{\alpha}, \alpha^*, \bar{\alpha}} |(\alpha \log(\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}]) + \alpha\rho) - (\alpha \log(\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]) + \alpha\rho)| \leq \tau/2. \tag{74}$$

It further implies that

$$\begin{aligned}
&\max_{\alpha \in [\underline{\alpha}, \bar{\alpha}]} \{(-\alpha \log(\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}]) - \alpha\rho)\} \\
&\stackrel{(i)}{\geq} -\alpha^* \log(\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha^*}}]) - \alpha^*\rho \\
&\stackrel{(ii)}{\geq} -\alpha^* \log(\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha^*}}]) - \alpha^*\rho - \tau/2 \\
&\stackrel{(iii)}{\geq} \max\{-\underline{\alpha} \log(\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\underline{\alpha}}}] - \underline{\alpha}\rho, -\bar{\alpha} \log(\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\bar{\alpha}}}] - \bar{\alpha}\rho) + \tau/2 \\
&\stackrel{(iv)}{\geq} \max\{-\underline{\alpha} \log(\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\underline{\alpha}}}] - \underline{\alpha}\rho, -\bar{\alpha} \log(\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\bar{\alpha}}}] - \bar{\alpha}\rho)\}. \tag{75}
\end{aligned}$$

where (i) follows from the fact that  $\alpha^* \in [\underline{\alpha}, \bar{\alpha}]$ , (ii) follows from Equation (74), (iii) follows from the definition of  $\tau$  in Equation (71) and (iv) again follows from Equation (74).

Thus  $\hat{\alpha}_n^* \in [\underline{\alpha}, \bar{\alpha}]$  follows from Equation (75) and concavity of  $-\alpha \log(\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]) - \alpha\rho$  w.r.t.  $\alpha$ . Note that  $\alpha^*$  also belongs in this set. We bound  $\max_{\alpha \geq 0} \{-\alpha \log(\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]) - \alpha\rho\} - \max_{\alpha \geq 0} \{-\alpha \log(\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}]) - \alpha\rho\}$  only between  $\alpha \in [\underline{\alpha}, \bar{\alpha}]$  instead of all  $\alpha > 0$ . As a result, it holds that

$$\begin{aligned}
&\left| \max_{\alpha \in [\underline{\alpha}, \bar{\alpha}]} \{-\alpha \log(\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]) - \alpha\rho\} - \max_{\alpha \in [\underline{\alpha}, \bar{\alpha}]} \{-\alpha \log(\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}]) - \alpha\rho\} \right| \tag{76} \\
&\leq \max_{\alpha \in [\underline{\alpha}, \bar{\alpha}]} \left| \{-\alpha \log(\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]) - \alpha\rho\} - \{-\alpha \log(\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}]) - \alpha\rho\} \right| \\
&= \max_{\alpha \in [\underline{\alpha}, \bar{\alpha}]} \alpha \left| \log\left(1 + \frac{\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}] - \mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]}{\mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]}\right) \right| \\
&\leq \max_{\alpha \in [\underline{\alpha}, \bar{\alpha}]} 2\alpha e^{\frac{M}{\alpha}} |\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}] - \mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]| \\
&\leq 2\frac{M}{\rho} e^{\frac{M}{\alpha}} \max_{\alpha \in [\underline{\alpha}, \bar{\alpha}]} |\mathbb{E}_{s' \sim P_{\hat{f}_n}}[e^{-\frac{V(s')}{\alpha}}] - \mathbb{E}_{s' \sim P_f}[e^{-\frac{V(s')}{\alpha}}]|,
\end{aligned}$$

where the first inequality follows from Equation (71) and second inequality follows from the bounds of  $\alpha$ . Taking a maximum over all  $(s, a)$  gets the desired result.  $\blacksquare$

**Lemma 15** (*Bound by difference between estimated model  $\hat{f}_n$  and true  $f$* ) For any value function  $V(s') : \mathcal{S} \rightarrow [0, 1/(1 - \gamma)]$  and any  $\alpha > 0$ , it holds that

$$|\mathbb{E}_{s' \sim P_{\hat{f}_n}(s, a)}[e^{-\frac{V(s')}{\alpha}}] - \mathbb{E}_{s' \sim P_f(s, a)}[e^{-\frac{V(s')}{\alpha}}]| \leq \sigma^{-1} \|f(s, a) - \hat{f}_n(s, a)\|,$$

where  $P_{\hat{f}_n}(s, a) = \mathcal{N}(\hat{f}_n(s, a), \sigma^2 I)$  and  $P_f(s, a) = \mathcal{N}(f(s, a), \sigma^2 I)$ .

**Proof**

$$\begin{aligned} \left| \mathbb{E}_{s' \sim P_{\hat{f}_n}(s, a)}[e^{-\frac{V(s')}{\alpha}}] - \mathbb{E}_{s' \sim P_f(s, a)}[e^{-\frac{V(s')}{\alpha}}] \right| &= \left| \int_{\mathbb{R}^d} \frac{1}{\sqrt{(2\pi\sigma^2)^d}} e^{-\frac{V(s')}{\alpha}} \left( e^{-\frac{\|x-f(s, a)\|^2}{2\sigma^2}} - e^{-\frac{\|x-\hat{f}_n(s, a)\|^2}{2\sigma^2}} \right) \right| \\ &\leq \int_{\mathbb{R}^d} \frac{1}{\sqrt{(2\pi\sigma^2)^d}} e^{-\frac{V(s')}{\alpha}} \left| e^{-\frac{\|x-f(s, a)\|^2}{2\sigma^2}} - e^{-\frac{\|x-\hat{f}_n(s, a)\|^2}{2\sigma^2}} \right| \\ &\leq \int_{\mathbb{R}^d} \frac{1}{\sqrt{(2\pi\sigma^2)^d}} \left| e^{-\frac{\|x-f(s, a)\|^2}{2\sigma^2}} - e^{-\frac{\|x-\hat{f}_n(s, a)\|^2}{2\sigma^2}} \right| \\ &\stackrel{(i)}{\leq} 2 \cdot \text{TV}(P_{\hat{f}_n}(s, a), P_f(s, a)) \\ &\stackrel{(ii)}{\leq} 2\sqrt{\text{KL}(P_{\hat{f}_n}(s, a), P_f(s, a))/2} \\ &\stackrel{(iii)}{\leq} 2\sqrt{\|f(s, a) - \hat{f}_n(s, a)\|^2/4\sigma^2} \\ &\leq \|f(s, a) - \hat{f}_n(s, a)\|/\sigma, \end{aligned}$$

where (i) follows from the definition of Total Variation (TV) distance between any two multivariate Gaussians, (ii) uses the Pinsker's inequality, (iii) uses the formula for KL-divergence between multivariate Gaussian distributions.  $\blacksquare$

**Lemma 16** (*Proposition-2 in (26)*) For any function  $V(\cdot) : \mathcal{S} \rightarrow [0, 1/(1 - \gamma)]$  and random variable  $Y = V(s')$  for  $s' \sim P_f(s, a)$ , we have

$$\lim_{\alpha \rightarrow 0} [-\alpha \log(\mathbb{E}_{s' \sim P_f(s, a)}[e^{-\frac{V(s')}{\alpha}}]) - \alpha \rho] = \text{ESI}(Y), \quad (77)$$

where  $\text{ESI}(Y) = \sup\{t \in \mathbb{R} : \mathbb{P}\{Y < t\} = 0\}$  (*essential infimum*).

**Proof** Consider the case when  $M > \text{ESI}(Y)$ . Let  $\kappa_M = \mathbb{P}(V(s') \leq M) = \int_{s'} 1(V(s') \leq M) e^{-\frac{\|s'-f(s, a)\|^2}{\sigma^2}}$ . It holds that

$$-\alpha \log \left( \mathbb{E}_{s' \sim P_f(s, a)}[e^{-\frac{V(s')}{\alpha}}] \right)$$



$$\begin{aligned}
&= -\alpha \log \left( \mathbb{E}_{s' \sim P_f(s,a)} [1(V(s') \leq M)e^{-\frac{V(s')}{\alpha}} + 1(V(s') > M)e^{-\frac{V(s')}{\alpha}}] \right) \\
&\leq -\alpha \log \left( \mathbb{E}_{s' \sim P_f(s,a)} [1(V(s') \leq M)e^{-\frac{V(s')}{\alpha}}] \right) \\
&\leq -\alpha \log \left( \mathbb{E}_{s' \sim P_f(s,a)} [1(V(s') \leq M)e^{-\frac{M}{\alpha}}] \right) \\
&\leq -\alpha \log \left( \kappa_M e^{-\frac{M}{\alpha}} \right) \\
&= M - \alpha \log(\kappa_M).
\end{aligned} \tag{78}$$

Thus for any  $M > \text{ESI}(Y)$ , we have

$$\lim_{\alpha \rightarrow 0} [\{-\alpha \log(\mathbb{E}_{s' \sim P_f(s,a)} [e^{-\frac{V(s')}{\alpha}}]) - \alpha \rho\} \leq M.$$

Combining with the fact that  $\lim_{\alpha \rightarrow 0} [\{-\alpha \log(\mathbb{E}_{s' \sim P_f(s,a)} [e^{-\frac{V(s')}{\alpha}}]) - \alpha \rho\} \geq \text{ESI}(Y)$ , we get the desired result.  $\blacksquare$

**Lemma 17** ( $\zeta$ -cover construction) For  $\mathcal{V}$  denoting the set of value functions from  $\mathcal{S} \rightarrow [0, 1/(1-\gamma)]$ ,  $\bar{\alpha} = M/\rho$ ,  $\underline{\alpha}$  as defined in Lemma 14 we have with probability at least  $1 - \delta$ ,

$$\begin{aligned}
\max_{V \in \mathcal{V}} \max_{s,a} 2\bar{\alpha} e^{\frac{M}{\alpha}} \max_{\alpha \in [\underline{\alpha}, \bar{\alpha}]} |\mathbb{E}_{s' \sim P_{\hat{f}_n}(s,a)} [e^{-\frac{V(s')}{\alpha}}] - \mathbb{E}_{s' \sim P_f(s,a)} [e^{-\frac{V(s')}{\alpha}}]| \\
\leq \mathcal{O} \left( 2 \left( \frac{M}{\rho} \right) e^{\frac{M}{\alpha_{kl}}} e^{\frac{\zeta}{\alpha_{kl}}} \frac{\beta_n(\delta) \sqrt{2ed^2 \Gamma_{nd}}}{\sqrt{n}} \right).
\end{aligned}$$

**Proof** Let  $\mathcal{N}_{\mathcal{V}}(\zeta)$  be the  $\zeta$ -cover of the set  $\mathcal{V}$ . By definition, there exists  $V' \in \mathcal{N}_{\mathcal{V}}(\zeta)$  such that  $\|V' - V\| \leq \zeta$  for every  $V \in \mathcal{V}$ .

$$\begin{aligned}
&|\mathbb{E}_{s' \sim P_{\hat{f}_n}(s,a)} [e^{-V(s')/\alpha}] - \mathbb{E}_{s' \sim P_f(s,a)} [e^{-V(s')/\alpha}]| \\
&\leq \left| \int_{\mathbb{R}^d} \frac{1}{\sqrt{(2\pi\sigma^2)^d}} e^{-\frac{V(s')}{\alpha}} \left( e^{-\frac{\|s' - f(s,a)\|^2}{\sigma^2}} - e^{-\frac{\|s' - \hat{f}_n(s,a)\|^2}{\sigma^2}} \right) \right| \\
&\leq \int_{\mathbb{R}^d} \frac{1}{\sqrt{(2\pi\sigma^2)^d}} e^{-\frac{V(s')}{\alpha}} \left| e^{-\frac{\|s' - f(s,a)\|^2}{\sigma^2}} - e^{-\frac{\|s' - \hat{f}_n(s,a)\|^2}{\sigma^2}} \right| \\
&\leq \int_{\mathbb{R}^d} \frac{1}{\sqrt{(2\pi\sigma^2)^d}} e^{-\frac{V(s') + V'(s')}{\alpha}} e^{-\frac{V'(s')}{\alpha}} \left| e^{-\frac{\|s' - f(s,a)\|^2}{\sigma^2}} - e^{-\frac{\|s' - \hat{f}_n(s,a)\|^2}{\sigma^2}} \right| \\
&\stackrel{(i)}{\leq} e^{\frac{\zeta}{\alpha_{kl}}} \int_{\mathbb{R}^d} \frac{1}{\sqrt{(2\pi\sigma^2)^d}} e^{-\frac{V'(s')}{\alpha}} \left| e^{-\frac{\|s' - f(s,a)\|^2}{\sigma^2}} - e^{-\frac{\|s' - \hat{f}_n(s,a)\|^2}{\sigma^2}} \right| \\
&\leq \max_{V' \in \mathcal{N}_{\mathcal{V}}(\zeta)} \max_{s,a} \max_{\alpha \in [\alpha_{kl}, \bar{\alpha}]} e^{\frac{\zeta}{\alpha_{kl}}} \int_{\mathbb{R}^d} \frac{1}{\sqrt{(2\pi\sigma^2)^d}} e^{-\frac{V'(s')}{\alpha}} \left| e^{-\frac{\|s' - f(s,a)\|^2}{\sigma^2}} - e^{-\frac{\|s' - \hat{f}_n(s,a)\|^2}{\sigma^2}} \right|. \tag{79}
\end{aligned}$$

