

# STATISTICAL GUARANTEES FOR OFFLINE DOMAIN RANDOMIZATION

Arnaud Fickinger\*<sup>1</sup>  
UC Berkeley

Abderrahim Bendahi\*<sup>†2</sup>  
École Polytechnique

Stuart Russell<sup>3</sup>  
UC Berkeley

## ABSTRACT

Reinforcement-learning (RL) agents often struggle when deployed from simulation to the real-world. A dominant strategy for reducing the sim-to-real gap is domain randomization (DR) which trains the policy across many simulators produced by sampling dynamics parameters, but standard DR ignores offline data already available from the real system. We study offline domain randomization (ODR), which first fits a distribution over simulator parameters to an offline dataset. While a growing body of empirical work reports substantial gains with algorithms such as DROPO (Tiboni et al., 2023), the theoretical foundations of ODR remain largely unexplored. In this work, we cast ODR as a maximum-likelihood estimation over a parametric simulator family and provide statistical guarantees: under mild regularity and identifiability conditions, the estimator is weakly consistent (it converges in probability to the true dynamics as data grows), and it becomes strongly consistent (i.e., it converges almost surely to the true dynamics) when an additional uniform Lipschitz continuity assumption holds. We examine the practicality of these assumptions and outline relaxations that justify ODR’s applicability across a broader range of settings. Taken together, our results place ODR on a principled footing and clarify when offline data can soundly guide the choice of a randomization distribution for downstream offline RL.

## 1 INTRODUCTION

In recent years, RL has achieved many empirical successes, attaining human-level performance in tasks such as games (Mnih et al., 2013; Silver et al., 2016), robotics (Kalashnikov et al., 2018; Schulman et al., 2015), and recommender systems (Afsar et al., 2021; Chen et al., 2021). Yet, RL algorithms often require vast amounts of training data to learn effective policies, which severely limits their applicability in real world settings where data collection is expensive, time-consuming, or unsafe (Levine et al., 2020; Kiran et al., 2020).

*Sim-to-real transfer* tackles this problem by learning in simulation and transferring the resulting policy to the real world (Sadeghi & Levine, 2016; Tan et al., 2018; Zhao et al., 2020). However, although simulation provides fast and safe data collection, inevitable discrepancies between the simulated dynamics and the real world, commonly termed the *sim-to-real gap*, typically induce a drop in performance upon deployment.

One of the most widely-used approaches to bridge this gap is *domain randomization* (DR). Rather than training on a single fixed simulator, DR defines a family of simulators parameterized by physical factors (e.g., masses, friction coefficients, sensor noise) and at the start of each episode *randomly samples* one instance from this family for training. DR has enabled zero-shot transfer in robotic control (Tobin et al., 2017; Sadeghi & Levine, 2016), dexterous manipulation (OpenAI et al., 2018) and agile locomotion (Peng et al., 2017).

---

\*Equal contribution.

<sup>†</sup>Work done during internship at UC Berkeley.

<sup>1</sup>arnaud.fickinger@berkeley.edu

<sup>2</sup>abderrahim.bendahi@polytechnique.edu

<sup>3</sup>russell@berkeley.edu

Despite this empirical track record, the choice of *how* to randomize is a fundamental challenge. In the original form of DR (Tobin et al., 2017; Sadeghi & Levine, 2016), broad *uniform* ranges that look reasonable for every parameter are chosen. While recent theoretical work (Chen et al., 2022) shows that such *uniform DR* (UDR) can indeed bound the sim-to-real gap, the bound unfavorably scales in  $O(N^3 \log(N))$  with respect to the number of candidate simulators, in part because UDR ignores any data already available from the target system.

In contrast, *Offline Domain Randomization* exploits a static dataset from the real environment before policy training to fit a sampling distribution that concentrates on plausible dynamics while remaining stochastic. Empirically, ODR variants such as DROID (Tsai et al., 2021) or DROPO (Tiboni et al., 2023) recover parameter distributions that explain the data and yield stronger zero-shot transfer than hand-tuned UDR. Yet, to the best of our knowledge, ODR lacks a principled foundation: we do not know (i) whether the fitted distribution converges to the real dynamics as data grows, nor (ii) how much it actually reduces the sim-to-real gap compared with UDR.

### Our Contributions:

- **Weak consistency (Section 4).** We formalize ODR as maximum-likelihood estimation over a parametric simulator family and prove *weak consistency*: under mild regularity, positivity, and identifiability assumptions, empirical maximizers converge in probability to the population maximizers.
- **Strong consistency (Section 5).** Adding a single *uniform Lipschitz continuity* assumption on the likelihood, we upgrade convergence to *strong consistency*: the ODR estimator converges almost surely to the true parameter when it is uniquely identified.
- **Assumptions in practice: discussion and relaxations (Section 6).** We analyze when the assumptions hold and provide drop-in relaxations and diagnostics: replacing i.i.d. by strict stationarity and ergodicity for the, weakening mixture positivity via a logarithmic tail condition, and giving simple sufficient conditions that imply the uniform Lipschitz requirement.

## 2 RELATED WORKS

**Sim-to-real transfer** The *sim-to-real gap* has led to extensive research in sim-to-real transfer. Early works exploited system identification or progressive networks to adapt controllers online (Floreato et al., 2008; Kober et al., 2013), while more recent efforts have focused on purely offline training in high-fidelity simulators. Although zero-shot transfer has been demonstrated for specific settings such as legged locomotion (Peng et al., 2017), dexterous manipulation (Chebotar et al., 2018; OpenAI et al., 2018) and visuomotor control (Rusu et al., 2016) a noticeable performance gap persists in unstructured environments. Similar ideas have been explored in autonomous driving (Pouyanfar et al., 2019; Niu et al., 2021).

**Domain randomization** Domain randomization (DR) varies environment parameters at every training episode with the goal of producing policies that generalize across the induced simulator family. Vision-based DR first showed zero-shot transfer for quadrotor flight from purely synthetic images (Sadeghi & Levine, 2016), and dynamics randomization extended this success to legged robots and manipulation (OpenAI et al., 2018). To avoid manual tuning of randomization ranges, online methods adapt the DR distribution using real-world feedback. Ensemble-based robust optimization and Bayesian optimization techniques refine parameters via real rollouts (Rajeswaran et al., 2016; Muratore et al., 2020), while meta RL further accelerates adaptation (Clavera et al., 2018; Arndt et al., 2019). However, these require repeated—and potentially unsafe—hardware interactions during training.

**Offline domain randomization** A growing line of work aims to find the best strategy to perform domain randomization from a fixed offline dataset, obviating any further real-world trials. DROID (Tsai et al., 2021) tunes simulator parameters using CMA-ES (Hansen & Ostermeier, 2001; Hansen, 2006) with the  $L^2$  distance between a single human demonstration and its simulated counterpart as objective function. BayesSim (Ramos et al., 2019) trains a conditional density estimator to predict a posterior over simulator parameters given offline off-policy rollouts. Most recently, DROPO (Tiboni

et al., 2023) introduces a likelihood-based framework that fits both the mean and covariance of a Gaussian parameter distribution by maximizing the log-likelihood of the offline data under a mixture simulator. This approach recovers rich uncertainty estimates, handles non-differentiable black-box simulators via gradient-free optimizers, and outperforms DROID, BayesSim and uniform DR in zero-shot transfer on standard benchmarks without any on-policy real-world interaction.

**Theoretical analyses** Let  $M$  be the number of candidate simulators and  $H$  the horizon length. Chen et al. (2022) modeled uniform DR as a *latent MDP* and proved that the performance gap between the optimal policy in the true system and the policy trained with DR scales as  $O(M^3 \log(MH))^1$  in the case where the simulator class is finite and separated and  $O(\sqrt{M^3 H \log(MH)})$  in the finite non-separated simulator class case. Other works have studied the information-theoretical limit of sim-to-real transfer (Jiang, 2018), PAC-style guarantees via approximate simulators (Feng et al., 2019) and generalization in rich-observation MDPs (Zhong et al., 2019; Krishnamurthy et al., 2016). But none address the statistical benefits of offline DR. Our work bridges this gap by providing the first consistency proofs and finite-sample gap bounds for offline DR, thereby unifying empirical successes and theoretical understanding in a single framework.

### 3 PROBLEM SETUP AND ODR FORMULATION

**Episodic MDPs** We consider the episodic RL setting where each MDP corresponds to  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, P_{\mathcal{M}}, R, H, s_1)$ .  $\mathcal{S}$  is the set of states,  $\mathcal{A}$  is the set of actions,  $P_{\mathcal{M}}: \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{A})$  is the transition probability matrix,  $R: \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  is the reward function,  $H$  is the number of steps of each episode, and  $s_1$  is the initial state at step  $h = 1$ ; we assume w.l.o.g. that the agent starts from the same state in each episode.

At step  $h \in [H]$ , the agent observes the current state  $s_h \in \mathcal{S}$ , takes action  $a_h \in \mathcal{A}$ , receives reward  $R(s_h, a_h)$ , and moves to state  $s_{h+1}$  with probability  $P_{\mathcal{M}}(s_{h+1} | s_h, a_h)$ . The episode ends when state  $s_{H+1}$  is reached.

A policy  $\pi$  is a sequence  $\{\pi_h\}_{h=1}^H$  where each  $\pi_h$  maps histories  $\text{traj}_h = \{(s_1, a_1, \dots, s_h)\}$  to action distributions. Denote by  $\Pi$  the set of all such history-dependent policies. We denote by  $V_{\mathcal{M}, h}^{\pi}: \mathcal{S} \rightarrow \mathbb{R}$  the value function at step  $h$  under policy  $\pi$  on MDP  $\mathcal{M}$ , i.e.,  $V_{\mathcal{M}, h}^{\pi}(s) := \mathbb{E}_{\mathcal{M}, \pi} \left[ \sum_{t=h}^H R(s_t, a_t) \mid s_h = s \right]^2$ . We use  $\pi_{\mathcal{M}}^*$  to denote the optimal policy for the MDP  $\mathcal{M}$ , and  $V_{\mathcal{M}, h}^*$  to denote the optimal value under the optimal policy at step  $h$ .

We fix a *simulator class*  $\mathcal{U} = \{\mathcal{M}_{\xi} : \xi \in \Xi \subset \mathbb{R}^d\}$  of candidate MDPs that share  $(\mathcal{S}, \mathcal{A}, R, H, s_1)$  but can differ in  $P_{\mathcal{M}}$  via the physical parameter vector  $\xi$ . The unknown *real-world* environment is  $\mathcal{M}^* = \mathcal{M}_{\xi^*} \in \mathcal{U}$ . We assume full observability and that the learner can interact freely with any  $\mathcal{M} \in \mathcal{U}$  in simulation, but never observes  $\xi^*$  directly.

**Sim-to-real Transfer Problem** Given access to the simulators in  $\mathcal{U}$ , the goal is to output a policy  $\pi$  that attains high return when executed in the real-world MDP  $\mathcal{M}^*$ . We quantify performance via the *sim-to-real gap* which is defined as the difference between the value of the learned policy  $\pi$  during the simulation phase (or training phase), and the value of an optimal policy for the real world, i.e.

$$\text{Gap}(\pi) := V_{\mathcal{M}^*, 1}^*(s_1) - V_{\mathcal{M}^*, 1}^{\pi}(s_1).$$

**Domain Randomization** Domain randomization specifies a prior distribution  $\nu$  over parameters  $\Xi$  and thus over  $\mathcal{U}$ . Sampling  $\xi \sim \nu$  at the start of every episode induces a *latent MDP* (LMDP) whose optimal Bayes policy is

$$\pi_{\text{DR}}^* := \arg \max_{\pi \in \Pi} \mathbb{E}_{\xi \sim \nu} [V_{\mathcal{M}_{\xi}, 1}^{\pi}(s_1)].$$

In practice we approximate  $\pi_{\text{DR}}^*$  with any RL algorithm that trains in the simulator while resampling  $\xi \sim \nu$  each episode.

<sup>1</sup>The original paper derived a looser bound, see Section A.1 for a tighter derivation.

<sup>2</sup>Since the policy  $\pi$  is allowed to be non Markovian, this quantity can be defined using the history  $H_h = \{s_1, \dots, s_h\}$  as follows:  $\mathbb{E}_{\mathcal{M}, \pi} \left[ \sum_{t=h}^H R(s_t, a_t) \mid s_h = s \right] = \mathbb{E}_{H_h | s_h = s} \mathbb{E}_{\mathcal{M}, \pi} \left[ \sum_{t=h}^H R(s_t, a_t) \mid H_h \right]$ .

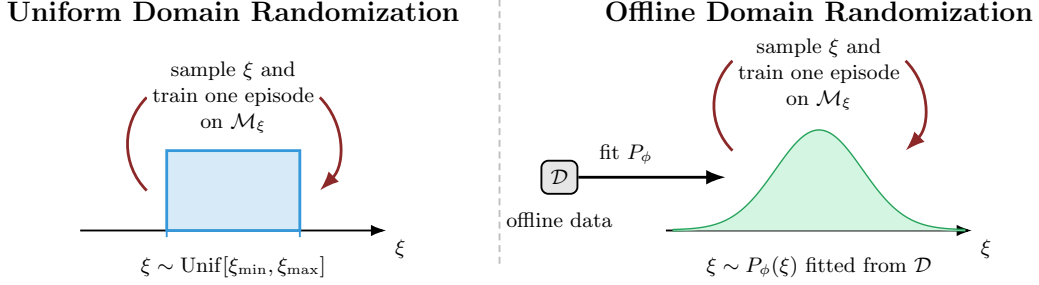


Figure 1: Conceptual comparison between Uniform Domain Randomization (left) and Offline Domain Randomization (right).

**Offline Domain Randomization** ODR assumes an offline data set  $\mathcal{D} = \{(s_i, a_i, s'_i)\}_{i=1}^N$  of i.i.d. transitions collected in the real system  $\mathcal{M}^*$  under some unknown behavior policy. The aim is to estimate a distribution  $p^*$  over  $\Xi$  that explains the data and can later be used for policy training. We restrict  $p_\phi(\xi) = \mathcal{N}(\mu, \Sigma)$ <sup>3</sup> and learn  $\phi$  by maximum likelihood:

$$p^*(\xi) = \arg \max_{p_\phi(\xi)} \prod_{(s_t, a_t, s_{t+1}) \in \mathcal{D}} \mathbb{E}_{\xi \sim p_\phi(\xi)} [P_\xi(s_{t+1} | s_t, a_t)] \quad (1)$$

$$= \arg \max_{p_\phi(\xi)} \sum_{(s_t, a_t, s_{t+1}) \in \mathcal{D}} \log [\mathbb{E}_{\xi \sim p_\phi(\xi)} [P_\xi(s_{t+1} | s_t, a_t)]] . \quad (2)$$

We justify that this formulation is well-motivated in wSection A.2.

Finally, we train a policy with the learned distribution:

$$\pi_{\text{ODR}}^* := \arg \max_{\pi \in \Pi} \mathbb{E}_{\xi \sim p^*} [V_{\mathcal{M}_{\xi,1}}^\pi(s_1)],$$

expecting  $\pi_{\text{ODR}}^*$  to transfer with lower gap thanks to the data-informed parameter distribution.

A conceptual comparison between Uniform Domain Randomization and Offline Domain Randomization is illustrated in Figure 1.

## 4 WEAK CONSISTENCY OF THE ODR ESTIMATOR

### 4.1 TECHNICAL ASSUMPTIONS

Before stating the theoretical guarantees for ODR, we introduce some mild assumptions of regularity and identifiability that will be useful for our proofs.

The following assumption assures that  $P_\xi$  is regular in the following sense.

