

---

# A Theoretical Justification for Asymmetric Actor-Critic Algorithms

---

Gaspard Lambrechts<sup>1\*</sup> Damien Ernst<sup>1</sup> Aditya Mahajan<sup>2</sup>

## Abstract

In reinforcement learning for partially observable environments, many successful algorithms have been developed within the asymmetric learning paradigm. This paradigm leverages additional state information available at training time for faster learning. Although the proposed learning objectives are usually theoretically sound, these methods still lack a precise theoretical justification for their potential benefits. We propose such a justification for asymmetric actor-critic algorithms with linear function approximators by adapting a finite-time convergence analysis to this setting. The resulting finite-time bound reveals that the asymmetric critic eliminates error terms arising from aliasing in the agent state.

## 1. Introduction

Reinforcement learning (RL) is an appealing framework for solving decision making problems, notably because it makes very few assumptions about the problem at hand. In its purest form, the promise of an RL algorithm is to learn an optimal behavior from interaction with an environment whose dynamics are unknown. More formally, an RL algorithm aims to learn a policy – which is defined as a mapping from observations to actions – from interaction samples, in order to maximize a reward signal. While RL has obtained empirical successes for a plethora of challenging problems ranging from games to robotics (Mnih et al., 2015; Schrittwieser et al., 2020; Levine et al., 2015; Akkaya et al., 2019), most of these achievements have assumed full state observability. A more realistic assumption is partial state observability, where only a partial observation of the state of the environment is available for taking actions. In this setting, the optimal action generally depends on the complete history of past observations and actions. Tradi-

tional RL approaches have thus been adapted by considering history-dependent policies, usually with a recurrent neural network to process histories (Bakker, 2001; Wierstra et al., 2007; Hausknecht & Stone, 2015; Heess et al., 2015; Zhang et al., 2016; Zhu et al., 2017). Given the difficulty of learning effective history-dependent policies, various auxiliary representation learning objectives have been proposed to compress the history into useful representations (Igl et al., 2018; Buesing et al., 2018; Guo et al., 2018; Gregor et al., 2019; Han et al., 2019; Guo et al., 2020; Lee et al., 2020; Subramanian et al., 2022; Ni et al., 2024). Such methods usually seek to learn history representations that encode the belief, defined as the posterior distributions over the states given the history, which is a sufficient statistic of the history for optimal control.

While these methods are theoretically able to learn optimal history-dependent policies, they usually learn solely from the partial state observations, which can be restrictive. Indeed, assuming the same partial observability at training time and execution time can be too pessimistic for many environments, notably for those that are simulated. This motivated the asymmetric learning paradigm, where additional state information available at training time is leveraged during the process of learning a history-dependent policy. Although the optimal policies obtained by asymmetric learning are theoretically equivalent to those learned by symmetric learning, the promise of asymmetric learning is to improve the convergence speed. Early approaches proposed to imitate a privileged policy conditioned on the state (Choudhury et al., 2018), or to use an asymmetric critic conditioned on the state (Pinto et al., 2018). These heuristic methods initially lacked a theoretical framework, and a recent line of work has focused on proposing theoretically grounded asymmetric learning objectives. First, imitation learning of a privileged policy was known to be suboptimal, and it was addressed by constraining the privileged policy so that its imitation results in an optimal policy for the partially observable environment (Warrington et al., 2021). Similarly, asymmetric actor-critic approaches were proven to provide biased gradients, and an unbiased actor-critic approach was proposed by introducing the history-state value function (Baisero & Amato, 2022). In model-based RL, several works proposed world model objectives that are proved to provide sufficient statistics of the history, by leveraging

---

<sup>\*</sup>Work done at McGill University and Mila Québec. <sup>1</sup>Montefiore Institute, University of Liège <sup>2</sup>Department of Electrical and Computer Engineering, McGill University. Correspondence to: Gaspard Lambrechts <gaspard.lambrechts@uliege.be>.

the state (Avalos et al., 2024) or arbitrary state information (Lambrechts et al., 2024). Finally, asymmetric representation learning approaches were proposed to learn sufficient statistics from state samples (Wang et al., 2023; Sinha & Mahajan, 2023). It is worth noting that many recent successful applications of RL have greatly benefited from asymmetric learning, usually through an asymmetric critic (Degraeve et al., 2022; Kaufmann et al., 2023; Vasco et al., 2024).

Despite these methods being theoretically grounded, in the sense that policies satisfying these objectives are optimal policies, they still lack a theoretical justification for their potential benefit. In particular, there is no theoretical justification for the improved convergence speed of asymmetric learning. In this work, we propose such a justification for an asymmetric actor-critic algorithm, using agent-state policies and linear function approximators. Agent-state policies rely on an internal state, which is updated recurrently based on successive actions and observations, from which the next action is selected. This agent state can introduce aliasing, a phenomenon in which an agent state may correspond to two different beliefs. Our argument relies on the comparison of two analogous finite-time bounds: one for a symmetric natural actor-critic algorithm (Cayci et al., 2024), and its adaptation to the asymmetric setting that we derive in this paper. This comparison reveals that asymmetric learning eliminates error terms arising from aliasing in the agent state in symmetric learning. These aliasing terms are given by the difference between the true belief (i.e., the posterior distribution over the states given the history) and the approximate belief (i.e., the posterior distribution over the states given the agent state). This suggests that asymmetric learning may be particularly useful when aliasing is high.

A recent related work proposed a model-based asymmetric actor-critic algorithm relying on belief approximation, and proved its sample efficiency (Cai et al., 2024). It also considered agent-state policies, and studied the finite-time performance by providing a probably approximately correct (PAC) bound, instead of an expectation bound as here. While the algorithm was restricted to finite horizon and discrete spaces, notably for implementing count-based exploration strategies, it tackled the online exploration setting and its performance bound did not present a concentrability coefficient. This related analysis thus provides a promising framework for future works in a more challenging setting. However, it did not study the existing asymmetric actor-critic algorithm, and did not provide a direct comparison with symmetric learning. In contrast, we focus on providing comparable bounds for the existing model-free asymmetric actor-critic algorithm and its symmetric counterpart.

In Section 2, we formalize the environments, policies, and Q-functions that are considered. In Section 3, we introduce the asymmetric and symmetric actor-critic algorithms that

are studied. In Section 4, we provide the finite-time bounds for the asymmetric and symmetric actor-critic algorithms. Finally, in Section 5, we conclude by summarizing the contributions and providing avenues for future works.

## 2. Background

In Subsection 2.1, we introduce the decision processes and agent-state policies that are considered. Then, we introduce the asymmetric and symmetric Q-function for such policies, in Subsection 2.2 and Subsection 2.3, respectively.

### 2.1. Partially Observable Markov Decision Process

A partially observable Markov decision process (POMDP) is a tuple  $\mathcal{P} = (\mathcal{S}, \mathcal{A}, \mathcal{O}, P, T, R, O, \gamma)$ , with discrete state space  $\mathcal{S}$ , discrete action space  $\mathcal{A}$ , and discrete observation space  $\mathcal{O}$ . The initial state distribution  $P$  gives the probability  $P(s_0)$  of  $s_0 \in \mathcal{S}$  being the initial state of the decision process. The dynamics are described by the transition distribution  $T$  that gives the probability  $T(s_{t+1}|s_t, a_t)$  of  $s_{t+1} \in \mathcal{S}$  being the state resulting from action  $a_t \in \mathcal{A}$  in state  $s_t \in \mathcal{S}$ . The reward function  $R$  gives the immediate reward  $r_t = R(s_t, a_t, s_{t+1})$  of the reward  $r_t \in [0, 1]$  resulting from this transition. The observation distribution  $O$  gives the probability  $O(o_t|s_t)$  to get observation  $o_t \in \mathcal{O}$  in state  $s_t \in \mathcal{S}$ . Finally, the discount factor  $\gamma \in [0, 1]$  weights the relative importance of future rewards. Taking a sequence of  $t$  actions in the POMDP conditions its execution and provides the history  $h_t = (o_0, a_0, \dots, o_t) \in \mathcal{H}$ , where  $\mathcal{H}$  is the set of histories of arbitrary length. In general, the optimal policy in a POMDP depends on the complete history.

However, in practice it is infeasible to learn a policy conditioned on the full history, since the latter grows unboundedly with time. We consider an agent-state policy  $\pi \in \Pi_{\mathcal{M}}$  that uses an agent-state process  $\mathcal{M} = (\mathcal{Z}, U)$ , in order to take actions (Dong et al., 2022; Sinha & Mahajan, 2024). More formally, we consider a discrete agent state space  $\mathcal{Z}$ , and an update distribution  $U$  that gives the probability  $U(z_{t+1}|z_t, a_t, o_{t+1})$  of  $z_{t+1} \in \mathcal{Z}$  being the state resulting from action  $a_t \in \mathcal{A}$  and observation  $o_{t+1} \in \mathcal{O}$  in agent state  $z_t \in \mathcal{Z}$ . Note that the update distribution  $U$  also describes the initial agent state distribution with  $z_{-1} \notin \mathcal{Z}$  the null agent state and  $a_{-1} \notin \mathcal{A}$  the null action. Some examples of agent states that are often used are a sliding window of past observations, or a belief filter. Aliasing may occur when the agent state does not summarize all information from the history about the state of the environment, see Appendix A for an example. Given the agent state  $z_t$ , the policy  $\pi$  samples actions according to  $a_t \sim \pi(\cdot|z_t)$ . An agent-state policy  $\pi^* \in \Pi_{\mathcal{M}}$  is said to be optimal for an agent-state process  $\mathcal{M}$  if it maximizes the expected discounted sum of rewards:  $\pi^* \in \arg \max_{\pi \in \Pi_{\mathcal{M}}} J(\pi)$  with  $J(\pi) = \mathbb{E}^{\pi}[\sum_{t=0}^{\infty} \gamma^t R_t]$ .

In the following, we denote by  $S_t, O_t, Z_t, A_t$  and  $R_t$  the random variables induced by the POMDP  $\mathcal{P}$ . Given a POMDP  $\mathcal{P}$  and an agent-state process  $\mathcal{M}$ , the initial environment-agent state distribution  $P$  is given by,

$$P(s_0, z_0) = P(s_0) \sum_{o_0 \in \mathcal{O}} O(o_0|s_0) U(z_0|z_{-1}, a_{-1}, o_0). \quad (1)$$

Furthermore, given an agent-state policy  $\pi \in \Pi_{\mathcal{M}}$ , we define the discounted visitation distribution as,

$$d^\pi(s, z) = (1 - \gamma) \sum_{s_0, z_0} P(s_0, z_0) \times \sum_{t=0}^{\infty} \gamma^t \Pr(S_t = s, Z_t = z | S_0 = s_0, Z_0 = z_0). \quad (2)$$

Finally, we define the visitation distribution  $m$  steps from the discounted visitation distribution as,

$$d_m^\pi(s, z) = \sum_{s_0, z_0} d^\pi(s_0, z_0) \times \Pr(S_m = s, Z_m = z | S_0 = s_0, Z_0 = z_0). \quad (3)$$

In the following, we define the various value functions for the policies that we defined. Note that we use calligraphic letters  $\mathcal{Q}^\pi, \mathcal{V}^\pi$  and  $\mathcal{A}^\pi$  for the asymmetric functions, and regular letters  $Q^\pi, V^\pi$  and  $A^\pi$  for the symmetric ones.

## 2.2. Asymmetric Q-function

Similarly to the asymmetric Q-function of [Baisero & Amato \(2022\)](#), which is conditioned on  $(s, h, a)$ , we define an asymmetric Q-function that we condition on  $(s, z, a)$ , where  $z$  is the agent state resulting from history  $h$ . The asymmetric Q-function  $\mathcal{Q}^\pi$  of an agent-state policy  $\pi \in \Pi_{\mathcal{M}}$  is defined as the expected discounted sum of rewards, starting from environment state  $s$ , agent state  $z$ , and action  $a$ , and using policy  $\pi$  afterwards,

$$\mathcal{Q}^\pi(s, z, a) = \mathbb{E}^\pi \left[ \sum_{t=0}^{\infty} \gamma^t R_t \mid S_0 = s, Z_0 = z, A_0 = a \right]. \quad (4)$$

The asymmetric value function  $\mathcal{V}^\pi$  of an agent-state policy  $\pi \in \Pi_{\mathcal{M}}$  is defined as  $\mathcal{V}^\pi(s, z) = \sum_{a \in \mathcal{A}} \pi(a|z) \mathcal{Q}^\pi(s, z, a)$ . We also define the asymmetric advantage function  $\mathcal{A}^\pi(s, z, a) = \mathcal{Q}^\pi(s, z, a) - \mathcal{V}^\pi(s, z)$ .

Let us define the  $m$ -step asymmetric Bellman operator as,

$$\tilde{\mathcal{Q}}^\pi(s, z, a) = \mathbb{E}^\pi \left[ \sum_{t=0}^{m-1} \gamma^t R_t + \gamma^m \tilde{\mathcal{Q}}^\pi(S_m, Z_m, A_m) \mid S_0 = s, Z_0 = z, A_0 = a \right]. \quad (5)$$

Since this  $m$ -step asymmetric Bellman operator is  $\gamma^m$ -contractive, equation (5) has a unique fixed point  $\tilde{\mathcal{Q}}^\pi$ . Notice that, when using an agent-state policy, the environment

state and agent state  $(S_t, Z_t)$  are Markovian. Therefore, it can be shown that the fixed point  $\tilde{\mathcal{Q}}^\pi$  is the same as the asymmetric Q-function  $\mathcal{Q}^\pi$ .

## 2.3. Symmetric Q-function

The symmetric Q-function  $Q^\pi$  of an agent-state policy  $\pi \in \Pi_{\mathcal{M}}$  in a POMDP  $\mathcal{P}$  is defined as the expected discounted sum of rewards, starting from agent state  $z$  and action  $a$ , and using policy  $\pi$  afterwards,

$$Q^\pi(z, a) = \mathbb{E}^\pi \left[ \sum_{t=0}^{\infty} \gamma^t R_t \mid Z_0 = z, A_0 = a \right]. \quad (6)$$

The symmetric value function  $V^\pi$  of an agent-state policy  $\pi \in \Pi_{\mathcal{M}}$  is defined as  $V^\pi(z) = \sum_{a \in \mathcal{A}} \pi(a|z) Q^\pi(z, a)$ . We also define the symmetric advantage function  $A^\pi(z, a) = Q^\pi(z, a) - V^\pi(z)$ .

Let us define the  $m$ -step symmetric Bellman operator as,

$$\tilde{Q}^\pi(z, a) = \mathbb{E}^\pi \left[ \sum_{t=0}^{m-1} \gamma^t R_t + \gamma^m \tilde{Q}^\pi(Z_m, A_m) \mid Z_0 = z, A_0 = a \right]. \quad (7)$$

It can be verified that the  $m$ -step symmetric Bellman operator is  $\gamma^m$ -contractive. Therefore, equation (7) has a unique fixed point  $\tilde{Q}^\pi$ . However, because the agent state is not necessarily Markovian, in general  $Q^\pi \neq \tilde{Q}^\pi$ .

## 3. Natural Actor-Critic Algorithms

In this section, we present the asymmetric and symmetric natural actor-critic algorithms, which make use of an actor, or policy, and a critic, or Q-function. The asymmetric variant will use an asymmetric critic, learned using asymmetric temporal difference learning, while the symmetric variant will use a symmetric critic, learned using symmetric temporal difference learning. These temporal difference learning algorithms are presented in [Subsection 3.1](#) and [Subsection 3.2](#), respectively. Then, [Subsection 3.3](#) presents the complete natural actor-critic algorithm that uses a temporal difference learning algorithm as a subroutine.