Here (i) is obtained using the fact that  $\|V' - V\| \leq \zeta$  and  $\alpha_{kl}$  is the minimum value of  $\underline{\alpha}$  as defined in Lemma 14. Using Equation (79), we bound uniformly over all  $V \in \mathcal{V}$ , we have

$$\max_{V \in \mathcal{V}} \max_{s,a} 2\bar{\alpha} e^{\frac{M}{\alpha_{kl}}} \max_{\alpha \in [\alpha_{kl}, \bar{\alpha}]} |\mathbb{E}_{s' \sim P_{\hat{f}_n}(s,a)} [e^{-\frac{V(s')}{\alpha}}] - \mathbb{E}_{s' \sim P_f(s,a)} [e^{-\frac{V(s')}{\alpha}}]|$$

$$\begin{aligned}
&\leq \max_{V' \in \mathcal{N}_V(\zeta)} \max_{s,a} \max_{\alpha \in [\alpha_{kl}, \bar{\alpha}]} 2\bar{\alpha} e^{\frac{M}{\alpha_{kl}}} e^{\frac{\zeta}{\alpha_{kl}}} \int_{\mathbb{R}^d} \frac{1}{\sqrt{(2\pi\sigma^2)^d}} e^{-\frac{V'(s')}{\alpha}} \left| e^{-\frac{\|s'-f(s,a)\|^2}{\sigma^2}} - e^{-\frac{\|s'-\hat{f}_n(s,a)\|^2}{\sigma^2}} \right| \\
&\leq \max_{V' \in \mathcal{N}_V(\zeta)} \max_{s,a} \max_{\alpha \in [\alpha_{kl}, \bar{\alpha}]} 2\bar{\alpha} e^{\frac{M}{\alpha_{kl}}} e^{\frac{\zeta}{\alpha_{kl}}} \int_{\mathbb{R}^d} \frac{1}{\sqrt{(2\pi\sigma^2)^d}} \left| e^{-\frac{\|s'-f(s,a)\|^2}{\sigma^2}} - e^{-\frac{\|s'-\hat{f}_n(s,a)\|^2}{\sigma^2}} \right| \quad (80) \\
&\stackrel{(i)}{\leq} \max_{s,a} 4\bar{\alpha}\sigma^{-1} e^{\frac{M}{\alpha_{kl}}} e^{\frac{\zeta}{\alpha_{kl}}} \|f(s,a) - \hat{f}_n(s,a)\| \\
&\stackrel{(ii)}{\leq} \mathcal{O}\left(2\left(\frac{M}{\rho}\right) e^{\frac{M}{\alpha_{kl}}} e^{\frac{\zeta}{\alpha_{kl}}} \frac{\beta_n(\delta) \sqrt{2ed^2\Gamma_{nd}}}{\sigma\sqrt{n}}\right)
\end{aligned}$$

Here (i) follows from Lemma 15 and by the fact that none of the remaining terms inside max depend on  $V'$  or  $\alpha$ . And (ii) follows from  $\bar{\alpha} = \frac{M}{\rho}$  and Equation (21). ■

## Appendix E. Other Uncertainty Sets

### E.1 $\chi^2$ Uncertainty Set

The f-divergence ((3; 16)) between probability measures  $P$  and  $P_0$  defined over  $\mathcal{X}$  for a convex function  $f : \mathbb{R} \rightarrow \mathbb{R}_+ = \mathbb{R}_+ \cup \{\infty\}$  satisfying  $f(1) = 0$  and  $f(t) = \infty$  for any  $t < 0$  is defined as follows:

$$D_f(P||P_0) = \int f\left(\frac{dP}{dP_0}\right) dP_0. \quad (81)$$

Specifically (22) considers the Cressie-Read family of f-divergences ((15), see Appendix E.1) which includes  $\chi^2$  divergence ( $k = 2$ ), etc. This family of f-divergences can be parametrized by  $k \in (-\infty, \infty) \setminus \{0, 1\}$  with  $f_k(t) := \frac{t^k - kt + k - 1}{k(k-1)}$ . Using this, we state the reformulation result from (22, Lemma-1).

**Lemma 18** For  $k \in (1, \infty)$ ,  $k_* = k/k - 1$ , any  $\rho > 0$  and  $c_k(\rho) = (1 + k(k-1)\rho)^{\frac{1}{k}}$  and  $X \sim P_0$  where  $P_0$  is any probability distribution over  $\mathcal{X}$  with  $H : \mathcal{X} \rightarrow \mathbb{R}$ , we have

$$\sup_{P: D_{f_k}(P||P_0) \leq \rho} \mathbb{E}_P[H(X)] = \inf_{\eta \in \mathbb{R}} \{c_k(\rho) (\mathbb{E}_{P_0}[(H(X) - \eta)_+^{k_*}])^{\frac{1}{k_*}} + \eta\}. \quad (82)$$

**Theorem 19** (Sample Complexity under  $\chi^2$  uncertainty set) Consider a robust MDP (see Section 2) with nominal transition dynamics  $f$  and uncertainty set defined as in Equation (2) w.r.t.  $\chi^2$  divergence. For  $\pi^*$  denoting the robust optimal policy w.r.t. nominal transition dynamics  $f$  and  $\pi_N^*$  denoting the robust optimal policy w.r.t. learned nominal transition dynamics  $\hat{f}_N$  via Algorithm 1, and  $\delta \in (0, 1)$ ,  $\epsilon \in (0, \frac{1}{1-\gamma})$ , it holds that  $\max_s |V_{\pi_N^*, f}^R(s) - V_{\pi^*, f}^R(s)| \leq \epsilon$  with probability at least  $1 - \delta$  for any  $N \geq N_{\chi^2}$ , where

$$N_{\chi^2} = \mathcal{O}\left(\left(\frac{1 + 2\rho}{\sqrt{1 + 2\rho} - 1}\right)^4 \frac{\gamma^4 \beta_n(\delta)^2 d^2 \gamma_{nd}}{\epsilon^4 (1 - \gamma)^8}\right). \quad (83)$$

**Proof Step (i):** As detailed in the proof outline of Section 4, in order to bound  $V_{\hat{\pi}_n, f}^R(s) - V_{\pi^*, f}^R(s)$ , we begin by adding and subtracting  $V_{\hat{\pi}_n, \hat{f}_n}^R(s)$  which is the robust value function

w.r.t. the nominal transition dynamics  $\hat{f}_n$  and its corresponding optimal policy  $\hat{\pi}_n$ . Then, we split the difference into two terms as follows:

$$V_{\hat{\pi}_n, f}^R(s) - V_{\pi^*, f}^R(s) = \underbrace{V_{\hat{\pi}_n, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s)}_{(i)} + \underbrace{V_{\hat{\pi}_n, \hat{f}_n}^R(s) - V_{\pi^*, f}^R(s)}_{(ii)}. \quad (84)$$

In order to not disturb the flow of the proof we bound (i) and (ii) separately Lemma 11 and Lemma 12 respectively. From Lemma 11, we obtain that

$$\begin{aligned} (i) &\leq \max_s \left| V_{\hat{\pi}_n, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s) \right| \\ &\leq \frac{\gamma}{1-\gamma} \max_s \left| \inf_{\chi^2(p||P_{f(s, \hat{\pi}_n(s))}) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] - \inf_{\chi^2(p||P_{\hat{f}_n(s, \hat{\pi}_n(s))}) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] \right|. \end{aligned} \quad (85)$$

And from Lemma 12, we obtain that

$$\begin{aligned} (ii) &\leq \max_s \left| V_{\hat{\pi}_n, \hat{f}_n}^R(s) - V_{\pi^*, f}^R(s) \right| \\ &\leq \frac{\gamma}{1-\gamma} \max_s \left| \inf_{\chi^2(p||P_{\hat{f}_n(s, \hat{\pi}_n(s))}) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\pi^*, f}^R(s') \right] - \inf_{\chi^2(p||P_{f(s, \hat{\pi}_n(s))}) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\pi^*, f}^R(s') \right] \right|. \end{aligned} \quad (86)$$

Note that both these terms in Equations (85) and (86) are of the form mentioned in the **Step (i)** of Section 4.

**Step (ii):** Next, corresponding to **Step (ii)** of the proof outline in Section 4, we use Lemma 18 to bound Equations (85) and (86). Denote  $M := \frac{1}{1-\gamma} \geq \max_s V_{\pi^*}^R(s)$  and  $c_2(\rho) := \sqrt{1+2\rho}$  for convenience. Using Equation (85) and Lemma 20 (internally using Lemma 18), it holds that

$$\begin{aligned} (i) &\leq \max_s \left| V_{\hat{\pi}_n, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s) \right| \\ &\leq \frac{1}{1-\gamma} \max_s \left| \gamma \inf_{\chi^2(p||P_{f(s, \hat{\pi}_n(s))}) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] - \gamma \inf_{\chi^2(p||P_{\hat{f}_n(s, \hat{\pi}_n(s))}) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] \right| \\ &\leq \max_{s,a} \left( \frac{\gamma \sqrt{1+2\rho}}{1-\gamma} \sup_{\eta \in [0, \frac{c_2(\rho)M}{c_2(\rho)-1}]} \left\{ \left| \mathbb{E}_{P_{f(s,a)}} [(-V_{\hat{\pi}_n, f}^R(s') + \eta)_+]^2 - \mathbb{E}_{P_{\hat{f}_n(s,a)}} [(-V_{\hat{\pi}_n, f}^R(s') + \eta)_+]^2 \right|^{\frac{1}{2}} \right\} \right). \end{aligned} \quad (87)$$

$$\leq \max_{V(\cdot) \in \mathcal{V}} \max_{s,a} \left( \frac{\gamma \sqrt{1+2\rho}}{1-\gamma} \sup_{\eta \in [0, \frac{c_2(\rho)M}{c_2(\rho)-1}]} \left\{ \left| \mathbb{E}_{P_{f(s,a)}} [(-V(s') + \eta)_+]^2 - \mathbb{E}_{P_{\hat{f}_n(s,a)}} [(-V(s') + \eta)_+]^2 \right|^{\frac{1}{2}} \right\} \right). \quad (88)$$

We can bound (ii) similarly.

$$(ii) \leq \max_s \left| V_{\hat{\pi}_n, \hat{f}_n}^R(s) - V_{\pi^*, f}^R(s) \right| \quad (89)$$

$$\leq \max_{V(\cdot) \in \mathcal{V}} \max_{s,a} \left( \frac{\gamma\sqrt{1+2\rho}}{1-\gamma} \sup_{\eta \in [0, \frac{c_2(\rho)M}{c_2(\rho)-1}]} \left\{ \left| \mathbb{E}_{P_f(s,a)}[(-V(s') + \eta)_+]^2] - \mathbb{E}_{P_{\hat{f}_n}(s,a)}[(-V(s') + \eta)_+]^2 \right|^{\frac{1}{2}} \right\} \right). \quad (90)$$

**Step (iii):** Next, we want to utilize the learning error bound (Equation (21)) that bounds the difference between the means of true nominal transition dynamics  $P_f$  and learned nominal transition dynamics  $P_{\hat{f}_n}$  to bound Equations (88) and (90).