**Assumption 1** (Simulator Regularity). *There exists a  $\sigma$ -finite measure  $\lambda$  on  $\mathcal{S}$  and a constant  $K < \infty$  such that for all  $\xi \in \Xi$  and  $(s, a, s')$*

$$P_\xi(ds' | s, a) = p_\xi(s' | s, a) \lambda(ds'), \quad 0 \leq p_\xi(s' | s, a) \leq K, \quad (3)$$

and  $\xi \mapsto p_\xi(s' | s, a)$  is continuous.

Notice that when  $\mathcal{S}$  is finite, and  $\lambda$  is the counting measure on  $\mathcal{S}$ , then the first assumption clearly holds with  $K = 1$  because  $p_\xi(s' | s, a) = P_\xi(\{s'\} | s, a) \leq 1$ . In this case, it suffices for the mass probability to depend continuously on  $\xi$  in order to verify Assumption 1. Another case where this continuity holds is the Gaussian case  $p_\xi(s'|s, a) = \mathcal{N}(s'; A(\xi)s + B(\xi)a, C(\xi))$ , where  $A(\xi), B(\xi), C(\xi)$  are matrices that vary continuously in  $\xi$ .

<sup>3</sup>The Gaussian parameterization over  $\xi$  is only a modeling choice, not a mathematical necessity. Any other parametric family for  $P_\phi$  that satisfies our upcoming assumptions could be substituted without changing the arguments.

**Assumption 2** (Parameter-Space Compactness). *We fit  $\phi = (\mu, \Sigma)$  in  $\Phi = \{\mu \in \tilde{\Xi} : 0 \preceq \Sigma \preceq \sigma_{\max} I\}$  where  $\tilde{\Xi}$  is compact, hence  $\Phi$  is compact in the product topology.*

This is a natural assumption, since in practice one always has prior bounds on each physical parameter, yielding a known compact search region.

Furthermore, we assume that all the transitions that appear in our dataset correspond to positive mixture probability. More formally,

**Assumption 3** (Mixture Positivity). *There exists some constant  $c > 0$  such that the induced kernel*

$$q_\phi(s' | s, a) := \mathbb{E}_{\xi \sim P_\phi(\xi)} [p_\xi(s' | s, a)] = \int p_\xi(s' | s, a) P_\phi(d\xi), \quad (4)$$

satisfies  $q_\phi(s' | s, a) \geq c > 0$  for every  $(s, a, s') \in \mathcal{D}$  and every  $\phi \in \Phi$ .

This guarantees that every transition in the dataset lies within the support of the simulator under the learned domain randomization distribution, so the log-likelihood is always well defined.

Furthermore, we assume that the only mixture distribution which exactly recovers the true transition kernel is the degenerate distribution concentrated at the true parameters  $\xi^*$ .

**Assumption 4** (Identifiability). *Let  $\mu$  be the dataset's distribution. If for  $\mu$ -almost every  $(s, a)$   $q_\phi(\cdot | s, a) = p_{\xi^*}(\cdot | s, a)$ , then  $\phi = (\xi^*, 0)$ .*

## 4.2 NOTATION FOR ODR

Throughout this work, we use a capital letter,  $P$ , to denote a probability distribution, and the corresponding lowercase letter,  $p$ , to denote its probability density (or mass) function.

We define the empirical and population log-likelihoods by

$$L_N(\phi) := \frac{1}{N} \sum_{i=1}^N a(X_i, \phi), \quad L(\phi) := \mathbb{E}_{X \sim P_{\xi^*}} [a(X, \phi)], \quad (5)$$

where  $X_i = (s_i, a_i, s'_i)$  is the  $i$ -th transition in  $\mathcal{D}$ , and  $X = (s, a, s')$  is a generic transition. The function  $a$  is defined by

$$a(x, \phi) := \log q_\phi(s' | s, a) = \log \int_\xi p_\xi(s' | s, a) p_\phi(\xi) d\xi. \quad (6)$$

## 4.3 MAIN THEOREM

The first lemma proves the uniqueness of the maximizer of the population log-likelihood  $L$ . A detailed proof of this lemma can be found in Section B.

**Lemma 1** (Uniqueness of the Population Maximizer). *Under assumptions 1, 3 and 4, the population log-likelihood*

$$L(\phi) = \mathbb{E}_{(s, a, s') \sim P_{\xi^*}} [\log q_\phi(s' | s, a)]$$

where  $q_\phi(s' | s, a) = \int P_\xi(s' | s, a) P_\phi(d\xi)$ , has the unique maximizer  $\phi^* = (\mu^*, \Sigma^*) = (\xi^*, 0)$ .

We now state our first consistency result for ODR.

**Theorem 1** (Weak Consistency of ODR). *Under Assumptions 1, 2, 3 and 4, any measurable maximizer  $\hat{\phi}_N \in \arg \max_{\phi \in \Phi} L_N(\phi)$  satisfies  $\hat{\phi}_N \xrightarrow[N \rightarrow \infty]{P} \phi^*$ .*

Theorem 1 guarantees that with a sufficiently large offline dataset, ODR recovers a distribution arbitrarily close to the true parameter  $\xi^*$ .

The following lemma is particularly strong: it establishes uniform convergence in probability of  $L_N$ .

**Lemma 2.** *The function  $\phi \mapsto L(\phi)$  is uniformly continuous on  $\Phi$ , and furthermore*

$$\sup_{\phi \in \Phi} |L_N(\phi) - L(\phi)| \xrightarrow[N \rightarrow \infty]{P} 0. \quad (7)$$

The proof of this lemma relies on a *uniform law of large numbers* (ULLN) -in particular the ULLN for *Glivenko-Cantelli* classes from Newey & McFadden (1994)- and is deferred to Section B. In contrast, the ordinary law of large numbers only guarantees that for each *fixed*  $\phi$  one has  $L_N(\phi) \rightarrow L(\phi)$  in probability, i.e.,  $|L_N(\phi) - L(\phi)| \rightarrow 0$  for that particular  $\phi$ . This pointwise convergence does *not* imply that  $\sup_{\phi \in \Phi} |L_N(\phi) - L(\phi)| \rightarrow 0$ , which is exactly what the ULLN provides. Uniform convergence over all  $\phi \in \Phi$  is crucial to control the behavior of the empirical maximizers and hence to establish the consistency of our estimator.

The following lemma formalizes a uniform separation property: any parameter  $\phi$  lying outside an  $\epsilon$ -neighborhood of the true maximizer  $\phi^*$  must have its population log-likelihood at least  $\eta > 0$  below  $L(\phi^*)$ .

**Lemma 3.** *Let  $\phi^*$  be the unique maximizer of  $L$ . We have*

$$\forall \epsilon > 0, \exists \eta(\epsilon) > 0, \forall \phi \in \Phi, \|\phi^* - \phi\| \geq \epsilon \implies L(\phi^*) - L(\phi) \geq \eta(\epsilon) > 0. \quad (8)$$

The proof of Lemma 3 is deferred to Section B.

*Proof of Theorem 1.* We consider a sequence of measurable maximizers  $\hat{\phi}_N \in \arg \max_{\phi \in \Phi} L_N(\phi)$ . Let  $\epsilon > 0$  be a fixed positive real number. Our goal is to prove that

$$P\left(\|\hat{\phi}_N - \phi^*\| \geq \epsilon\right) \xrightarrow{N \rightarrow \infty} 0. \quad (9)$$

Using Lemma 3, we conclude that there exists some  $\eta > 0$  such that  $\forall \phi \in \Phi$  if  $\|\phi^* - \phi\| \geq \epsilon$  then  $L(\phi^*) - L(\phi) \geq \eta > 0$ . Now, let  $E_\eta$  be the event

$$E_\eta = \left\{ \sup_{\phi \in \Phi} |L_N(\phi) - L(\phi)| < \eta/3 \right\} \quad (10)$$

then under  $E_\eta$ , if  $\|\phi^* - \phi\| \geq \epsilon$  we have

$$L_N(\phi^*) = L_N(\phi^*) - L(\phi^*) + L(\phi^*) \geq -|L_N(\phi^*) - L(\phi^*)| + L(\phi^*) \geq -\eta/3 + L(\phi^*), \quad (11)$$

since under  $E_\eta$ ,  $-|L_N(\phi) - L(\phi)| \geq -\eta/3$ , similarly

$$L(\phi^*) \geq L(\phi) + \eta = L(\phi) - L_N(\phi) + L_N(\phi) + \eta \geq -|L_N(\phi) - L(\phi)| + L_N(\phi) + \eta. \quad (12)$$

and combining these two inequalities gives  $L_N(\phi^*) \geq L_N(\phi) + \eta/3$ . This proves that, under  $E_\eta$ ,  $\hat{\phi}_N \in \mathcal{B}(\phi^*, \epsilon) := \{\phi \in \Phi : \|\phi - \phi^*\| < \epsilon\}$  thus  $\{\|\hat{\phi}_N - \phi^*\| \geq \epsilon\} \subset E_\eta^c$ , which yields

$$P(\|\hat{\phi}_N - \phi^*\| \geq \epsilon) \leq P\left(\sup_{\phi \in \Phi} |L_N(\phi) - L(\phi)| \geq \eta/3\right) \xrightarrow[n \rightarrow \infty]{\text{By Lemma 2}} 0. \quad (13)$$

□

The result is a *weak* consistency statement ( $\hat{\phi}_N \rightarrow \phi^*$  in probability). In Section 5 we strengthen this to almost-sure convergence by adding a Lipschitz regularity assumption.

## 5 STRONG CONSISTENCY UNDER UNIFORM LIPSCHITZ CONDITIONS

While Theorem 1 guarantees that the ODR estimate converges in probability to the true parameter distribution, in many practical settings one desires a stronger, almost sure guarantee. Intuitively, *strong consistency* asserts that, with probability one, the estimated distribution will converge exactly to the true one as more offline data is observed. In this section we show that, under an additional Lipschitz continuity assumption on the log-likelihood function, ODR enjoys this almost-sure convergence property.

### 5.1 ADDITIONAL ASSUMPTION

The key extra ingredient is a uniform control over how rapidly the single step log-likelihood  $a(x, \phi)$  can change as we vary the distributional parameter  $\phi = (\mu, \Sigma)$ . Formally:

**Assumption 5** (Uniform Lipschitz Continuity). *There exists a constant  $L < \infty$  such that for every transition  $x = (s, a, s')$  and all  $\phi, \psi \in \Phi$ , we have  $|a(x, \phi) - a(x, \psi)| \leq L \|\phi - \psi\|_2$ .*

This condition ensures that the family  $\{a(\cdot, \phi) : \phi \in \Phi\}$  is *equi-Lipschitz*, which -together with compactness of  $\Phi$ - yields a *uniform strong law of large numbers*. In turn, this uniform convergence of the empirical log-likelihood to its population counterpart underpins the almost sure convergence of the maximizers.

## 5.2 MAIN RESULT

We can now state our strong consistency result:

**Theorem 2** (Strong Consistency of ODR). *Under Assumptions 1 to 5, let  $\hat{\phi}_N \in \arg \max_{\phi \in \Phi} L_N(\phi)$  be any measurable maximizer of the empirical log-likelihood, then*

$$\hat{\phi}_N \xrightarrow[N \rightarrow \infty]{\text{a.s.}} \phi^* = (\xi^*, 0), \quad (14)$$

*i.e., almost surely the estimated distribution collapses exactly onto the true simulator parameters.*

The distinction between convergence *in probability* and *almost surely* is subtle but meaningful: almost-sure consistency implies that, except on a set of histories of measure zero, as soon as enough data is collected the optimizer will *never* stray from the true maximum again. In contrast, convergence in probability only assures that large deviations become increasingly unlikely.

The heart of the proof is the following uniform strong law, which follows from empirical process arguments once we have the Lipschitz control:

**Lemma 4** (Uniform Strong Law). *Under Assumptions 1 to 5, the empirical and population log-likelihoods satisfy  $\sup_{\phi \in \Phi} |L_N(\phi) - L(\phi)| \xrightarrow[N \rightarrow \infty]{\text{a.s.}} 0$ .*

Lemma 4 tells us that with probability one the worst-case difference between the finite-sample objective and its ideal limit vanishes. Once this uniform convergence is in hand, classical arguments on continuity and compactness show that the maximizers converge almost surely.

*Proof (Sketch).* We first show  $\sup_{\phi \in \Phi} |L_N(\phi) - L(\phi)| \xrightarrow[N \rightarrow \infty]{\text{a.s.}} 0$  by verifying for each  $\epsilon > 0$  that  $\sum_N P(\sup_{\phi \in \Phi} |L_N(\phi) - L(\phi)| > 2\epsilon) < \infty$ . By compactness of  $\Phi$  there is a finite  $\epsilon/L$ -net  $\{\phi_1, \dots, \phi_K\}$  so that Lipschitz continuity gives  $|L_N(\phi) - L_N(\phi_i)| + |L(\phi) - L(\phi_i)| \leq \epsilon$  whenever  $\|\phi - \phi_i\| \leq \epsilon/L$ . Hence

$$\left\{ \sup_{\phi} |L_N(\phi) - L(\phi)| > 2\epsilon \right\} \subset \bigcup_{i=1}^K \{|L_N(\phi_i) - L(\phi_i)| > \epsilon\}, \quad (15)$$

and *Hoeffding's inequality* yields

$$P(|L_N(\phi_i) - L(\phi_i)| > \epsilon) \leq 2 \exp\left(-\frac{N\epsilon^2}{2\widetilde{M}^2}\right), \quad (16)$$

where  $\widetilde{M} := \max\{|\log K|, |\log c|\}$ . So  $P(\sup_{\phi} |L_N(\phi) - L(\phi)| > 2\epsilon) \leq 2K \exp(-cN\epsilon^2)$ , which is summable in  $N$ . *Borel-Cantelli lemma* then gives uniform almost sure convergence. Finally, on the event of uniform convergence one repeats the identification neighborhood argument of Theorem 1 to conclude  $\hat{\phi}_N \rightarrow \phi^*$  almost surely.  $\square$

Full details of the proof are deferred to Section C, but the key takeaway is that the Lipschitz assumption upgrades our earlier *in probability* consistency to the far stronger *almost sure* statement, giving robust guarantees for ODR even in worst case data realizations.

### 5.3 A NOTION OF $\alpha$ -INFORMATIVENESS

The strong consistency yields the following.

**Lemma 5.** *Let  $\epsilon > 0$ . If  $\hat{\phi}_N = (\mu_N, \Sigma_N) \xrightarrow{\text{a.s.}} (\xi^*, 0)$  then almost surely there is  $N_0$  so that for all  $N \geq N_0$ ,  $P_{\hat{\phi}_N}(\mathbb{B}(\xi^*, \epsilon)) > \frac{1}{2}$ .*

*Proof of Lemma 5.* Fix  $\epsilon > 0$  and let  $Z_N \sim \mathcal{N}(\mu_N, \Sigma_N)$ . Then  $P(\|Z_N - \xi^*\| \geq \epsilon) \leq P(\|Z_N - \mu_N\| \geq \frac{\epsilon}{2}) + P(\|\mu_N - \xi^*\| \geq \frac{\epsilon}{2})$ . By Chebyshev’s inequality,  $P(\|Z_N - \mu_N\| \geq \frac{\epsilon}{2}) \leq \frac{\mathbb{E}\|Z_N - \mu_N\|^2}{(\epsilon/2)^2} = \frac{\text{tr}(\Sigma_N)}{(\epsilon/2)^2}$ . Hence  $P_{\hat{\phi}_N}(\mathbb{B}(\xi^*, \epsilon)) = 1 - P(\|Z_N - \xi^*\| \geq \epsilon) \geq 1 - \frac{4\text{tr}(\Sigma_N)}{\epsilon^2} - P(\|\mu_N - \xi^*\| \geq \epsilon/2)$ . As  $(\mu_N, \Sigma_N) \rightarrow (\xi^*, 0)$  a.s., we have  $\|\mu_N - \xi^*\| \rightarrow 0$  and  $\text{tr}(\Sigma_N) \rightarrow 0$ , so the right hand side tends to 1 almost surely. Hence  $P_{\hat{\phi}_N}(\mathbb{B}(\xi^*, \epsilon)) \rightarrow 1$  almost surely.  $\square$

The lemma states that when the estimator  $(\mu_N, \Sigma_N)$  converges almost surely to the true mean with vanishing covariance, the Gaussian distribution fitted by ODR eventually assigns *more than half of its probability mass* to any fixed  $\epsilon$ -ball around  $\xi^*$ . In other words, ODR is so informative that the learned randomization concentrates near the real world. This observation motivates a general, model-agnostic notion of “informativeness” for ODR, applicable beyond the Gaussian setting.