For any Euclidean space  $\mathcal{X}$ , let  $\mathcal{B}_2(0, B)$  be the  $\ell_2$ -ball centered at the origin with radius  $B > 0$ , and let  $\Gamma_{\mathcal{C}} : \mathcal{X} \rightarrow \mathcal{C}$  be a projection operator into the closed and convex set  $\mathcal{C} \subseteq \mathcal{X}$  in  $\ell_2$ -norm:  $\Gamma_{\mathcal{C}}(x) \in \arg \min_{c \in \mathcal{C}} \|c - x\|_2^2 \subseteq \mathcal{C}, \forall x \in \mathcal{X}$ . Finally, let us define the  $\mu$ -weighted  $\ell_2$ -norm, for any probability measures  $\mu \in \Delta(\mathcal{X})$  as,

$$\|f\|_\mu = \sqrt{\sum_{x \in \mathcal{X}} \mu(x) |f(x)|^2}. \quad (8)$$

In the algorithms, we implicitly assume to be able to directly sample from the discounted visitation measure  $d^\pi$ . When this assumption is unrealistic, it is still possible to sample from  $d^\pi$  by sampling an initial timestep  $t_0 \sim \text{Geom}(1 - \gamma)$  from a geometric distribution with success rate  $1 - \gamma$ , and then taking  $t_0 - 1$  actions in the POMDP. The resulting sample  $(s_{t_0}, z_{t_0})$  follows the distribution  $d^\pi$ .

### 3.1. Asymmetric Critic

Suppose we are given features  $\phi: \mathcal{S} \times \mathcal{Z} \times \mathcal{A} \rightarrow \mathbb{R}^{d_\phi}$ . Without loss of generality, we assume  $\sup_{s,z,a} \|\phi(s, z, a)\|_2 \leq 1$ . Given a weight vector  $\beta \in \mathbb{R}^{d_\phi}$ , let  $\hat{Q}_\beta^\pi$  denote the linear approximation of the asymmetric Q-function  $Q^\pi$  that uses features  $\phi$  with weight  $\beta$ ,

$$\hat{Q}_\beta^\pi(s, z, a) = \langle \beta, \phi(s, z, a) \rangle. \quad (9)$$

Given an arbitrary projection radius  $B > 0$ , we define the hypothesis space as,

$$\mathcal{F}_\phi^B = \{(s, z, a) \mapsto \langle \beta, \phi(s, z, a) \rangle : \beta \in \mathcal{B}_2(0, B)\}. \quad (10)$$

We denote the optimal parameter of the asymmetric critic approximation by  $\beta_*^\pi \in \arg \min_{\beta \in \mathcal{B}_2(0, B)} \|\langle \beta, \phi(\cdot) \rangle - Q^\pi(\cdot)\|_d$ , and denote the corresponding approximation by  $\hat{Q}_*^\pi(\cdot) = \langle \beta_*^\pi, \phi(\cdot) \rangle$ . The corresponding error is,

$$\varepsilon_{\text{app}} = \min_{f \in \mathcal{F}_\phi^B} \|f - Q^\pi\|_d = \|\hat{Q}_*^\pi - Q^\pi\|_d, \quad (11)$$

with  $d(s, z, a) = d^\pi(s, z) \pi(a|z)$  the sampling distribution.

In Algorithm 1, we present the  $m$ -step temporal difference learning algorithm for approximating the asymmetric Q-function  $Q^\pi$  of an arbitrary agent-state policy  $\pi \in \Pi_{\mathcal{M}}$ . At each step  $k$ , the algorithm obtains one sample  $(s_{k,0}, z_{k,0}) \sim d^\pi$  from the discounted visitation distribution. Then,  $m$  actions are selected according to policy  $\pi$  to provide samples  $(a_{k,t}, r_{k,t}, s_{k,t+1}, o_{k,t+1}, z_{k,t+1})$  for  $0 \leq t < m$ . Next, the temporal difference  $\delta_k$  and semi-gradient  $g_k$  are computed, based on a last action  $a_{k,m} \sim \pi(\cdot|z_{k,m})$ ,

$$\delta_k = \sum_{i=0}^{m-1} \gamma^i r_{k,i} + \gamma^m \hat{Q}_{\beta_k}^\pi(s_{k,m}, z_{k,m}, a_{k,m}) - \hat{Q}_{\beta_k}^\pi(s_{k,0}, z_{k,0}, a_{k,0}), \quad (12)$$

$$g_k = \delta_k \nabla_\beta \hat{Q}_{\beta_k}^\pi(s_{k,0}, z_{k,0}, a_{k,0}). \quad (13)$$

Then, the semi-gradient update is performed with  $\beta_{k+1}^- = \beta_k + \alpha g_k$  and the parameters are projected onto the ball of radius  $B$ :  $\beta_{k+1} = \Gamma_{\mathcal{B}_2(0, B)}(\beta_{k+1}^-)$ . At the end, the algorithm computes the average parameter  $\bar{\beta} = \frac{1}{K} \sum_{k=0}^{K-1} \beta_k$  and returns the average approximation  $\bar{Q}^\pi = \hat{Q}_{\bar{\beta}}^\pi$ .

---

#### Algorithm 1 $m$ -step temporal difference learning algorithm

---

**input:** policy  $\pi \in \Pi_{\mathcal{M}}$ , bootstrap timestep  $m$ , step size  $\alpha$ , number of updates  $K$ , projection radius  $B$ .  
**for**  $k = 0 \dots K - 1$  **do**  
 Initialize  $(s_{k,0}, z_{k,0}) \sim d^\pi$ .  
**for**  $i = 0 \dots m - 1$  **do**  
 Select action  $a_{k,i} \sim \pi(\cdot|z_{k,i})$ .  
 Get environment state  $s_{k,i+1} \sim T(\cdot|s_{k,i}, a_{k,i})$ .  
 Get reward  $r_{k,i} = R(s_{k,i}, a_{k,i}, s_{k,i+1})$ .  
 Get observation  $o_{k,i+1} \sim O(\cdot|s_{k,i+1})$ .  
 Update agent state  $z_{k,i+1} \sim U(\cdot|z_{k,i}, a_{k,i}, o_{k,i+1})$ .  
**end for**  
 Sample last action  $a_{k,m} \sim \pi(\cdot|z_{k,m})$ .  
 Compute semi-gradient  $g_k$  according to equation (13) or equation (17).  
 Update  $\beta_{k+1} = \Gamma_{\mathcal{B}_2(0, B)}(\beta_k + \alpha g_k)$ .  
**end for**  
**return:** average estimate  $\bar{Q}^\pi(\cdot) = \hat{Q}_{\bar{\beta}}^\pi(\cdot) = \langle \bar{\beta}, \phi(\cdot) \rangle$  or  $\bar{Q}^\pi(\cdot) = \hat{Q}_{\bar{\beta}}^\pi(\cdot) = \langle \bar{\beta}, \chi(\cdot) \rangle$  with  $\bar{\beta} = \frac{1}{K} \sum_{k=0}^{K-1} \beta_k$ .

---

### 3.2. Symmetric Critic

Similarly, we suppose that we are given features  $\chi: \mathcal{Z} \times \mathcal{A} \rightarrow \mathbb{R}^{d_\chi}$ . Without loss of generality, we assume  $\sup_{z,a} \|\chi(z, a)\|_2 \leq 1$ . Given a weight vector  $\beta \in \mathbb{R}^{d_\chi}$ , let  $\hat{Q}_\beta^\pi$  denote the linear approximation of the symmetric Q-function  $Q^\pi$  that uses features  $\chi$  with weight  $\beta$ ,

$$\hat{Q}_\beta^\pi(z, a) = \langle \beta, \chi(z, a) \rangle. \quad (14)$$

The corresponding hypothesis space for an arbitrary projection radius  $B > 0$  is denoted with  $\mathcal{F}_\chi^B$ . The optimal parameter is also denoted by  $\beta_*^\pi \in \arg \min_{\beta \in \mathcal{B}_2(0, B)} \|\langle \beta, \chi(\cdot) \rangle - Q^\pi(\cdot)\|_d$ , the corresponding optimal approximation is  $\hat{Q}_*^\pi = \langle \beta_*^\pi, \chi(\cdot) \rangle$ , and the corresponding error is,

$$\varepsilon_{\text{app}} = \min_{f \in \mathcal{F}_\chi^B} \|f - Q^\pi\|_d = \|\hat{Q}_*^\pi - Q^\pi\|_d, \quad (15)$$

with  $d(z, a) = \sum_{s \in \mathcal{S}} d^\pi(s, z) \pi(a|z)$  the sampling distribution.

Algorithm 1 also presents the  $m$ -step temporal difference learning algorithm for approximating the symmetric Q-function. The latter is identical to that of the asymmetric Q-function except that states are not exploited, such that the temporal difference  $\delta_k$  and semi-gradient  $g_k$  are given by,

$$\delta_k = \sum_{i=0}^{m-1} \gamma^i r_{k,i} + \gamma^m \hat{Q}_{\beta_k}^\pi(z_{k,m}, a_{k,m}) - \hat{Q}_{\beta_k}^\pi(z_{k,0}, a_{k,0}), \quad (16)$$

$$g_k = \delta_k \nabla_\beta \hat{Q}_{\beta_k}^\pi(z_{k,0}, a_{k,0}). \quad (17)$$

At the end, the algorithm returns the average symmetric approximation  $\bar{Q}^\pi = \hat{Q}_{\bar{\beta}}^\pi$ . Note that this symmetric critic



approximation and temporal difference learning algorithm corresponds to the one proposed by Cayci et al. (2024).

### 3.3. Natural Actor-Critic Algorithms

For both the asymmetric and symmetric actor-critic algorithms, we consider a log-linear agent-state policy  $\pi_\theta \in \Pi_{\mathcal{M}}$ . More precisely, the policy uses features  $\psi: \mathcal{Z} \times \mathcal{A} \rightarrow \mathbb{R}^{d_\psi}$ , with  $\sup_{z,a} \|\psi(z,a)\|_2 \leq 1$  without loss of generality, and a softmax readout,

$$\pi_\theta(a_t|z_t) = \frac{\exp(\langle \theta, \psi(z_t, a_t) \rangle)}{\sum_{a \in \mathcal{A}} \exp(\langle \theta, \psi(z_t, a) \rangle)}. \quad (18)$$

In this work, we consider natural policy gradients, which are less sensitive to policy parametrization (Kakade, 2001). Instead of computing the policy gradient in the original metric space, the idea is to compute the policy gradient on a statistical manifold, defined by the expected Fisher information metric. The natural policy gradient is thus given by the standard policy gradient multiplied by a preconditioner Fisher information matrix. Natural policy gradients are at the core of many effective modern policy-gradient methods (Schulman et al., 2015).

The natural policy gradient of policy  $\pi_\theta \in \Pi_{\mathcal{M}}$  is defined as follows (Kakade, 2001),

$$w_*^{\pi_\theta} = (1 - \gamma) F_{\pi_\theta}^\dagger \nabla_\theta J(\pi_\theta), \quad (19)$$

where  $F_{\pi_\theta}^\dagger$  is the pseudoinverse of the Fisher information matrix, which is defined as the outer product of the score of the policy,

$$F_{\pi_\theta} = \mathbb{E}^{d^{\pi_\theta}} [\nabla_\theta \log \pi_\theta(A|Z) \otimes \nabla_\theta \log \pi_\theta(A|Z)]. \quad (20)$$

As shown in Theorem 1, the natural policy gradient  $w_*^{\pi_\theta}$  is the minimizer of the asymmetric objective (22).

**Theorem 1** (Asymmetric Natural Policy Gradient). For any POMDP  $\mathcal{P}$  and any agent-state policy  $\pi_\theta \in \Pi_{\mathcal{M}}$ , we have,

$$w_*^{\pi_\theta} = (1 - \gamma) F_{\pi_\theta}^\dagger \nabla_\theta J(\pi_\theta) \in \arg \min_{w \in \mathbb{R}^{d_\psi}} \mathcal{L}(w), \quad (21)$$

with,

$$\mathcal{L}(w) = \mathbb{E}^{d^{\pi_\theta}} \left[ (\langle \nabla_\theta \log \pi_\theta(A|Z), w \rangle - \mathcal{A}^{\pi_\theta}(S, Z, A))^2 \right]. \quad (22)$$

The proof is given in Appendix B. In practice, since the asymmetric advantage function is unknown, the algorithm estimates the natural policy gradient by stochastic gradient descent of  $\mathcal{L}(w)$  using the approximation  $\bar{\mathcal{A}}^{\pi_\theta}(S, Z, A) = \bar{\mathcal{Q}}^{\pi_\theta}(S, Z, A) - \bar{\mathcal{V}}^{\pi_\theta}(S, Z)$  with  $\bar{\mathcal{V}}^{\pi_\theta} = \sum_{a \in \mathcal{A}} \pi_\theta(a|Z) \bar{\mathcal{Q}}(S, Z, a)$ .

Our natural actor-critic algorithm generalizes the one of Cayci et al. (2024) to the asymmetric setting and is detailed in Algorithm 2. For each policy gradient step  $0 \leq t < T$ , the natural policy gradient  $w_*^{\pi_t}$  is first estimated using  $N$  steps of stochastic gradient descent. At each natural policy gradient estimation step  $0 \leq n < N$ , the algorithm samples an initial state  $(s_{t,n}, z_{t,n}) \sim d^{\pi_t}$  from the discounted distribution  $d^{\pi_t}$  and an action  $a_{t,n} \sim \pi_t(\cdot|z_{t,n})$  according to the policy  $\pi_t = \pi_{\theta_t}$ . Then, the gradient  $v_{t,n}$  of the natural policy gradient estimate  $w_{t,n}$  is computed with,

$$v_{t,n} = \nabla_w \left( \langle \nabla_\theta \log \pi_\theta(a_{t,n}|z_{t,n}), w_{t,n} \rangle - \bar{\mathcal{A}}^{\pi_\theta}(s_{t,n}, z_{t,n}, a_{t,n}) \right)^2, \quad (23)$$

The gradient step is performed with  $w_{t,n+1}^- = w_{t,n} - \zeta v_{t,n}$  and the parameters are projected onto the ball of radius  $B$ :  $w_{t,n+1} = \Gamma_{B_2(0,B)}(w_{t,n+1}^-)$ . Finally, the algorithm computes the average parameter  $\bar{w}_t = \frac{1}{N} \sum_{n=0}^{N-1} w_{t,n}$  and performs the policy gradient step:  $\theta_{t+1} = \theta_t + \eta \bar{w}_t$ . After all policy gradient steps, the final policy is returned.

---

#### Algorithm 2 Natural actor-critic algorithm

---

**input:** number of updates  $T$ , number of steps  $N$ , step sizes  $\zeta, \eta$ , projection radius  $B$ .  
 Initialize  $\theta_0 = 0$ .  
**for**  $t = 0 \dots T - 1$  **do**  
   Obtain  $\bar{\mathcal{Q}}^{\pi_t}$  or  $\bar{\mathcal{Q}}^{\pi_t}$  using Algorithm 1.  
   Initialize  $w_{t,0} = 0$   
   **for**  $n = 0 \dots N - 1$  **do**  
     Initialize  $(s_{t,n}, z_{t,n}) \sim d^{\pi_t}$ .  
     Sample  $a_{t,n} \sim \pi_{\theta_t}(\cdot|z_{t,n})$ .  
     Compute the gradient  $v_{t,n}$  of the policy gradient using equation (23) or equation (26).  
     Update  $w_{t,n+1}^- = w_{t,n} - \zeta v_{t,n}$ .  
     Project  $w_{t,n+1} = \Gamma_{B_2(0,B)}(w_{t,n+1}^-)$ .  
   **end for**  
   Update  $\theta_{t+1} = \theta_t + \eta \frac{1}{N} \sum_{n=0}^{N-1} w_{t,n}$ .  
**end for**  
**return:** final policy  $\pi_T = \pi_{\theta_T}$ .

---

As shown in Theorem 2, the natural policy gradient  $w_*^{\pi_\theta}$  is also the minimizer of the symmetric objective (25).