We begin by bounding the difference  $\left| \mathbb{E}_{P_f(s,a)}[(-V(s') + \eta)_+]^2] - \mathbb{E}_{P_{\hat{f}_n}(s,a)}[(-V(s') + \eta)_+]^2 \right|$ , by the difference in means of  $P_f$  and  $P_{\hat{f}_n}$  in Lemma 23. Since Equation (88) has a max over all value functions, we introduce a covering number argument in Lemma 21 to reform it to a max over the functions in the  $\zeta$ -covering set. We then use Lemma 23 to obtain bounds in terms of maximum information gain  $\Gamma_{Nd}$  (Equation (18)) and  $\zeta$ . Further details regarding the covering number argument are deferred to Lemma 21. Then, we apply the result of Lemma 21 with  $\zeta = 1$  (defined in Lemma 21) on Equation (88). Then, it holds that

$$(i) \leq \max_s \left| V_{\hat{\pi}_n, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s) \right| = \mathcal{O} \left( \left( \frac{\gamma(c_2(\rho))^2 M^2}{c_2(\rho) - 1} \right) \left( \frac{\beta_n(\delta) \sqrt{2ed^2 \gamma_{nd}}}{\sigma \sqrt{n}} \right)^{\frac{1}{2}} \right). \quad (91)$$

Note that  $\beta_n$ , which appears in Lemma 2, has a logarithmic dependence on  $n$ . Similarly, from Equation (90), and Lemmas 23, 21, we obtain

$$(ii) \leq \max_s \left| V_{\hat{\pi}_n, \hat{f}_n}^R(s) - V_{\pi^*, f}^R(s) \right| = \mathcal{O} \left( \left( \frac{\gamma(c_2(\rho))^2 M^2}{c_2(\rho) - 1} \right) \left( \frac{\beta_n(\delta) \sqrt{2ed^2 \gamma_{nd}}}{\sigma \sqrt{n}} \right)^{\frac{1}{2}} \right). \quad (92)$$

Note that we want to bound  $V_{\hat{\pi}_n, f}^R(s) - V_{\pi^*, f}^R(s) = (i) + (ii)$  over all  $s \in \mathcal{S}$ . Using  $\max_s \left| V_{\hat{\pi}_n, f}^R(s) - V_{\pi^*, f}^R(s) \right| \leq \max_s \left| V_{\hat{\pi}_n, \hat{f}_n}^R(s) - V_{\pi^*, f}^R(s) \right| + \max_s \left| V_{\hat{\pi}_n, \hat{f}_n}^R(s) - V_{\pi^*, f}^R(s) \right|$  and substituting  $M$  by  $1/(1-\gamma)$ , we obtain from Equation (91) and Equation (92)

$$\max_s \left| V_{\hat{\pi}_n, f}^R(s) - V_{\pi^*, f}^R(s) \right| = \mathcal{O} \left( \left( \frac{\gamma(c_2(\rho))^2 M^2}{c_2(\rho) - 1} \right) \left( \frac{\beta_n(\delta) \sqrt{2ed^2 \gamma_{nd}}}{\sigma \sqrt{n}} \right)^{\frac{1}{2}} \right).$$

Finally, to ensure that  $\max_s |V_{\hat{\pi}_n, f}^R(s) - V_{\pi^*, f}^R(s)| \leq \epsilon$ , it suffices to have

$$\max_s \left| V_{\hat{\pi}_n, f}^R(s) - V_{\pi^*, f}^R(s) \right| = \mathcal{O} \left( \left( \frac{\gamma(c_2(\rho))^2 M^2}{c_2(\rho) - 1} \right) \left( \frac{\beta_n(\delta) \sqrt{2ed^2 \gamma_{nd}}}{\sigma \sqrt{n}} \right)^{\frac{1}{2}} \right) = \epsilon.$$

Moving  $\sqrt{n}$  and  $\epsilon$  to opposite sides and squaring both sides twice, we obtain

$$n = \mathcal{O} \left( \left( \frac{1+2\rho}{\sqrt{1+2\rho}-1} \right)^4 \frac{\gamma^4 \beta_n(\delta)^2 d^2 \gamma_{nd}}{\sigma^2 \epsilon^4 (1-\gamma)^8} \right).$$

■

**Lemma 20** (*Simplification using Lemma 18 reformulation*) For any value function  $V$  from  $\mathcal{S} \rightarrow [0, 1/(1 - \gamma)]$ , it holds that

$$\begin{aligned} \max_s \left| \inf_{\chi^2(p \| P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] - \inf_{\chi^2(p \| P_f(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] \right| \leq \\ \max_{s, a} c_2(\rho) \sup_{\eta \in [0, \frac{c_2(\rho)M}{c_2(\rho)-1}]} \{ |\mathbb{E}_{P_f(s, a)} [(-V(s') + \eta)_+^2] - \mathbb{E}_{P_{\hat{f}_n}(s, a)} [(-V(s') + \eta)_+^2] |^{\frac{1}{2}} \}, \end{aligned} \quad (93)$$

where  $c_2(\rho) = \sqrt{1 + 2\rho}$  and  $M = 1/(1 - \gamma)$ .

**Proof** First note that,

$$\begin{aligned} \max_s \left| \inf_{\chi^2(p \| P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] - \inf_{\chi^2(p \| P_f(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] \right| \leq \\ \max_{s, a} \left| \inf_{\chi^2(p \| P_{\hat{f}_n}(s, a)) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] - \inf_{\chi^2(p \| P_f(s, a)) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] \right|. \end{aligned} \quad (94)$$

Using Lemma 18 and focusing to bound right side of Equation (94) for one particular  $(s, a)$  state-action pair, we obtain

$$\begin{aligned} \left| \inf_{\chi^2(p \| P_{\hat{f}_n}(s, a)) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] - \inf_{\chi^2(p \| P_f(s, a)) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] \right| = \\ \left| \sup_{\eta \in \mathbb{R}} \{ -c_2(\rho) (\mathbb{E}_{P_f(s, a)} [(-V(s') - \eta)_+^2])^{\frac{1}{2}} - \eta \} - \sup_{\eta \in \mathbb{R}} \{ -c_2(\rho) (\mathbb{E}_{P_{\hat{f}_n}(s, a)} [(-V(s') - \eta)_+^2])^{\frac{1}{2}} - \eta \} \right| \\ \stackrel{(i)}{=} \left| \sup_{\eta \in \mathbb{R}} \{ -c_2(\rho) (\mathbb{E}_{P_f(s, a)} [(-V(s') + \eta)_+^2])^{\frac{1}{2}} + \eta \} - \sup_{\eta \in \mathbb{R}} \{ -c_2(\rho) (\mathbb{E}_{P_{\hat{f}_n}(s, a)} [(-V(s') + \eta)_+^2])^{\frac{1}{2}} + \eta \} \right|, \end{aligned} \quad (95)$$

where (i) is obtained by replacing  $\eta$  with  $-\eta$ .

Let  $g_{\chi^2}(\eta, P_f(s, a)) := \left( -c_2(\rho) (\mathbb{E}_{P_f(s, a)} [(-V(s') + \eta)_+^2])^{\frac{1}{2}} + \eta \right)$ . Note that  $g_{\chi^2}(\eta, P_f(s, a))$  satisfies the following: For  $\eta \leq 0$  (implying  $(-V(s') + \eta) \leq 0$  and  $(-V(s') + \eta)_+ = 0$ ),

$$g_{\chi^2}(\eta, P_f(s, a)) = \eta \leq 0. \quad (96)$$

And for  $\eta = \frac{c_2(\rho)M}{c_2(\rho)-1} > 0$ ,

$$\begin{aligned} g_{\chi^2}(\frac{c_2(\rho)M}{c_2(\rho)-1}, P_f(s, a)) &= -c_2(\rho) (\mathbb{E}_{P_f(s, a)} [(-V(s') + \frac{c_2(\rho)M}{c_2(\rho)-1})_+^2])^{\frac{1}{2}} + \frac{c_2(\rho)M}{c_2(\rho)-1} \\ &\stackrel{(i)}{\leq} \frac{c_2(\rho)M}{c_2(\rho)-1} - c_2(\rho) (\mathbb{E}_{P_f(s, a)} [(-M + \frac{c_2(\rho)M}{c_2(\rho)-1})_+^2])^{\frac{1}{2}} \\ &\leq \frac{c_2(\rho)M}{c_2(\rho)-1} - c_2(\rho) (\mathbb{E}_{P_f(s, a)} [(\frac{M}{c_2(\rho)-1})_+^2])^{\frac{1}{2}} \\ &\leq \frac{c_2(\rho)M}{c_2(\rho)-1} - \frac{c_2(\rho)M}{c_2(\rho)-1} \\ &= 0, \end{aligned} \quad (97)$$

where (i) follows from the fact that the random variable  $V(s')$  is bounded by  $M = 1/1 - \gamma$ . A similar result can be shown for  $g_{\chi^2}(\eta, P_{\hat{f}_n}(s, a))$  (or for any  $P$ ). Along with the convexity of  $\eta \rightarrow g_{\chi}(\eta, P)$  ((22)), and  $\inf_{\chi^2(p||P) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] \geq 0$ , Equation (96) and Equation (97) imply that the sup is attained between  $[0, \frac{c_2(\rho)M}{c_2(\rho)-1}]$  for both  $\sup_{\eta \in \mathbb{R}} g_{\chi}(\eta, P_f(s, a))$  and  $\sup_{\eta \in \mathbb{R}} g_{\chi}(\eta, P_{\hat{f}_n}(s, a))$ . Using this in Equation (95) we have,

$$\left| \sup_{\eta \in \mathbb{R}} \{g_{\chi}(\eta, P_f(s, a))\} - \sup_{\eta \in \mathbb{R}} \{g_{\chi}(\eta, P_{\hat{f}_n}(s, a))\} \right| \quad (98)$$

$$= \left| \sup_{\eta \in [0, \frac{c_2(\rho)M}{c_2(\rho)-1}]} \{g_{\chi}(\eta, P_f(s, a))\} - \sup_{\eta \in [0, \frac{c_2(\rho)M}{c_2(\rho)-1}]} \{g_{\chi}(\eta, P_{\hat{f}_n}(s, a))\} \right| \quad (99)$$

$$\leq \sup_{\eta \in [0, \frac{c_2(\rho)M}{c_2(\rho)-1}]} \{|g_{\chi}(\eta, P_f(s, a)) - g_{\chi}(\eta, P_{\hat{f}_n}(s, a))|\} \quad (100)$$

$$\leq \sup_{\eta \in [0, \frac{c_2(\rho)M}{c_2(\rho)-1}]} \{[c_2(\rho)(\mathbb{E}_{P_f(s, a)}[(-V(s') + \eta)_+^2])^{\frac{1}{2}} - c_2(\rho)\mathbb{E}_{P_{\hat{f}_n}(s, a)}[(-V(s') + \eta)_+^2])^{\frac{1}{2}}]\} \quad (101)$$

$$\leq c_2(\rho) \sup_{\eta \in [0, \frac{c_2(\rho)M}{c_2(\rho)-1}]} \{|\mathbb{E}_{P_f(s, a)}[(-V(s') + \eta)_+^2] - \mathbb{E}_{P_{\hat{f}_n}(s, a)}[(-V(s') + \eta)_+^2]|\}^{\frac{1}{2}}. \quad (102)$$

The last step is obtained using the basic inequality  $|\sqrt{a} - \sqrt{b}| \leq \sqrt{|a - b|}$ . ■

**Lemma 21** ( $\zeta$ -cover construction) For  $\mathcal{V}$  denoting the set of value functions from  $\mathcal{S} \rightarrow [0, 1/(1 - \gamma)]$  it holds with probability at least  $1 - \delta$ ,

$$\max_{V \in \mathcal{V}} \max_{s, a} \sup_{\eta \in [0, \frac{c_2(\rho)M}{c_2(\rho)-1}]} \{|\mathbb{E}_{P_f(s, a)}[(-V(s') + \eta)_+^2] - \mathbb{E}_{P_{\hat{f}_n}(s, a)}[(-V(s') + \eta)_+^2]|\}^{\frac{1}{2}} \leq \mathcal{O}\left(\left(\frac{c_2(\rho)M}{c_2(\rho) - 1}\right)\left(\frac{\beta_n(\delta)\sqrt{2ed^2\gamma na}}{\sigma\sqrt{n}}\right)^{\frac{1}{2}}\right), \quad (103)$$

where  $c_2(\rho) = \sqrt{1 + 2\rho}$ ,  $M = 1/(1 - \gamma)$ .