**Definition 1** ( $\alpha, \epsilon$ -Informativeness of an ODR Algorithm  $\mathcal{A}$ ). *Let  $\alpha \in (0, 1)$  and  $\epsilon > 0$ , an algorithm  $\mathcal{A}$  is  $\alpha, \epsilon$ -informative if there exists almost surely  $N_0 \geq 1$  such that for all  $N \geq N_0$ , running  $\mathcal{A}$  on any collection  $\mathcal{D} = \{(s_i, a_i, s'_i)\}_{i=1}^N$  of i.i.d. transitions (from the real system) produces an ODR distribution  $\hat{\phi}_N$  such that*

$$P_{\hat{\phi}_N}(\mathbb{B}(\xi^*, \epsilon)) \geq \alpha.$$

We say algorithm  $\mathcal{A}$  is  $\alpha$ -informative if  $\mathcal{A}$  is  $\alpha, \epsilon$ -informative for any  $\epsilon > 0$ .

Under this language, Lemma 5 states that the Gaussian ODR procedure from Section 3 is  $\alpha$ -informative for every  $\alpha < 1$ . When the simulator class  $\Xi$  is finite,  $\alpha$ -informativeness is equivalent to the almost-sure existence of an index  $N_0$  such that, for all  $N \geq N_0$ , the fitted distribution assigns at least  $\alpha$  mass to the singleton  $\{\xi^*\}$ , that is,  $P_{\hat{\phi}_N}(\xi^*) \geq \alpha$ .

## 6 ASSUMPTIONS: PRACTICALITY, VIOLATIONS, AND RELAXATIONS

### 6.1 THE I.I.D. ASSUMPTION

The i.i.d. assumption on the offline dataset  $\mathcal{D}$  holds whenever the offline dataset is collected using a *fixed, stationary* behavior policy  $\pi(\cdot | s)$ . This assumption is stronger than needed for our weak consistency result: we invoke it only to apply a uniform law of large numbers at the end of the proof of Lemma 2. As noted after Lemma 2.4 in Newey & McFadden (1994), the same conclusion holds (even for dependent data) for ergodic and strictly stationary sequences  $\{X_i = (s_i, a_i, s'_i)\}$  which means that the joint distribution of the vector  $(X_i, \dots, X_{i+m})$  does not depend on  $i$  for any  $m$ . This is much weaker than the i.i.d. assumption and is satisfied whenever the offline dataset is collected by a *fixed* behavior policy (not necessarily a stationary policy). In practice, weak consistency should therefore hold broadly.

### 6.2 THE MIXTURE POSITIVITY ASSUMPTION

Assumption 3 is a strong requirement: it holds if and only if  $\inf_x \inf_{\phi} q_{\phi}(x) > 0$ , i.e., the density is uniformly bounded away from zero over both  $x$  and  $\phi$ . This excludes common light-tailed families (e.g., Gaussian-like), for which  $\inf_x q_{\phi}(x) = 0$ . For *weak consistency*, however, Assumption 3 can be relaxed:

**Lemma 6** (Relaxation of Assumption 3). *Weak consistency of ODR still holds if Assumption 3 is replaced by the following tail condition: there exists  $\epsilon_0 > 0$  such that for all  $\epsilon \in (0, \epsilon_0]$ ,*

$$P\left(\inf_{\phi} q_{\phi}(X) \leq \epsilon\right) \leq \frac{1}{\log(1/\epsilon)^2}. \quad (17)$$

This assumption is strictly weaker than uniform positivity. The key point is that, to apply the uniform law of large numbers from Newey & McFadden (1994) in the weak-consistency proof, it suffices to have an *integrable envelope*  $d(x)$  with  $a(x, \phi) \leq d(x)$  for all  $\phi$ , rather than a uniform bound in  $(x, \phi)$ , the above tail control yields such an envelope. The proof is deferred to Section D.1.

### 6.3 THE UNIFORM LIPSCHITZ CONTINUITY

Assumption 5 is not immediately interpretable. We give a simple sufficient condition under which it holds:

**Lemma 7** (Sufficient Condition for the Uniform Lipschitz Continuity Assumption). *Suppose the following holds for every  $x = (s, a, s')$*

1. *The function  $\xi \mapsto p_\xi(s' | s, a)$  is twice continuously differentiable (of class  $C^2$ ),*
2. *There exists two constants  $G_1 > 0$  and  $G_2 > 0$  such that  $|\nabla_\xi p_\xi(s' | s, a)| \leq G_1$  and  $|\nabla_\xi^2 p_\xi(s' | s, a)| \leq G_2$ ,*

*then Assumption 5 holds with  $L = \frac{G_1 + G_2/2}{c}$ .*

A complete proof appears in Section D.2. This sufficient condition is easy to interpret because it depends only on the simulator’s transition kernel  $p_\xi$ . In practice, it is satisfied whenever the simulators are governed by smooth physics.

### 6.4 THE IDENTIFIABILITY ASSUMPTION

Assumption 3 is a coverage condition on the dataset: it requires that any mixing Gaussian distribution that reproduces the transition kernel on the state–action pairs observed in  $\mathcal{D}$  must equal the degenerate Dirac mass at the true parameter. Intuitively, the dataset must visit state–action pairs that are informative about  $\xi$ . This is information-theoretically minimal: no method can distinguish parameters that are observationally identical on  $\text{supp}(\mu)$ .

In the case of partial coverage, we naturally define the *identified set under coverage*  $\mu$  as follows:

$$\mathcal{Q}_\mu^* := \{\phi \in \Phi : q_\phi(\cdot | s, a) = p_{\xi^*}(\cdot | s, a) \text{ for } \mu - \text{a.e. } (s, a)\}. \quad (18)$$

It follows from this definition and the proof of Lemma 1 that:

**Lemma 8.** *The following holds:*

$$\mathcal{Q}_\mu^* = \arg \max_{\phi} L(\phi). \quad (19)$$

Using this notion of identified set, we can generalize Theorem 1 when we relax Assumption 4 as follows:

**Theorem 3.** *Under Assumptions 1, 2 and 3, the following holds, Any measurable maximizer  $\hat{\phi}_N \in \arg \max_{\phi \in \Phi} L_N(\phi)$  satisfies  $\text{dist}(\hat{\phi}_N, \mathcal{Q}_\mu^*) \xrightarrow[N \rightarrow \infty]{P} 0$ <sup>4</sup>.*

This theorem states that under partial coverage, our estimator does not select a single parameter but converges to an *identified set* of parameters that are observationally indistinguishable on the state–action pairs visited by the data. The proof is deferred to Section D.3. The proof of this theorem is very general. In particular, even in the misspecified case where  $\mathcal{M}^* \notin \mathcal{U}$ , we still have  $\hat{\phi}_N \rightarrow \phi^\dagger \in \arg \max_{\phi} L(\phi)$ . A more detailed discussion of this case is deferred to Section D.4.

Without any additional assumptions, the only structural result that we can derive on the identified set is:

<sup>4</sup>where  $\text{dist}$  is the distance to a set defined by  $\text{dist}(\phi, \mathcal{Q}) := \inf_{\psi \in \mathcal{Q}} \|\phi - \psi\|$ .

**Lemma 9** (Upper Hemicontinuity of  $\mathcal{Q}_\mu^*$ ). *Under Assumptions 1, 2 and 3 The identified set  $\mathcal{Q}_\mu^*$  is non-empty and compact and the correspondence  $\mu \mapsto \mathcal{Q}_\mu^*$  is upper hemicontinuous<sup>5</sup> with respect to total variation.*

The proof of this lemma uses *Berge’s Maximum Theorem* and is deferred to Section D.3.

In short, this lemma says if we perturb the dataset’s coverage only slightly (in total-variation distance), the set of maximizers cannot “jump” to a faraway region: any limit of maximizers for the perturbed coverages remains a maximizer at the limit coverage (upper hemicontinuity). Intuitively, modestly adding or reweighting offline data will not create spurious, distant optima, it keeps the solution set nearby, and, as coverage includes more informative state–action pairs, typically makes it tighter.

The main limitation is that, *without additional assumptions*, we cannot provide a quantitative radius for this set or a Lipschitz-type bound on how much it can move when coverage changes.

## 7 CONCLUSION

In this paper, we present a rigorous framework for ODR, bridging the gap between empirical success and theoretical understanding in sim-to-real transfer. By casting ODR as maximum likelihood estimation over a parametric family of simulator distributions, we proved that, under mild regularity conditions, the learned distribution is weakly consistent, concentrating on the true dynamics as the offline dataset grows. With the addition of a uniform Lipschitz continuity assumption, we further established strong consistency. Beyond these core results, we scrutinized the practicality of the assumptions and provided diagnostics and relaxations—replacing i.i.d. with stationarity/ergodicity for the ULLN, weakening mixture positivity via a logarithmic tail condition, and giving checkable smoothness criteria that imply the uniform Lipschitz requirement—thereby justifying ODR’s applicability across a broader range of settings. By demonstrating that offline logs are not merely passive datasets but a powerful tool for principled domain randomization, we hope our formulation and analysis can provide insight that paves the way for safer, more data-efficient sim-to-real pipelines in robotics, autonomous vehicles, and beyond.

---

<sup>5</sup>A set-valued map  $F$  is upper hemicontinuous at  $x_0$  if, whenever  $x_n \rightarrow x_0$  and  $y_n \in F(x_n)$  with  $y_n \rightarrow y$ , then  $y \in F(x_0)$ . Equivalently: for every open  $U$  with  $F(x_0) \subseteq U$ , there exists a neighborhood  $V$  of  $x_0$  such that  $F(x) \subseteq U$  for all  $x \in V$ .

## REFERENCES

- Mohammad Mehdi Afsar, Trafford Crump, and Behrouz H. Far. Reinforcement learning based recommender systems: A survey. *CoRR*, abs/2101.06286, 2021. URL <https://arxiv.org/abs/2101.06286>.
- Karol Arndt, Murtaza Hazara, Ali Ghadirzadeh, and Ville Kyrki. Meta reinforcement learning for sim-to-real domain adaptation. *CoRR*, abs/1909.12906, 2019. URL <http://arxiv.org/abs/1909.12906>.
- Pierre C Bellec and Cun-Hui Zhang. Second order stein: Sure for sure and other applications in high-dimensional inference, 2020. URL <https://arxiv.org/abs/1811.04121>.
- Claude Berge. *Topological Spaces: Including a Treatment of Multi-Valued Functions, Vector Spaces, and Convexity*. Oliver and Boyd, Edinburgh, 1963. Translated from the French by E.M. Patterson.
- Yevgen Chebotar, Ankur Handa, Viktor Makoviychuk, Miles Macklin, Jan Issac, Nathan D. Ratliff, and Dieter Fox. Closing the sim-to-real loop: Adapting simulation randomization with real world experience. *CoRR*, abs/1810.05687, 2018. URL <http://arxiv.org/abs/1810.05687>.
- Xiacong Chen, Lina Yao, Julian J. McAuley, Guanglin Zhou, and Xianzhi Wang. A survey of deep reinforcement learning in recommender systems: A systematic review and future directions. *CoRR*, abs/2109.03540, 2021. URL <https://arxiv.org/abs/2109.03540>.
- Xiaoyu Chen, Jiachen Hu, Chi Jin, Lihong Li, and Liwei Wang. Understanding domain randomization for sim-to-real transfer, 2022. URL <https://arxiv.org/abs/2110.03239>.
- Ignasi Clavera, Anusha Nagabandi, Ronald S. Fearing, Pieter Abbeel, Sergey Levine, and Chelsea Finn. Learning to adapt: Meta-learning for model-based control. *CoRR*, abs/1803.11347, 2018. URL <http://arxiv.org/abs/1803.11347>.
- Fei Feng, Wotao Yin, and Lin F. Yang. Does knowledge transfer always help to learn a better policy? *CoRR*, abs/1912.02986, 2019. URL <http://arxiv.org/abs/1912.02986>.
- Dario Floreano, Phil Husbands, and Stefano Nolfi. *Evolutionary Robotics*, pp. 1423–1451. 01 2008. ISBN 978-3-540-23957-4. doi: 10.1007/978-3-540-30301-5\_62.
- Nikolaus Hansen. *The CMA Evolution Strategy: A Comparing Review*, pp. 75–102. Springer Berlin Heidelberg, Berlin, Heidelberg, 2006. ISBN 978-3-540-32494-2. doi: 10.1007/3-540-32494-1\_4. URL [https://doi.org/10.1007/3-540-32494-1\\_4](https://doi.org/10.1007/3-540-32494-1_4).
- Nikolaus Hansen and Andreas Ostermeier. Completely derandomized self-adaptation in evolution strategies. *Evolutionary Computation*, 9(2):159–195, 2001. doi: 10.1162/106365601750190398.
- Nan Jiang. Pac reinforcement learning with an imperfect model. *Proceedings of the AAAI Conference on Artificial Intelligence*, 32, 04 2018. doi: 10.1609/aaai.v32i1.11594.
- Dmitry Kalashnikov, Alex Irpan, Peter Pastor, Julian Ibarz, Alexander Herzog, Eric Jang, Deirdre Quillen, Ethan Holly, Mrinal Kalakrishnan, Vincent Vanhoucke, and Sergey Levine. Qt-opt: Scalable deep reinforcement learning for vision-based robotic manipulation. *CoRR*, abs/1806.10293, 2018. URL <http://arxiv.org/abs/1806.10293>.
- Bangalore Ravi Kiran, Ibrahim Sobh, Victor Talpaert, Patrick Mannion, Ahmad A. Al Sallab, Senthil Kumar Yogamani, and Patrick Pérez. Deep reinforcement learning for autonomous driving: A survey. *CoRR*, abs/2002.00444, 2020. URL <https://arxiv.org/abs/2002.00444>.
- Jens Kober, J. Bagnell, and Jan Peters. Reinforcement learning in robotics: A survey. *The International Journal of Robotics Research*, 32:1238–1274, 09 2013. doi: 10.1177/0278364913495721.
- Akshay Krishnamurthy, Alekh Agarwal, and John Langford. Pac reinforcement learning with rich observations. *Advances in Neural Information Processing Systems*, 29, 2016.
- S. Kullback and R. A. Leibler. On Information and Sufficiency. *The Annals of Mathematical Statistics*, 22(1):79 – 86, 1951. doi: 10.1214/aoms/1177729694. URL <https://doi.org/10.1214/aoms/1177729694>.

- Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *CoRR*, abs/2005.01643, 2020. URL <https://arxiv.org/abs/2005.01643>.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin A. Riedmiller. Playing atari with deep reinforcement learning. *CoRR*, abs/1312.5602, 2013. URL <http://arxiv.org/abs/1312.5602>.
- Fabio Muratore, Christian Eilers, Michael Gienger, and Jan Peters. Bayesian domain randomization for sim-to-real transfer. *CoRR*, abs/2003.02471, 2020. URL <https://arxiv.org/abs/2003.02471>.
- Whitney K. Newey and Daniel McFadden. Chapter 36 large sample estimation and hypothesis testing. volume 4 of *Handbook of Econometrics*, pp. 2111–2245. Elsevier, 1994. doi: [https://doi.org/10.1016/S1573-4412\(05\)80005-4](https://doi.org/10.1016/S1573-4412(05)80005-4). URL <https://www.sciencedirect.com/science/article/pii/S1573441205800054>.
- Haoyi Niu, Jianming Hu, Zheyu Cui, and Yi Zhang. DR2L: surfacing corner cases to robustify autonomous driving via domain randomization reinforcement learning. *CoRR*, abs/2107.11762, 2021. URL <https://arxiv.org/abs/2107.11762>.
- OpenAI, Marcin Andrychowicz, Bowen Baker, Maciek Chociej, Rafal Józefowicz, Bob McGrew, Jakub Pachocki, Arthur Petron, Matthias Plappert, Glenn Powell, Alex Ray, Jonas Schneider, Szymon Sidor, Josh Tobin, Peter Welinder, Lilian Weng, and Wojciech Zaremba. Learning dexterous in-hand manipulation. *CoRR*, abs/1808.00177, 2018. URL <http://arxiv.org/abs/1808.00177>.
- Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Sim-to-real transfer of robotic control with dynamics randomization. *CoRR*, abs/1710.06537, 2017. URL <http://arxiv.org/abs/1710.06537>.
- Samira Pouyanfar, Muneeb Saleem, Nikhil George, and Shu-Ching Chen. Roads: Randomization for obstacle avoidance and driving in simulation. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 1267–1276, 2019. doi: 10.1109/CVPRW.2019.00166.
- Aravind Rajeswaran, Sarvjeet Ghotra, Sergey Levine, and Balaraman Ravindran. Epopt: Learning robust neural network policies using model ensembles. *CoRR*, abs/1610.01283, 2016. URL <http://arxiv.org/abs/1610.01283>.
- Fabio Ramos, Rafael Carvalhaes Possas, and Dieter Fox. Bayessim: adaptive domain randomization via probabilistic inference for robotics simulators. *CoRR*, abs/1906.01728, 2019. URL <http://arxiv.org/abs/1906.01728>.
- Andrei A. Rusu, Matej Vecerík, Thomas Rothörl, Nicolas Manfred Otto Heess, Razvan Pascanu, and Raia Hadsell. Sim-to-real robot learning from pixels with progressive nets. *ArXiv*, abs/1610.04286, 2016. URL <https://api.semanticscholar.org/CorpusID:876231>.
- Fereshteh Sadeghi and Sergey Levine. (cad)<sup>2</sup>rl: Real single-image flight without a single real image. *CoRR*, abs/1611.04201, 2016. URL <http://arxiv.org/abs/1611.04201>.
- John Schulman, Sergey Levine, Philipp Moritz, Michael I. Jordan, and Pieter Abbeel. Trust region policy optimization. *CoRR*, abs/1502.05477, 2015. URL <http://arxiv.org/abs/1502.05477>.
- David Silver, Aja Huang, Christopher Maddison, Arthur Guez, Laurent Sifre, George Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap, Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the game of go with deep neural networks and tree search. *Nature*, 529:484–489, 01 2016. doi: 10.1038/nature16961.
- Charles M. Stein. Estimation of the Mean of a Multivariate Normal Distribution. *The Annals of Statistics*, 9(6):1135 – 1151, 1981. doi: 10.1214/aos/1176345632. URL <https://doi.org/10.1214/aos/1176345632>.

Jie Tan, Tingnan Zhang, Erwin Coumans, Atil Iscen, Yunfei Bai, Danijar Hafner, Steven Bohez, and Vincent Vanhoucke. Sim-to-real: Learning agile locomotion for quadruped robots. *CoRR*, abs/1804.10332, 2018. URL <http://arxiv.org/abs/1804.10332>.

George Tauchen. Diagnostic testing and evaluation of maximum likelihood models. *Journal of Econometrics*, 30(1):415–443, 1985. ISSN 0304-4076. doi: [https://doi.org/10.1016/0304-4076\(85\)90149-6](https://doi.org/10.1016/0304-4076(85)90149-6). URL <https://www.sciencedirect.com/science/article/pii/0304407685901496>.

Gabriele Tiboni, Karol Arndt, and Ville Kyrki. Dropo: Sim-to-real transfer with offline domain randomization, 2023. URL <https://arxiv.org/abs/2201.08434>.

Joshua Tobin, Rachel Fong, Alex Ray, Jonas Schneider, Wojciech Zaremba, and Pieter Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. *CoRR*, abs/1703.06907, 2017. URL <http://arxiv.org/abs/1703.06907>.

Ya-Yen Tsai, Hui Xu, Zihan Ding, Chong Zhang, Edward Johns, and Bidan Huang. DROID: minimizing the reality gap using single-shot human demonstration. *CoRR*, abs/2102.11003, 2021. URL <https://arxiv.org/abs/2102.11003>.

David Williams. *Probability with Martingales*. Cambridge University Press, 1991.

Wenshuai Zhao, Jorge Peña Queraltá, and Tomi Westerlund. Sim-to-real transfer in deep reinforcement learning for robotics: a survey. *CoRR*, abs/2009.13303, 2020. URL <https://arxiv.org/abs/2009.13303>.

Yuren Zhong, Aniket Anand Deshmukh, and Clayton Scott. PAC reinforcement learning without real-world feedback. *CoRR*, abs/1909.10449, 2019. URL <http://arxiv.org/abs/1909.10449>.

## A ADDITIONAL PRELIMINARIES

### A.1 REFINED ANALYSIS OF THE UNIFORM DR SIM-TO-REAL GAP

In this section, we tighten the worst-case sim-to-real gap bound in the finite,  $\delta$ -separable setting originally proved by Chen et al. (2022).

In the proof of Lemma 5 of the paper, Inequality (47) yields with probability at least  $1 - \delta_0$ ,

$$\sum_{s' \in \mathcal{H}} \ln \left( \frac{P_{\mathcal{M}^*}(s' | s_0, a_0)}{P_{\mathcal{M}_1}(s' | s_0, a_0)} \right) \geq \frac{n_0 \delta^2}{2} - \log(1/\alpha) \sqrt{2n_0 \log(2/\delta_0)} - \sqrt{n_0 \log(2/\delta_0)/c} - 2\alpha S n_0. \quad (20)$$

The objective is to find a setting of parameters that guarantee with probability at least  $1 - \frac{1}{MH}$ ,

$$\sum_{s' \in \mathcal{H}} \ln \left( \frac{P_{\mathcal{M}^*}(s' | s_0, a_0)}{P_{\mathcal{M}_1}(s' | s_0, a_0)} \right) > 0.$$

It is sufficient to have the right term positive in Equation (20), i.e.,

$$\frac{n_0 \delta^2}{2} - \log(1/\alpha) \sqrt{2n_0 \log(2/\delta_0)} - \sqrt{n_0 \log(2/\delta_0)/c} - 2\alpha S n_0 > 0. \quad (21)$$

Setting  $\alpha = \frac{\delta^2}{8S}$ ,  $\delta_0 = \frac{1}{MH}$  (the same values as in the paper), this term becomes

$$\frac{n_0 \delta^2}{2} - \log(1/\alpha) \sqrt{2n_0 \log(2/\delta_0)} - \sqrt{n_0 \log(2/\delta_0)/c} - 2\alpha S n_0 \quad (22)$$

$$= \frac{n_0 \delta^2}{2} - \log\left(\frac{8S}{\delta^2}\right) \sqrt{2n_0 \log(2MH)} - \sqrt{n_0 \log(2MH)/c} - \frac{\delta^2}{4} n_0 \quad (23)$$

$$= \frac{n_0 \delta^2}{4} - \log\left(\frac{8S}{\delta^2}\right) \sqrt{2n_0 \log(2MH)} - \sqrt{n_0 \log(2MH)/c} \quad (24)$$

$$= \sqrt{n_0} \frac{\delta^2}{4} \left[ \sqrt{n_0} - \frac{4}{\delta^2} \left( \log\left(\frac{8S}{\delta^2}\right) \sqrt{2 \log(2MH)} - \sqrt{\log(2MH)/c} \right) \right] \quad (25)$$

hence the condition 21 becomes equivalent to

$$\sqrt{n_0} > \frac{4}{\delta^2} \sqrt{\log(2MH)} \left( \sqrt{2} \log\left(\frac{8S}{\delta^2}\right) + \frac{1}{\sqrt{c}} \right), \quad (26)$$

or, equivalently,

$$n_0 > \frac{16}{\delta^4} \log(2MH) \left( \sqrt{2} \log\left(\frac{8S}{\delta^2}\right) + \frac{1}{\sqrt{c}} \right)^2. \quad (27)$$

Thus, there exists a valid setting that satisfies condition 21 which can be expressed as

$$\alpha = \frac{\delta^2}{8S}, \quad \delta_0 = \frac{1}{MH}, \quad n_0 = \frac{c_0 \log(MH) \log^2(S/\delta^2)}{\delta^4}, \quad (28)$$

for some constant  $c_0 > 0$  sufficiently large.

With this new setting, the result of the Lemma 7 of the paper becomes

$$\mathbb{E}[h_0] \leq O\left(\frac{DM^2 \log(MH) \log^2(S/\delta^2)}{\delta^4}\right). \quad (29)$$

The proof of Theorem 5 of the paper is not affected by the new expression of  $n_0$  and gives

$$V_{\mathcal{M}^*,1}^*(s_1) - V_{\mathcal{M}^*,1}^{\hat{\pi}}(s_1) \leq O(\mathbb{E}[h_0] + D) = O\left(\frac{DM^2 \log(MH) \log^2(S/\delta^2)}{\delta^4}\right). \quad (30)$$

Combining this result with Lemma 1 of the paper leads to

$$\text{Gap}(\pi_{DR}^*, \mathcal{U}) = O\left(\frac{DM^3 \log(MH) \log^2(S/\delta^2)}{\delta^4}\right). \quad (31)$$

This shows that in the regime where  $H$  and  $M$  are relatively large, the  $O(M^3 \log^3(MH))$  bound of Chen et al. (2022) can be tightened to  $O(M^3 \log(MH))$ .

## A.2 INSIGHTS INTO THE ODR OBJECTIVE

In this section, we explain why the formal ODR problem in Equation 2 corresponds exactly to fitting the simulator parameter distribution that maximizes the likelihood of our offline dataset.

We seek the parameter  $\phi$  of the distribution  $P_\phi(\xi)$  that maximizes the probability of observing the triples  $(s_i, a_i, s'_i)$  of our dataset, i.e., to solve

$$\phi^* = \arg \max_{\phi} P\left(\{(s_i, a_i, s'_i)\}_{i=1}^N \mid \phi\right), \quad (32)$$

This probability corresponds to  $P(\cap_{i=1}^N \{(s_i, a_i, s'_i)\} \mid \phi)$  and since the data is i.i.d.,  $\phi^*$  can be rewritten as follows

$$\phi^* = \arg \max_{\phi} \prod_{i=1}^N P(\{(s_i, a_i, s'_i)\} \mid \phi). \quad (33)$$

Now, we can approximate  $P(\{(s_i, a_i, s'_i)\} \mid \phi)$  by the expected transition probability over all  $\xi \sim P_\phi(\xi)$ , i.e.,

$$\phi^* \approx \arg \max_{\phi} \prod_{i=1}^N \mathbb{E}_{\xi \sim P_\phi(\xi)} [P_\xi(s'_i \mid s_i, a_i)]. \quad (34)$$

Since the logarithm is increasing, this is equivalent to

$$\phi^* \approx \arg \max_{\phi} \sum_{i=1}^N \log \mathbb{E}_{\xi \sim P_\phi(\xi)} [P_\xi(s'_i \mid s_i, a_i)], \quad (35)$$

which recovers exactly the empirical log-likelihood objective stated in Equation 2.

## B OMITTED PROOFS IN SECTION 4

*Proof of Lemma 1.* We have

$$L(\phi) = \mathbb{E}_{(s,a)} \mathbb{E}_{s' \sim p_{\xi^*}(\cdot \mid s,a)} \left[ \log \mathbb{E}_{\xi \sim p_\phi(\xi)} [p_\xi(s' \mid s, a)] \right]. \quad (36)$$

We rewrite the inner expectation as follows

$$\mathbb{E}_{s' \sim p_{\xi^*}(\cdot \mid s,a)} \left[ \log q_\phi(s' \mid s, a) \right] = \mathbb{E}_{s' \sim p_{\xi^*}(\cdot \mid s,a)} \left[ \log \frac{q_\phi(s' \mid s, a)}{p_{\xi^*}(s' \mid s, a)} + \log p_{\xi^*}(s' \mid s, a) \right]. \quad (37)$$

Notice that

$$\int_{s' \in \mathcal{S}} q_\phi(s' \mid s, a) \lambda(ds') = \int_{s' \in \mathcal{S}} \int_{\xi \in \Xi} p_\xi(s' \mid s, a) p_\phi(d\xi) \lambda(ds'), \quad (38)$$

and using *Fubini-Tonelli's theorem*, it follows,

$$\int_{s' \in \mathcal{S}} \int_{\xi \in \Xi} p_\xi(s' \mid s, a) p_\phi(d\xi) \lambda(ds') = \int_{\xi \in \Xi} p_\phi(d\xi) \int_{s' \in \mathcal{S}} p_\xi(s' \mid s, a) \lambda(ds'). \quad (39)$$

Since  $p_\xi(\cdot | s, a)$  and  $p_\phi$  are probability densities, their total mass is 1, which yields

$$\int_{s' \in \mathcal{S}} q_\phi(s' | s, a) \lambda(ds') = \int_{\xi \in \Xi} p_\phi(d\xi) = 1. \quad (40)$$

Hence  $q_\phi(\cdot | s, a)$  is a probability density, and one can rewrite  $L(\phi)$  using *Kullback-Leibler (KL) divergence* (defined in Kullback & Leibler (1951)) as follows

$$L(\phi) = \mathbb{E}_{(s,a)} \left[ -D_{KL}(p_{\xi^*}(\cdot | s, a) \| q_\phi(\cdot | s, a)) + \mathbb{E}_{s' \sim p_{\xi^*}(\cdot | s, a)} [\log p_{\xi^*}(s' | s, a)] \right] \quad (41)$$

$$= \mathbb{E}_{(s,a)} [-D_{KL}(p_{\xi^*}(\cdot | s, a) \| q_\phi(\cdot | s, a))] + H(\xi^*), \quad (42)$$

where  $H(\xi^*) = \mathbb{E}_{(s,a)} \mathbb{E}_{s' \sim p_{\xi^*}(\cdot | s, a)} [\log p_{\xi^*}(s' | s, a)]$  is independent of  $\phi$ , and for a fixed  $(s, a)$ ,  $D_{KL}(p_{\xi^*}(\cdot | s, a) \| q_\phi(\cdot | s, a)) \geq 0$  with equality if and only if  $p_{\xi^*}(\cdot | s, a) = q_\phi(\cdot | s, a)$ .

Hence, for all  $\phi \in \Phi$ ,  $L(\phi) \leq H(\xi^*)$ , and

$$L(\phi) = H(\xi^*) \iff \mathbb{E}_{(s,a)} [-D_{KL}(p_{\xi^*}(\cdot | s, a) \| q_\phi(\cdot | s, a))] = 0 \quad (43)$$

$$\iff \text{For almost every } (s, a), D_{KL}(p_{\xi^*}(\cdot | s, a) \| q_\phi(\cdot | s, a)) = 0 \quad (44)$$

$$\iff \text{For almost every } (s, a), p_{\xi^*}(\cdot | s, a) = q_\phi(\cdot | s, a) \quad (45)$$

$$\iff \phi = (\xi^*, 0), \quad (46)$$

where the last equivalence follows from Assumption 4. This concludes the proof.  $\square$

*Proof of Lemma 2.* We begin by stating and proving a few intermediate lemmas that will simplify the proof.