**Theorem 2** (Symmetric Natural Policy Gradient). For any POMDP  $\mathcal{P}$  and any agent-state policy  $\pi_\theta \in \Pi_{\mathcal{M}}$ , we have,

$$w_*^{\pi_\theta} = (1 - \gamma) F_{\pi_\theta}^\dagger \nabla_\theta J(\pi_\theta) \in \arg \min_{w \in \mathbb{R}^{d_\psi}} L(w), \quad (24)$$

with,

$$L(w) = \mathbb{E}^{d^{\pi_\theta}} \left[ (\langle \nabla_\theta \log \pi_\theta(A|Z), w \rangle - \mathcal{A}^{\pi_\theta}(Z, A))^2 \right]. \quad (25)$$

The proof is given in [Appendix B](#). As in the asymmetric case, the symmetric advantage function is unknown, and the algorithm estimates the natural gradient by stochastic gradient descent of equation (25) using the approximation  $\bar{A}^{\pi_\theta}(Z, A) = \bar{Q}^{\pi_\theta}(Z, A) - \bar{V}^{\pi_\theta}(Z)$  with  $\bar{V}^{\pi_\theta} = \sum_{a \in \mathcal{A}} \pi_\theta(a|Z) \bar{Q}^{\pi_\theta}(Z, a)$ .

[Algorithm 2](#) also presents the symmetric natural actor-critic algorithm, initially proposed by [Cayci et al. \(2024\)](#). The latter is similar to the asymmetric algorithm except that it uses the symmetric advantage function, such that the gradient of the policy gradient is given by,

$$v_{t,n} = \nabla_w \left( \langle \nabla_\theta \log \pi_\theta(a_{t,n}|z_{t,n}), w_{t,n} \rangle - \bar{A}^{\pi_\theta}(z_{t,n}, a_{t,n}) \right)^2. \quad (26)$$

While [Theorem 1](#) and [Theorem 2](#) show that  $w_*^{\pi_\theta}$  is the minimizer of both the asymmetric and the symmetric objectives, the next section establishes the benefit of using the asymmetric loss. More precisely, asymmetric learning is shown to improve the estimation of the critic and thus the advantage function, which in turn results in a better estimation of the natural policy gradient.

## 4. Finite-Time Analysis

In this section, we give the finite-time bounds of the previous algorithms in both the asymmetric and symmetric cases. The bounds of the asymmetric and symmetric temporal difference learning algorithms are presented in [Subsection 4.1](#) and [Subsection 4.2](#), respectively. In [Subsection 4.3](#), the bounds of the asymmetric and symmetric natural actor-critic algorithms are given.

We use  $\|\mu - \nu\|_{\text{TV}}$  to denote the total variation between two probability measures  $\mu, \nu \in \Delta(\mathcal{X})$  over a discrete space  $\mathcal{X}$ ,

$$\|\mu - \nu\|_{\text{TV}} = \sup_{A \subseteq \mathcal{X}} |\mu(A) - \nu(A)| \quad (27)$$

$$= \frac{1}{2} \sum_{x \in \mathcal{X}} |\mu(x) - \nu(x)|. \quad (28)$$

### 4.1. Finite-Time Bound for the Asymmetric Critic

Our main result is to establish the following finite-time bound for the Q-function approximation resulting from the asymmetric temporal difference learning algorithm detailed in [Algorithm 1](#).

**Theorem 3** (Finite-time bound for asymmetric  $m$ -step temporal difference learning). For any agent-state policy  $\pi \in \Pi_{\mathcal{M}}$ , and any  $m \in \mathbb{N}$ , we have for [Algorithm 1](#) with  $\alpha = \frac{1}{\sqrt{K}}$  and arbitrary  $B > 0$ ,

$$\sqrt{\mathbb{E} [\|Q^\pi - \bar{Q}^\pi\|_d^2]} \leq \varepsilon_{\text{td}} + \varepsilon_{\text{app}} + \varepsilon_{\text{shift}}, \quad (29)$$

where the temporal difference learning, function approximation, and distribution shift terms are given by,

$$\varepsilon_{\text{td}} = \sqrt{\frac{4B^2 + \left(\frac{1}{1-\gamma} + 2B\right)^2}{2\sqrt{K}(1-\gamma^m)}} \quad (30)$$

$$\varepsilon_{\text{app}} = \frac{1+\gamma^m}{1-\gamma^m} \min_{f \in \mathcal{F}_\phi^B} \|f - Q^\pi\|_d \quad (31)$$

$$\varepsilon_{\text{shift}} = \left(B + \frac{1}{1-\gamma}\right) \sqrt{\frac{2\gamma^m}{1-\gamma^m}} \sqrt{\|d_m - d\|_{\text{TV}}}, \quad (32)$$

with  $d(s, z, a) = d^\pi(s, z) \pi(a|z)$  the sampling distribution, and  $d_m(s, z, a) = d_m^\pi(s, z) \pi(a|z)$  the bootstrapping distribution.

The proof is given in [Appendix C](#), and adapts the proof of [Cayci et al. \(2024\)](#) to the asymmetric setting. The first term  $\varepsilon_{\text{td}}$  is the usual temporal difference error term, decreasing in  $K^{-1/4}$ . The second term  $\varepsilon_{\text{app}}$  results from the use of linear function approximators. The third term  $\varepsilon_{\text{shift}}$  arises from the distribution shift between the sampling distribution  $d^\pi \otimes \pi$  (i.e., the discounted visitation measure) and the bootstrapping distribution  $d_m^\pi \otimes \pi$  (i.e., the distribution  $m$  steps from the discounted visitation measure). It is a consequence of not assuming the existence of a stationary distribution nor assuming to sample from the stationary distribution.

### 4.2. Finite-Time Bound for the Symmetric Critic

Given a history  $h_t = (o_0, a_0, \dots, o_t)$ , the belief is defined as,

$$b_t(s_t|h_t) = \Pr(S_t = s_t | H_t = h_t). \quad (33)$$

Given an agent state  $z_t$ , the approximate belief is defined as,

$$\hat{b}_t(s_t|z_t) = \Pr(S_t = s_t | Z_t = z_t). \quad (34)$$

We obtain the following finite-time bound for the Q-function approximation resulting from the symmetric temporal difference learning algorithm detailed in [Algorithm 1](#).

**Theorem 4** (Finite-time bound for symmetric  $m$ -step temporal difference learning ([Cayci et al., 2024](#))). For any agent-state policy  $\pi \in \Pi_{\mathcal{M}}$ , and any  $m \in \mathbb{N}$ , we have for [Algorithm 1](#) with  $\alpha = \frac{1}{\sqrt{K}}$ , and arbitrary  $B > 0$ ,

$$\sqrt{\mathbb{E} [\|Q^\pi - \bar{Q}^\pi\|_d^2]} \leq \varepsilon_{\text{td}} + \varepsilon_{\text{app}} + \varepsilon_{\text{shift}} + \varepsilon_{\text{alias}}, \quad (35)$$

where the temporal difference learning, function approximation, distribution shift, and aliasing terms are given by,

$$\varepsilon_{\text{td}} = \sqrt{\frac{4B^2 + \left(\frac{1}{1-\gamma} + 2B\right)^2}{2\sqrt{K}(1-\gamma^m)}} \quad (36)$$

$$\varepsilon_{\text{app}} = \frac{1 + \gamma^m}{1 - \gamma^m} \min_{f \in \mathcal{F}_x^B} \|f - Q^\pi\|_d \quad (37)$$

$$\varepsilon_{\text{shift}} = \left( B + \frac{1}{1 - \gamma} \right) \sqrt{\frac{2\gamma^m}{1 - \gamma^m}} \sqrt{\|d_m - d\|_{\text{TV}}} \quad (38)$$

$$\varepsilon_{\text{alias}} = \frac{2}{1 - \gamma} \left\| \mathbb{E}^\pi \left[ \sum_{k=0}^{\infty} \gamma^{km} \left\| \hat{b}_{km} - b_{km} \right\|_{\text{TV}} \middle| Z_0 = \cdot \right] \right\|_d \quad (39)$$

with  $d(z, a) = \sum_{s \in \mathcal{S}} d^\pi(s, z) \pi(a|z)$  the sampling distribution, and  $d_m(z, a) = \sum_{s \in \mathcal{S}} d_m^\pi(s, z) \pi(a|z)$  the bootstrap distribution.

The first three terms are identical or analogous to the asymmetric case. The fourth term  $\varepsilon_{\text{alias}}$  results from the difference between the fixed point  $\tilde{Q}^\pi$  of the symmetric Bellman operator (7) and the true Q-function  $Q^\pi$ .

We note some minor differences with respect to the original result of Cayci et al. (2024) that appear to be typos and minor mistakes in the original proof.<sup>1</sup> We provide the corrected proof in Appendix D.

The results of Theorem 3 and Theorem 4 can be straightforwardly generalized to any other sampling distribution. However, obtaining bounds in term of  $d^\pi \otimes \pi$  is useful for bounding the performance of the actor-critic algorithm.

### 4.3. Finite-Time Bound for the Natural Actor-Critic

Following Cayci et al. (2024), we assume that there exists a concentrability coefficient  $\bar{C}_\infty < \infty$  such that  $\sup_{0 \leq t < T} \mathbb{E}[C_t] \leq \bar{C}_\infty$  with,

$$C_t = \sup_{s, z, a} \left| \frac{d^{\pi^*}(s, z) \pi^*(a|z)}{d^{\pi_{\theta_t}}(s, z) \pi_{\theta_t}(a|z)} \right|. \quad (40)$$

Roughly speaking, this assumption means that all successive policies should visit every agent states and actions visited by the optimal policy with nonzero probability. It motivates the log-linear policy parametrization in equation (18) and the initialization to the maximum entropy policy in Algorithm 2. We obtain the following finite-time bound for the suboptimality of the policy resulting from Algorithm 2.

**Theorem 5** (Finite-time bound for asymmetric and symmetric natural actor-critic algorithm). For any agent-state process  $\mathcal{M} = (\mathcal{Z}, U)$ , we have for Algorithm 2 with  $\alpha = \frac{1}{\sqrt{K}}$ ,  $\zeta = \frac{B\sqrt{1-\gamma}}{\sqrt{2N}}$ ,  $\eta = \frac{1}{\sqrt{T}}$  and arbitrary  $B > 0$ ,

$$(1 - \gamma) \min_{0 \leq t < T} \mathbb{E}[J(\pi^*) - J(\pi_t)] \leq \varepsilon_{\text{nac}} + 2\varepsilon_{\text{inf}} + \bar{C}_\infty \left( \varepsilon_{\text{actor}} + 2\varepsilon_{\text{grad}} + 2\sqrt{6} \frac{1}{T} \sum_{t=0}^{T-1} \varepsilon_{\text{critic}}^{\pi_t} \right), \quad (41)$$

<sup>1</sup>The authors notably wrongly bound the distance  $\|\tilde{Q}^* - \tilde{Q}^\pi\|_d$  by  $\varepsilon_{\text{app}}$  at one point, which nevertheless yields a similar result.

where the different terms may differ for asymmetric and symmetric critics,

$$\varepsilon_{\text{nac}} = \frac{B^2 + 2 \log |\mathcal{A}|}{2\sqrt{T}} \quad (42)$$

$$\varepsilon_{\text{actor}} = \sqrt{\frac{(2 - \gamma)B}{(1 - \gamma)\sqrt{N}}} \quad (43)$$

$$\varepsilon_{\text{inf,asym}} = 0 \quad (44)$$

$$\varepsilon_{\text{inf,sym}} = \mathbb{E}^{\pi^*} \left[ \sum_{k=0}^{\infty} \gamma^k \left\| \hat{b}_k - b_k \right\|_{\text{TV}} \right] \quad (45)$$

$$\varepsilon_{\text{grad,asym}} = \sup_{0 \leq t < T} \sqrt{\min_w \mathcal{L}_t(w)} \quad (46)$$

$$\varepsilon_{\text{grad,sym}} = \sup_{0 \leq t < T} \sqrt{\min_w L_t(w)} \quad (47)$$

and  $\varepsilon_{\text{critic}}^{\pi_t}$  is given in Theorem 3 and Theorem 4.

The first term  $\varepsilon_{\text{nac}}$  is the usual natural actor-critic term decreasing in  $T^{-1/2}$  (Agarwal et al., 2021). The second term  $\varepsilon_{\text{inf}}$  is the inference error resulting from use of an agent state in a POMDP (Cayci et al., 2024). This term is zero for the asymmetric algorithm. The third term  $\varepsilon_{\text{actor}}$  is the error resulting from the estimation of the natural policy gradient by stochastic gradient descent. The fourth term  $\varepsilon_{\text{grad}}$  is the error resulting from the use of a linear function approximator with features  $\nabla_\theta \log \pi_t(a|z)$  for the natural policy gradient. Finally, the fifth term  $\frac{1}{T} \sum_{t=0}^{T-1} \varepsilon_{\text{critic}}^{\pi_t}$  is the error arising from the successive critic approximations. Inside of each  $\varepsilon_{\text{critic}}^{\pi_t}$  terms, the aliasing term is thus zero for the asymmetric algorithm. The proof, generalizing that of Cayci et al. (2024) to the asymmetric setting, is available in Appendix E.

### 4.4. Discussion

As can be seen from Theorem 3 and Theorem 4, compared to the symmetric temporal difference learning algorithm, the asymmetric one eliminates a term arising from aliasing in the agent state, in the sense of equation (39). In other words, even for an aliased agent-state process, leveraging the state to learn the asymmetric Q-function instead of the symmetric Q-function does not suffer from aliasing, while still providing a valid critic for the policy gradient algorithm. That said, these bounds are given in expectation, and future works may want to study the variance of the error of such Q-function approximations.

From Theorem 5, we notice that the inference term (45) in the suboptimality bound vanishes in the asymmetric setting. Moreover, the average error  $\frac{1}{T} \sum_{t=0}^{T-1} \varepsilon_{\text{critic}}^{\pi_t}$  made in the evaluation of all policies  $\pi_0, \dots, \pi_{t-1}$  appears in the finite-time bound that we obtain for the suboptimality of the policy. Thus, the suboptimality bound for the actor also improves in the asymmetric setting by eliminating the aliasing terms with respect to the symmetric setting.

By diving into the proof of [Theorem 5](#) at equations (236) and (237), we understand that the Q-function error impacts the suboptimality bound through the estimation of the natural policy gradient (19). Indeed, this error term in the suboptimality bound directly results from the error on the advantage function estimation used in the target of the natural policy gradient estimation loss of equations (23) and (26). This advantage function estimation is derived from the estimation of the Q-function, such that the error on the latter directly impacts the error on the former, as detailed in equations (236) and (237). This improvement in the average critic error unfortunately comes at the expense of a different residual error  $\varepsilon_{\text{grad}}$  on the natural policy gradient loss. Indeed, as can be seen in equation (47), we obtain a residual error  $\varepsilon_{\text{grad,asym}}$  using the best approximation of the asymmetric advantage  $\mathcal{A}^{\pi_t}(s, z, a)$ , instead of a residual error  $\varepsilon_{\text{grad,sym}}$  using the best approximation of the symmetric critic  $\mathcal{A}^{\pi_t}(z, a)$ . Since both natural policy gradients are obtained through a linear regression with features  $\nabla_{\theta} \log \pi_t(a|z)$ , it is clear that the asymmetric residual error may be higher than the symmetric residual error, even in the tabular case.

We conclude that the effectiveness of asymmetric actor-critic algorithms notably results from a better approximation of the Q-function by eliminating the aliasing bias, which in turn provides a better estimate of the policy gradient.