**Proof** Let  $\mathcal{N}_{\mathcal{V}}(\zeta)$  be the  $\zeta$ -cover of the set  $\mathcal{V}$ . By definition, there exists  $V' \in \mathcal{N}_{\mathcal{V}}(\zeta)$  such that  $\|V' - V\| \leq \zeta$  for every  $V \in \mathcal{V}$ .

$$\begin{aligned} & |\mathbb{E}_{P_f(s, a)}[(-V(s') + \eta)_+^2] - \mathbb{E}_{P_{\hat{f}_n}(s, a)}[(-V(s') + \eta)_+^2]| \\ & \leq |\mathbb{E}_{P_f(s, a)}[(-V(s') + \eta)_+^2] - \mathbb{E}_{P_f(s, a)}[(-V'(s') + \eta)_+^2]| \\ & \quad + |\mathbb{E}_{P_f(s, a)}[(-V'(s') + \eta)_+^2] - \mathbb{E}_{P_{\hat{f}_n}(s, a)}[(-V'(s') + \eta)_+^2]| \end{aligned} \quad (104)$$

$$\begin{aligned} & + |\mathbb{E}_{P_{\hat{f}_n}(s, a)}[(-V'(s') + \eta)_+^2] - \mathbb{E}_{P_{\hat{f}_n}(s, a)}[(-V(s') + \eta)_+^2]|. \\ & \stackrel{(i)}{\leq} 4\|V' - V\|^2 + 4\eta\|V' - V\| + |\mathbb{E}_{P_f(s, a)}[(-V'(s') + \eta)_+^2] - \mathbb{E}_{P_{\hat{f}_n}(s, a)}[(-V'(s') + \eta)_+^2]|, \end{aligned} \quad (105)$$

where (i) follows from Lemma 22. Using Equation (105) we bound uniformly over all  $V \in \mathcal{V}$ ,

$$\max_{V \in \mathcal{V}} \max_{s,a} \sup_{\eta \in [0, \frac{c_2(\rho)M}{c_2(\rho)-1}]} \{ |\mathbb{E}_{P_f(s,a)}[(-V(s') + \eta)_+]^2] - \mathbb{E}_{P_{\hat{f}_n}(s,a)}[(-V(s') + \eta)_+]^2] |^{\frac{1}{2}} \} \quad (106)$$

$$\begin{aligned} &\leq \max_{V' \in \mathcal{N}_{\mathcal{V}}(\zeta)} \max_{s,a} \sup_{\eta \in [0, \frac{c_2(\rho)M}{c_2(\rho)-1}]} \left\{ \left( 4\|V' - V\|^2 + 4\eta\|V' - V\| + |\mathbb{E}_{P_f(s,a)}[(-V'(s') + \eta)_+]^2] \right. \right. \\ &\quad \left. \left. - \mathbb{E}_{P_{\hat{f}_n}(s,a)}[(-V'(s') + \eta)_+]^2] \right)^{\frac{1}{2}} \right\} \\ &\stackrel{(ii)}{\leq} \max_{V' \in \mathcal{N}_{\mathcal{V}}(\zeta)} \max_{s,a} \sup_{\eta \in [0, \frac{c_2(\rho)M}{c_2(\rho)-1}]} \left\{ \left( \mathbb{E}_{P_f(s,a)}[(-V'(s') + \eta)_+]^2] - \mathbb{E}_{P_{\hat{f}_n}(s,a)}[(-V'(s') + \eta)_+]^2] \right)^{\frac{1}{2}} \right\} \\ &\quad + \sqrt{4\zeta^2 + 4\zeta \frac{c_2(\rho)M}{c_2(\rho)-1}} \\ &\stackrel{(iii)}{\leq} \max_{V' \in \mathcal{N}_{\mathcal{V}}(\zeta)} \max_{s,a} \sup_{\eta \in [0, \frac{c_2(\rho)M}{c_2(\rho)-1}]} \left\{ \left( \frac{c_2(\rho)M}{c_2(\rho)-1} \sqrt{2\sigma^{-1}\|f(s,a) - \hat{f}_n(s,a)\|} \right) \right\} + \sqrt{4\zeta^2 + 4\zeta \frac{c_2(\rho)M}{c_2(\rho)-1}} \\ &\stackrel{(iv)}{\leq} \mathcal{O} \left( \left( \frac{c_2(\rho)M}{c_2(\rho)-1} \right) \left( \frac{\beta_n(\delta)\sqrt{2ed^2\gamma_{nd}}}{\sigma\sqrt{n}} \right)^{\frac{1}{2}} \right) + \sqrt{4\zeta^2 + 4\zeta \frac{c_2(\rho)M}{c_2(\rho)-1}} \quad (107) \end{aligned}$$

$$\stackrel{(v)}{\leq} \mathcal{O} \left( \left( \frac{c_2(\rho)M}{c_2(\rho)-1} \right) \left( \frac{\beta_n(\delta)\sqrt{2ed^2\gamma_{nd}}}{\sigma\sqrt{n}} \right)^{\frac{1}{2}} \right), \quad (108)$$

where (ii) follows from  $\|V' - V\| \leq \zeta$  and  $\eta \leq \frac{c_2(\rho)M}{c_2(\rho)-1}$ , (iii) follows from Lemma 23, (iv) follows from Equation (21), and (v) follows from substituing  $\zeta = 1$  (or any constant).  $\blacksquare$

**Lemma 22** For any two value functions  $V, V'$  from  $\mathcal{S} \rightarrow [0, 1/(1-\gamma)]$ , it holds that

$$\left| \mathbb{E}_{P_f(s,a)}[(-V'(s') + \eta)_+]^2] - \mathbb{E}_{P_f(s,a)}[(-V(s') + \eta)_+]^2] \right| \leq 2\|V' - V\|^2 + 2\eta\|V' - V\|. \quad (109)$$

**Proof** Let  $p_{P_f(s,a)}(\cdot)$  denote the probability density function of  $P_f(s,a)$ . Then,

$$\begin{aligned} &\mathbb{E}_{P_f(s,a)}[(-V'(s') + \eta)_+]^2] - \mathbb{E}_{P_f(s,a)}[(-V(s') + \eta)_+]^2] \\ &\leq \int_{s' \sim P_f(s,a)} \left( 1(V'(s') < \eta)(-V'(s') + \eta)^2 - 1(V(s') < \eta)(-V(s') + \eta)^2 \right) p_{P_f(s,a)}(s') ds'. \\ &\leq \underbrace{\int_{s' \sim P_f(s,a)} \left( 1(V'(s') < \eta) - 1(V(s') < \eta) \right) (-V'(s') + \eta)^2 p_{P_f(s,a)}(s') ds'}_{(i)} \\ &\quad + \underbrace{\int_{s' \sim P_f(s,a)} 1(V(s') < \eta) \left( (-V'(s') + \eta)^2 - (-V(s') + \eta)^2 \right) p_{P_f(s,a)}(s') ds'}_{(ii)}. \quad (110) \end{aligned}$$

where the last inequality is obtained by adding and subtracting  $1(V(s') < \eta)(-V'(s') + \eta)^2$ .

We begin by bounding (ii). We have,

$$\begin{aligned}
(ii) &= \int_{s' \sim P_f(s,a)} 1(V(s') < \eta) \left( (-V'(s') + \eta)^2 - (-V(s') + \eta)^2 \right) p_{P_f(s,a)}(s') ds' \\
&= \int_{s' \sim P_f(s,a)} 1(V(s') < \eta) \left( -V'(s') + V(s') \right) \left( -V'(s') - V(s') + 2\eta \right) p_{P_f(s,a)}(s') ds' \\
&\leq \int_{s' \sim P_f(s,a)} 1(V(s') < \eta) \left( 1(V'(s') < \eta) + 1(V'(s') \geq \eta) \right) \left( -V'(s') + V(s') \right) \\
&\quad \left( -V'(s') - V(s') + 2\eta \right) p_{P_f(s,a)}(s') ds' \\
&\leq \underbrace{\int 1(V(s'), V'(s') < \eta) (-V'(s') + V(s')) (-V'(s') - V(s') + 2\eta) p_{P_f(s,a)}(s') ds'}_{(ii-a)} \\
&\quad + \underbrace{\int 1(V(s') < \eta \leq V'(s')) (-V'(s') + V(s')) (-V'(s') - V(s') + 2\eta) p_{P_f(s,a)}(s') ds'}_{(ii-b)}.
\end{aligned} \tag{111}$$

Bounding (ii - a) first, we have,

$$\begin{aligned}
(ii-a) &= \int 1(V(s'), V'(s') < \eta) (-V'(s') + V(s')) (-V'(s') - V(s') + 2\eta) p_{P_f(s,a)}(s') ds' \\
&\stackrel{(a)}{\leq} \int 1(V(s'), V'(s') < \eta) \left| -V'(s') + V(s') \right| \left( -V'(s') - V(s') + 2\eta \right) p_{P_f(s,a)}(s') ds' \\
&\stackrel{(b)}{\leq} \int_{s' \sim P_f(s,a)} 1(V(s'), V'(s') < \eta) \left| -V'(s') + V(s') \right| \left( 2\eta \right) p_{P_f(s,a)}(s') ds' \\
&\leq 2\eta \|V' - V\|,
\end{aligned} \tag{112}$$

where (a) and (b) follows from  $(-V'(s') - V(s') + 2\eta) > 0$  as  $V(s'), V'(s') < \eta$ . And (ii - b) can be bounded as,

$$\begin{aligned}
(ii-b) &= \int 1(V(s') < \eta \leq V'(s')) (-V'(s') + V(s')) (-V'(s') - V(s') + 2\eta) p_{P_f(s,a)}(s') ds' \\
&\leq \int 1(V(s') < \eta \leq V'(s')) \left| -V'(s') + V(s') \right| \left| -V'(s') - V(s') + 2\eta \right| p_{P_f(s,a)}(s') ds' \\
&\stackrel{(c)}{\leq} \int 1(V(s') < \eta \leq V'(s')) \left| -V'(s') + V(s') \right| \left| -V(s') + V'(s') \right| p_{P_f(s,a)}(s') ds' \\
&\leq \int_{s' \sim P_f(s,a)} 1(V(s') < \eta \leq V'(s')) \left| -V'(s') + V(s') \right|^2 p_{P_f(s,a)}(s') ds' \\
&\leq \|V' - V\|^2,
\end{aligned} \tag{113}$$



where (c) follows from  $\eta \leq V'(s')$ . Bounding (i) similarly,

$$\begin{aligned}
i &= \int_{s' \sim P_f(s,a)} \left(1(V'(s') < \eta) - 1(V(s') < \eta)\right) (-V'(s') + \eta)^2 p_{P_f(s,a)}(s') ds' \\
&\leq \int_{s' \sim P_f(s,a)} \left(1(V'(s') < \eta \leq V(s'))\right) (-V'(s') + \eta)^2 p_{P_f(s,a)}(s') ds' \\
&\leq \int_{s' \sim P_f(s,a)} \left(1(V'(s') < \eta \leq V(s'))\right) (-V'(s') + V(s'))^2 p_{P_f(s,a)}(s') ds' \\
&\leq \|V' - V\|^2.
\end{aligned} \tag{114}$$

Using Equations (110) to (114) we get the desired result.  $\blacksquare$

**Lemma 23** (Bound by difference between estimated model  $\hat{f}_n$  and true  $f$ ) For any value function  $V(s') : \mathcal{S} \rightarrow [0, 1/(1 - \gamma)]$  and any  $\alpha > 0$ , it holds that

$$|\mathbb{E}_{P_f(s,a)}[(-V(s') + \eta)_+^2] - \mathbb{E}_{P_{\hat{f}_n}(s,a)}[(-V(s') + \eta)_+^2]| \leq 2\sigma^{-1} \left(\frac{c_2(\rho)M}{c_2(\rho) - 1}\right)^2 \|f(s, a) - \hat{f}_n(s, a)\|,$$

where  $P_{\hat{f}_n}(s, a) = \mathcal{N}(\hat{f}_n(s, a), \sigma^2 I)$  and  $P_f(s, a) = \mathcal{N}(f(s, a), \sigma^2 I)$ ,  $\eta \in [0, \frac{c_2(\rho)M}{c_2(\rho)-1}]$ ,  $c_2(\rho) = \sqrt{1 + 2\bar{\rho}}$  and  $M = 1/(1 - \gamma)$ .