The following lemma states that convergence of  $\phi$  implies convergence in distribution of  $P_\phi$ .

**Lemma 10.** *Let  $\{\phi_n\} := \{(\mu_n, \Sigma_n)\} \in \Phi^{\mathbb{N}}$  a sequence that converges to  $\phi := (\mu, \Sigma)$  (i.e.  $\|\mu_n - \mu\| \rightarrow 0$  and  $\|\Sigma_n - \Sigma\|_{op} \rightarrow 0$ ). Then  $P_{\phi_n}$  converges weakly to  $P_\phi$  ( $P_{\phi_n} \implies P_\phi$ ).*

*Proof of Lemma 10.* We denote

$$G_n = \mathcal{N}(\mu_n, \Sigma_n), \quad G = \mathcal{N}(\mu, \Sigma). \quad (47)$$

The characteristic function of  $G_n$  is

$$\varphi_{G_n}(t) = \exp\left(it^\top \mu_n - \frac{1}{2}t^\top \Sigma_n t\right), \quad t \in \mathbb{R}^d. \quad (48)$$

For every fixed  $t \in \mathbb{R}^d$ , we have

$$\varphi_{G_n}(t) \xrightarrow{n \rightarrow \infty} \exp\left(it^\top \mu - \frac{1}{2}t^\top \Sigma t\right) = \varphi_G(t). \quad (49)$$

By *Lévy's continuity theorem* (see Williams (1991)), we have  $P_{\phi_n} \implies P_\phi$ .  $\square$

Notice that the result holds also in the case where  $\Sigma = 0$ . In that case,  $\varphi_G(t) = \exp(it^\top \mu)$  which is the characteristic function of the degenerate distribution  $\delta_\mu = \mathcal{N}(\mu, 0)$ .

This result will be used to derive the continuity of the function  $\phi \mapsto a(x, \phi)$  in the following lemma.

**Lemma 11.** *For some fixed  $x = (s, a, s')$  and  $\phi \in \Phi$ , the function*

$$\phi \mapsto a(x, \phi) := \log \int_{\xi} p_\xi(s' | s, a) p_\phi(\xi) d\xi$$

*is continuous on  $\Phi$ .*

*Proof of Lemma 11.* For  $\xi \in \Xi$ , we denote  $h_x(\xi) := p_\xi(s' | s, a)$ .

$h_x$  is continuous on  $\Xi$  (by Assumption 1) and bounded on  $\Xi$ , because

$$\forall \xi \in \Xi, |h_x(\xi)| = |p_\xi(s' | s, a)| \leq M \quad (\text{again by Assumption 1}).$$

Let  $\{\phi_n\} := \{(\mu_n, \Sigma_n)\} \in \Phi^{\mathbb{N}}$  a sequence that converges to  $\phi := (\mu, \Sigma)$ . Notice that

$$\int_{\xi} p_\xi(s' | s, a) p_{\phi_n}(\xi) d\xi = \mathbb{E}_{P_{\phi_n}} [h_x], \quad (50)$$

and since  $P_{\phi_n} \implies P_\phi$  (from Lemma 10), then  $\mathbb{E}_{P_{\phi_n}} [h_x] \xrightarrow{n \rightarrow \infty} \mathbb{E}_{P_\phi} [h_x]$ .

We then compose by the logarithm function which is continuous on  $(0, \infty)$ . This yields  $\log \mathbb{E}_{P_{\phi_n}} [h_x] \xrightarrow{n \rightarrow \infty} \log \mathbb{E}_{P_\phi} [h_x]$ . Equivalently,

$$a(x, \phi_n) \xrightarrow{n \rightarrow \infty} a(x, \phi). \quad (51)$$

This concludes the proof by the *sequential characterization of continuity*.  $\square$

Now we prove Lemma 2:

We have  $L_N(\phi) = \frac{1}{N} \sum_{i=1}^N a(X_i, \phi)$ , where  $X_i = (s_i, a_i, s'_i) \stackrel{\text{iid}}{\sim} P_{\xi^*}$ .

$\Phi$  is compact (by Assumption 2), and by Lemma 11, for each  $x$ ,  $\phi \mapsto a(x, \phi)$  is continuous on  $\Phi$ .

Additionally, the following holds for any  $\phi \in \Phi$ ,

$$|a(x, \phi)| = \left| \log \int_{\xi} p_\xi(s' | s, a) p_\phi(\xi) d\xi \right| \quad (52)$$

By Assumptions 1 and 4, we have  $c \leq \int_{\xi} p_\xi(s' | s, a) p_\phi(\xi) d\xi \leq K$ . Hence

$$|a(x, \phi)| \leq \tilde{M} := \max\{|\log c|, |\log K|\}. \quad (53)$$

Since  $L(\phi) = \mathbb{E}_{X \sim P_{\xi^*}} [a(X, \phi)]$ , this implies (by Lemma 2.4 from Newey & McFadden (1994) which is implied by Lemma 1 from Tauchen (1985)) that  $L$  is continuous on  $\Phi$  and thus uniformly continuous since  $\Phi$  is compact by *Heine-Cantor theorem*. Furthermore,

$$\sup_{\phi \in \Phi} |L_N(\phi) - L(\phi)| \xrightarrow{N \rightarrow \infty} 0. \quad (54)$$

$\square$

*Proof of Lemma 3.* Let  $\epsilon > 0$ . We consider the set defined as follows

$$C_{\phi^*, \epsilon} := \{\phi \in \Phi \mid \|\phi - \phi^*\| \geq \epsilon\}. \quad (55)$$

$C_{\phi^*, \epsilon}$  is compact because it can be written as the intersection of a compact set

$$C_{\phi^*, \epsilon} = \Phi \cap f_{\phi^*}^{-1}([\epsilon, \infty)), \quad (56)$$

where we denote  $f_{\phi^*} : \phi \mapsto \|\phi - \phi^*\|$ . Indeed,  $\Phi$  is compact (by Assumption 2) and  $f_{\phi^*}^{-1}([\epsilon, \infty))$  is closed as the inverse image of the closed set  $[\epsilon, \infty)$  by the continuous function  $f_{\phi^*}$ .

The function  $g : \phi \mapsto L(\phi^*) - L(\phi)$  is continuous (by Lemma 2) on the compact set  $C_{\phi^*, \epsilon}$ , hence by the *extreme value theorem*,  $g$  attains its minimum on  $C_{\phi^*, \epsilon}$  in some  $\tilde{\phi} \in \Phi$ .

Thus

$$\forall \phi \in C_{\phi^*, \epsilon}, L(\phi^*) - L(\phi) \geq g(\tilde{\phi}). \quad (57)$$

By Lemma 1,  $g \geq 0$  on  $\Phi$  and

$$g(\phi) = 0 \iff \phi = \phi^*. \quad (58)$$

Since  $\tilde{\phi} \neq \phi^*$  (because  $\tilde{\phi} \in C_{\phi^*, \epsilon}$ ), we have  $g(\tilde{\phi}) > 0$ . Thus, the lemma holds with the choice of  $\eta(\epsilon) = g(\tilde{\phi}) > 0$ .  $\square$

## C OMITTED PROOFS IN SECTION 5

Before proving Lemma 4, we state and prove a few preliminary lemmas.

**Notation for Strong Consistency** We define the *diameter* of  $\Phi$  by

$$\text{Diam}(\Phi) := \sup_{\phi, \psi \in \Phi} \|\phi - \psi\|. \quad (59)$$

We begin with the following technical lemma, which gives an upper bound on the number of closed balls of radius  $r = \epsilon/L$  needed to cover  $\Phi$ .

**Lemma 12.** *Let  $0 < \epsilon < 2 \text{Diam}(\Phi) L$ , and let  $N_\epsilon$  be the minimum number of closed balls of radius  $r = \frac{\epsilon}{L}$  required to cover  $\Phi$ . Then*

$$N_\epsilon \leq 4^d \left( \frac{\text{Diam}(\Phi) L}{\epsilon} \right)^d. \quad (60)$$

*Proof of Lemma 12.* We construct a sequence  $\phi_1, \phi_2, \dots$  in  $\Phi$  satisfying

$$\forall i \neq j, \quad \|\phi_i - \phi_j\| > r. \quad (61)$$

This process must terminate after finitely many steps; denote the final index by  $K$ . Indeed, if it were infinite, then compactness of  $\Phi$  would yield a convergent subsequence of  $\{\phi_n\}$ , contradicting equation 61.

By construction,

$$\Phi \subset \bigcup_{k=1}^K \text{B}(\phi_k, r), \quad (62)$$

for otherwise we could pick some  $\phi \notin \bigcup_{k=1}^K \text{B}(\phi_k, r)$  to continue the process, contradicting the definition of  $K$ . Hence  $N_\epsilon \leq K$ .

Next, observe that the closed balls  $\text{B}(\phi_k, r/2)$ ,  $k = 1, \dots, K$ , are pairwise disjoint: if there were  $\phi \in \text{B}(\phi_i, r/2) \cap \text{B}(\phi_j, r/2)$  with  $i \neq j$ , then

$$\|\phi_i - \phi_j\| \leq \|\phi_i - \phi\| + \|\phi - \phi_j\| \leq r/2 + r/2 = r, \quad (63)$$

contradicting equation 61.

Moreover, for each  $k$ ,

$$\text{B}(\phi_k, r/2) \subset \text{B}(\phi_1, \text{Diam}(\Phi) + r/2), \quad (64)$$

since if  $\|\phi - \phi_k\| \leq r/2$  then  $\|\phi - \phi_1\| \leq \|\phi - \phi_k\| + \|\phi_k - \phi_1\| \leq r/2 + \text{Diam}(\Phi)$ .

Thus

$$\bigcup_{k=1}^K \text{B}(\phi_k, r/2) \subset \text{B}(\phi_1, \text{Diam}(\Phi) + r/2), \quad (65)$$

and by comparing volumes of disjoint balls in  $\mathbb{R}^d$  we get

$$K \frac{(r/2)^d \pi^{d/2}}{\Gamma(\frac{d}{2} + 1)} \leq \frac{(\text{Diam}(\Phi) + r/2)^d \pi^{d/2}}{\Gamma(\frac{d}{2} + 1)}. \quad (66)$$

Hence

$$N_\epsilon \leq K \leq \left( 1 + \frac{2 \text{Diam}(\Phi)}{r} \right)^d \leq \left( \frac{4 \text{Diam}(\Phi)}{r} \right)^d, \quad (67)$$

where the final inequality uses  $\epsilon < 2 \text{Diam}(\Phi) L$ .  $\square$

In the following two lemmas establish a sufficient condition for the almost sure convergence.

**Lemma 13.** Let  $(A_\ell)_{\ell \geq 1}$  be a sequence of events. We have

$$P\left(\bigcup_{\ell \geq 1} A_\ell\right) = 0 \iff \forall \ell \geq 1, \quad P(A_\ell) = 0. \quad (68)$$

*Proof of Lemma 13.* If  $P\left(\bigcup_{\ell \geq 1} A_\ell\right) = 1$ , then for all  $\ell \geq 1$ , we have clearly  $P(A_\ell) \leq P\left(\bigcup_{\ell \geq 1} A_\ell\right) = 1$  and so  $P(A_\ell) = 0$ .

If  $P(A_\ell) = 0$  for every  $\ell \geq 0$ , then we have by *Boole's inequality*,

$$P\left(\bigcup_{\ell \geq 1} A_\ell\right) \leq \sum_{\ell \geq 0} P(A_\ell) = 0. \quad (69)$$

□

**Lemma 14.** Let  $\{Z_n\}_n$  be a sequence of random variables. We have

$$\forall \epsilon > 0, \quad \sum_{n \geq 1} P(|Z_n| \geq \epsilon) < \infty \implies Z_n \xrightarrow[n \rightarrow \infty]{\text{a.s.}} 0. \quad (70)$$

*Proof of Lemma 14.* We have by definition of the almost sure convergence,  $Z_n \xrightarrow[n \rightarrow \infty]{\text{a.s.}} 0$  if and only if  $P\left(\lim_{n \rightarrow \infty} Z_n = 0\right) = 1$ . Equivalently,

$$P(\forall \epsilon > 0, \exists n \geq 1, \forall m \geq n, |Z_n| < \epsilon) = 1, \quad (71)$$

and since we can replace  $\epsilon$  by any sequence of positive real numbers that converges to 0, the previous condition is equivalent to

$$P\left(\bigcap_{\ell \geq 1} \bigcup_{n \geq 1} \bigcap_{m \geq n} \left\{|Z_n| < \frac{1}{\ell}\right\}\right) = 1. \quad (72)$$

Considering the complementary event, this is equivalent to

$$P\left(\bigcup_{\ell \geq 1} \bigcap_{n \geq 1} \bigcup_{m \geq n} \left\{|Z_n| \geq \frac{1}{\ell}\right\}\right) = 0. \quad (73)$$

Using Lemma 13, in order to have the almost sure convergence of  $Z_n$  to 0, it is sufficient to prove that

$$\forall \ell \geq 1, \quad P\left(\bigcap_{n \geq 1} \bigcup_{m \geq n} \left\{|Z_n| \geq \frac{1}{\ell}\right\}\right) = 0. \quad (74)$$

Now suppose that for all  $\epsilon > 0$ ,  $\sum_{n \geq 1} P(|Z_n| \geq \epsilon) < \infty$ . This implies that for all  $\ell \geq 1$ , we have

$$\sum_{n \geq 1} P\left(|Z_n| \geq \frac{1}{\ell}\right) < \infty. \quad (75)$$

Using *Borel-Cantelli lemma*, this implies that

$$\forall \ell \geq 1, \quad P\left(\bigcap_{n \geq 1} \bigcup_{m \geq n} \left\{|Z_n| \geq \frac{1}{\ell}\right\}\right) = 0. \quad (76)$$

This concludes the proof. □

**Lemma 15.** For any fixed  $\phi \in \Phi$ , and  $\epsilon > 0$ , we have

$$P(|L_N(\phi) - L(\phi)| \geq \epsilon) \leq 2 \exp\left(-\frac{N\epsilon^2}{2\tilde{M}^2}\right). \quad (77)$$

*Proof of Lemma 15.* We have

$$L_N(\phi) = \frac{1}{N} \sum_{i=1}^N a(X_i, \phi), \quad L(\phi) = \mathbb{E}_{X \sim P^*} [a(X, \phi)], \quad (78)$$

where  $X, X_1, \dots, X_N \stackrel{\text{iid}}{\sim} P_{\xi^*}$ .

We already establish that  $|a(x, \phi)| \leq \tilde{M}$  (see 53), hence

$$P(|L_N(\phi) - L(\phi)| \geq \epsilon) = P\left(\left|\sum_{i=1}^N (a(X_i, \phi) - \mathbb{E}_{X \sim P^*} [a(X, \phi)])\right| \geq N\epsilon\right) \quad (79)$$

$$\leq 2 \exp\left(-\frac{2(N\epsilon)^2}{\sum_{i=1}^N (2\tilde{M})^2}\right) \quad (80)$$

$$\leq 2 \exp\left(-\frac{N\epsilon^2}{2\tilde{M}^2}\right), \quad (81)$$

where Inequality 80 results from *Hoeffding's inequality*.  $\square$

*Proof of Lemma 4.* Let  $0 < \epsilon < 2DL$ . We cover  $\Phi$  by  $N_\epsilon$  closed balls of radius  $r = \epsilon/L$ , i.e.,

$$\Phi \subset \bigcup_{k=1}^{N_\epsilon} \mathbb{B}(\phi_k, r),$$

for some  $\phi_1, \dots, \phi_{N_\epsilon} \in \Phi$ , where  $N_\epsilon \leq 4^d \left(\frac{DL}{\epsilon}\right)^d$  by Lemma 12.