## 5. Conclusion

In this work, we extended the unbiased asymmetric actor-critic algorithm to agent-state policies. Then, we adapted a finite-time analysis for natural actor-critic to the asymmetric setting. This analysis highlighted that on the contrary to symmetric learning, asymmetric learning is less sensitive to aliasing in the agent state. While this analysis assumed a fixed agent-state process, we argue that it is useful to interpret the causes of effectiveness of asymmetric learning with learnable agent-state processes. Indeed, aliasing can be present in the agent-state process throughout learning, and in particular at initialization. Moreover, it should be noted that this analysis can be straightforwardly generalized to learnable agent-state processes by extending the action space to select future agent states. More formally, we would extend the action space to  $\mathcal{A}^+ = \mathcal{A} \times \Delta(\mathcal{Z})$  with  $a_t^+ = (a_t, a_t^z)$ , the agent state space to  $\mathcal{Z}^+ = \mathcal{Z} \times \mathcal{O}$  with  $z_t^+ = (z_t, z_t^o)$ , and the agent-state process to  $U(z_{t+1}^+ | z_t^+, a_t, o_{t+1}) \propto \exp(a_t^{z_{t+1}^+}) \delta_{z_{t+1}^o, o_{t+1}}$ . This alternative to backpropagation through time would nevertheless still not reflect the common setting of recurrent actor-critic algorithms. We consider this as a future work that could build on recent advances in finite-time bound for recurrent actor-critic algorithms (Cayci & Eryilmaz, 2024a;b). Alternatively, generalizing this analysis to nonlinear approximators may include recurrent neural networks, which can

be seen as nonlinear approximators with a sliding window as agent state. Our analysis also motivates future work studying other asymmetric learning approaches that consider representation losses to reduce the aliasing bias (Sinha & Mahajan, 2023; Lambrechts et al., 2022; 2024).

## Acknowledgements

Gaspard Lambrechts acknowledges the financial support of the *Wallonia-Brussels Federation* for his FRIA grant.

## Impact Statement

This paper presents work whose goal is to advance the field of machine learning. There are many potential societal consequences of our work, none which we feel must be specifically highlighted here.

## References

- Agarwal, A., Kakade, S. M., Lee, J. D., and Mahajan, G. On the Theory of Policy Gradient Methods: Optimality, Approximation, and Distribution Shift. *Journal of Machine Learning Research*, 2021.
- Akkaya, I., Andrychowicz, M., Chociej, M., teusz Litwin, M., McGrew, B., Petron, A., Paino, A., Plappert, M., Powell, G., Ribas, R., Schneider, J., Tezak, N., Tworek, J., Welinder, P., Weng, L., Yuan, Q., Zaremba, W., and Zhang, L. Solving Rubik’s Cube with a Robot Hand. *arXiv:1910.07113*, 2019.
- Avalos, R., Delgrange, F., Nowe, A., Perez, G., and Roijers, D. M. The Wasserstein Believer: Learning Belief Updates for Partially Observable Environments through Reliable Latent Space Models. *The Twelfth International Conference on Learning Representations*, 2024.
- Baisero, A. and Amato, C. Unbiased Asymmetric Reinforcement Learning under Partial Observability. *International Conference on Autonomous Agents and Multiagent Systems*, 2022.
- Bakker, B. Reinforcement Learning with Long Short-Term Memory. *Advances in Neural Information Processing Systems*, 2001.
- Buesing, L., Weber, T., Racanière, S., Eslami, S. M. A., Rezende, D. J., Reichert, D. P., Viola, F., Besse, F., Gregor, K., Hassabis, D., and Wierstra, D. Learning and Querying Fast Generative Models for Reinforcement Learning. *arXiv:1802.03006*, 2018.
- Cai, Y., Liu, X., Oikonomou, A., and Zhang, K. Provable Partially Observable Reinforcement Learning with Privileged Information. *Advances in Neural Information Processing Systems*, 2024.



- Cayci, S. and Eryilmaz, A. Convergence of Gradient Descent for Recurrent Neural Networks: A Nonasymptotic Analysis. *arXiv:2402.12241*, 2024a.
- Cayci, S. and Eryilmaz, A. Recurrent Natural Policy Gradient for POMDPs. *ICML Workshop on the Foundations of Reinforcement Learning and Control*, 2024b.
- Cayci, S., He, N., and Srikant, R. Finite-Time Analysis of Natural Actor-Critic for POMDPs. *SIAM Journal on Mathematics of Data Science*, 2024.
- Choudhury, S., Bhardwaj, M., Arora, S., Kapoor, A., Ranade, G., Scherer, S., and Dey, D. Data-Driven Planning via Imitation Learning. *The International Journal of Robotics Research*, 2018.
- Degrave, J., Felici, F., Buchli, J., Neunert, M., Tracey, B. D., Carpanese, F., Ewalds, T., Hafner, R., Abdolmaleki, A., de Las Casas, D., Donner, C., Fritz, L., Galperti, C., Huber, A., Keeling, J., Tsimpoukelli, M., Kay, J., Merle, A., Moret, J.-M., Noury, S., Pesamosca, F., Pfau, D., Sauter, O., Sommariva, C., Coda, S., Duval, B., Fasoli, A., Kohli, P., Kavukcuoglu, K., Hassabis, D., and Riedmiller, M. A. Magnetic Control of Tokamak Plasmas through Deep Reinforcement Learning. *Nature*, 2022.
- Dong, S., Roy, B. V., and Zhou, Z. Simple Agent, Complex Environment: Efficient Reinforcement Learning with Agent States. *Journal of Machine Learning Research*, 2022.
- Gregor, K., Rezende, D. J., Besse, F., Wu, Y., Merzic, H., and van den Oord, A. Shaping Belief States with Generative Environment Models for RL. *Advances in Neural Information Processing Systems*, 2019.
- Guo, Z. D., Azar, M. G., Piot, B., Pires, B. A., and Ré Munosmi. Neural Predictive Belief Representations. *arXiv:1811.06407*, 2018.
- Guo, Z. D., Pires, B. A., Piot, B., Grill, J.-B., Althé, F., Munos, R., and Azar, M. G. Bootstrap Latent-Predictive Representations for Multitask Reinforcement Learning. *International Conference on Machine Learning*, 2020.
- Han, D., Doya, K., and Tani, J. Variational Recurrent Models for Solving Partially Observable Control Tasks. *International Conference on Learning Representations*, 2019.
- Hausknecht, M. and Stone, P. Deep Recurrent Q-learning for Partially Observable MDPs. *AAAI Fall Symposium Series*, 2015.
- Heess, N., Hunt, J. J., Lillicrap, T. P., and Silver, D. Memory-Based Control with Recurrent Neural Networks. *arXiv:1512.04455*, 2015.
- Igl, M., Zintgraf, L., Le, T. A., Wood, F., and Whiteson, S. Deep Variational Reinforcement Learning for POMDPs. *International Conference on Machine Learning*, 2018.
- Kaelbling, L. P., Littman, M. L., and Cassandra, A. R. Planning and Acting in Partially Observable Stochastic Domains. *Artificial Intelligence*, 1998.
- Kakade, S. and Langford, J. Approximately optimal approximate reinforcement learning. *International Conference on Machine Learning*, 2002.
- Kakade, S. M. A Natural Policy Gradient. *Advances in Neural Information Processing Systems*, 2001.
- Kaufmann, E., Bauersfeld, L., Loquercio, A., Müller, M., Koltun, V., and Scaramuzza, D. Champion-Level Drone Racing using Deep Reinforcement Learning. *Nature*, 2023.
- Lambrechts, G., Bolland, A., and Ernst, D. Recurrent Networks, Hidden States and Beliefs in Partially Observable Environments. *Transactions on Machine Learning Research*, 2022.
- Lambrechts, G., Bolland, A., and Ernst, D. Informed POMDP: Leveraging Additional Information in Model-Based RL. *Reinforcement Learning Journal*, 2024.
- Lee, A. X., Nagabandi, A., Abbeel, P., and Levine, S. Stochastic Latent Actor-Critic: Deep Reinforcement Learning with a Latent Variable Model. *Advances in Neural Information Processing Systems*, 2020.
- Levine, S., Finn, C., Darrell, T., and Abbeel, P. End-to-End Training of Deep Visuomotor Policies. *Journal of Machine Learning Research*, 2015.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M. A., Fidjeland, A. K., Ostrovski, G., Petersen, S., Beattie, C., Sadik, A., Antonoglou, I., King, H., Kumaran, D., Wierstra, D., Legg, S., and Hassabis, D. Human-Level Control through Deep Reinforcement Learning. *Nature*, 2015.
- Ni, T., Eysenbach, B., SeyedSalehi, E., Ma, M., Gehring, C., Mahajan, A., and Bacon, P.-L. Bridging State and History Representations: Understanding Self-Predictive RL. *International Conference on Learning Representations*, 2024.
- Pinto, L., Andrychowicz, M., Welinder, P., Zaremba, W., and Abbeel, P. Asymmetric Actor Critic for Image-Based Robot Learning. *Robotics: Science and Systems*, 2018.
- Schrittwieser, J., Antonoglou, I., Hubert, T., Simonyan, K., Sifre, L., Schmitt, S., Guez, A., Lockhart, E., Hassabis, D., Graepel, T., Lillicrap, T. P., and Silver, D. Mastering

- Atari, Go, Chess and Shogi by Planning with a Learned Model. *Nature*, 2020.
- Schulman, J., Levine, S., Abbeel, P., Jordan, M., and Moritz, P. Trust Region Policy Optimization. *International Conference on Machine Learning*, 2015.
- Shalev-Shwartz, S. and Ben-David, S. *Understanding Machine Learning: From Theory to Algorithms*. Cambridge University Press, 2014.
- Sinha, A. and Mahajan, A. Asymmetric Actor-Critic with Approximate Information State. *IEEE Conference on Decision and Control*, 2023.
- Sinha, A. and Mahajan, A. Agent-State Based Policies in POMDPs: Beyond Belief-State MDPs. *arXiv: 2409.15703*, 2024.
- Subramanian, J., Sinha, A., Seraj, R., and Mahajan, A. Approximate Information State for Approximate Planning and Reinforcement Learning in Partially Observed Systems. *Journal of Machine Learning Research*, 2022.
- Sutton, R. S., McAllester, D., Singh, S., and Mansour, Y. Policy Gradient Methods for Reinforcement Learning with Function Approximation. *Advances in Neural Information Processing Systems*, 1999.
- Vasco, M., Seno, T., Kawamoto, K., Subramanian, K., Wurman, P. R., and Stone, P. A Super-Human Vision-Based Reinforcement Learning Agent for Autonomous Racing in Gran Turismo. *Reinforcement Learning Journal*, 2024.
- Wang, A., Li, A. C., Klassen, T. Q., Icarte, R. T., and McIlraith, S. A. Learning Belief Representations for Partially Observable Deep RL. *International Conference on Machine Learning*, 2023.
- Warrington, A., Lavington, J. W., Scibior, A., Schmidt, M., and Wood, F. Robust Asymmetric Learning in POMDPs. *International Conference on Machine Learning*, 2021.
- Wierstra, D., Förster, A., Peters, J., and Schmidhuber, J. Solving Deep Memory POMDPs with Recurrent Policy Gradients. *International Conference on Artificial Neural Networks*, 2007.
- Zhang, M., McCarthy, Z., Finn, C., Levine, S., and Abbeel, P. Learning Deep Neural Network Policies with Continuous Memory States. *IEEE International Conference on Robotics and Automation*, 2016.
- Zhu, P., Li, X., Poupart, P., and Miao, G. On Improving Deep Reinforcement Learning for POMDPs. *arXiv: 1704.07978*, 2017.

## A. Agent State Aliasing

In this section, we provide an example of aliased agent state, and discuss the corresponding aliasing bias. For this purpose, we introduce a slightly modified version of the Tiger POMDP (Kaelbling et al., 1998), see Figure 1. In this POMDP, there are two doors: one opening on a room with a treasure on the left, and another opening on a room with a tiger on the right. There are four states for this POMDP: being in the treasure room (Treasure), being in the tiger room (Tiger), being in front of the treasure door (Left) or being in front of the tiger door (Right). The rooms are labeled outside (Left or Right), but inside it is completely dark (Dark), such that we do not observe in which room we are. When outside of the rooms, the agent can switch to the other door (Swap) or it can open the door and enter the room (Enter). Once in a room (Treasure or Tiger), the agent stays locked forever, and gets a positive reward (+1) if it is in the treasure room (Treasure) whatever the action taken (Swap or Enter). We consider the agent state to be simply the last observation (Left, Right, or Dark). Notice that the optimal agent-state policy conditioned on this agent state is also an optimal history-dependent policy. In other words, the current observation is a sufficient statistic for optimal control in this POMDP. We consider a uniform initial distributions over the four states.

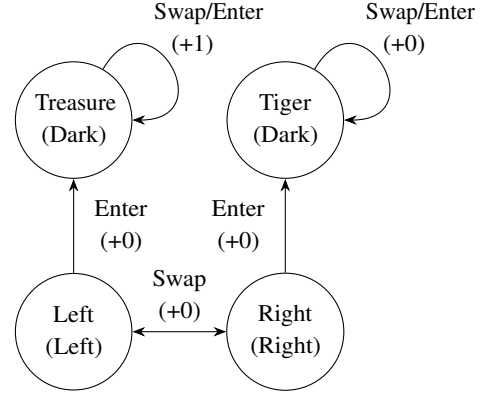


Figure 1: Aliased Tiger POMDP.

For a given agent state (Dark), there exist two different underlying states (Treasure or Tiger). We call this phenomenon aliasing. Now, let us consider a simple policy  $\pi$  that always takes the same action (Enter). It is clear that the symmetric value function defined according to equation (6) is given by  $V^\pi(z = \text{Dark}) = \frac{1}{2(1-\gamma)}$ ,  $V^\pi(z = \text{Left}) = \frac{\gamma}{1-\gamma}$ , and  $V^\pi(z = \text{Right}) = 0$ . However, when considering the unique fixed point of the aliased Bellman operator of equation (7) with  $m = 1$ , we have instead  $\tilde{V}^\pi(z = \text{Dark}) = \frac{1}{2(1-\gamma)}$ ,  $\tilde{V}^\pi(z = \text{Left}) = \frac{\gamma}{2(1-\gamma)}$ , and  $\tilde{V}^\pi(z = \text{Right}) = \frac{\gamma}{2(1-\gamma)}$ . We refer to the distance between  $V^\pi$  and  $\tilde{V}^\pi$ , or similarly  $Q^\pi$  and  $\tilde{Q}^\pi$ , as the aliasing bias. In the analysis of this paper, this distance appears as the weighted  $\ell_2$ -norm  $\|Q^\pi - \tilde{Q}^\pi\|_d$  where  $d(s, z, a) = d^\pi(s, z)\pi(a|z)$ . In the analysis, we also define the aliasing term  $\varepsilon_{\text{alias}}$  as an upper bound on this aliasing bias, see Lemma D.1 for a detailed definition.

## B. Proof of the Natural Policy Gradients

In this section, we prove that the natural policy gradient is the minimizer of analogous asymmetric and symmetric losses.

### B.1. Proof of the Asymmetric Natural Policy Gradient

In this section, we prove that the natural policy gradient is the minimizer of an asymmetric loss.