**Proof**

$$\begin{aligned}
&\left| \mathbb{E}_{P_f(s,a)}[(-V(s') + \eta)_+^2] - \mathbb{E}_{P_{\hat{f}_n}(s,a)}[(-V(s') + \eta)_+^2] \right| \\
&= \left| \int_{\mathbb{R}^d} \frac{1}{\sqrt{(2\pi\sigma^2)^d}} (-V(s') + \eta)_+^2 \left( e^{-\frac{\|x-f(s,a)\|^2}{2\sigma^2}} - e^{-\frac{\|x-\hat{f}_n(s,a)\|^2}{2\sigma^2}} \right) \right| \\
&\leq \int_{\mathbb{R}^d} \frac{1}{\sqrt{(2\pi\sigma^2)^d}} (-V(s') + \eta)_+^2 \left| e^{-\frac{\|x-f(s,a)\|^2}{2\sigma^2}} - e^{-\frac{\|x-\hat{f}_n(s,a)\|^2}{2\sigma^2}} \right| \\
&\stackrel{(i)}{\leq} \left(\frac{c_2(\rho)M}{c_2(\rho) - 1}\right)^2 \int_{\mathbb{R}^d} \frac{1}{\sqrt{(2\pi\sigma^2)^d}} \left| e^{-\frac{\|x-f(s,a)\|^2}{2\sigma^2}} - e^{-\frac{\|x-\hat{f}_n(s,a)\|^2}{2\sigma^2}} \right| \\
&\stackrel{(ii)}{\leq} 2 \left(\frac{c_2(\rho)M}{c_2(\rho) - 1}\right)^2 \cdot \text{TV}(P_{\hat{f}_n}(s, a), P_f(s, a)) \\
&\stackrel{(iii)}{\leq} 2 \left(\frac{c_2(\rho)M}{c_2(\rho) - 1}\right)^2 \sqrt{\text{KL}(P_{\hat{f}_n}(s, a), P_f(s, a))/2} \\
&\stackrel{(iv)}{\leq} 2 \left(\frac{c_2(\rho)M}{c_2(\rho) - 1}\right)^2 \sqrt{\|f(s, a) - \hat{f}_n(s, a)\|^2/4\sigma^2} \\
&\leq \left(\frac{c_2(\rho)M}{c_2(\rho) - 1}\right)^2 \|f(s, a) - \hat{f}_n(s, a)\|/\sigma,
\end{aligned}$$

where (i) follows from  $(-V(s') + \eta)_+^2 \leq \left(\frac{c_2(\rho)M}{c_2(\rho)-1}\right)^2$  as  $\eta \leq \left(\frac{c_2(\rho)M}{c_2(\rho)-1}\right)$ , (ii) follows from the definition of Total Variation (TV) distance between any two multivariate Gaussians, (iii) uses

the Pinsker's inequality, and (iv) uses the formula for KL-divergence between multivariate Gaussian distributions. ■

## E.2 Total Variation Distance

Similar to Lemma 18, we want a similar convex reformulation for the variation distance. We derive such a reformulation starting from the dual reformulation from (50) and (6) stated as Proposition-1 in (22).

**Lemma 24** *For  $X \sim P_0$  where  $P_0$  is any probability distribution over  $\mathcal{X}$  with  $H : \mathcal{X} \rightarrow \mathbb{R}$ ,  $\rho > 0$  and,  $D_f(P||P_0)$  defined as in Equation (81), it holds that*

$$\sup_{P: D_f(P||P_0) \leq \rho} \mathbb{E}_P[H(X)] = \inf_{\lambda \geq 0, \eta \in \mathbb{R}} \left\{ \mathbb{E}_{P_0} \left[ \lambda f^* \left( \frac{H(X) - \eta}{\lambda} \right) \right] + \lambda \rho + \eta \right\}. \quad (115)$$

Note that the total variation distance between two probability distributions  $P$  and  $P_0$  is attained by substituting  $f_{\text{TV}}(t) = |t - 1|$  in  $D_f(P||P_0) = \int f \left( \frac{dP}{dP_0} \right) dP_0$ . The corresponding Fenchel conjugate  $f_{\text{TV}}^*(s)$  for  $f_{\text{TV}}(t) = |t - 1|$  would be

$$f_{\text{TV}}^*(s) = \begin{cases} -1, & s \leq -1 \\ s, & s \in [-1, 1] \\ \infty, & s > 1 \end{cases} \quad (116)$$

As we require  $\inf_{P: \text{TV}(P||P_0) \leq \rho} \mathbb{E}_P[H(X)]$ , using Equation (115) and replacing  $\eta$  with  $-\eta$ , we have

$$\inf_{P: \text{TV}(P||P_0) \leq \rho} \mathbb{E}_P[H(X)] = \sup_{\lambda \geq 0, \eta \in \mathbb{R}} \left\{ -\mathbb{E}_{P_0} \left[ \lambda f_{\text{TV}}^* \left( \frac{-H(X) + \eta}{\lambda} \right) \right] - \lambda \rho + \eta \right\}. \quad (117)$$

Using Equation (117), we derive a convex reformulation in Lemma 25

**Lemma 25** *(Reformulation for total variation distance based on (59)) For  $\rho > 0$  and  $X \sim P_0$  where  $P_0$  is any probability distribution over  $\mathcal{X}$  with  $H : \mathcal{X} \rightarrow \mathbb{R}$ , for  $0 \leq H(x) \leq \frac{1}{1-\gamma}$  and  $ESI(Y) = \sup\{t \in \mathbb{R} : \mathbb{P}\{Y < t\} = 0\}$  (essential infimum), it holds that*

$$\inf_{P: \text{TV}(P||P_0) \leq \rho} \mathbb{E}_P[H(X)] = \sup_{\eta \in [0, \frac{(2+\rho)}{\rho(1-\gamma)}]} \left\{ -\mathbb{E}_{P_0}[-H(X) + \eta]_+ - \frac{(-ESI(H(x)) + \eta)_+}{2} \rho + \eta \right\}. \quad (118)$$

where TV denotes the total variation distance.

### Proof

Substituting Equation (116) in Equation (117) to obtain the reformulation for total variation distance, we have

$$\inf_{P: \text{TV}(P||P_0) \leq \rho} \mathbb{E}_P[H(X)] \quad (119)$$

$$= \sup_{\lambda \geq 0, \eta \in \mathbb{R}, \frac{-H(x)+\eta}{\lambda} \leq 1} \{-\mathbb{E}_{P_0} \left[ \lambda \max \left\{ \frac{-H(X) + \eta}{\lambda}, -1 \right\} \right] - \lambda \rho + \eta\} \quad (120)$$

$$= \sup_{\lambda \geq 0, \eta \in \mathbb{R}, \frac{-H(x)+\eta}{\lambda} \leq 1} \{-\mathbb{E}_{P_0} \left[ \max \left\{ -H(X) + \eta, -\lambda \right\} \right] - \lambda \rho + \eta\} \quad (121)$$

$$= \sup_{\lambda \geq 0, \eta \in \mathbb{R}, \frac{-H(x)+\eta}{\lambda} \leq 2} \{-\mathbb{E}_{P_0} \left[ \max \left\{ -H(X) + \eta - \lambda, -\lambda \right\} \right] - \lambda \rho + \eta - \lambda\} \quad (122)$$

$$= \sup_{\lambda \geq 0, \eta \in \mathbb{R}, \frac{-H(x)+\eta}{\lambda} \leq 2} \{-\mathbb{E}_{P_0} \left[ \max \left\{ -H(X) + \eta, 0 \right\} \right] - \lambda \rho + \eta\} \quad (123)$$

$$= \sup_{\lambda \geq 0, \eta \in \mathbb{R}, \frac{-H(x)+\eta}{\lambda} \leq 2} \{-\mathbb{E}_{P_0} \left[ -H(X) + \eta \right]_+ - \lambda \rho + \eta\}. \quad (124)$$

Here Equation (122) is obtained by substituting  $\eta$  with  $\eta - \lambda$ . In order to optimize over  $\lambda$ , we need to choose the minimum  $\lambda$  satisfying the constraints. We require  $\lambda \geq \frac{-H(x)+\eta}{2}$  which translates to  $\lambda \geq \frac{-ESI(H(x))+\eta}{2}$  (as this constraint originates inside the expectation, points with zero mass,  $\{t \in \mathbb{R} : \mathbb{P}\{Y < t\} = 0\}$ , will have no effect). Substituting this, we have

$$\inf_{P: \text{TV}(P||P_0) \leq \rho} \mathbb{E}_P[H(X)] = \sup_{\eta \in \mathbb{R}} \{-\mathbb{E}_{P_0} \left[ -H(X) + \eta \right]_+ - \frac{(-ESI(H(x)) + \eta)_+}{2} \rho + \eta\}. \quad (125)$$

Denote the inner function in Equation (125), as

$$g_{\text{TV}}(\eta, P_0) = -\mathbb{E}_{P_0} \left[ -H(X) + \eta \right]_+ - \frac{(-ESI(H(x)) + \eta)_+}{2} \rho + \eta. \quad (126)$$

Note that for  $\eta \leq 0$ , the first two terms in  $g_{\text{TV}}(\eta, P_0)$  will be 0 if  $H(x) > 0$  for all  $x$ . This implies

$$g_{\text{TV}}(\eta, P_0) = \eta \leq 0 \quad \forall \quad \eta \leq 0. \quad (127)$$

Also, as  $H(x) \leq \frac{1}{1-\gamma}$ , we substitute  $\eta = \frac{2+\rho}{\rho(1-\gamma)}$  in  $g_{\text{TV}}(\eta, P_0)$ , and bound it as follows:

$$g_{\text{TV}}\left(\frac{2+\rho}{\rho(1-\gamma)}, P_0\right) = -\mathbb{E}_{P_0} \left[ -H(X) + \frac{(2+\rho)}{\rho(1-\gamma)} \right]_+ - \frac{(-ESI(H(x)) + \frac{(2+\rho)}{\rho(1-\gamma)})_+}{2} \rho + \frac{(2+\rho)}{\rho(1-\gamma)} \quad (128)$$

$$= \mathbb{E}_{P_0} \left[ H(X) \right] - \frac{(2+\rho)}{\rho(1-\gamma)} - \frac{(-ESI(H(x)) + \frac{(2+\rho)}{\rho(1-\gamma)})_+}{2} \rho + \frac{(2+\rho)}{\rho(1-\gamma)} \quad (129)$$

$$= \mathbb{E}_{P_0} \left[ H(X) \right] - \frac{(-ESI(H(x)) + \frac{(2+\rho)}{\rho(1-\gamma)})_+}{2} \rho \quad (130)$$

$$= \mathbb{E}_{P_0} \left[ H(X) \right] - \frac{(-ESI(H(x)) + \frac{(2+\rho)}{\rho(1-\gamma)})}{2} \rho \quad (131)$$

$$= \mathbb{E}_{P_0} \left[ H(X) - \frac{1}{1-\gamma} \right] + \frac{\rho ESI(H(x))}{2} - \frac{\rho}{2(1-\gamma)} \quad (132)$$

$$= \mathbb{E}_{P_0} \left[ H(X) - \frac{1}{1-\gamma} \right] + \frac{\rho}{2} \left( \text{ESI}(H(x)) - \frac{1}{(1-\gamma)} \right) \quad (133)$$

$$\leq 0. \quad (134)$$

Here Equation (129), Equation (131) and Equation (134) are obtained from the fact that that  $H(x) \leq \frac{1}{1-\gamma}$  ( $-H(x) + \frac{(2+\rho)}{\rho(1-\gamma)} > 0$ ) and  $\text{ESI}(H(x)) \leq \frac{1}{1-\gamma}$  ( $-\text{ESI}(H(x)) + \frac{(2+\rho)}{\rho(1-\gamma)} > 0$ ). Along with the convexity of  $g_{\text{TV}}(\eta, P_0)$ , Equation (127) and Equation (134) imply that the  $\sup_{\eta \in \mathbb{R}} \{g_{\text{TV}}(\eta, P_0)\}$  is attained in the  $\eta$  range  $[0, \frac{(2+\rho)}{\rho(1-\gamma)}]$ .  $\blacksquare$

**Theorem 26** (*Sample Complexity under TV uncertainty set*) Consider a robust MDP (see Section 2) with nominal transition dynamics  $f$  and uncertainty set defined as in Equation (2) w.r.t. TV distance. For  $\pi^*$  denoting the robust optimal policy w.r.t. nominal transition dynamics  $f$  and  $\pi_N^*$  denoting the robust optimal policy w.r.t. learned nominal transition dynamics  $\hat{f}_N$  via Algorithm 1, and  $\delta \in (0, 1)$ ,  $\epsilon \in (0, \frac{1}{1-\gamma})$ , it holds that  $\max_s |V_{\pi_N^*, f}^R(s) - V_{\pi^*, f}^R(s)| \leq \epsilon$  with probability at least  $1 - \delta$  for any  $N \geq N_{\text{TV}}$ , where

$$N_{\text{TV}} = \mathcal{O} \left( \frac{(2+\rho)^2 \gamma^2 \beta_n(\delta)^2 d^2 \gamma n d}{\rho^2 (1-\gamma)^4 \epsilon^2} \right). \quad (135)$$