For all  $\phi \in \Phi$ , there exists an integer  $1 \leq k(\phi) \leq N_\epsilon$  such that  $\|\phi - \phi_{k(\phi)}\| \leq r$ , hence it follows from Assumption 5 that

$$\forall x, \quad \|a(x, \phi) - a(x, \phi_{k(\phi)})\| \leq L \|\phi - \phi_{k(\phi)}\| \leq Lr = \epsilon. \quad (82)$$

We have

$$|L_N(\phi) - L(\phi)| \leq |L_N(\phi) - L_N(\phi_{k(\phi)})| + |L_N(\phi_{k(\phi)}) - L(\phi_{k(\phi)})| + |L(\phi_{k(\phi)}) - L(\phi)|.$$

The first term can be bounded using Inequality 82 as follows,

$$|L_N(\phi) - L_N(\phi_{k(\phi)})| \leq \frac{1}{N} \sum_{i=1}^N |a(X_i, \phi) - a(X_i, \phi_{k(\phi)})| \leq \epsilon. \quad (83)$$

Similarly, the third term satisfies

$$|L(\phi_{k(\phi)}) - L(\phi)| = |\mathbb{E}_X a(X, \phi_{k(\phi)}) - \mathbb{E}_X a(X, \phi)| \leq \mathbb{E}_X |a(X, \phi_{k(\phi)}) - a(X, \phi)| \leq \epsilon,$$

where the first equality holds from *Jensen's inequality*.

Putting these inequalities together yields

$$\sup_{\phi \in \Phi} |L_N(\phi) - L(\phi)| \leq \max_{i=1, \dots, N_\epsilon} |L_N(\phi_i) - L(\phi_i)| + 2\epsilon. \quad (84)$$

This implies that

$$P\left(\sup_{\phi \in \Phi} |L_N(\phi) - L(\phi)| \geq 3\epsilon\right) \leq P\left(\max_{i=1, \dots, N_\epsilon} |L_N(\phi_i) - L(\phi_i)| \geq \epsilon\right) \quad (85)$$

$$\leq \sum_{i=1}^{N_\epsilon} P(|L_N(\phi_i) - L(\phi_i)| \geq \epsilon) \quad (86)$$

$$\leq \sum_{i=1}^{N_\epsilon} 2 \exp\left(-\frac{N\epsilon^2}{2M^2}\right) \quad (87)$$

$$= 2N_\epsilon \exp\left(-\frac{N\epsilon^2}{2M^2}\right) \quad (88)$$

$$\leq 2 \cdot 4^d \left(\frac{DL}{\epsilon}\right)^d \exp\left(-\frac{N\epsilon^2}{2M^2}\right) \quad (89)$$

where Equation (86) uses union bound, Equation (87) follows from Lemma 15 and the last inequality follows from Lemma 12.

This yields when  $N \rightarrow \infty$

$$P\left(\sup_{\phi \in \Phi} |L_N(\phi) - L(\phi)| \geq 3\epsilon\right) = o\left(\frac{1}{N^2}\right). \quad (90)$$

This assures that  $\sum_{N \geq 1} P(\sup_{\phi \in \Phi} |L_N(\phi) - L(\phi)| \geq 3\epsilon) < \infty$ , which gives by Lemma 14:

$$\sup_{\phi \in \Phi} |L_N(\phi) - L(\phi)| \xrightarrow[N \rightarrow \infty]{\text{a.s.}} 0. \quad (91)$$

□

*Proof of Theorem 2.* By the preceding lemma we have the event

$$P\left(\Omega_0 := \left\{\omega : \sup_{\phi \in \Phi} |L_N(\phi, \omega) - L(\phi)| \xrightarrow[N \rightarrow \infty]{} 0\right\}\right) = 1. \quad (92)$$

Fix  $\omega \in \Omega_0$  and, let  $\epsilon > 0$ . From Lemma 3, there exists  $\eta > 0$  such that

$$\forall \phi \in \Phi, \quad \|\phi^* - \phi\| \geq \epsilon \implies L(\phi^*) - L(\phi) \geq \eta > 0. \quad (93)$$

Since  $\omega \in \Omega_0$ , there exists a random index  $N_0(\omega, \eta)$  with

$$\sup_{\phi \in \Phi} |L_N(\phi) - L(\phi)| < \eta/3 \quad \forall N \geq N_0(\omega, \eta). \quad (94)$$

Take  $N \geq N_0(\omega, \eta)$  and suppose, towards a contradiction, that  $\|\widehat{\phi}_N(\omega) - \phi^*\| \geq \epsilon$ . Then, using equation 93–equation 94,

$$L_N(\widehat{\phi}_N(\omega)) \leq L(\widehat{\phi}_N(\omega)) + \eta/3 \leq L(\phi^*) - \eta + \eta/3 = L(\phi^*) - 2\eta/3 \leq L_N(\phi^*) - \eta/3 < L_N(\phi^*), \quad (95)$$

which contradicts the maximality of  $\widehat{\phi}_N(\omega)$ . Hence, for all  $N \geq N_0(\omega, \eta)$ ,

$$\|\widehat{\phi}_N(\omega) - \phi^*\| < \epsilon. \quad (96)$$

This implies that  $\Omega_0 \subset \left\{\omega : \widehat{\phi}_N(\omega) \xrightarrow[N \rightarrow \infty]{} \phi^*\right\}$ . Since  $P(\Omega_0) = 1$ , we conclude

$$\widehat{\phi}_N \xrightarrow[N \rightarrow \infty]{\text{a.s.}} \phi^*. \quad (97)$$

□

## D OMITTED PROOFS IN SECTION 6

### D.1 RELAXATION OF THE MIXTURE POSITIVITY ASSUMPTION

**Lemma 6.** *The weak consistency of ODR still holds if we replace Assumption 3 with the following (weaker) assumption:*

$$P\left(\inf_{\phi} q_{\phi}(x) \leq \epsilon\right) \leq \frac{1}{(\log(\epsilon))^2} \quad \text{for } \epsilon \text{ sufficiently small.} \quad (98)$$

*Proof of Lemma 6.* We start by proving these two elementary lemmas.

**Lemma 16.** *For any almost surely non-negative random variable  $Z$ , i.e.,  $P(Z \geq 0) = 1$ , we have*

$$\mathbb{E}[Z] = \int_0^{\infty} P(Z \geq \alpha) d\alpha. \quad (99)$$

*Proof of Lemma 16.* We have

$$\int_0^{\infty} P(Z \geq \alpha) d\alpha = \int_0^{\infty} \mathbb{E}[\mathbb{1}_{Z \geq \alpha}] d\alpha \quad (100)$$

$$= \int_{\alpha=0}^{\infty} \int_{z=0}^{\infty} \mathbb{1}_{z \geq \alpha} dP(z) d\alpha \quad (101)$$

$$= \int_{z=0}^{\infty} \left[ \int_{\alpha=0}^{\infty} \mathbb{1}_{z \geq \alpha} d\alpha \right] dP(z) \quad (102)$$

$$= \int_{z=0}^{\infty} \left[ \int_{\alpha=0}^z 1 d\alpha \right] dP(z) \quad (103)$$

$$= \int_{z=0}^{\infty} z dP(z) \quad (104)$$

$$= \mathbb{E}[Z], \quad (105)$$

where Equality equation 102 follows from *Fubini-Tonelli's theorem*, and Equality equation 105 follows from the non-negativity of the random variable  $Z$ .  $\square$

**Lemma 17.** *For any positive function  $f : I \rightarrow (0, \infty)$  defined on some interval  $I \subset \mathbb{R}$ , we have*

$$\sup_x \log f(x) = \log \sup_x f(x). \quad (106)$$

*Proof of Lemma 17.* For any  $x \in I$  we have by monotonicity of the logarithm function

$$\log f(x) \leq \log \sup_x f(x), \quad (107)$$

hence,  $\sup_x \log f(x) \leq \log \sup_x f(x)$ . Furthermore,

$$f(x) = e^{\log f(x)} \leq e^{\sup_x \log f(x)}, \quad (108)$$

and taking the supremum over  $x \in I$  yields  $\sup_x f(x) \leq e^{\sup_x \log f(x)}$ , thus

$$\log \sup_x f(x) \leq \sup_x \log f(x), \quad (109)$$

which concludes the proof.  $\square$

Note that the only passage of the proof of Theorem 1 in which we use Assumption 3 is when we derive a uniform bound on the function  $a$  in Inequality 53. More precisely, we proved that

$$\forall x, \forall \phi \in \Phi, |a(x, \phi)| \leq \tilde{M} := \max\{|\log(c)|, |\log(M)|\}. \quad (110)$$

While this is sufficient to apply Lemma 2.4 from Newey & McFadden (1994), this lemma only require to bound  $a(x, \phi)$  by some quantity  $d(x)$  that is independent of  $\phi$  and integrable in  $x$ .

We have

$$|a(x, \phi)| = |\log q_\phi(x)| \quad (111)$$

$$= (\log q_\phi(x))^+ + (\log q_\phi(x))^- \quad (112)$$

where  $z^+$  and  $z^-$  denote respectively the positive and negative parts of  $z$ .

We have

$$(\log q_\phi(x))^+ = \max(0, \log q_\phi(x)) \quad (113)$$

$$= \max\left(0, \log \int_{\xi} p_{\xi}(s' | s, a) p_{\phi}(\xi) d\xi\right), \quad (114)$$

and by Assumption 1,  $p_{\xi}(s' | s, a) \leq M$ , hence  $(\log q_\phi(x))^+ \leq |\log(M)|$ . Thus, the first term of equation 112 is bounded by  $|\log(M)|$  which is independent of  $\phi$  and integrable in  $x$ .

Furthermore,

$$(\log q_\phi(x))^- = \max(0, -\log q_\phi(x)) \quad (115)$$

$$= \max\left(0, \log \frac{1}{q_\phi(x)}\right) \quad (116)$$

$$\leq \max\left(0, \sup_{\phi} \log \frac{1}{q_\phi(x)}\right) \quad (117)$$

$$= \max\left(0, \log \sup_{\phi} \frac{1}{q_\phi(x)}\right) \quad (118)$$

$$= \max\left(0, \log \frac{1}{\inf_{\phi} q_\phi(x)}\right), \quad (119)$$

where Equality equation 118 follows from Lemma 17. The last quantity is independent of  $\phi$ , so we only need it to be integrable in order for the weak consistency result to hold.

Since this quantity is non-negative, Lemma 16 yields

$$\mathbb{E} \left[ \max\left(0, \log \frac{1}{\inf_{\phi} q_\phi(x)}\right) \right] = \int_0^{\infty} P\left(\max\left(0, \log \frac{1}{\inf_{\phi} q_\phi(x)}\right) \geq \alpha\right) d\alpha \quad (120)$$

$$= \int_0^{\infty} P\left(\log \frac{1}{\inf_{\phi} q_\phi(x)} \geq \alpha\right) d\alpha \quad (121)$$

$$= \int_0^{\infty} P\left(\inf_{\phi} q_\phi(x) \leq e^{-\alpha}\right) d\alpha, \quad (122)$$

and hence we only need to have the convergence of this integral. The integrand is bounded (between 0 and 1), so the integral is always convergent on  $(0, 1]$ . Hence, it is sufficient to have the convergence of the integral on  $[1, \infty)$ , e.g., one sufficient condition might be

$$P\left(\inf_{\phi} q_\phi(x) \leq e^{-\alpha}\right) \leq \frac{1}{\alpha^2} \text{ for } \alpha \text{ sufficiently large,} \quad (123)$$

equivalently,

$$P\left(\inf_{\phi} q_\phi(x) \leq \epsilon\right) \leq \frac{1}{(\log(\epsilon))^2} \text{ for } \epsilon \text{ sufficiently small.} \quad (124)$$

Notice that Assumption 3 implies this condition, since it implies that  $\inf_{\phi} q_\phi(x) > 0$  and hence for sufficiently small  $\epsilon > 0$  we have

$$P\left(\inf_{\phi} q_\phi(x) \leq \epsilon\right) = 0 \leq \frac{1}{(\log(\epsilon))^2}. \quad (125)$$

□

## D.2 SUFFICIENT CONDITION FOR THE UNIFORM LIPSCHITZ CONTINUITY ASSUMPTION

In this section, we prove a practical sufficient condition for Assumption 5. More formally, the following holds:

**Lemma 7** (Sufficient Condition for the Uniform Lipschitz Continuity Assumption). *Suppose the following holds for every  $x = (s, a, s')$*

1. *The function  $\xi \mapsto p_\xi(s' | s, a)$  is twice continuously differentiable (of class  $C^2$ ),*
2. *There exists two constants  $G_1 > 0$  and  $G_2 > 0$  such that  $|\nabla_\xi p_\xi(s' | s, a)| \leq G_1$  and  $|\nabla_\xi^2 p_\xi(s' | s, a)| \leq G_2$ ,*

*then Assumption 5 holds with  $L = \frac{G_1 + G_2/2}{c}$ .*

Before proving this result, state and prove a technical lemma that we use in our proof.

**Lemma 18.** *For any  $c > 0$ , the logarithm function  $\log$  is  $\frac{1}{c}$ -Lipschitz on  $[c, \infty)$ .*

*Proof of Lemma 18.* Let  $x$  and  $y$  be two real numbers such that  $c \leq x < y$ . We have

$$|\log(y) - \log(x)| = \log(y) - \log(x) = \log\left(\frac{y}{x}\right) = \log\left(1 + \frac{y-x}{x}\right) \leq \frac{y-x}{x}, \quad (126)$$

and since  $x \geq c$ , it follows

$$|\log(y) - \log(x)| \leq \frac{1}{c}(y-x) = \frac{1}{c}|y-x|. \quad (127)$$

□

Notice that this result can also be proved using the *mean value inequality*.

*Proof of Lemma 7.* Our goal is to prove that under the two assumptions of Lemma 7, we have

$$\forall \phi := (\mu, \Sigma), \phi' := (\mu', \Sigma') \in \Phi, \forall x, |a(x, \phi) - a(x, \phi')| \leq L \|\phi - \phi'\|_2. \quad (128)$$

First, notice that using Lemma 18 and Assumption 3, we have

$$|a(x, \phi) - a(x, \phi')| = |\log(f_x(\phi)) - \log(f_x(\phi'))| \leq \frac{1}{c}|f_x(\phi) - f_x(\phi')|, \quad (129)$$

where we used the notation  $f_x(\phi) := q_\phi(s' | s, a) = \mathbb{E}_{\xi \sim P_\phi}[p_\xi(s' | s, a)]$ . Hence, it is sufficient to prove that  $|f_x(\phi) - f_x(\phi')| \leq \tilde{L} \|\phi - \phi'\|$  for every  $x$  for some constant  $\tilde{L} > 0$ .

We start by treating the case where  $\Sigma$  and  $\Sigma'$  are non-singular.