**Theorem 1** (Asymmetric Natural Policy Gradient). For any POMDP  $\mathcal{P}$  and any agent-state policy  $\pi_\theta \in \Pi_{\mathcal{M}}$ , we have,

$$w_*^{\pi_\theta} = (1 - \gamma) F_{\pi_\theta}^\dagger \nabla_\theta J(\pi_\theta) \in \arg \min_{w \in \mathbb{R}^{d_\psi}} \mathcal{L}(w), \quad (21)$$

with,

$$\mathcal{L}(w) = \mathbb{E}^{d^{\pi_\theta}} \left[ \left( \langle \nabla_\theta \log \pi_\theta(A|Z), w \rangle - \mathcal{A}^{\pi_\theta}(S, Z, A) \right)^2 \right]. \quad (22)$$

*Proof.* Let us note that,

$$\nabla_w \mathcal{L}(w) = 2 \mathbb{E}^{d^{\pi_\theta}} [\nabla_\theta \log \pi_\theta(A|Z) (\langle \nabla_\theta \log \pi_\theta(A|Z), w \rangle - \mathcal{A}^{\pi_\theta}(S, Z, A))] \quad (48)$$

Therefore, for any  $w_*^{\pi_\theta} \in \mathbb{R}^{d_\psi}$  minimizing  $\mathcal{L}(w)$ , we have  $\nabla_w \mathcal{L}(w) = 0$ , such that,

$$\mathbb{E}^{d^{\pi_\theta}} [\nabla_\theta \log \pi_\theta(A|Z) \mathcal{A}^{\pi_\theta}(S, Z, A)] = \mathbb{E}^{d^{\pi_\theta}} [\nabla_\theta \log \pi_\theta(A|Z) \langle \nabla_\theta \log \pi_\theta(A|Z), w_*^{\pi_\theta} \rangle] \quad (49)$$

$$= \mathbb{E}^{d^{\pi_\theta}} [(\nabla_\theta \log \pi_\theta(A|Z) \otimes \nabla_\theta \log \pi_\theta(A|Z)) w_*^{\pi_\theta}] \quad (50)$$

$$= \mathbb{E}^{d^{\pi_\theta}} [\nabla_\theta \log \pi_\theta(A|Z) \otimes \nabla_\theta \log \pi_\theta(A|Z)] w_*^{\pi_\theta} \quad (51)$$

$$= F_{\pi_\theta} w_*^{\pi_\theta}. \quad (52)$$

which follows from the definition of the Fisher information matrix  $F_{\pi_\theta}$  in equation (20). Now, let us define the policy  $\pi_\theta^+(A|S, Z) = \pi_\theta(A|Z)$ , which ignores the state  $S$ . From there, we have,

$$F_{\pi_\theta} w_*^{\pi_\theta} = \mathbb{E}^{d^{\pi_\theta}} [\nabla_\theta \log \pi_\theta(A|Z) \mathcal{A}(S, Z, A)] \quad (53)$$

$$= \mathbb{E}^{d^{\pi_\theta^+}} [\nabla_\theta \log \pi_\theta^+(A|S, Z) \mathcal{A}(S, Z, A)] \quad (54)$$

$$= \mathbb{E}^{d^{\pi_\theta^+}} [\nabla_\theta \log \pi_\theta^+(A|S, Z) (\mathcal{A}(S, Z, A) + \mathcal{V}(S, Z) - \mathcal{V}(S, Z))] \quad (55)$$

$$= \mathbb{E}^{d^{\pi_\theta^+}} [\nabla_\theta \log \pi_\theta^+(A|S, Z) \mathcal{Q}(S, Z, A)] - \mathbb{E}^{d^{\pi_\theta^+}} [\nabla_\theta \log \pi_\theta^+(A|S, Z) \mathcal{V}(S, Z)] \quad (56)$$

$$= \mathbb{E}^{d^{\pi_\theta^+}} [\nabla_\theta \log \pi_\theta^+(A|S, Z) \mathcal{Q}(S, Z, A)] - \mathbb{E}^{d^{\pi_\theta^+}} \left[ \mathcal{V}(S, Z) \sum_{a \in \mathcal{A}} \pi_\theta^+(a|S, Z) \nabla_\theta \log \pi_\theta^+(a|S, Z) \right] \quad (57)$$

$$= \mathbb{E}^{d^{\pi_\theta^+}} [\nabla_\theta \log \pi_\theta^+(A|S, Z) \mathcal{Q}(S, Z, A)] - \mathbb{E}^{d^{\pi_\theta^+}} \left[ \mathcal{V}(S, Z) \sum_{a \in \mathcal{A}} \nabla_\theta \pi_\theta^+(a|S, Z) \right] \quad (58)$$

$$= \mathbb{E}^{d^{\pi_\theta^+}} [\nabla_\theta \log \pi_\theta^+(A|S, Z) \mathcal{Q}(S, Z, A)] - \mathbb{E}^{d^{\pi_\theta^+}} \left[ \mathcal{V}(S, Z) \nabla_\theta \sum_{a \in \mathcal{A}} \pi_\theta^+(a|S, Z) \right] \quad (59)$$

$$= \mathbb{E}^{d^{\pi_\theta^+}} [\nabla_\theta \log \pi_\theta^+(A|S, Z) \mathcal{Q}(S, Z, A)] - \mathbb{E}^{d^{\pi_\theta^+}} [\mathcal{V}(S, Z) \nabla_\theta 1] \quad (60)$$

$$= \mathbb{E}^{d^{\pi_\theta^+}} [\nabla_\theta \log \pi_\theta^+(A|S, Z) \mathcal{Q}(S, Z, A)]. \quad (61)$$

Using the policy gradient theorem (Sutton et al., 1999) and equation (61),

$$F_{\pi_\theta} w_*^{\pi_\theta} = (1 - \gamma) \nabla_\theta J(\pi_\theta^+), \quad (62)$$

From there, we obtain using the definition of  $\pi_\theta^+$ ,

$$F_{\pi_\theta} w_*^{\pi_\theta} = (1 - \gamma) \nabla_\theta J(\pi_\theta^+) \quad (63)$$

$$= (1 - \gamma) \nabla_\theta J(\pi_\theta). \quad (64)$$

This concludes the proof.  $\square$

## B.2. Proof of the Symmetric Natural Policy Gradient

In this section, we prove that the natural policy gradient is the minimizer of an asymmetric loss.

**Theorem 2** (Symmetric Natural Policy Gradient). For any POMDP  $\mathcal{P}$  and any agent-state policy  $\pi_\theta \in \Pi_{\mathcal{M}}$ , we have,

$$w_*^{\pi_\theta} = (1 - \gamma) F_{\pi_\theta}^\dagger \nabla_\theta J(\pi_\theta) \in \arg \min_{w \in \mathbb{R}^{d_\psi}} L(w), \quad (24)$$

with,

$$L(w) = \mathbb{E}^{d^{\pi_\theta}} \left[ (\langle \nabla_\theta \log \pi_\theta(A|Z), w \rangle - A^{\pi_\theta}(Z, A))^2 \right]. \quad (25)$$

*Proof.* Similarly to the asymmetric setting, for any  $w_*^{\pi_\theta}$  minimizing  $L(w)$ , we have  $\nabla_w L(w) = 0$ , such that,

$$\mathbb{E}^{d^{\pi_\theta}} [\nabla_\theta \log \pi_\theta(A|Z) A(Z, A)] = \mathbb{E}^{d^{\pi_\theta}} [\nabla_\theta \log \pi_\theta(A|Z) \langle \nabla_\theta \log \pi_\theta(A|Z), w_*^{\pi_\theta} \rangle] \quad (65)$$

$$= \mathbb{E}^{d^{\pi_\theta}} [(\nabla_\theta \log \pi_\theta(A|Z) \otimes \nabla_\theta \log \pi_\theta(A|Z)) w_*^{\pi_\theta}] \quad (66)$$

$$= \mathbb{E}^{d^{\pi_\theta}} [\nabla_\theta \log \pi_\theta(A|Z) \otimes \nabla_\theta \log \pi_\theta(A|Z)] w_*^{\pi_\theta} \quad (67)$$

$$= F_{\pi_\theta} w_*^{\pi_\theta}, \quad (68)$$



which follows from the definition of the Fisher information matrix  $F_{\pi_\theta}$  in equation (20). From there, we have,

$$F_{\pi_\theta} w_*^{\pi_\theta} = \mathbb{E}^{d^{\pi_\theta}} [\nabla_\theta \log \pi_\theta(A|Z) \mathcal{A}(Z, A)] \quad (69)$$

$$F_{\pi_\theta} w_*^{\pi_\theta} = \mathbb{E}^{d^{\pi_\theta}} [\nabla_\theta \log \pi_\theta(A|Z) \mathbb{E}^{d^{\pi_\theta}} [\mathcal{A}(S, Z, A)|Z, A]] \quad (70)$$

$$F_{\pi_\theta} w_*^{\pi_\theta} = \mathbb{E}^{d^{\pi_\theta}} [\mathbb{E}^{d^{\pi_\theta}} [\nabla_\theta \log \pi_\theta(A|Z) \mathcal{A}(S, Z, A)|Z, A]] \quad (71)$$

$$F_{\pi_\theta} w_*^{\pi_\theta} = \mathbb{E}^{d^{\pi_\theta}} [\nabla_\theta \log \pi_\theta(A|Z) \mathcal{A}(S, Z, A)], \quad (72)$$

which follows from the law of total probability. From there, by following the same steps as in the asymmetric case (see Subsection B.1), we obtain,

$$F_{\pi_\theta} w_*^{\pi_\theta} = (1 - \gamma) \nabla_\theta J(\pi_\theta). \quad (73)$$

This concludes the proof.  $\square$

### C. Proof of the Finite-Time Bound for the Asymmetric Critic

In this section, we prove Theorem 3, that is recalled below.

**Theorem 3** (Finite-time bound for asymmetric  $m$ -step temporal difference learning). For any agent-state policy  $\pi \in \Pi_{\mathcal{M}}$ , and any  $m \in \mathbb{N}$ , we have for Algorithm 1 with  $\alpha = \frac{1}{\sqrt{K}}$  and arbitrary  $B > 0$ ,

$$\sqrt{\mathbb{E} [\|\mathcal{Q}^\pi - \bar{\mathcal{Q}}^\pi\|_d^2]} \leq \varepsilon_{\text{td}} + \varepsilon_{\text{app}} + \varepsilon_{\text{shift}}, \quad (29)$$

where the temporal difference learning, function approximation, and distribution shift terms are given by,

$$\varepsilon_{\text{td}} = \sqrt{\frac{4B^2 + \left(\frac{1}{1-\gamma} + 2B\right)^2}{2\sqrt{K}(1-\gamma^m)}} \quad (30)$$

$$\varepsilon_{\text{app}} = \frac{1+\gamma^m}{1-\gamma^m} \min_{f \in \mathcal{F}_\phi^B} \|f - \mathcal{Q}^\pi\|_d \quad (31)$$

$$\varepsilon_{\text{shift}} = \left(B + \frac{1}{1-\gamma}\right) \sqrt{\frac{2\gamma^m}{1-\gamma^m}} \sqrt{\|d_m - d\|_{\text{TV}}}, \quad (32)$$

with  $d(s, z, a) = d^\pi(s, z)\pi(a|z)$  the sampling distribution, and  $d_m(s, z, a) = d_m^\pi(s, z)\pi(a|z)$  the bootstrapping distribution.

*Proof.* To simplify notation, we drop the dependence on  $\pi$  and  $\beta$  and use  $\mathcal{Q}$  as a shorthand for  $\mathcal{Q}^\pi$ ,  $\hat{\mathcal{Q}}^*$  as a shorthand for  $\hat{\mathcal{Q}}_*^\pi$ ,  $\bar{\mathcal{Q}}$  as a shorthand for  $\bar{\mathcal{Q}}^\pi$  and  $\hat{\mathcal{Q}}_k$  as a shorthand for  $\hat{\mathcal{Q}}_{\beta_k}^\pi$ , where the subscripts and superscripts remain implicit but are assumed clear from context. When evaluating the Q-functions, we go one step further by using  $\mathcal{Q}_{k,i}$  to denote  $\mathcal{Q}(S_{k,i}, Z_{k,i}, A_{k,i})$ ,  $\hat{\mathcal{Q}}_{k,i}^*$  to denote  $\hat{\mathcal{Q}}^*(Z_{k,i}, A_{k,i})$  or  $\hat{\mathcal{Q}}_{k,i}$  to denote  $\hat{\mathcal{Q}}_k(S_{k,i}, Z_{k,i}, A_{k,i})$ , and  $\phi_{k,i}$  to denote  $\phi(S_{k,i}, Z_{k,i}, A_{k,i})$ . In addition, we define  $d$  as a shorthand for  $d^\pi \otimes \pi$ , such that  $d(s, z, a) = d^\pi(s, z)\pi(a|z)$ , and  $d_m$  as a shorthand for  $d_m^\pi \otimes \pi$ , such that  $d_m(s, z, a) = d_m^\pi(s, z)\pi(a|z)$ .

First, let us define  $\Delta_k$  as,

$$\Delta_k = \sqrt{\mathbb{E} [\|\mathcal{Q} - \hat{\mathcal{Q}}_k\|_d^2]} = \sqrt{\mathbb{E} [\|\mathcal{Q}(\cdot) - \langle \beta_k, \phi(\cdot) \rangle\|_d^2]}. \quad (74)$$

Using the linearity of  $\bar{\mathcal{Q}}$  in  $\beta_1, \dots, \beta_{K-1}$ , the triangle inequality, the subadditivity of the square root, and Jensen's inequality, we have,

$$\sqrt{\mathbb{E} [\|\mathcal{Q} - \bar{\mathcal{Q}}\|_d^2]} = \sqrt{\mathbb{E} \left[ \left\| \mathcal{Q}(\cdot) - \left\langle \frac{1}{K} \sum_{k=0}^{K-1} \beta_k, \phi(\cdot) \right\rangle \right\|_d^2 \right]} \quad (75)$$

$$= \sqrt{\mathbb{E} \left[ \left\| \frac{1}{K} \sum_{k=0}^{K-1} (\mathcal{Q}(\cdot) - \langle \beta_k, \phi(\cdot) \rangle) \right\|_d^2 \right]} \quad (76)$$

$$= \sqrt{\mathbb{E} \left[ \left\| \sum_{k=0}^{K-1} \frac{1}{K} (\mathcal{Q}(\cdot) - \langle \beta_k, \phi(\cdot) \rangle) \right\|_d^2 \right]} \quad (77)$$

$$\leq \sqrt{\mathbb{E} \left[ \sum_{k=0}^{K-1} \frac{1}{K^2} \|\mathcal{Q}(\cdot) - \langle \beta_k, \phi(\cdot) \rangle\|_d^2 \right]} \quad (78)$$

$$= \sqrt{\frac{1}{K^2} \sum_{k=0}^{K-1} \mathbb{E} [\|\mathcal{Q}(\cdot) - \langle \beta_k, \phi(\cdot) \rangle\|_d^2]} \quad (79)$$

$$= \frac{1}{K} \sqrt{\sum_{k=0}^{K-1} \Delta_k^2} \quad (80)$$

$$\leq \frac{1}{K} \sum_{k=0}^{K-1} \sqrt{\Delta_k^2} \quad (81)$$

$$= \frac{1}{K} \sum_{k=0}^{K-1} \Delta_k \quad (82)$$

$$= \frac{1}{K} \sum_{k=0}^{K-1} (\Delta_k - l) + l \quad (83)$$

$$\leq \sqrt{\left( \frac{1}{K} \sum_{k=0}^{K-1} (\Delta_k - l) \right)^2} + l \quad (84)$$

$$\leq \sqrt{\frac{1}{K} \sum_{k=0}^{K-1} (\Delta_k - l)^2} + l, \quad (85)$$

where  $l$  is arbitrary.