**Proof Step (i):** As detailed in the proof outline of Section 4, in order to bound  $V_{\pi_N^*, f}^R(s) - V_{\pi^*, f}^R(s)$ , we begin by adding and subtracting  $V_{\hat{\pi}_n, \hat{f}_n}^R(s)$  which is the robust value function w.r.t. the nominal transition dynamics  $\hat{f}_n$  and its corresponding optimal policy  $\hat{\pi}_n$ . Then, we split the difference into two terms as follows:

$$V_{\pi_N^*, f}^R(s) - V_{\pi^*, f}^R(s) = \underbrace{V_{\hat{\pi}_n, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s)}_{(i)} + \underbrace{V_{\hat{\pi}_n, \hat{f}_n}^R(s) - V_{\pi^*, f}^R(s)}_{(ii)}. \quad (136)$$

In order to not disturb the flow of the proof we bound (i) and (ii) separately Lemma 11 and Lemma 12 respectively. From Lemma 11, we obtain that

$$\begin{aligned} (i) &\leq \max_s \left| V_{\hat{\pi}_n, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s) \right| \\ &\leq \frac{\gamma}{1-\gamma} \max_s \left| \inf_{\text{TV}(p||P_f(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] - \inf_{\text{TV}(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] \right|. \end{aligned} \quad (137)$$

And from Lemma 12, we obtain that

$$\begin{aligned} (ii) &\leq \max_s \left| V_{\hat{\pi}_n, \hat{f}_n}^R(s) - V_{\pi^*, f}^R(s) \right| \\ &\leq \frac{\gamma}{1-\gamma} \max_s \left| \inf_{\text{TV}(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\pi^*, f}^R(s') \right] - \inf_{\text{TV}(p||P_f(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\pi^*, f}^R(s') \right] \right|. \end{aligned} \quad (138)$$

Note that both these terms in Equations (137) and (138) are of the form mentioned in the **Step (i)** of Section 4.

**Step (ii):** Next, corresponding to **step (ii)** of the proof outline in Section 4, we use Lemma 25 to bound Equations (137) and (138). Denote  $M := \frac{1}{1-\gamma} \geq \max_s V_\pi^R(s)$  for convenience. Using Equation (137) and Lemma 27 (internally using Lemma 25), it holds that

$$\begin{aligned}
(i) &\leq \max_s \left| V_{\hat{\pi}_n, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s) \right| \\
&\leq \frac{1}{1-\gamma} \max_s \left| \gamma \inf_{\text{TV}(p \| P_f(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] - \gamma \inf_{\text{TV}(p \| P_{\hat{f}_n}(s, \hat{\pi}_n(s))) \leq \rho} \mathbb{E}_{s' \sim p} \left[ V_{\hat{\pi}_n, f}^R(s') \right] \right| \\
&\leq \frac{\gamma}{1-\gamma} \max_{s, a} \left( \sup_{\eta \in [0, \frac{(2+\rho)}{\rho(1-\gamma)}]} \left\{ \left| \mathbb{E}_{P_f(s, a)}[(-V_{\hat{\pi}_n, \hat{f}_n}^R(s') + \eta)_+] \right| - \mathbb{E}_{P_{\hat{f}_n}(s, a)}[(-V_{\hat{\pi}_n, \hat{f}_n}^R(s') + \eta)_+] \right\} \right) \tag{139}
\end{aligned}$$

$$\leq \frac{\gamma}{1-\gamma} \max_{V(\cdot) \in \mathcal{V}} \max_{s, a} \left( \sup_{\eta \in [0, \frac{(2+\rho)}{\rho(1-\gamma)}]} \left\{ \left| \mathbb{E}_{P_f(s, a)}[(-V(s') + \eta)_+] \right| - \mathbb{E}_{P_{\hat{f}_n}(s, a)}[(-V(s') + \eta)_+] \right\} \right). \tag{140}$$

We can bound (ii) similarly.

$$(ii) \leq \max_s \left| V_{\hat{\pi}_n, \hat{f}_n}^R(s) - V_{\pi^*, f}^R(s) \right| \tag{141}$$

$$\leq \frac{\gamma}{1-\gamma} \max_{V(\cdot) \in \mathcal{V}} \max_{s, a} \left( \sup_{\eta \in [0, \frac{(2+\rho)}{\rho(1-\gamma)}]} \left\{ \left| \mathbb{E}_{P_f(s, a)}[(-V(s') + \eta)_+] \right| - \mathbb{E}_{P_{\hat{f}_n}(s, a)}[(-V(s') + \eta)_+] \right\} \right). \tag{142}$$

**Step (iii):** Next, we want to utilize the learning error bound (Equation (21)) that bounds the difference between the means of true nominal transition dynamics  $P_f$  and learned nominal transition dynamics  $P_{\hat{f}_n}$  to bound Equations (140) and (142).

We begin by bounding the difference  $\left| \mathbb{E}_{P_f(s, a)}[(-V(s') + \eta)_+] - \mathbb{E}_{P_{\hat{f}_n}(s, a)}[(-V(s') + \eta)_+] \right|$ , by the difference in means of  $P_f$  and  $P_{\hat{f}_n}$  in Lemma 28. Since Equation (140) has a max over all value functions, we introduce a covering number argument in Lemma 29 to reform it to a max over the functions in the  $\zeta$ -covering set. We then use Lemma 28 to obtain bounds in terms of maximum information gain  $\Gamma_{Nd}$  (Equation (18)) and  $\zeta$ . Further details regarding the covering number argument are deferred to Lemma 29. Then, we apply the result of Lemma 29 with  $\zeta = 1$  (defined in Lemma 29) on Equation (140). Then, it holds that

$$(i) \leq \max_s \left| V_{\hat{\pi}_n, f}^R(s) - V_{\hat{\pi}_n, \hat{f}_n}^R(s) \right| = \mathcal{O} \left( \left( \frac{(2+\rho)\gamma}{\rho(1-\gamma)^2} \right) \left( \frac{\beta_n(\delta)\sqrt{2ed^2\gamma_{nd}}}{\sigma\sqrt{n}} \right) \right). \tag{143}$$

Note that  $\beta_n$ , which appears in Lemma 2, has a logarithmic dependence on  $n$ . Similarly, from Equation (142), and Lemmas 28, 29, we obtain

$$(ii) \leq \max_s \left| V_{\hat{\pi}_n, \hat{f}_n}^R(s) - V_{\pi^*, f}^R(s) \right| = \mathcal{O} \left( \left( \frac{(2+\rho)\gamma}{\rho(1-\gamma)^2} \right) \left( \frac{\beta_n(\delta)\sqrt{2ed^2\gamma_{nd}}}{\sigma\sqrt{n}} \right) \right). \tag{144}$$

Note that we want to bound  $V_{\hat{\pi}_n, f}^R(s) - V_{\pi^*, f}^R(s) = (i) + (ii)$  over all  $s \in \mathcal{S}$ . Using  $\max_s |V_{\hat{\pi}_n, f}^R(s) - V_{\pi^*, f}^R(s)| \leq \max_s |V_{\hat{\pi}_n, \hat{f}_n}^R(s) - V_{\pi^*, f}^R(s)| + \max_s |V_{\hat{\pi}_n, \hat{f}_n}^R(s) - V_{\pi^*, f}^R(s)|$  and substituting  $M$  by  $1/(1-\gamma)$ , we obtain from Equation (143) and Equation (144)

$$\max_s |V_{\hat{\pi}_n, f}^R(s) - V_{\pi^*, f}^R(s)| = \mathcal{O}\left(\left(\frac{(2+\rho)\gamma}{\rho(1-\gamma)^2}\right)\left(\frac{\beta_n(\delta)\sqrt{2ed^2\gamma nd}}{\sigma\sqrt{n}}\right)\right).$$

Finally, to ensure that  $\max_s |V_{\hat{\pi}_n, f}^R(s) - V_{\pi^*, f}^R(s)| \leq \epsilon$ , it suffices to have

$$\max_s |V_{\hat{\pi}_n, f}^R(s) - V_{\pi^*, f}^R(s)| = \mathcal{O}\left(\left(\frac{(2+\rho)\gamma}{\rho(1-\gamma)^2}\right)\left(\frac{\beta_n(\delta)\sqrt{2ed^2\gamma nd}}{\sigma\sqrt{n}}\right)\right) = \epsilon.$$

Moving  $\sqrt{n}$  and  $\epsilon$  to opposite sides and squaring both sides, we obtain

$$n = \mathcal{O}\left(\left(\frac{(2+\rho)^2\gamma^2}{\rho^2(1-\gamma)^4}\right)\left(\frac{\beta_n(\delta)^2 2ed^2\gamma nd}{\sigma^2\epsilon^2}\right)\right).$$

■

**Lemma 27** (Simplification using Lemma 25 reformulation) *Let  $V$  be a value function from  $\mathcal{S} \rightarrow [0, 1/(1-\gamma)]$ . Then, it holds that*

$$\begin{aligned} \max_s \left| \inf_{\text{TV}(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s)) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] - \inf_{\text{TV}(p||P_f(s, \hat{\pi}_n(s)) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] \right| \leq \\ \max_{s, a} \sup_{\eta \in [0, \frac{(2+\rho)}{\rho(1-\gamma)}]} \{ |(\mathbb{E}_{P_f(s, a)}[(-V(s') + \eta)_+] - \mathbb{E}_{P_{\hat{f}_n}(s, a)}[(-V(s') + \eta)_+])| \}. \end{aligned}$$

**Proof** First note that,

$$\begin{aligned} \max_s \left| \inf_{\text{TV}(p||P_{\hat{f}_n}(s, \hat{\pi}_n(s)) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] - \inf_{\text{TV}(p||P_f(s, \hat{\pi}_n(s)) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] \right| \leq \\ \max_{s, a} \left| \inf_{\text{TV}(p||P_{\hat{f}_n}(s, a) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] - \inf_{\text{TV}(p||P_f(s, a) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] \right| \quad (145) \end{aligned}$$

Using Lemma 25 and focusing to bound right side of Equation (145) for one particular  $(s, a)$  state action pair, we obtain

$$\begin{aligned} \left| \inf_{\text{TV}(p||P_{\hat{f}_n}(s, a) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] - \inf_{\text{TV}(p||P_f(s, a) \leq \rho} \mathbb{E}_{s' \sim p} [V(s')] \right| \\ = \left| \sup_{\eta \in [0, \frac{(2+\rho)}{\rho(1-\gamma)}]} \left\{ -\mathbb{E}_{P_f(s, a)} \left[ -V(s') + \eta \right]_+ - \frac{(-ESI_{P_f(s, a)}(V(s')) + \eta)_+}{2} \rho + \eta \right\} - \right. \\ \left. \sup_{\eta \in [0, \frac{(2+\rho)}{\rho(1-\gamma)}]} \left\{ -\mathbb{E}_{P_{\hat{f}_n}(s, a)} \left[ -V(s') + \eta \right]_+ - \frac{(-ESI_{P_{\hat{f}_n}(s, a)}(V(s')) + \eta)_+}{2} \rho + \eta \right\} \right| \quad (146) \end{aligned}$$

$$\leq \sup_{\eta \in [0, \frac{(2+\rho)}{\rho(1-\gamma)}]} \{|\mathbb{E}_{P_f(s,a)}[(-V(s') + \eta)_+] - \mathbb{E}_{P_{\hat{f}_n}(s,a)}[(-V(s') + \eta)_+]\}|. \quad (147)$$

Here, Equation (147) is obtained using  $ESI_{P_f(s,a)}(V(s')) = ESI_{P_{\hat{f}_n}(s,a)}(V(s'))$  as shown in proof of Lemma 14 (Case-1). ■

**Lemma 28** (*Bound by difference between estimated model  $\hat{f}_n$  and true  $f$* ) Let  $V$  be a value function from  $\mathcal{S} \rightarrow [0, 1/(1-\gamma)]$ . Then, it holds that

$$|\mathbb{E}_{P_f(s,a)}[(-V(s') + \eta)_+] - \mathbb{E}_{P_{\hat{f}_n}(s,a)}[(-V(s') + \eta)_+]| \leq \left(\frac{(2+\rho)}{\rho(1-\gamma)}\right) \sigma^{-1} \|f(s,a) - \hat{f}_n(s,a)\|, \quad (148)$$

where  $P_{\hat{f}_n}(s,a) = \mathcal{N}(\hat{f}_n(s,a), \sigma^2 I)$  and  $P_f(s,a) = \mathcal{N}(f(s,a), \sigma^2 I)$  and  $\eta \in [0, \frac{(2+\rho)}{\rho(1-\gamma)}]$ .