**Case 1: non-singular covariance matrices.** In the case where  $\Sigma$  is non-singular,

$$f_x(\phi) = \int_{\xi} h_x(\xi) \mathcal{N}(\xi; \mu, \Sigma) d\xi, \quad (130)$$

where  $h_x(\xi) := p_\xi(s' | s, a)$  and  $\mathcal{N}(\xi; \mu, \Sigma) := (2\pi)^{-d/2} \det(\Sigma)^{-1/2} \exp\left(-\frac{1}{2}(\xi - \mu)^\top \Sigma^{-1}(\xi - \mu)\right)$ . Since  $\mu \mapsto \mu^\top$  and  $\Sigma \mapsto \Sigma^{-1}$  are continuously differentiable respectively on  $\mathbb{R}^d$  and  $\text{GL}_d(\mathbb{R})$ , then the function  $\phi \mapsto \mathcal{N}(\xi; \mu, \Sigma)$  is  $C^1$  as long as  $\Sigma \succ 0$  with

$$\boxed{\nabla_\mu \mathcal{N}(\xi; \mu, \Sigma) = \Sigma^{-1}(\xi - \mu) \mathcal{N}(\xi; \mu, \Sigma)}. \quad (131)$$

and using the matrix-calculus identities

$$d \log \det \Sigma = \text{tr}(\Sigma^{-1} d\Sigma), \quad d(\Sigma^{-1}) = -\Sigma^{-1}(d\Sigma)\Sigma^{-1}, \quad (132)$$

we compute

$$d \log \mathcal{N} = d \left[ -\frac{1}{2} \log \det \Sigma - \frac{1}{2} (\xi - \mu)^\top \Sigma^{-1} (\xi - \mu) \right] \quad (133)$$

$$= -\frac{1}{2} \operatorname{tr}(\Sigma^{-1} d\Sigma) - \frac{1}{2} (\xi - \mu)^\top d(\Sigma^{-1})(\xi - \mu) \quad (134)$$

$$= -\frac{1}{2} \operatorname{tr}(\Sigma^{-1} d\Sigma) + \frac{1}{2} (\xi - \mu)^\top [\Sigma^{-1} (d\Sigma) \Sigma^{-1}] (\xi - \mu). \quad (135)$$

Since  $d\mathcal{N} = \mathcal{N} d \log \mathcal{N}$ , we get

$$d\mathcal{N} = \frac{1}{2} \mathcal{N} \left[ (\xi - \mu)^\top \Sigma^{-1} (d\Sigma) \Sigma^{-1} (\xi - \mu) - \operatorname{tr}(\Sigma^{-1} d\Sigma) \right]. \quad (136)$$

Rewriting in Frobenius inner product form,

$$d\mathcal{N} = \operatorname{tr} \left[ \left( \frac{1}{2} \mathcal{N} [\Sigma^{-1} (\xi - \mu) (\xi - \mu)^\top \Sigma^{-1} - \Sigma^{-1}] \right)^\top d\Sigma \right]. \quad (137)$$

Thus the gradient is

$$\boxed{\nabla_{\Sigma} \mathcal{N}(\xi; \mu, \Sigma) = \frac{1}{2} \mathcal{N}(\xi; \mu, \Sigma) \left[ \Sigma^{-1} (\xi - \mu) (\xi - \mu)^\top \Sigma^{-1} - \Sigma^{-1} \right]}. \quad (138)$$

On each compact subset  $K$  of  $\Phi \cap \{(\xi, \Sigma) : \Sigma \succ 0\}$ , we have by the sub-multiplicativity of the norm

$$\|h_x(\xi) \nabla_{\mu} \mathcal{N}(\xi; \mu, \Sigma)\|_2 \leq M \|\Sigma^{-1}\|_2 \|\xi - \mu\|_2 \mathcal{N}(\xi; \mu, \Sigma), \quad (139)$$

and since the function  $\phi \mapsto \|\Sigma^{-1}\|_2 \|\xi - \mu\|_2 \mathcal{N}(\xi; \mu, \Sigma)$  is continuous on  $K$ , it attains its maximum in some point of  $K$ , hence, there exists some  $\mu_0$  and  $\Sigma_0 \succ 0$  such that for all  $\phi \in K$ ,

$$\|h_x(\xi) \nabla_{\mu} \mathcal{N}(\xi; \mu, \Sigma)\| \leq M \|\Sigma_0^{-1}\| \|\xi - \mu_0\| \mathcal{N}(\xi; \mu_0, \Sigma_0), \quad (140)$$

where the right term is integrable in  $\xi$  since  $\mathbb{E}_{X \sim \mathcal{N}(\xi; \mu_0, \Sigma_0)}[\|X - \mu_0\|] < \infty$ . Furthermore,

$$\|h_x(\xi) \nabla_{\Sigma} \mathcal{N}(\xi; \mu, \Sigma)\|_F \leq \frac{1}{2} M \mathcal{N}(\xi; \mu, \Sigma) \left( \|\Sigma^{-1} (\xi - \mu) (\xi - \mu)^\top \Sigma^{-1}\|_F + \|\Sigma^{-1}\|_F \right), \quad (141)$$

and  $\|\Sigma^{-1} (\xi - \mu) (\xi - \mu)^\top \Sigma^{-1}\|_F \leq \|\Sigma^{-1}\|_F \|(\xi - \mu) (\xi - \mu)^\top\|_F \|\Sigma^{-1}\|_F$ . The middle factor can be rewritten as follows

$$\|(\xi - \mu) (\xi - \mu)^\top\|_F = \operatorname{tr} [(\xi - \mu) (\xi - \mu)^\top (\xi - \mu) (\xi - \mu)^\top] \quad (142)$$

$$= \operatorname{tr} [(\xi - \mu)^\top (\xi - \mu) (\xi - \mu)^\top (\xi - \mu)] \quad (143)$$

$$= \|\xi - \mu\|_2^2, \quad (144)$$

which yields

$$\|h_x(\xi) \nabla_{\Sigma} \mathcal{N}(\xi; \mu, \Sigma)\|_F \leq \frac{1}{2} M \mathcal{N}(\xi; \mu, \Sigma) \left( \|\Sigma^{-1}\|_F^2 \|\xi - \mu\|_2^2 + \|\Sigma^{-1}\|_F \right). \quad (145)$$

Again, the function  $\phi \mapsto \left( \|\Sigma^{-1}\|_F^2 \|\xi - \mu\|_2^2 + \|\Sigma^{-1}\|_F \right) \mathcal{N}(\xi; \mu, \Sigma)$  is continuous on  $K$ , it attains its maximum in some point of  $K$ , hence, there exists some  $\mu_1$  and  $\Sigma_1 \succ 0$  such that for all  $\phi \in K$ ,

$$\|h_x(\xi) \nabla_{\Sigma} \mathcal{N}(\xi; \mu, \Sigma)\|_F \leq \frac{1}{2} M \left( \|\Sigma_1^{-1}\|_F^2 \|\xi - \mu_1\|_2^2 + \|\Sigma_1^{-1}\|_F \right) \mathcal{N}(\xi; \mu_1, \Sigma_1), \quad (146)$$

where the right term is integrable in  $\xi$  since the Gaussian distribution has finite second order moment.

Using *Leibniz integral rule*, the function  $\phi \mapsto f_x(\phi)$  is  $C^1$  and we may interchange differentiation and integration to get

$$\nabla_{\mu} f_x(\phi) = \int_{\xi} h_x(\xi) \nabla_{\mu} \mathcal{N}(\xi; \mu, \Sigma) d\xi \quad (147)$$

$$= \int_{\xi} h_x(\xi) \Sigma^{-1} (\xi - \mu) \mathcal{N}(\xi; \mu, \Sigma) d\xi \quad (148)$$

$$= \int_{\xi} h_x(\xi) [-\nabla_{\xi} \mathcal{N}(\xi; \mu, \Sigma)] d\xi \quad (149)$$

$$= \int_{\xi} \nabla_{\xi} h_x(\xi) \mathcal{N}(\xi; \mu, \Sigma) d\xi \quad (150)$$

$$= \mathbb{E}_{\xi \sim \mathcal{N}(\mu, \Sigma)} [\nabla_{\xi} h_x(\xi)], \quad (151)$$

where Equation (150) follows from an *integration by part*.<sup>6</sup> Furthermore,

$$\nabla_{\Sigma} f_x(\phi) = \int_{\xi} h_x(\xi) \nabla_{\Sigma} \mathcal{N}(\xi; \mu, \Sigma) d\xi \quad (152)$$

$$= \int_{\xi} h_x(\xi) \frac{1}{2} \left[ \Sigma^{-1} (\xi - \mu) (\xi - \mu)^{\top} \Sigma^{-1} - \Sigma^{-1} \right] \mathcal{N}(\xi; \mu, \Sigma) d\xi \quad (153)$$

$$= \frac{1}{2} \Sigma^{-1} \mathbb{E}_{\xi \sim \mathcal{N}(\mu, \Sigma)} \left[ h_x(\xi) \left[ (\xi - \mu) (\xi - \mu)^{\top} - \Sigma \right] \right] \Sigma^{-1} \quad (154)$$

$$= \frac{1}{2} \Sigma^{-1/2} \mathbb{E}_{\xi \sim \mathcal{N}(\mu, \Sigma)} \left[ h_x(\xi) \left[ \Sigma^{-1/2} (\xi - \mu) (\Sigma^{-1/2} (\xi - \mu))^{\top} - \mathbf{I}_d \right] \right] \Sigma^{-1/2} \quad (155)$$

$$= \frac{1}{2} \Sigma^{-1/2} \mathbb{E}_{\xi \sim \mathcal{N}(\mu, \Sigma)} \left[ g(\mathbf{z}) (\mathbf{z} \mathbf{z}^{\top} - \mathbf{I}_d) \right] \Sigma^{-1/2}, \quad (156)$$

where  $\Sigma^{-1/2}$  is the unique positive definite square root of  $\Sigma^{-1}$ ,  $\mathbf{z} := \Sigma^{-1/2} (\xi - \mu)$  and  $g(\mathbf{z}) := h_x(\xi) = h_x(\Sigma^{1/2} \mathbf{z} + \mu)$ . Using the *Iterated Stein formula* (Bellec & Zhang, 2020; Stein, 1981) we have

$$\mathbb{E}_{\xi \sim \mathcal{N}(\mu, \Sigma)} \left[ g(\mathbf{z}) (\mathbf{z} \mathbf{z}^{\top} - \mathbf{I}_d) \right] = \mathbb{E}_{\xi \sim \mathcal{N}(\mu, \Sigma)} \left[ \nabla_{\mathbf{z}}^2 g(\mathbf{z}) \right] \quad (157)$$

$$= \mathbb{E}_{\xi \sim \mathcal{N}(\mu, \Sigma)} \left[ \Sigma \nabla_{\xi}^2 h_x(\xi) \right] \quad (158)$$

$$= \Sigma \mathbb{E}_{\xi \sim \mathcal{N}(\mu, \Sigma)} \left[ \nabla_{\xi}^2 h_x(\xi) \right]. \quad (159)$$

Combining this equation with Equation (156) yields

$$\nabla_{\Sigma} f_x(\phi) = \frac{1}{2} \mathbb{E}_{\xi \sim \mathcal{N}(\mu, \Sigma)} \left[ \nabla_{\xi}^2 h_x(\xi) \right]. \quad (160)$$

Since  $f_x$  is  $C^1$  when  $\Sigma \succ 0$ , for any two points  $\phi, \phi' \in \Phi$  such that  $\Sigma \succ 0$  and  $\Sigma' \succ 0$ , there is  $\tilde{\phi}$  on the segment joining them (and thus  $\tilde{\Sigma} \succ 0$ )<sup>7</sup> so that by the *mean-value theorem*

$$f_x(\phi) - f_x(\phi') = \langle \nabla_{\phi} f_x(\tilde{\phi}), \phi - \phi' \rangle. \quad (161)$$

In particular

$$|f_x(\phi) - f_x(\phi')| \leq \|\nabla_{\phi} f_x(\tilde{\phi})\| \|\phi - \phi'\|. \quad (162)$$

By assumption (ii),  $\|\nabla_{\xi} h_x\| \leq G_1$  and  $\|\nabla_{\xi}^2 h_x\| \leq G_2$ . Hence

$$\|\nabla_{\mu} f_x(\phi)\| = \|\mathbb{E}[\nabla_{\xi} h_x(\xi)]\| \leq G_1, \quad \|\nabla_{\Sigma} f_x(\phi)\| = \frac{1}{2} \|\mathbb{E}[\nabla_{\xi}^2 h_x(\xi)]\| \leq \frac{G_2}{2}. \quad (163)$$

Assembling the two blocks,

$$|f_x(\phi) - f_x(\phi')| \leq (G_1 + \frac{G_2}{2}) \|\phi - \phi'\|. \quad (164)$$

Therefore  $f_x$  is Lipschitz in  $\phi$ , with constant  $L' = G_1 + G_2/2$ , and by Lemma 18 so is  $a(x, \phi) = \log f_x(\phi)$  with constant  $L = \frac{G_1 + G_2/2}{c}$ .

**General case.** For the case where we no longer suppose that  $\Sigma$  and  $\Sigma'$  are non-singular, we use the density of the set of invertible matrices in  $M_d(\mathbb{R})$ . More precisely, there exists two sequences of non-singular matrices  $\{\Sigma_N\}_N$  and  $\{\Sigma'_N\}_N$  such that  $\Sigma_N \rightarrow \Sigma$  and  $\Sigma'_N \rightarrow \Sigma'$  when  $N \rightarrow \infty$ . We denote  $\phi_N := (\mu, \Sigma_N)$  and  $\phi'_N := (\mu, \Sigma'_N)$ . The previous result yields

$$\forall N \geq 0, \forall x, |a(x, \phi_N) - a(x, \phi'_N)| \leq L \|\phi_N - \phi'_N\|, \quad (165)$$

and thus, when  $N \rightarrow \infty$  we get

$$\forall x, |a(x, \phi) - a(x, \phi')| \leq L \|\phi - \phi'\|, \quad (166)$$

where we used the continuity of the function  $\phi \mapsto a(x, \phi)$  on  $\Phi$  (Lemma 11). This concludes the proof.  $\square$

<sup>6</sup>The first term of the integration by part vanishes since  $|h_x(\xi) \mathcal{N}(\xi; \mu, \Sigma)| \leq M \mathcal{N}(\xi; \mu, \Sigma) \xrightarrow{\|\xi\| \rightarrow \infty} 0$ .

<sup>7</sup>Indeed, there exists  $t \in [0, 1]$  such that  $\tilde{\Sigma} = t\Sigma + (1-t)\Sigma'$  where  $\Sigma \succ 0$  and  $\Sigma' \succ 0$ , thus for any  $z \in \mathbb{R}^d$ ,  $z^{\top} \tilde{\Sigma} z = tz^{\top} \Sigma z + (1-t)z^{\top} \Sigma' z > 0$ .

Notice that even if in many robotic systems have strongly non-smooth dynamics in  $(s, a)$  due to hard contacts and friction. However, this non-smoothness concerns the map  $(s, a) \mapsto p_\xi(s' | s, a)$  (for example, if the state encodes the position and velocity of a robot arm, near a hard contact the probability of next-step positive velocity can change discontinuously). However, the map  $\xi \mapsto p_\xi(s' | s, a)$  typically remains smooth with respect to the physical parameters. Intuitively, if we slightly perturb masses, friction coefficients, or gains, we expect the transition probabilities to change only slightly, even though the contact dynamics in  $(s, a)$  are themselves non-smooth.

### D.3 WEAK CONSISTENCY UNDER PARTIAL COVERAGE

**Theorem 3.** *Under Assumptions 1, 2 and 3, the following holds, Any measurable maximizer  $\hat{\phi}_N \in \arg \max_{\phi \in \Phi} L_N(\phi)$  satisfies  $\text{dist}(\hat{\phi}_N, \mathcal{Q}_\mu^*) \xrightarrow[N \rightarrow \infty]{P} 0$ <sup>8</sup>.*

*Proof of Theorem 3.* As in Theorem 1, the uniform law of large numbers holds:

$$\sup_{\phi \in \Phi} |L_N(\phi) - L(\phi)| \xrightarrow{P} 0. \quad (167)$$

Lemma 9 proves that  $\mathcal{Q}_\mu^*$  is nonempty and compact.