Now, we consider the Lyapounov function  $\mathcal{L}(\beta) = \|\beta_* - \beta\|_2^2$  in order to find a bound on  $\frac{1}{K} \sum_{k=0}^{K-1} (\Delta_k - l)^2$ . Since  $\beta_* \in \mathcal{B}_2(0, B)$ , with  $\mathcal{B}_2(0, B)$  a convex subset of  $\mathbb{R}^{d_\phi}$ , and the projection  $\Gamma_{\mathcal{C}}$  is non-expansive for closed and convex  $\mathcal{C}$ , we have for all  $k \geq 0$ ,

$$\mathcal{L}(\beta_{k+1}) = \|\beta_* - \beta_{k+1}\|_2^2 \quad (86)$$

$$\leq \|\beta_* - \beta_{k+1}^-\|_2^2 \quad (87)$$

$$= \|\beta_* - (\beta_k + \alpha g_k)\|_2^2 \quad (88)$$

$$= \|(\beta_* - \beta_k) - \alpha g_k\|_2^2 \quad (89)$$

$$= \langle (\beta_* - \beta_k) - \alpha g_k, (\beta_* - \beta_k) - \alpha g_k \rangle \quad (90)$$

$$= \langle \beta_* - \beta_k, \beta_* - \beta_k \rangle - 2\alpha \langle \beta_* - \beta_k, g_k \rangle + \alpha^2 \langle g_k, g_k \rangle \quad (91)$$

$$= \mathcal{L}(\beta_k) - 2\alpha \langle \beta_* - \beta_k, g_k \rangle + \alpha^2 \|g_k\|_2^2 \quad (92)$$

$$= \mathcal{L}(\beta_k) + 2\alpha \langle \beta_k - \beta_*, g_k \rangle + \alpha^2 \|g_k\|_2^2. \quad (93)$$

Let us consider the Lyapounov drift  $\mathbb{E}[\mathcal{L}(\beta_{k+1}) - \mathcal{L}(\beta_k)]$ , and exploit the fact that environments samples used to compute  $g_k$  are independent and identically distributed. Formally, we define  $\mathfrak{G}_k = \sigma(S_{i,j}, Z_{i,j}, A_{i,j}, i \leq k, j \leq m)$  and  $\mathfrak{F}_k = \sigma(S_{k,0}, Z_{k,0}, A_{k,0})$ , where  $\sigma(X_i : i \in \mathcal{I})$  denotes the  $\sigma$ -algebra generated by a collection  $\{X_i : i \in \mathcal{I}\}$  of random

variables. We can write, using to the law of total expectation,

$$\mathbb{E}[\mathcal{L}(\beta_{k+1}) - \mathcal{L}(\beta_k)] = \mathbb{E}[\mathbb{E}[\mathcal{L}(\beta_{k+1}) - \mathcal{L}(\beta_k) | \mathfrak{G}_{k-1}]] \quad (94)$$

$$\leq 2\alpha \mathbb{E}[\mathbb{E}[\langle \beta_k - \beta_*, g_k \rangle | \mathfrak{G}_{k-1}]] + \alpha^2 \mathbb{E}[\mathbb{E}[\|g_k\|_2^2 | \mathfrak{G}_{k-1}]]]. \quad (95)$$

Let us focus on the first term of equation (95) with  $\mathbb{E}[\langle g_k, \beta_k - \beta_* \rangle | \mathfrak{G}_{k-1}]$ . First, since  $\nabla_{\beta} \hat{\mathcal{Q}}_{k,0} = \phi_{k,0}$ , the semi-gradient  $g_k$  is given by (see equation (13)),

$$g_k = \left( \sum_{t=0}^{m-1} \gamma^t R_{k,t} + \gamma^m \hat{\mathcal{Q}}_{k,m} - \hat{\mathcal{Q}}_{k,0} \right) \phi_{k,0}. \quad (96)$$

By conditioning on the sigma-fields  $\mathfrak{G}_{k-1}$  and  $\mathfrak{F}_k$ , we have,

$$\mathbb{E}[\langle \beta_k - \beta_*, g_k \rangle | \mathfrak{F}_k, \mathfrak{G}_{k-1}] = \left( \mathbb{E} \left[ \sum_{t=0}^{m-1} \gamma^t R_{k,t} + \gamma^m \hat{\mathcal{Q}}_{k,m} \middle| \mathfrak{F}_k, \mathfrak{G}_{k-1} \right] - \hat{\mathcal{Q}}_{k,0} \right) \langle \beta_k - \beta_*, \phi_{k,0} \rangle \quad (97)$$

$$= \left( \mathbb{E} \left[ \sum_{t=0}^{m-1} \gamma^t R_{k,t} + \gamma^m \hat{\mathcal{Q}}_{k,m} \middle| \mathfrak{F}_k, \mathfrak{G}_{k-1} \right] - \hat{\mathcal{Q}}_{k,0} \right) (\hat{\mathcal{Q}}_{k,0} - \hat{\mathcal{Q}}_{k,0}^*). \quad (98)$$

Note that according to the Bellman operator (5) we have,

$$\mathbb{E} \left[ \sum_{t=0}^{m-1} \gamma^t R_{k,t} \middle| \mathfrak{F}_k, \mathfrak{G}_{k-1} \right] = \mathcal{Q}_{k,0} - \gamma^m \mathbb{E}[\mathcal{Q}_{k,m} | \mathfrak{F}_k, \mathfrak{G}_{k-1}]. \quad (99)$$

By substituting equation (99) in equation (98), we obtain,

$$\begin{aligned} & \mathbb{E}[\langle \beta_k - \beta_*, g_k \rangle | \mathfrak{F}_k, \mathfrak{G}_{k-1}] \\ &= \left( \mathbb{E} \left[ \sum_{t=0}^{m-1} \gamma^t R_{k,t} \middle| \mathfrak{F}_k, \mathfrak{G}_{k-1} \right] + \gamma^m \mathbb{E}[\hat{\mathcal{Q}}_{k,m} | \mathfrak{F}_k, \mathfrak{G}_{k-1}] - \hat{\mathcal{Q}}_{k,0} \right) (\hat{\mathcal{Q}}_{k,0} - \hat{\mathcal{Q}}_{k,0}^*) \end{aligned} \quad (100)$$

$$= (\mathcal{Q}_{k,0} - \gamma^m \mathbb{E}[\mathcal{Q}_{k,m} | \mathfrak{F}_k, \mathfrak{G}_{k-1}] + \gamma^m \mathbb{E}[\hat{\mathcal{Q}}_{k,m} | \mathfrak{F}_k, \mathfrak{G}_{k-1}] - \hat{\mathcal{Q}}_{k,0}) (\hat{\mathcal{Q}}_{k,0} - \hat{\mathcal{Q}}_{k,0}^*) \quad (101)$$

$$= ((\mathcal{Q}_{k,0} - \hat{\mathcal{Q}}_{k,0}) - \gamma^m \mathbb{E}[\mathcal{Q}_{k,m} - \hat{\mathcal{Q}}_{k,m} | \mathfrak{F}_k, \mathfrak{G}_{k-1}]) ((\hat{\mathcal{Q}}_{k,0} - \mathcal{Q}_{k,0}) + (\mathcal{Q}_{k,0} - \hat{\mathcal{Q}}_{k,0}^*)) \quad (102)$$

$$\begin{aligned} &= -(\mathcal{Q}_{k,0} - \hat{\mathcal{Q}}_{k,0})^2 + (\mathcal{Q}_{k,0} - \hat{\mathcal{Q}}_{k,0})(\mathcal{Q}_{k,0} - \hat{\mathcal{Q}}_{k,0}^*) \\ &\quad + \gamma^m \mathbb{E}[\hat{\mathcal{Q}}_{k,m} - \mathcal{Q}_{k,m} | \mathfrak{F}_k, \mathfrak{G}_{k-1}] (\hat{\mathcal{Q}}_{k,0} - \mathcal{Q}_{k,0}) + \gamma^m \mathbb{E}[\hat{\mathcal{Q}}_{k,m} - \mathcal{Q}_{k,m} | \mathfrak{F}_k, \mathfrak{G}_{k-1}] (\mathcal{Q}_{k,0} - \hat{\mathcal{Q}}_{k,0}^*). \end{aligned} \quad (103)$$

Let us now take the expectation of (103) over  $\mathfrak{F}_k$  given  $\mathfrak{G}_{k-1}$ , for each term separately,

- For the first term, we have,

$$\mathbb{E}[-(\mathcal{Q}_{k,0} - \hat{\mathcal{Q}}_{k,0})^2 | \mathfrak{G}_{k-1}] = -\|\mathcal{Q} - \hat{\mathcal{Q}}_k\|_d^2. \quad (104)$$

- For the second term, we have, using the Cauchy-Schwarz inequality,

$$\mathbb{E}[(\mathcal{Q}_{k,0} - \hat{\mathcal{Q}}_{k,0})(\mathcal{Q}_{k,0} - \hat{\mathcal{Q}}_{k,0}^*) | \mathfrak{G}_{k-1}] = \|(\mathcal{Q} - \hat{\mathcal{Q}}_k)(\mathcal{Q} - \hat{\mathcal{Q}}^*)\|_d \quad (105)$$

$$\leq \|\mathcal{Q} - \hat{\mathcal{Q}}_k\|_d \|\mathcal{Q} - \hat{\mathcal{Q}}^*\|_d. \quad (106)$$

Before proceeding to the third and fourth terms, let us notice that,

$$\mathbb{E}[\hat{\mathcal{Q}}_{k,m} - \mathcal{Q}_{k,m} | \mathfrak{G}_{k-1}] = \sum_{s,z,a} d_m(s, z, a) (\hat{\mathcal{Q}}_k(s, z, a) - \mathcal{Q}(s, z, a)) \quad (107)$$

$$= \sum_{s,z,a} (d(s,z,a) + d_m(s,z,a) - d(s,z,a)) \left( \hat{\mathcal{Q}}_k(s,z,a) - \mathcal{Q}(s,z,a) \right). \quad (108)$$

Remembering that  $\sup_{s,z,a} \hat{\mathcal{Q}}_k(s,z,a) \leq B$  and  $\sup_{s,z,a} \mathcal{Q}(s,z,a) \leq \frac{1}{1-\gamma}$ , we have,

$$\mathbb{E} \left[ \left( \hat{\mathcal{Q}}_{k,m} - \mathcal{Q}_{k,m} \right)^2 \middle| \mathfrak{G}_{k-1} \right] = \sum_{s,z,a} (d(s,z,a) + d_m(s,z,a) - d(s,z,a)) \left( \hat{\mathcal{Q}}_k(s,z,a) - \mathcal{Q}(s,z,a) \right)^2 \quad (109)$$

$$= \left\| \hat{\mathcal{Q}}_k - \mathcal{Q} \right\|_d^2 + \sum_{s,z,a} (d_m(s,z,a) - d(s,z,a)) \left( \hat{\mathcal{Q}}_k(s,z,a) - \mathcal{Q}(s,z,a) \right)^2 \quad (110)$$

$$\leq \left\| \hat{\mathcal{Q}}_k - \mathcal{Q} \right\|_d^2 + \|d_m - d\|_{\text{TV}} \sup_{s,z,a} \left( \hat{\mathcal{Q}}_k(s,z,a) - \mathcal{Q}(s,z,a) \right)^2 \quad (111)$$

$$\leq \left\| \hat{\mathcal{Q}}_k - \mathcal{Q} \right\|_d^2 + \|d_m - d\|_{\text{TV}} \left( B + \frac{1}{1-\gamma} \right)^2, \quad (112)$$

where  $\left( B + \frac{1}{1-\gamma} \right)$  is an upper bound on  $\sup_{s,z,a} \left| \hat{\mathcal{Q}}_k(s,z,a) - \mathcal{Q}(s,z,a) \right|$ . Now, using Jensen's inequality and the subadditivity of the square root, we have,

$$\mathbb{E} \left[ \left| \hat{\mathcal{Q}}_{k,m} - \mathcal{Q}_{k,m} \right| \middle| \mathfrak{G}_{k-1} \right] \leq \mathbb{E} \left[ \sqrt{\left( \hat{\mathcal{Q}}_{k,m} - \mathcal{Q}_{k,m} \right)^2} \middle| \mathfrak{G}_{k-1} \right] \quad (113)$$

$$\leq \sqrt{\mathbb{E} \left[ \left( \hat{\mathcal{Q}}_{k,m} - \mathcal{Q}_{k,m} \right)^2 \middle| \mathfrak{G}_{k-1} \right]} \quad (114)$$

$$\leq \left\| \hat{\mathcal{Q}}_k - \mathcal{Q} \right\|_d + \left( B + \frac{1}{1-\gamma} \right) \sqrt{\|d_m - d\|_{\text{TV}}}. \quad (115)$$

With this, we proceed to the third and fourth terms (without the multiplier  $\gamma^m$ ) and show the following.

- For the third term, we have by upper bounding  $|\hat{\mathcal{Q}}_{k,0} - \mathcal{Q}_{k,0}|$  by  $B + \frac{1}{1-\gamma}$ ,

$$\mathbb{E} \left[ (\hat{\mathcal{Q}}_{k,m} - \mathcal{Q}_{k,m})(\hat{\mathcal{Q}}_{k,0} - \mathcal{Q}_{k,0}) \middle| \mathfrak{G}_{k-1} \right] \leq \left\| \hat{\mathcal{Q}}_k - \mathcal{Q} \right\|_d^2 + \left( B + \frac{1}{1-\gamma} \right)^2 \sqrt{\|d_m - d\|_{\text{TV}}}. \quad (116)$$

- For the fourth term, we have by upper bounding  $|\mathcal{Q}_{k,0} - \hat{\mathcal{Q}}_{k,0}^*|$  by  $\frac{1}{1-\gamma} + B$ ,

$$\mathbb{E} \left[ (\hat{\mathcal{Q}}_{k,m} - \mathcal{Q}_{k,m})(\mathcal{Q}_{k,0} - \hat{\mathcal{Q}}_{k,0}^*) \middle| \mathfrak{G}_{k-1} \right] \leq \left\| \hat{\mathcal{Q}}_k - \mathcal{Q} \right\|_d \left\| \mathcal{Q} - \hat{\mathcal{Q}}^* \right\|_d + \left( B + \frac{1}{1-\gamma} \right)^2 \sqrt{\|d_m - d\|_{\text{TV}}}. \quad (117)$$

By taking expectation over  $\mathfrak{G}_{k-1}$  of the four terms and using the previous upper bounds, we obtain,

$$\mathbb{E} [\langle \beta_k - \beta_*, g_k \rangle] = \mathbb{E} [\mathbb{E} [\langle \beta_k - \beta_*, g_k \rangle \middle| \mathfrak{G}_{k-1}]] \quad (118)$$

$$\begin{aligned} &\leq -(1-\gamma^m) \mathbb{E} \left[ \left\| \hat{\mathcal{Q}}_k - \mathcal{Q} \right\|_d^2 \right] + (1+\gamma^m) \mathbb{E} \left[ \left\| \hat{\mathcal{Q}}_k - \mathcal{Q} \right\|_d \right] \left\| \hat{\mathcal{Q}}^* - \mathcal{Q} \right\|_d \\ &\quad + 2\gamma^m \left( B + \frac{1}{1-\gamma} \right)^2 \sqrt{\|d_m - d\|_{\text{TV}}} \end{aligned} \quad (119)$$

$$= -(1-\gamma^m) \Delta_k^2 + (1+\gamma^m) \Delta_k \left\| \hat{\mathcal{Q}}^* - \mathcal{Q} \right\|_d + 2\gamma^m \left( B + \frac{1}{1-\gamma} \right)^2 \sqrt{\|d_m - d\|_{\text{TV}}}. \quad (120)$$

Let us now focus on the second term of equation (95) with  $\mathbb{E} \left[ \|g_k\|_2^2 \middle| \mathfrak{G}_{k-1} \right]$ . Since  $\sup_{s,z,a} \|\phi(s,z,a)\|_2 \leq 1$  and  $\|\beta_k\|_2 \leq B$  for all  $k \geq 0$ , and  $r_{k,i} \leq 1$  for all  $k \geq 0$  and for all  $i < m-1$ , the norm of the gradient (96) is bounded as follows,

$$\sup_{k \geq 0} \|g_k\|_2 \leq \frac{1-\gamma^m}{1-\gamma} + (1+\gamma^m)B \leq \frac{1}{1-\gamma} + 2B. \quad (121)$$