**Proof**

$$\begin{aligned} & \left| \mathbb{E}_{P_f(s,a)}[(-V(s') + \eta)_+] - \mathbb{E}_{P_{\hat{f}_n}(s,a)}[(-V(s') + \eta)_+] \right| \\ &= \left| \int_{\mathbb{R}^d} \frac{1}{\sqrt{(2\pi\sigma^2)^d}} (-V(s') + \eta)_+ \left( e^{-\frac{\|x-f(s,a)\|^2}{2\sigma^2}} - e^{-\frac{\|x-\hat{f}_n(s,a)\|^2}{2\sigma^2}} \right) \right| \\ &\leq \int_{\mathbb{R}^d} \frac{1}{\sqrt{(2\pi\sigma^2)^d}} (-V(s') + \eta)_+ \left| e^{-\frac{\|x-f(s,a)\|^2}{2\sigma^2}} - e^{-\frac{\|x-\hat{f}_n(s,a)\|^2}{2\sigma^2}} \right| \\ &\stackrel{(i)}{\leq} \frac{(2+\rho)}{\rho(1-\gamma)} \int_{\mathbb{R}^d} \frac{1}{\sqrt{(2\pi\sigma^2)^d}} \left| e^{-\frac{\|x-f(s,a)\|^2}{2\sigma^2}} - e^{-\frac{\|x-\hat{f}_n(s,a)\|^2}{2\sigma^2}} \right| \\ &\stackrel{(ii)}{\leq} 2 \frac{(2+\rho)}{\rho(1-\gamma)} \cdot \text{TV}(P_{\hat{f}_n}(s,a), P_f(s,a)) \\ &\stackrel{(iii)}{\leq} 2 \frac{(2+\rho)}{\rho(1-\gamma)} \sqrt{\text{KL}(P_{\hat{f}_n}(s,a), P_f(s,a))/2} \\ &\stackrel{(iv)}{\leq} 2 \frac{(2+\rho)}{\rho(1-\gamma)} \sqrt{\|f(s,a) - \hat{f}_n(s,a)\|^2/4\sigma^2} \\ &\leq \frac{(2+\rho)}{\rho(1-\gamma)} \|f(s,a) - \hat{f}_n(s,a)\|/\sigma, \end{aligned}$$

where (i) follows from  $(-V(s') + \eta)_+^2 \leq \frac{(2+\rho)}{\rho(1-\gamma)}$  as  $\eta \leq \frac{(2+\rho)}{\rho(1-\gamma)}$ , (ii) follows from the definition of Total Variation (TV) distance between any two multivariate Gaussians, (iii) uses the Pinsker's inequality, and (iv) uses the formula for KL-divergence between multivariate Gaussian distributions. ■

**Lemma 29** ( $\zeta$ -cover construction) For  $\mathcal{V}$  denoting the set of value functions from  $\mathcal{S} \rightarrow [0, 1/(1-\gamma)]$ , with probability at least  $1 - \delta$  it holds that

$$\begin{aligned} \max_{V \in \mathcal{V}} \max_{s,a} \sup_{\eta \in [0, \frac{(2+\rho)}{\rho(1-\gamma)}]} \{ & |\mathbb{E}_{P_f(s,a)}[(-V(s') + \eta)_+] - \mathbb{E}_{P_{\hat{f}_n}(s,a)}[(-V(s') + \eta)_+]| \} \\ & \leq \mathcal{O} \left( \left( \frac{(2+\rho)}{\rho(1-\gamma)} \right) \left( \frac{\beta_n(\delta) \sqrt{2ed^2 \gamma_{nd}}}{\sigma \sqrt{n}} \right) \right), \end{aligned} \quad (149)$$

where  $\mathcal{N}_{\mathcal{V}}(\zeta)$  is the  $\zeta$ -cover for  $\mathcal{V}$ .

**Proof** Let  $\mathcal{N}_{\mathcal{V}}(\zeta)$  be the  $\zeta$ -cover of the set  $\mathcal{V}$ . By definition, there exists  $V' \in \mathcal{N}_{\mathcal{V}}(\zeta)$  such that  $\|V' - V\| \leq \zeta$  for every  $V \in \mathcal{V}$ .

$$\begin{aligned} & |\mathbb{E}_{P_f(s,a)}[(-V(s') + \eta)_+] - \mathbb{E}_{P_{\hat{f}_n}(s,a)}[(-V(s') + \eta)_+]| \\ & \leq |\mathbb{E}_{P_f(s,a)}[(-V(s') + \eta)_+] - \mathbb{E}_{P_f(s,a)}[(-V'(s') + \eta)_+]| \\ & \quad + |\mathbb{E}_{P_f(s,a)}[(-V'(s') + \eta)_+] - \mathbb{E}_{P_{\hat{f}_n}(s,a)}[(-V'(s') + \eta)_+]| \\ & \quad + |\mathbb{E}_{P_{\hat{f}_n}(s,a)}[(-V'(s') + \eta)_+] - \mathbb{E}_{P_{\hat{f}_n}(s,a)}[(-V(s') + \eta)_+]|. \end{aligned} \quad (150)$$

$$\stackrel{(i)}{\leq} 2\|V' - V\| + |\mathbb{E}_{P_f(s,a)}[(-V'(s') + \eta)_+] - \mathbb{E}_{P_{\hat{f}_n}(s,a)}[(-V'(s') + \eta)_+]|, \quad (151)$$

where (i) follows from Lemma 30. Using Equation (151), we bound uniformly over all  $V \in \mathcal{V}$ . Using Equation (151) we bound uniformly over all  $V \in \mathcal{V}$ ,

$$\max_{V \in \mathcal{V}} \max_{s,a} \sup_{\eta \in [0, \frac{(2+\rho)}{\rho(1-\gamma)}]} \{ |\mathbb{E}_{P_f(s,a)}[(-V(s') + \eta)_+] - \mathbb{E}_{P_{\hat{f}_n}(s,a)}[(-V(s') + \eta)_+]| \} \quad (152)$$

$$\leq \max_{V' \in \mathcal{N}_{\mathcal{V}}(\zeta)} \max_{s,a} \sup_{\eta \in [0, \frac{(2+\rho)}{\rho(1-\gamma)}]} \left\{ \left| 2\|V' - V\| + |\mathbb{E}_{P_f(s,a)}[(-V'(s') + \eta)_+] - \mathbb{E}_{P_{\hat{f}_n}(s,a)}[(-V'(s') + \eta)_+]| \right| \right\}$$

$$\stackrel{(ii)}{\leq} \max_{V' \in \mathcal{N}_{\mathcal{V}}(\zeta)} \max_{s,a} \sup_{\eta \in [0, \frac{(2+\rho)}{\rho(1-\gamma)}]} \left\{ \left| \mathbb{E}_{P_f(s,a)}[(-V'(s') + \eta)_+] - \mathbb{E}_{P_{\hat{f}_n}(s,a)}[(-V'(s') + \eta)_+] \right| \right\} + 2\zeta$$

$$\stackrel{(iii)}{\leq} \max_{V' \in \mathcal{N}_{\mathcal{V}}(\zeta)} \max_{s,a} \sup_{\eta \in [0, \frac{(2+\rho)}{\rho(1-\gamma)}]} \left\{ \left( \frac{(2+\rho)}{\rho(1-\gamma)} \right) \sigma^{-1} \|f(s,a) - \hat{f}_n(s,a)\| \right\} + 2\zeta$$

$$\stackrel{(iv)}{\leq} \mathcal{O} \left( \left( \frac{(2+\rho)}{\rho(1-\gamma)} \right) \left( \frac{\beta_n(\delta) \sqrt{2ed^2 \gamma_{nd}}}{\sigma \sqrt{n}} \right) \right) + 2\zeta \quad (153)$$

$$\stackrel{(v)}{\leq} \mathcal{O} \left( \left( \frac{(2+\rho)}{\rho(1-\gamma)} \right) \left( \frac{\beta_n(\delta) \sqrt{2ed^2 \gamma_{nd}}}{\sigma \sqrt{n}} \right) \right), \quad (154)$$

where (ii) follows from  $\|V' - V\| \leq \zeta$ , (iii) follows from Lemma 28, (iv) follows from Equation (21), and (v) follows from substituting  $\zeta = 1$  (or any constant). ■

**Lemma 30** For any two value functions  $V, V' : \mathcal{S} \rightarrow [0, \frac{1}{1-\gamma}]$ , it holds that

$$|\mathbb{E}_{P_f(s,a)}[(-V'(s') + \eta)_+] - \mathbb{E}_{P_f(s,a)}[(-V(s') + \eta)_+]| \leq \|V' - V\|. \quad (155)$$



**Proof** Noting that both the distributions are w.r.t. the same distribution  $P_f(s, a)$  we have,

$$\begin{aligned} & \mathbb{E}_{P_f(s,a)}[(-V'(s') + \eta)_+] - \mathbb{E}_{P_f(s,a)}[(-V(s') + \eta)_+] \\ & \leq \int_{s' \sim P_f(s,a)} \left( \mathbf{1}(V'(s') < \eta)(-V'(s') + \eta) - \mathbf{1}(V(s') < \eta)(-V(s') + \eta) \right) p_{P_f(s,a)}(s') ds'. \end{aligned} \quad (156)$$

Adding and subtracting  $\mathbf{1}(V(s') < \eta)(-V'(s') + \eta)$  to Equation (156), we obtain 2 terms,

$$i = \int_{s' \sim P_f(s,a)} \left( \mathbf{1}(V'(s') < \eta) - \mathbf{1}(V(s') < \eta) \right) (-V'(s') + \eta) p_{P_f(s,a)}(s') ds' \quad (157)$$

$$ii = \int_{s' \sim P_f(s,a)} \mathbf{1}(V(s') < \eta) \left( (-V'(s') + \eta) - (-V(s') + \eta) \right) p_{P_f(s,a)}(s') ds'. \quad (158)$$

Bounding i first,

$$i = \int_{s' \sim P_f(s,a)} \left( \mathbf{1}(V'(s') < \eta) - \mathbf{1}(V(s') < \eta) \right) (-V'(s') + \eta) p_{P_f(s,a)}(s') ds' \quad (159)$$

$$= \int_{s' \sim P_f(s,a)} \left( \mathbf{1}(V'(s') < \eta \leq V(s')) \right) (-V'(s') + \eta) p_{P_f(s,a)}(s') ds' \quad (160)$$

$$- \int_{s' \sim P_f(s,a)} \left( \mathbf{1}(V(s') < \eta < V'(s')) \right) (-V'(s') + \eta) p_{P_f(s,a)}(s') ds'$$

$$\leq \int_{s' \sim P_f(s,a)} \left( \mathbf{1}(V'(s') < \eta \leq V(s')) \right) (-V'(s') + V(s')) p_{P_f(s,a)}(s') ds' \quad (161)$$

$$- \int_{s' \sim P_f(s,a)} \left( \mathbf{1}(V(s') < \eta < V'(s')) \right) (-V'(s') + V(s')) p_{P_f(s,a)}(s') ds'$$

$$\leq \int_{s' \sim P_f(s,a)} \left( \mathbf{1}(V'(s') < \eta \leq V(s')) \right) (-V'(s') + V(s')) p_{P_f(s,a)}(s') ds' \quad (162)$$

$$+ \int_{s' \sim P_f(s,a)} \left( \mathbf{1}(V(s') < \eta < V'(s')) \right) (V'(s') - V(s')) p_{P_f(s,a)}(s') ds'$$

$$\leq \|V' - V\|. \quad (163)$$

Similarly bounding ii,

$$ii = \int_{s' \sim P_f(s,a)} \mathbf{1}(V(s') < \eta) \left( (-V'(s') + \eta) - (-V(s') + \eta) \right) p_{P_f(s,a)}(s') ds' \quad (164)$$

$$= \int_{s' \sim P_f(s,a)} \mathbf{1}(V(s') < \eta) \left( -V'(s') + V(s') \right) p_{P_f(s,a)}(s') ds' \quad (165)$$

$$\leq \int_{s' \sim P_f(s,a)} \mathbf{1}(V(s') < \eta) \left| -V'(s') + V(s') \right| p_{P_f(s,a)}(s') ds' \quad (166)$$

$$\leq \|V' - V\|. \quad (167)$$

Using Equations (163) and (167) we get the desired result.  $\blacksquare$

## Appendix F. Additional Experiments and Details

In this section, we report additional experiments and discuss further details of our experimental setup.