Fix  $\varepsilon > 0$  and define the separation (margin) outside the  $\varepsilon$ -neighborhood of  $\mathcal{Q}_\mu^*$ :

$$\eta(\varepsilon) := \inf \left\{ L(\phi^*) - L(\phi) : \phi^* \in \mathcal{Q}_\mu^*, \text{dist}(\phi, \mathcal{Q}_\mu^*) \geq \varepsilon \right\}.$$

Because  $L$  is continuous and  $\{\phi \in \Phi : \text{dist}(\phi, \mathcal{Q}_\mu^*) \geq \varepsilon\}$  is compact, we have  $\eta(\varepsilon) > 0$ .

By equation 167, there exists a sequence of events  $\mathcal{E}_N$  with  $P(\mathcal{E}_N) \rightarrow 1$  such that on  $\mathcal{E}_N$ ,

$$\sup_{\phi \in \Phi} |L_N(\phi) - L(\phi)| \leq \frac{1}{3} \eta(\varepsilon).$$

On  $\mathcal{E}_N$ , for any  $\phi$  with  $\text{dist}(\phi, \mathcal{Q}_\mu^*) \geq \varepsilon$  and any  $\phi^* \in \mathcal{Q}_\mu^*$ ,

$$L_N(\phi) \leq L(\phi) + \frac{1}{3} \eta(\varepsilon) \leq L(\phi^*) - \eta(\varepsilon) + \frac{1}{3} \eta(\varepsilon) = L(\phi^*) - \frac{2}{3} \eta(\varepsilon) < \sup_{\psi \in \mathcal{Q}_\mu^*} L_N(\psi),$$

where the last inequality uses  $L_N(\psi) \geq L(\psi) - \frac{1}{3} \eta(\varepsilon) = L(\phi^*) - \frac{1}{3} \eta(\varepsilon)$  for any  $\psi \in \mathcal{Q}_\mu^*$ . Therefore, no maximizer of  $L_N$  can lie outside the  $\varepsilon$ -neighborhood of  $\mathcal{Q}_\mu^*$  on  $\mathcal{E}_N$ . Equivalently,

$$\text{dist}(\hat{\phi}_N, \mathcal{Q}_\mu^*) < \varepsilon \quad \text{on } \mathcal{E}_N.$$

Since  $P(\mathcal{E}_N) \rightarrow 1$  and  $\varepsilon > 0$  is arbitrary, we conclude  $\text{dist}(\hat{\phi}_N, \mathcal{Q}_\mu^*) \xrightarrow{P} 0$ .  $\square$

**Lemma 9.** *Under Assumptions 1, 2 and 3 The identified set  $\mathcal{Q}_\mu^*$  is non-empty and compact and the correspondence  $\mu \mapsto \mathcal{Q}_\mu^*$  is upper hemicontinuous<sup>9</sup> with respect to total variation.*

*Proof of Lemma 9.* Write

$$L(\phi, \mu) = \mathbb{E}_{(S,A) \sim \mu} \mathbb{E}_{S' \sim p_{\xi^*}(\cdot | S, A)} [a((S, A, S'), \phi)] = \int_{\mathcal{S} \times \mathcal{A}} f_\phi(s, a) \mu(ds, da),$$

where

$$f_\phi(s, a) := \mathbb{E}_{S' | s, a} [a((s, a, S'), \phi)].$$

<sup>8</sup>where  $\text{dist}$  is the distance to a set defined by  $\text{dist}(\phi, \mathcal{Q}) := \inf_{\psi \in \mathcal{Q}} \|\phi - \psi\|$ .

<sup>9</sup>A set-valued map  $F$  is upper hemicontinuous at  $x_0$  if, whenever  $x_n \rightarrow x_0$  and  $y_n \in F(x_n)$  with  $y_n \rightarrow y$ , then  $y \in F(x_0)$ . Equivalently: for every open  $U$  with  $F(x_0) \subseteq U$ , there exists a neighborhood  $V$  of  $x_0$  such that  $F(x) \subseteq U$  for all  $x \in V$ .

**Step 1: Finite-valued and continuity in  $\phi$ .** We have

$$\sup_{\phi \in \Phi} |a(x, \phi)| \leq \widetilde{M} \quad \text{for all } x = (s, a, s').$$

Therefore  $|f_\phi(s, a)| \leq \widetilde{M}$  for all  $(s, a)$  and  $\phi$ , and  $L(\phi, \mu) \in \mathbb{R}$ . Moreover, Lemma 11 gives continuity of  $\phi \mapsto a(x, \phi)$  for each fixed  $x$ . By dominated convergence with the uniform bound  $\widetilde{M}$ , we obtain continuity (hence upper semicontinuity) of  $\phi \mapsto L(\phi, \mu)$  on  $\Phi$ .

**Step 2: Uniform TV–continuity in  $\mu$ .** Let  $\mu_n \rightarrow \mu$  in total variation. Then, for any  $\phi \in \Phi$ ,

$$|L(\phi, \mu_n) - L(\phi, \mu)| = \left| \int f_\phi(s, a) (\mu_n - \mu)(ds da) \right| \quad (168)$$

$$\leq \int |f_\phi(s, a)| |(\mu_n - \mu)|(ds da) \quad (169)$$

$$\leq \widetilde{M} \|\mu_n - \mu\|_{\text{TV}}. \quad (170)$$

Taking the supremum over  $\phi \in \Phi$  yields

$$\sup_{\phi \in \Phi} |L(\phi, \mu_n) - L(\phi, \mu)| \leq \widetilde{M} \|\mu_n - \mu\|_{\text{TV}} \xrightarrow{n \rightarrow \infty} 0. \quad (171)$$

**Step 3: Joint continuity of  $L$ .** Let  $(\phi_n, \mu_n) \rightarrow (\phi, \mu)$  with  $\phi_n \rightarrow \phi$  in  $\Phi$  and  $\mu_n \rightarrow \mu$  in TV. Then

$$|L(\phi_n, \mu_n) - L(\phi, \mu)| \leq |L(\phi_n, \mu_n) - L(\phi_n, \mu)| + |L(\phi_n, \mu) - L(\phi, \mu)|.$$

By uniform TV–continuity in  $\mu$  (from  $|a(x, \phi)| \leq \widetilde{M}$ ),

$$\sup_{\psi \in \Phi} |L(\psi, \mu_n) - L(\psi, \mu)| \leq \widetilde{M} \|\mu_n - \mu\|_{\text{TV}} \xrightarrow{n \rightarrow \infty} 0,$$

hence  $|L(\phi_n, \mu_n) - L(\phi_n, \mu)| \rightarrow 0$ . By continuity in  $\phi$  at fixed  $\mu$  (dominated convergence with the same bound),  $|L(\phi_n, \mu) - L(\phi, \mu)| \rightarrow 0$ . Therefore  $L(\phi_n, \mu_n) \rightarrow L(\phi, \mu)$ , i.e.,  $(\phi, \mu) \mapsto L(\phi, \mu)$  is jointly continuous.

Hence, by *Berge’s Maximum Theorem* (Berge, 1963), for each  $\mu$  the argmax set  $\mathcal{Q}_\mu^* = \arg \max_{\phi \in \Phi} L(\phi, \mu)$  is nonempty and compact, and the correspondence  $\mu \mapsto \mathcal{Q}_\mu^*$  is upper hemicontinuous (in total variation).  $\square$

#### D.4 MISSPECIFICATION AND REPRESENTATIVE MDPs

In this subsection, we study the case where we no longer suppose that the true dynamics  $\mathcal{M}^*$  belongs to the simulator class  $\mathcal{U}$ . Given a Markov kernel  $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ , we write  $\mathcal{M}(P)$  for the MDP  $(\mathcal{S}, \mathcal{A}, P, R, H, s_1)$  and we denote the value of a policy  $\pi$  at the initial state  $s_1$  by

$$V_{P,1}^\pi(s_1) := V_{\mathcal{M}(P),1}^\pi(s_1).$$

We use  $\text{TV}(p, q) := \frac{1}{2} \sum_{s' \in \mathcal{S}} |p(s') - q(s')|$  for the total variation distance between distributions  $p, q$  over  $\mathcal{S}$ .

**Lemma 19** (Value stability under kernel perturbations). *Let  $P, Q : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$  be two transition kernels defined on the same state–action space, with a common reward function  $R : \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$  and horizon  $H$ . Define*

$$\Delta(P, Q) := \sup_{(s,a) \in \mathcal{S} \times \mathcal{A}} \text{TV}(P(\cdot | s, a), Q(\cdot | s, a)),$$

where  $\text{TV}(p, q) := \frac{1}{2} \sum_{s' \in \mathcal{S}} |p(s') - q(s')|$  denotes the total variation distance between two probability distributions  $p, q$  on  $\mathcal{S}$ . Then, for any policy  $\pi$  and initial state  $s_1$ ,

$$|V_{P,1}^\pi(s_1) - V_{Q,1}^\pi(s_1)| \leq H^2 \Delta(P, Q).$$

*Proof.* For  $h \in \{1, \dots, H\}$ , let  $d_h^P$  and  $d_h^Q$  denote the distributions over state–action pairs  $(s_h, a_h)$  at step  $h$  when running policy  $\pi$  in the MDPs  $\mathcal{M}(P)$  and  $\mathcal{M}(Q)$  respectively, both initialized from the same state  $s_1$ . In particular,  $d_1^P = d_1^Q$ . We first control the evolution of the occupancy measures. By definition of the dynamics,

$$d_{h+1}^P(s', a') = \sum_{s,a} d_h^P(s, a) \pi_h(a' | \text{traj}_h) P(s' | s, a),$$

and analogously for  $Q$ . Hence

$$\begin{aligned} \|d_{h+1}^P - d_{h+1}^Q\|_1 &= \sum_{s', a'} \left| \sum_{s,a} d_h^P(s, a) \pi_h(a' | \text{traj}_h) P(s' | s, a) - d_h^Q(s, a) \pi_h(a' | \text{traj}_h) Q(s' | s, a) \right| \\ &\leq \sum_{s', a'} \sum_{s,a} |d_h^P(s, a) - d_h^Q(s, a)| \pi_h(a' | \text{traj}_h) P(s' | s, a) \\ &\quad + \sum_{s', a'} \sum_{s,a} d_h^Q(s, a) \pi_h(a' | \text{traj}_h) |P(s' | s, a) - Q(s' | s, a)| \\ &\leq \sum_{s,a} |d_h^P(s, a) - d_h^Q(s, a)| + \sup_{(s,a)} \sum_{s'} |P(s' | s, a) - Q(s' | s, a)| \\ &= \|d_h^P - d_h^Q\|_1 + 2 \sup_{(s,a)} \text{TV}(P(\cdot | s, a), Q(\cdot | s, a)) \\ &\leq \|d_h^P - d_h^Q\|_1 + 2\Delta(P, Q). \end{aligned}$$

Since  $\|d_1^P - d_1^Q\|_1 = 0$ , an induction on  $h$  yields

$$\|d_h^P - d_h^Q\|_1 \leq 2(h-1) \Delta(P, Q) \quad \forall h = 1, \dots, H.$$

Next, write the value of policy  $\pi$  under kernel  $P$  as

$$V_{P,1}^\pi(s_1) = \sum_{h=1}^H \mathbb{E}_{(s_h, a_h) \sim d_h^P} [R(s_h, a_h)],$$

and similarly  $V_{Q,1}^\pi(s_1)$  with  $d_h^Q$ . Since  $R(s, a) \in [0, 1]$ ,

$$|\mathbb{E}_{d_h^P}[R] - \mathbb{E}_{d_h^Q}[R]| \leq \|d_h^P - d_h^Q\|_1.$$

Therefore,

$$\begin{aligned} |V_{P,1}^\pi(s_1) - V_{Q,1}^\pi(s_1)| &\leq \sum_{h=1}^H |\mathbb{E}_{d_h^P}[R] - \mathbb{E}_{d_h^Q}[R]| \\ &\leq \sum_{h=1}^H \|d_h^P - d_h^Q\|_1 \\ &\leq \sum_{h=1}^H 2(h-1) \Delta(P, Q) \leq H^2 \Delta(P, Q), \end{aligned}$$

where the last inequality uses  $\sum_{h=1}^H (h-1) = H(H-1)/2 \leq H^2/2$ . This concludes the proof.  $\square$

**Theorem 4.** Let  $q_{\hat{\phi}_N}$  be the learned mixture kernel from offline data,  $\mathcal{M}_{\hat{\phi}_N} := \mathcal{M}(q_{\hat{\phi}_N})$  be the training MDP ODR uses, and  $\pi_N$  be the learned policy using any RL algorithm in  $\mathcal{M}_{\hat{\phi}_N}$ , then we have:

$$\text{Gap}_{\mathcal{M}^*}(\pi_N) \leq \text{Gap}_{\mathcal{M}_{\hat{\phi}_N}}(\pi_N) + 4H^2 \Delta(P^*, q_{\hat{\phi}_N}).$$

*Proof.* We have

$$\begin{aligned} &V_{\mathcal{M}^*,1}^*(s_1) - V_{\mathcal{M}^*,1}^{\pi_N}(s_1) \\ &= V_{\mathcal{M}^*,1}^*(s_1) - V_{\mathcal{M}_{\hat{\phi}_N},1}^*(s_1) + V_{\mathcal{M}_{\hat{\phi}_N},1}^*(s_1) - V_{\mathcal{M}_{\hat{\phi}_N},1}^{\pi_N}(s_1) + V_{\mathcal{M}_{\hat{\phi}_N},1}^{\pi_N}(s_1) - V_{\mathcal{M}^*,1}^{\pi_N}(s_1), \end{aligned}$$

By maximality of  $V_{\mathcal{M}_{\hat{\phi}_N},1}^*(s_1)$  we have:

$$V_{\mathcal{M}^*,1}^*(s_1) - V_{\mathcal{M}_{\hat{\phi}_N},1}^*(s_1) \leq V_{\mathcal{M}^*,1}^*(s_1) - V_{\mathcal{M}_{\hat{\phi}_N},1}^{\pi_{\mathcal{M}^*}}(s_1) \quad (172)$$

$$= V_{\mathcal{M}^*,1}^{\pi_{\mathcal{M}^*}}(s_1) - V_{\mathcal{M}_{\hat{\phi}_N},1}^{\pi_{\mathcal{M}^*}}(s_1) \quad (173)$$

$$\leq H^2 \Delta(P^*, q_{\hat{\phi}_N}), \quad (174)$$

where the last inequality follows from Lemma 19. Using the same lemma we have:

$$V_{\mathcal{M}_{\hat{\phi}_N},1}^{\pi_N}(s_1) - V_{\mathcal{M}^*,1}^{\pi_N}(s_1) \leq H^2 \Delta(P^*, q_{\hat{\phi}_N}).$$

Hence,

$$V_{\mathcal{M}^*,1}^*(s_1) - V_{\mathcal{M}^*,1}^{\pi_N}(s_1) \leq V_{\mathcal{M}_{\hat{\phi}_N},1}^*(s_1) - V_{\mathcal{M}_{\hat{\phi}_N},1}^{\pi_N}(s_1) + 2H^2 \Delta(P^*, q_{\hat{\phi}_N}).$$

□

The first term is the suboptimality of  $\pi_N$  in the *learned* mixture MDP (that ODR actually optimizes against), while the second term is a *closeness penalty* measuring how well the fitted mixture kernel  $q_{\hat{\phi}_N}$  approximates the real dynamics  $P^*$  in total variation. Under the well-specified and identifiable assumptions of our main results,  $q_{\hat{\phi}_N}$  converges to  $P^*$ , so the penalty term vanishes as  $N \rightarrow \infty$ . Under misspecification, the penalty converges to the best approximation error achievable within the simulator family.