We obtain, for the second term of equation (95),

$$\mathbb{E} [\|g_k\|_2^2] = \mathbb{E} [\mathbb{E} [\|g_k\|_2^2 | \mathfrak{G}_{k-1}]] \quad (122)$$

$$\leq \left( \frac{1}{1-\gamma} + 2B \right)^2. \quad (123)$$

By substituting equations (120) and (123) into the Lyapounov drift of equation (95), we obtain,

$$\begin{aligned} \mathbb{E} [\mathcal{L}(\beta_{k+1}) - \mathcal{L}(\beta_k)] &\leq -2\alpha(1-\gamma^m)\Delta_k^2 + 2\alpha(1+\gamma^m)\Delta_k \left\| \widehat{\mathcal{Q}}^* - \mathcal{Q} \right\|_d + \alpha^2 \left( \frac{1}{1-\gamma} + 2B \right)^2 \\ &\quad + 4\alpha\gamma^m \left( B + \frac{1}{1-\gamma} \right)^2 \sqrt{\|d_m - d\|_{\text{TV}}}. \end{aligned} \quad (124)$$

By setting  $l = \frac{1+\gamma^m}{2(1-\gamma^m)} \min_{f \in \mathcal{F}_\phi^B} \|f - \mathcal{Q}\|_d$ , we can write,

$$\begin{aligned} \mathbb{E} [\mathcal{L}(\beta_{k+1}) - \mathcal{L}(\beta_k)] &\leq -2\alpha(1-\gamma^m) (\Delta_k^2 - 2l\Delta_k) + \alpha^2 \left( \frac{1}{1-\gamma} + 2B \right)^2 \\ &\quad + 4\alpha\gamma^m \left( B + \frac{1}{1-\gamma} \right)^2 \sqrt{\|d_m - d\|_{\text{TV}}} \end{aligned} \quad (125)$$

$$\begin{aligned} &= -2\alpha(1-\gamma^m) (\Delta_k^2 - 2l\Delta_k + l^2) + 2\alpha(1-\gamma^m)l^2 + \alpha^2 \left( \frac{1}{1-\gamma} + 2B \right)^2 \\ &\quad + 4\alpha\gamma^m \left( B + \frac{1}{1-\gamma} \right)^2 \sqrt{\|d_m - d\|_{\text{TV}}} \end{aligned} \quad (126)$$

$$\begin{aligned} &= -2\alpha(1-\gamma^m) (\Delta_k - l)^2 + 2\alpha(1-\gamma^m)l^2 + \alpha^2 \left( \frac{1}{1-\gamma} + 2B \right)^2 \\ &\quad + 4\alpha\gamma^m \left( B + \frac{1}{1-\gamma} \right)^2 \sqrt{\|d_m - d\|_{\text{TV}}}. \end{aligned} \quad (127)$$

By summing all Lyapounov drifts  $\sum_{k=0}^{K-1} \mathbb{E} [\mathcal{L}(\beta_{k+1}) - \mathcal{L}(\beta_k)]$ , we get,

$$\begin{aligned} \mathbb{E} [\mathcal{L}(\beta_K) - \mathcal{L}(\beta_0)] &\leq -2\alpha(1-\gamma^m) \sum_{k=0}^{K-1} (\Delta_k - l)^2 + 2\alpha K(1-\gamma^m)l^2 + \alpha^2 K \left( \frac{1}{1-\gamma} + 2B \right)^2 \\ &\quad + 4\alpha K\gamma^m \left( B + \frac{1}{1-\gamma} \right)^2 \sqrt{\|d_m - d\|_{\text{TV}}}. \end{aligned} \quad (128)$$

By rearranging and dividing by  $2\alpha K(1-\gamma^m)$ , we obtain after neglecting  $\mathcal{L}(\beta_K) > 0$ ,

$$\begin{aligned} \frac{1}{K} \sum_{k=0}^{K-1} (\Delta_k - l)^2 &\leq \frac{\mathbb{E} [\mathcal{L}(\beta_0) - \mathcal{L}(\beta_K)]}{2\alpha K(1-\gamma^m)} + l^2 + \frac{\alpha}{2(1-\gamma^m)} \left( \frac{1}{1-\gamma} + 2B \right)^2 \\ &\quad + \frac{2\gamma^m}{1-\gamma^m} \left( B + \frac{1}{1-\gamma} \right)^2 \sqrt{\|d_m - d\|_{\text{TV}}} \end{aligned} \quad (129)$$

$$\begin{aligned} &\leq \frac{\|\beta_0 - \beta_*\|_2^2}{2\alpha K(1-\gamma^m)} + l^2 + \frac{\alpha}{2(1-\gamma^m)} \left( \frac{1}{1-\gamma} + 2B \right)^2 \\ &\quad + \frac{2\gamma^m}{1-\gamma^m} \left( B + \frac{1}{1-\gamma} \right)^2 \sqrt{\|d_m - d\|_{\text{TV}}}. \end{aligned} \quad (130)$$

The bound obtained through this Lyapounov drift summation can be used to further develop equation (85), using the subadditivity of the square root,

$$\sqrt{\mathbb{E}[\|\mathcal{Q} - \bar{\mathcal{Q}}\|_d^2]} \leq \sqrt{\frac{1}{K} \sum_{k=0}^{K-1} (\Delta_k - l)^2 + l} \quad (131)$$

$$\begin{aligned} &\leq \frac{\|\beta_0 - \beta_*\|_2}{\sqrt{2\alpha K(1-\gamma^m)}} + 2l + \sqrt{\frac{\alpha}{2(1-\gamma^m)}} \left( \frac{1}{1-\gamma} + 2B \right) \\ &\quad + \left( B + \frac{1}{1-\gamma} \right) \sqrt{\frac{2\gamma^m}{1-\gamma^m}} \sqrt{\|d_m - d\|_{\text{TV}}} \end{aligned} \quad (132)$$

$$\begin{aligned} &= \frac{\|\beta_0 - \beta_*\|_2}{\sqrt{2\alpha K(1-\gamma^m)}} + \frac{1+\gamma^m}{1-\gamma^m} \min_{f \in \mathcal{F}_\phi^B} \|f - \mathcal{Q}\|_d + \sqrt{\frac{\alpha}{2(1-\gamma^m)}} \left( \frac{1}{1-\gamma} + 2B \right) \\ &\quad + \left( B + \frac{1}{1-\gamma} \right) \sqrt{\frac{2\gamma^m}{1-\gamma^m}} \sqrt{\|d_m - d\|_{\text{TV}}}. \end{aligned} \quad (133)$$

By setting  $\alpha = \frac{1}{\sqrt{K}}$  and upper bounding  $\|\beta_0 - \beta_*\|$  by  $2B$ , we get,

$$\begin{aligned} \sqrt{\mathbb{E}[\|\mathcal{Q} - \bar{\mathcal{Q}}\|_d^2]} &\leq \frac{2B}{\sqrt{2\sqrt{K}(1-\gamma^m)}} + \frac{1+\gamma^m}{1-\gamma^m} \min_{f \in \mathcal{F}_\phi^B} \|f - \mathcal{Q}\|_d + \frac{1}{\sqrt{2\sqrt{K}(1-\gamma^m)}} \left( \frac{1}{1-\gamma} + 2B \right) \\ &\quad + \left( B + \frac{1}{1-\gamma} \right) \sqrt{\frac{2\gamma^m}{1-\gamma^m}} \sqrt{\|d_m - d\|_{\text{TV}}} \end{aligned} \quad (134)$$

$$\begin{aligned} &= \sqrt{\frac{4B^2 + \left( \frac{1}{1-\gamma} + 2B \right)^2}{2\sqrt{K}(1-\gamma^m)}} + \frac{1+\gamma^m}{1-\gamma^m} \min_{f \in \mathcal{F}_\phi^B} \|f - \mathcal{Q}\|_d \\ &\quad + \left( B + \frac{1}{1-\gamma} \right) \sqrt{\frac{2\gamma^m}{1-\gamma^m}} \sqrt{\|d_m - d\|_{\text{TV}}}. \end{aligned} \quad (135)$$

This concludes the proof.  $\square$

## D. Proof of the Finite-Time Bound for the Symmetric Critic

Let us first find an upper bound on the distance  $\|Q^\pi - \tilde{Q}^\pi\|_d^2$  between the Q-function  $Q^\pi$  and the fixed point  $\tilde{Q}^\pi$ .

**Lemma D.1** (Upper bound on the aliasing bias (Cayci et al., 2024)). For any agent-state policy  $\pi \in \Pi_{\mathcal{M}}$ , and any  $m \in \mathbb{N}$ , we have,

$$\|Q^\pi - \tilde{Q}^\pi\|_d \leq \frac{1-\gamma^m}{1-\gamma} \left\| \mathbb{E}^\pi \left[ \sum_{k=0}^{\infty} \gamma^{km} \left\| \hat{b}_{km} - b_{km} \right\|_{\text{TV}} \middle| Z_0 = \cdot \right] \right\|_d. \quad (136)$$

*Proof.* The proof is similar to the one of Cayci et al. (2024). Let us first define the expected  $m$ -step return,

$$\bar{r}_m(s, z, a) = \mathbb{E}^\pi \left[ \sum_{k=0}^{m-1} \gamma^k R_k \middle| S_0 = s, Z_0 = s, A_0 = a \right]. \quad (137)$$

Using the expected  $m$ -step return and the definition of the belief  $b$  in equation (33) and approximate belief  $\hat{b}$  in equation (34), it can be noted that,

$$Q^\pi(z, a) = \mathbb{E}^\pi \left[ \sum_{k=0}^{\infty} \gamma^{km} \sum_{s \in \mathcal{S}} b_{km}(s | H_{km}) \bar{r}_m(s, Z_{km}, A_{km}) \middle| Z_0 = z, A_0 = a \right] \quad (138)$$

$$\tilde{Q}^\pi(z, a) = \mathbb{E}^\pi \left[ \sum_{k=0}^{\infty} \gamma^{km} \sum_{s \in \mathcal{S}} \hat{b}_{km}(s|Z_{km}) \bar{r}_m(s, Z_{km}, A_{km}) \middle| Z_0 = z, A_0 = a \right]. \quad (139)$$

Indeed, bootstrapping at timestep  $m$  based on the agent state only is equivalent to considering the distribution of future states to be  $\hat{b}_m(\cdot|Z_m)$  instead of  $b_m(\cdot|H_m)$ . As a consequence, we have,

$$\left| Q^\pi(z, a) - \tilde{Q}^\pi(z, a) \right| = \mathbb{E}^\pi \left[ \sum_{k=0}^{\infty} \gamma^{km} \sum_{s \in \mathcal{S}} \left( b_{km}(s|H_{km}) - \hat{b}_{km}(s|Z_{km}) \right) \bar{r}_m(s, Z_{km}, A_{km}) \middle| Z_0 = z, A_0 = a \right] \quad (140)$$

$$\leq \mathbb{E}^\pi \left[ \sum_{k=0}^{\infty} \gamma^{km} \sup_{s \in \mathcal{S}} \left| b_{km}(s|H_{km}) - \hat{b}_{km}(s|Z_{km}) \right| \sup_{s \in \mathcal{S}} |\bar{r}_m(s, Z_{km}, A_{km})| \middle| Z_0 = z, A_0 = a \right] \quad (141)$$

$$\leq \mathbb{E}^\pi \left[ \sum_{k=0}^{\infty} \gamma^{km} \sup_{s \in \mathcal{S}} \left| b_{km}(s|H_{km}) - \hat{b}_{km}(s|Z_{km}) \right| \frac{1 - \gamma^m}{1 - \gamma} \middle| Z_0 = z, A_0 = a \right] \quad (142)$$

$$= \frac{1 - \gamma^m}{1 - \gamma} \mathbb{E}^\pi \left[ \sum_{k=0}^{\infty} \gamma^{km} \sup_{s \in \mathcal{S}} \left| b_{km}(s|H_{km}) - \hat{b}_{km}(s|Z_{km}) \right| \middle| Z_0 = z, A_0 = a \right] \quad (143)$$

$$\leq \frac{1 - \gamma^m}{1 - \gamma} \mathbb{E}^\pi \left[ \sum_{k=0}^{\infty} \gamma^{km} \left\| b_{km}(\cdot|H_{km}) - \hat{b}_{km}(\cdot|Z_{km}) \right\|_{\text{TV}} \middle| Z_0 = z, A_0 = a \right] \quad (144)$$

$$\leq \frac{1 - \gamma^m}{1 - \gamma} \mathbb{E}^\pi \left[ \sum_{k=0}^{\infty} \gamma^{km} \left\| b_{km} - \hat{b}_{km} \right\|_{\text{TV}} \middle| Z_0 = z, A_0 = a \right], \quad (145)$$

where we use  $b_{km}$  and  $\hat{b}_{km}$  to denote the random variables  $b_{km}(\cdot|H_{km})$  and  $\hat{b}_{km}(\cdot|Z_{km})$ , respectively. It illustrates that the aliasing bias can be bounded proportionally to the distance between the true belief and the approximate belief at the bootstrapping timesteps. Then, we obtain,

$$\left\| Q^\pi - \tilde{Q}^\pi \right\|_d \leq \frac{1 - \gamma^m}{1 - \gamma} \left\| \mathbb{E}^\pi \left[ \sum_{k=0}^{\infty} \gamma^{km} \left\| \hat{b}_{km} - b_{km} \right\|_{\text{TV}} \middle| Z_0 = \cdot \right] \right\|_d. \quad (146)$$

This concludes the proof.  $\square$

Using [Lemma D.1](#), we can prove [Theorem 4](#), that is recalled below. Note that some notations used in [Appendix C](#) will be reused with another meaning.

**Theorem 4** (Finite-time bound for symmetric  $m$ -step temporal difference learning ([Cayci et al., 2024](#))). For any agent-state policy  $\pi \in \Pi_{\mathcal{M}}$ , and any  $m \in \mathbb{N}$ , we have for [Algorithm 1](#) with  $\alpha = \frac{1}{\sqrt{K}}$ , and arbitrary  $B > 0$ ,

$$\sqrt{\mathbb{E} \left[ \left\| Q^\pi - \bar{Q}^\pi \right\|_d^2 \right]} \leq \varepsilon_{\text{td}} + \varepsilon_{\text{app}} + \varepsilon_{\text{shift}} + \varepsilon_{\text{alias}}, \quad (35)$$

where the temporal difference learning, function approximation, distribution shift, and aliasing terms are given by,

$$\varepsilon_{\text{td}} = \sqrt{\frac{4B^2 + \left( \frac{1}{1-\gamma} + 2B \right)^2}{2\sqrt{K}(1-\gamma^m)}} \quad (36)$$

$$\varepsilon_{\text{app}} = \frac{1 + \gamma^m}{1 - \gamma^m} \min_{f \in \mathcal{F}_x^B} \|f - Q^\pi\|_d \quad (37)$$

$$\varepsilon_{\text{shift}} = \left( B + \frac{1}{1-\gamma} \right) \sqrt{\frac{2\gamma^m}{1-\gamma^m}} \sqrt{\|d_m - d\|_{\text{TV}}} \quad (38)$$

$$\varepsilon_{\text{alias}} = \frac{2}{1-\gamma} \left\| \mathbb{E}^\pi \left[ \sum_{k=0}^{\infty} \gamma^{km} \left\| \hat{b}_{km} - b_{km} \right\|_{\text{TV}} \middle| Z_0 = \cdot \right] \right\|_d, \quad (39)$$

with  $d(z, a) = \sum_{s \in \mathcal{S}} d^\pi(s, z) \pi(a|z)$  the sampling distribution, and  $d_m(z, a) = \sum_{s \in \mathcal{S}} d_m^\pi(s, z) \pi(a|z)$  the bootstrapping distribution.