**Environments:** We consider the OpenAI’s gym (8) environments of swing-up Pendulum, Cartpole and Reacher, respectively. Pendulum has a 2-dimensional state space and scalar actions ((38)). For Cartpole, we consider a scalar continuous action space as done in (38), while states are 4-dimensional. Reacher, instead, consists of a 2DOF robot arm with 8-dimensional states. For each environment we test our approach against various perturbations as outlined below.

**Baselines:** We compare our approach, which we denote as MVR+RFQI, with the following baselines:

- MVR+FQI: This is a natural non-robust baseline that consists of computing a non-robust policy via the Fitted Q-Iteration (FQI) algorithm (23) on the same offline data used by MVR+RFQI,
- Soft Actor-Critic (SAC) (25), or Model Predictive Control (MPC) (9; 13), as model-free methods which compute non-robust policies *interacting with the nominal environment* (in case of MPC, the latter is used for planning),
- RFQI (42), which also requires the nominal environment and uses  $10^6$  offline data collected by SAC or MPC to train a robust policy,
- FQI (23), which trains a non-robust policy from the same data.

**Training:** Model-free methods are trained directly on the nominal environments. In particular, for Pendulum and Reacher we train SAC until convergence for  $10^4$  and  $10^6$  steps, respectively. On the continuous actions Cartpole, instead, we run MPC following the implementation of (45; 38) which requires a total of 2250 planning interactions to select the optimal action at each step. Depending on the environment, we utilize SAC or MPC mixed with an  $\epsilon$ -greedy rule to collect  $10^6$  offline data. These are used to train the offline methods RFQI and FQI as done in (42). For the model-based approaches, instead, we first run MVR for a sufficiently informative number of samples (60 for Pendulum, 150 for Cartpole and 2000 for Reacher) to obtain an estimated model  $\hat{f}_n$ . Then, we use SAC (trained against model  $\hat{f}_n$ ) or MPC to collect  $10^6$  offline data on such estimated environment. These data are then used to train MVR+RFQI and MVR+FQI.

**Evaluation:** For each environment, we evaluate the computed policy against different perturbation types and magnitudes. For Cartpole, we perturb the magnitude of the actuation force. Its nominal value is 10 and we perturb up to 300%. Also, we consider perturbations to gravity in the range of (-100%,100%) with the nominal value being 9.82. For the Pendulum, we consider perturbations to the length of the pendulum and action perturbations (where a random action is chosen with  $\epsilon$  probability). Finally, in the case of Reacher we consider perturbations to the joint’s stiffness (from 0 to 100) coupled with perturbations of the joint’s equilibrium position. Further details on the chosen perturbations and hyperparameters used are provided in Appendix F.

All experiments were run with GPU clusters: 10xNvidia 32Gb Tesla V100 with Intel(R) processors (2 cores, 2.50 GHz) and 256Gb RAM. For all the experiments, we use the environment implementations of (38) as done in <https://github.com/fusion-ml/trajectory-information-rl/tree/main>. Also, to learn the environment transition model, we use the same corresponding GP hyperparameters proposed by (38). For the offline RFQI/FQI algorithms we follow the implementation of (42; 10) in <https://github.com/zaiyan-x/RFQI>. We use the same default hyperparameters as used in their code except for training steps, batch size and robustness radius  $\rho$  (for RFQI) which we tune depending on the environment as outlined next. For SAC in Pendulum experiments, we use the implementation and hyperparameters of <https://github.com/DLR-RM/rl-baselines3-zoo>. Whereas, for SAC in Reacher experiments, we use the implementation and hyperparameters of <https://github.com/fusion-ml/bac-baselines>, <https://github.com/IanChar/rlkit2> (as done in (38)).

**Pendulum:** In Pendulum experiments, we construct the learned model using 60 samples from the true environment. Then, we train a SAC policy on such a model for  $2 * 10^4$  steps and use it (with the probability of choosing a random action being 0.3 or 0.5) to generate  $10^6$  offline data (these are used both for MVR+RFQI and MVR+FQI). For training steps and batch size we consider the following combinations:  $\{2000 - 100', 5000 - 100', 10000 - 100', 20000 - 100', 35000 - 100', 50000 - 100', 5000 - 500', 5000 - 1000'\}$ . We combine all these combinations with the following values of  $\rho - \{0.1, 0.2, 0.3, 0.5, 0.6, 0.7, 0.8, 0.9\}$ . For each algorithm, we pick the best-performing combination in terms of average reward over 20 episodes for all (or most) perturbation values. We do this separately for length perturbations and action perturbations. In the length perturbation, the pendulum’s length is changed from its nominal value to a new value depending on the perturbation percentage. In the action perturbation, a random action is chosen instead of the action chosen by the policy with various probabilities ranging from  $[0, 1]$ . We detail the optimal hyperparameters we realized for each algorithm in Table 2 for the length and action perturbation, respectively. Moreover, we plot the average performance (over 20 episodes) of the different baselines w.r.t. length and action perturbations in Figure 2. We notice that in the case of length perturbation, the robust algorithms (RFQI and MVR+RFQI) outperform the corresponding non-robust baselines. In the case of action perturbations, we observe all algorithms except for SAC achieve similar performance.

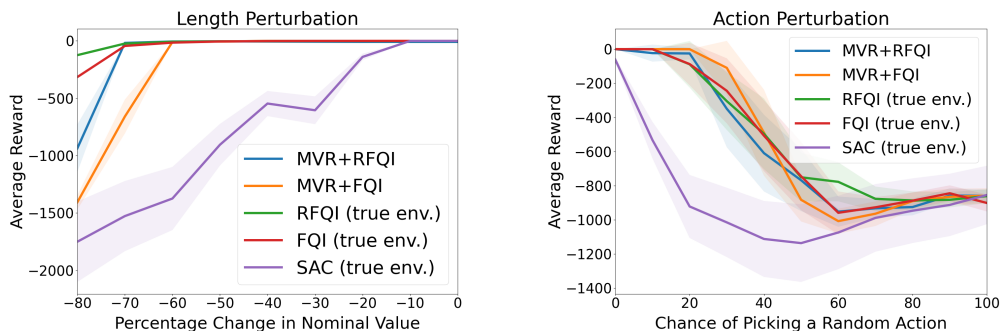


Figure 2: Pendulum experiments.

	TRAINING STEPS	BATCH-SIZE	$\rho$	RANDOM ACTION PROBABILITY (DATASET)
MVR+RFQI	5000	100	0.3	0.5
MVR+FQI	2000	100	-	0.5
RFQI	2000	100	0.9	0.5
FQI	5000	500	-	0.5

---

	TRAINING STEPS	BATCH-SIZE	$\rho$	RANDOM ACTION PROBABILITY (DATASET)
MVR+RFQI	20000	100	0.5	0.3
MVR+FQI	50000	100	-	0.3
RFQI	50000	100	0.1	0.5
FQI	5000	500	-	0.5

Table 2: Hyperparameters for Pendulum - length perturbation (top) and action perturbation (bottom).

**Cartpole:** In Cartpole experiments, we construct the learned model using 150 samples from the true environment. Then, we run MPC on such a model following the implementation and hyperparameters of (38; 45) requiring 2250 samples to calculate the optimal action at each step and use it (with the probability of choosing a random action being 0.3) to generate  $10^6$  offline data for MVR+RFQI and MVR+FQI. For training steps and batch size, we test the following combinations:  $\{2000 - 100', 5000 - 100', 10000 - 100', 20000 - 100', 35000 - 100', 50000 - 100', 5000 - 500', 5000 - 1000'\}$ , and consider radii  $\rho$  in  $\{0.1, 0.2, 0.3, 0.5, 0.6, 0.7, 0.8, 0.9\}$ . We consider perturbations of the force magnitude and the gravity, whereby the actuation force/gravity is changed from its nominal value to a new value depending on the perturbation percentage. We report the best-performing (average over 20 episodes) hyperparameters for each algorithm in Table 3. Such parameters were observed to be a good choice for both perturbation types. Finally, we plot the average performance (over 20 episodes) of the different baselines w.r.t. force magnitude and gravity perturbations in Figure 3. We notice that in both perturbations, the robust algorithms (RFQI and MVR+RFQI) outperform the corresponding non-robust baselines.

	TRAINING STEPS	BATCH-SIZE	$\rho$	RANDOM ACTION PROBABILITY (DATASET)
MVR+RFQI	5000	500	0.5	0.3
MVR+FQI	50000	100	-	0.3
RFQI	5000	100	0.3	0.3
FQI	10000	100	-	0.3

Table 3: Hyperparameters for Cartpole.

**Reacher:** In Reacher experiments, we construct the learned model using 2000 samples from the true environment. Then, we train a SAC policy on such a model for  $10^6$  steps and use it (with the probability of choosing a random action being 0.3) to generate  $10^6$  offline data for

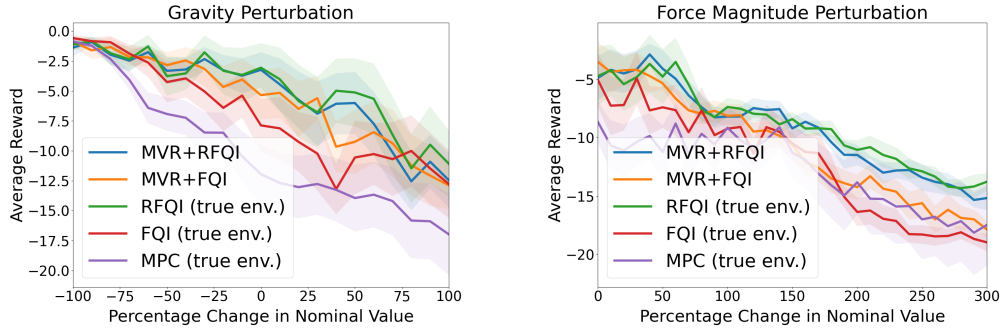


Figure 3: Cartpole experiments.

MVR+RFQI and MVR+FQI. For training steps and batch size, we consider the following combinations:  $\{10000 - 500', 20000 - 500', 40000 - 500', 80000 - 500', 160000 - 1000'\}$ , while we consider radii  $\rho$  in  $\{0.1, 0.3, 0.5, 0.7, 0.9\}$ . We consider perturbations of the joint stiffness subject to different equilibrium positions, the latter represented by the 'Springref' parameter which we take to be 50 or 100. In both perturbation types, the joint stiffness is changed from its nominal value of 0 to a new value depending on the perturbation magnitude. Best-performing hyperparameters' configurations are reported in Table 4. We plot the average performance (over 20 episodes) of the different baselines in Figure 4. Similar to the other environments, we observe the robust algorithms (RFQI and MVR+RFQI) outperform the corresponding non-robust baselines.

	TRAINING STEPS	BATCH-SIZE	$\rho$	RANDOM ACTION PROBABILITY (DATASET)
MVR+RFQI	10000	500	0.5	0.3
MVR+FQI	20000	500	-	0.3
RFQI	40000	500	0.1	0.3
FQI	20000	500	-	0.3

Table 4: Hyperparameters for Reacher.

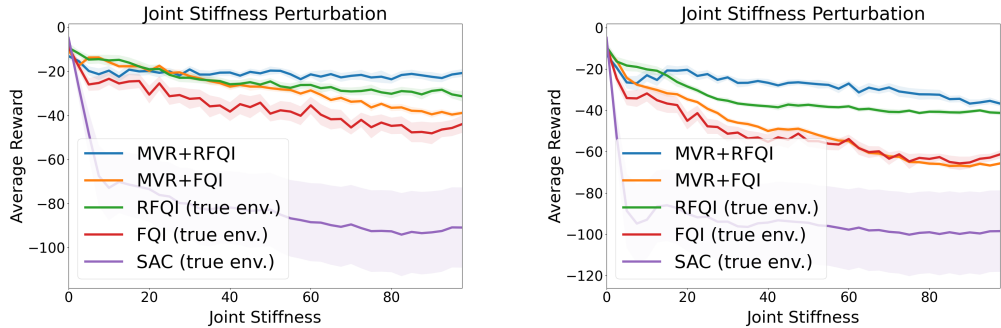


Figure 4: Reacher experiments with 'Springref' parameter set to 50 (left) or 100 (right).