*Proof.* To ease notation as for the proof of [Theorem 3](#) in [Appendix C](#), we use  $Q$  as a shorthand for  $Q^\pi$ ,  $\hat{Q}^*$  as a shorthand for  $\hat{Q}_*^\pi$ ,  $\tilde{Q}$  as a shorthand for  $\tilde{Q}^\pi$ ,  $\bar{Q}$  as a shorthand for  $\bar{Q}^\pi$  and  $\hat{Q}_k$  as a shorthand for  $\hat{Q}_{\beta_k}^\pi$ , where the subscripts and superscripts remain implicit but are assumed clear from context. When evaluating the Q-functions, we go one step further by using  $Q_{k,i}$  to denote  $Q(Z_{k,i}, A_{k,i})$ ,  $\hat{Q}_{k,i}^*$  to denote  $\hat{Q}_*^\pi(Z_{k,i}, A_{k,i})$ ,  $\tilde{Q}_{k,i}$  to denote  $\tilde{Q}(Z_{k,i}, A_{k,i})$  and  $\hat{Q}_{k,i}$  to denote  $\hat{Q}_k(Z_{k,i}, A_{k,i})$ , and  $\chi_{k,i}$  to denote  $\chi(Z_{k,i}, A_{k,i})$ . In addition, we define  $d$  as a shorthand for  $d^\pi \otimes \pi$ , such that  $d(z, a) = d^\pi(z)\pi(a|z)$ , and  $d_m$  as a shorthand for  $d_m^\pi \otimes \pi$ , such that  $d_m(z, a) = d_m^\pi(z)\pi(a|z)$ . Using the triangle inequality and the subadditivity of the square root, we have,

$$\sqrt{\mathbb{E}[\|Q - \bar{Q}\|_d^2]} \leq \sqrt{\mathbb{E}[\|Q - \tilde{Q}\|_d^2]} + \sqrt{\mathbb{E}[\|\tilde{Q} - \bar{Q}\|_d^2]} \quad (147)$$

$$\leq \sqrt{\mathbb{E}[\|Q - \tilde{Q}\|_d^2]} + \sqrt{\mathbb{E}[\|\tilde{Q} - \bar{Q}\|_d^2]} \quad (148)$$

$$\leq \|Q - \tilde{Q}\|_d + \sqrt{\mathbb{E}[\|\tilde{Q} - \bar{Q}\|_d^2]}. \quad (149)$$

We can bound the second term in equation (149) using similar steps as in the proof for the asymmetric finite-time bound (see [Appendix C](#)). We obtain,

$$\sqrt{\mathbb{E}[\|\tilde{Q} - \bar{Q}\|_d^2]} \leq \sqrt{\frac{1}{K} \sum_{k=0}^{K-1} (\Delta_k - l)^2} + l, \quad (150)$$

where  $l$  is arbitrary, and  $\Delta_k$  is defined as,

$$\Delta_k = \sqrt{\mathbb{E}[\|\tilde{Q} - \hat{Q}_k\|_d^2]} = \sqrt{\mathbb{E}[\|\tilde{Q}(\cdot) - \langle \beta_k, \chi(\cdot) \rangle\|_d^2]}. \quad (151)$$

Similarly to the asymmetric case (see [Appendix C](#)), we consider the Lyapounov function  $\mathcal{L}(\beta) = \|\beta_* - \beta\|_2^2$  in order to find a bound on  $\frac{1}{K} \sum_{k=0}^{K-1} (\Delta_k - l)^2$ . We define  $\mathfrak{G}_k = \sigma(Z_{i,j}, A_{i,j}, i \leq k, j \leq m)$  and  $\mathfrak{F}_k = \sigma(Z_{k,0}, A_{k,0})$ . As in the asymmetric case (see [Appendix C](#)), we obtain, using to the law of total expectation,

$$\mathbb{E}[\mathcal{L}(\beta_{k+1}) - \mathcal{L}(\beta_k)] \leq 2\alpha \mathbb{E}[\mathbb{E}[\langle \beta_k - \beta_*, g_k \rangle | \mathfrak{G}_{k-1}]] + \alpha^2 \mathbb{E}[\mathbb{E}[\|g_k\|_2^2 | \mathfrak{G}_{k-1}]]]. \quad (152)$$

Let us focus on the first term of equation (152) with  $\mathbb{E}[\langle \beta_k - \beta_*, g_k \rangle | \mathfrak{G}_{k-1}]$ . By conditioning on the sigma-fields  $\mathfrak{G}_{k-1}$  and  $\mathfrak{F}_k$ , we have,

$$\mathbb{E}[\langle \beta_k - \beta_*, g_k \rangle | \mathfrak{F}_k, \mathfrak{G}_{k-1}] = \left( \mathbb{E} \left[ \sum_{t=0}^{m-1} \gamma^t R_{k,t} + \gamma^m \hat{Q}_{k,m} \middle| \mathfrak{F}_k, \mathfrak{G}_{k-1} \right] - \hat{Q}_{k,0} \right) (\hat{Q}_{k,0} - \hat{Q}_{k,0}^*). \quad (153)$$

Note that, according to the Bellman operator (7), we have,

$$\mathbb{E} \left[ \sum_{t=0}^{m-1} \gamma^t R_{k,t} \middle| \mathfrak{F}_k, \mathfrak{G}_{k-1} \right] = \tilde{Q}_{k,0} - \gamma^m \mathbb{E}[\hat{Q}_{k,m} | \mathfrak{F}_k, \mathfrak{G}_{k-1}]. \quad (154)$$

It differs from the asymmetric case (see [Appendix C](#)) in that we do not necessarily have  $Q = \tilde{Q}$  here. By substituting equation (154) in equation (153), we obtain,

$$\begin{aligned} & \mathbb{E}[\langle \beta_k - \beta_*, g_k \rangle | \mathfrak{F}_k, \mathfrak{G}_{k-1}] \\ &= \left( \mathbb{E} \left[ \sum_{t=0}^{m-1} \gamma^t R_{k,t} \middle| \mathfrak{F}_k, \mathfrak{G}_{k-1} \right] + \gamma^m \mathbb{E}[\hat{Q}_{k,m} | \mathfrak{F}_k, \mathfrak{G}_{k-1}] - \hat{Q}_{k,0} \right) (\hat{Q}_{k,0} - \hat{Q}_{k,0}^*) \end{aligned} \quad (155)$$



$$= (\tilde{Q}_{k,0} - \gamma^m \mathbb{E} [\tilde{Q}_{k,m} | \mathfrak{F}_k, \mathfrak{G}_{k-1}] + \gamma^m \mathbb{E} [\hat{Q}_{k,m} | \mathfrak{F}_k, \mathfrak{G}_{k-1}] - \hat{Q}_{k,0}) (\hat{Q}_{k,0} - \hat{Q}_{k,0}^*) \quad (156)$$

$$= (\tilde{Q}_{k,0} - \gamma^m \mathbb{E} [\tilde{Q}_{k,m} - \hat{Q}_{k,m} | \mathfrak{F}_k, \mathfrak{G}_{k-1}] - \hat{Q}_{k,0}) (\hat{Q}_{k,0} - \tilde{Q}_{k,0} + \tilde{Q}_{k,0} - \hat{Q}_{k,0}^*) \quad (157)$$

$$= ((\tilde{Q}_{k,0} - \hat{Q}_{k,0}) - \gamma^m \mathbb{E} [\tilde{Q}_{k,m} - \hat{Q}_{k,m} | \mathfrak{F}_k, \mathfrak{G}_{k-1}]) ((\hat{Q}_{k,0} - \tilde{Q}_{k,0}) + (\tilde{Q}_{k,0} - \hat{Q}_{k,0}^*)) \quad (158)$$

$$= -(\tilde{Q}_{k,0} - \hat{Q}_{k,0})^2 + (\tilde{Q}_{k,0} - \hat{Q}_{k,0})(\tilde{Q}_{k,0} - \hat{Q}_{k,0}^*) + \gamma^m \mathbb{E} [\hat{Q}_{k,m} - \tilde{Q}_{k,m} | \mathfrak{F}_k, \mathfrak{G}_{k-1}] (\tilde{Q}_{k,0} - \hat{Q}_{k,0}^*). \quad (159)$$

We now follow the same technique as in the asymmetric case (see [Appendix C](#)) for each of the four terms. By taking the expectation over  $\mathfrak{F}_k$ , we get the following.

- For the first term, we have,

$$\mathbb{E} [-(\tilde{Q}_{k,0} - \hat{Q}_{k,0})^2 | \mathfrak{G}_{k-1}] = -\|\tilde{Q} - \hat{Q}\|_d^2. \quad (160)$$

- For the second term, we have,

$$\mathbb{E} [(\tilde{Q}_{k,0} - \hat{Q}_{k,0})(\tilde{Q}_{k,0} - \hat{Q}_{k,0}^*) | \mathfrak{G}_{k-1}] \leq \|\tilde{Q} - \hat{Q}\|_d \|\tilde{Q} - \hat{Q}^*\|_d. \quad (161)$$

- For the third term, we have,

$$\mathbb{E} [(\hat{Q}_{k,m} - \tilde{Q}_{k,m})(\hat{Q}_{k,0} - \tilde{Q}_{k,0}) | \mathfrak{G}_{k-1}] \leq \|\hat{Q}_k - \tilde{Q}\|_d^2 + \left(B + \frac{1}{1-\gamma}\right)^2 \sqrt{\|d_m - d\|_{\text{TV}}}. \quad (162)$$

- For the fourth term, we have,

$$\mathbb{E} [(\hat{Q}_{k,m} - \tilde{Q}_{k,m})(\tilde{Q}_{k,0} - \hat{Q}_{k,0}^*) | \mathfrak{G}_{k-1}] \leq \|\hat{Q}_k - \tilde{Q}\|_d \|\tilde{Q} - \hat{Q}^*\|_d + \left(B + \frac{1}{1-\gamma}\right)^2 \sqrt{\|d_m - d\|_{\text{TV}}}. \quad (163)$$

By taking expectation over  $\mathfrak{G}_{k-1}$  of the four terms and using the previous upper bounds, we obtain,

$$\mathbb{E} [\langle \beta_k - \beta_*, g_k \rangle] \leq -(1 - \gamma^m) \Delta_k^2 + (1 + \gamma^m) \Delta_k \|\hat{Q}^* - \tilde{Q}\|_d + 2\gamma^m \left(B + \frac{1}{1-\gamma}\right)^2 \sqrt{\|d_m - d\|_{\text{TV}}}. \quad (164)$$

The second term in equation (152) is treated similarly to the asymmetric case (see [Appendix C](#)), which yields,

$$\mathbb{E} [\|g_k\|_2^2] \leq \left(\frac{1}{1-\gamma} + 2B\right)^2. \quad (165)$$

By substituting equations (164) and (165) into the Lyapounov drift of equation (152), we obtain,

$$\begin{aligned} \mathbb{E} [\mathcal{L}(\beta_{k+1}) - \mathcal{L}(\beta_k)] &\leq -2\alpha(1 - \gamma^m) \Delta_k^2 + 2\alpha(1 + \gamma^m) \Delta_k \|\hat{Q}^* - \tilde{Q}\|_d + \alpha^2 \left(\frac{1}{1-\gamma} + 2B\right)^2 \\ &\quad + 4\alpha\gamma^m \left(B + \frac{1}{1-\gamma}\right)^2 \sqrt{\|d_m - d\|_{\text{TV}}}. \end{aligned} \quad (166)$$

We can upper bound  $\|\hat{Q}^* - \tilde{Q}\|_d$  as follows,

$$\|\hat{Q}^* - \tilde{Q}\|_d \leq \|\hat{Q}^* - Q\|_d + \|Q - \tilde{Q}\|_d. \quad (167)$$

By setting  $l = \frac{1+\gamma^m}{2(1-\gamma^m)} \left( \|\hat{Q}^* - Q\|_d + \|Q - \tilde{Q}\|_d \right)$ , we can write, following a similar strategy as in the asymmetric case (see [Appendix C](#)),

$$\begin{aligned} \mathbb{E}[\mathcal{L}(\beta_{k+1}) - \mathcal{L}(\beta_k)] &\leq -2\alpha(1-\gamma^m)(\Delta_k - l)^2 + 2\alpha(1-\gamma^m)l^2 + \alpha^2 \left( \frac{1}{1-\gamma} + 2B \right)^2 \\ &\quad + 4\alpha\gamma^m \left( B + \frac{1}{1-\gamma} \right)^2 \sqrt{\|d_m - d\|_{\text{TV}}}. \end{aligned} \quad (168)$$

By summing all drifts, rearranging, and dividing by  $2\alpha K(1-\gamma^m)$ , we obtain after neglecting  $\mathcal{L}(\beta_K) > 0$ ,

$$\begin{aligned} \frac{1}{K} \sum_{k=0}^{K-1} (\Delta_k - l)^2 &\leq \frac{\|\beta_0 - \beta_*\|_2^2}{2\alpha K(1-\gamma^m)} + l^2 + \frac{\alpha}{2(1-\gamma^m)} \left( \frac{1}{1-\gamma} + 2B \right)^2 \\ &\quad + \frac{2\gamma^m}{1-\gamma^m} \left( B + \frac{1}{1-\gamma} \right)^2 \sqrt{\|d_m - d\|_{\text{TV}}}. \end{aligned} \quad (169)$$

The bound obtained through this Lyapounov drift summation can be used to further develop equation (150), using the subadditivity of the square root,

$$\sqrt{\mathbb{E}[\|\tilde{Q} - \bar{Q}\|_d^2]} \leq \sqrt{\frac{1}{K} \sum_{k=0}^{K-1} (\Delta_k - l)^2 + l^2} \quad (170)$$

$$\begin{aligned} &\leq \frac{\|\beta_0 - \beta_*\|_2}{\sqrt{2\alpha K(1-\gamma^m)}} + 2l + \sqrt{\frac{\alpha}{2(1-\gamma^m)}} \left( \frac{1}{1-\gamma} + 2B \right) \\ &\quad + \left( B + \frac{1}{1-\gamma} \right) \sqrt{\frac{2\gamma^m}{1-\gamma^m}} \sqrt{\|d_m - d\|_{\text{TV}}}. \end{aligned} \quad (171)$$

$$\begin{aligned} &= \frac{\|\beta_0 - \beta_*\|_2}{\sqrt{2\alpha K(1-\gamma^m)}} + \frac{1+\gamma^m}{1-\gamma^m} \left( \frac{1}{1-\gamma} + B \right) + \sqrt{\frac{\alpha}{2(1-\gamma^m)}} \left( \frac{1}{1-\gamma} + 2B \right) \\ &\quad + \left( B + \frac{1}{1-\gamma} \right) \sqrt{\frac{2\gamma^m}{1-\gamma^m}} \sqrt{\|d_m - d\|_{\text{TV}}}. \end{aligned} \quad (172)$$

Plugging equation (172) into equation (149), and substituting back  $l$ , we finally have,

$$\begin{aligned} \sqrt{\mathbb{E}[\|Q - \bar{Q}\|_d^2]} &\leq \frac{\|\beta_0 - \beta_*\|_2}{\sqrt{2\alpha K(1-\gamma^m)}} + \frac{1+\gamma^m}{1-\gamma^m} \left( \|\hat{Q}^* - Q\|_d + \|Q - \tilde{Q}\|_d \right) + \sqrt{\frac{\alpha}{2(1-\gamma^m)}} \left( \frac{1}{1-\gamma} + 2B \right) \\ &\quad + \left( B + \frac{1}{1-\gamma} \right) \sqrt{\frac{2\gamma^m}{1-\gamma^m}} \sqrt{\|d_m - d\|_{\text{TV}}} + \|Q - \tilde{Q}\|_d \end{aligned} \quad (173)$$

$$\begin{aligned} &\leq \frac{\|\beta_0 - \beta_*\|_2}{\sqrt{2\alpha K(1-\gamma^m)}} + \frac{1+\gamma^m}{1-\gamma^m} \|\hat{Q}^* - Q\|_d + \sqrt{\frac{\alpha}{2(1-\gamma^m)}} \left( \frac{1}{1-\gamma} + 2B \right) \\ &\quad + \left( B + \frac{1}{1-\gamma} \right) \sqrt{\frac{2\gamma^m}{1-\gamma^m}} \sqrt{\|d_m - d\|_{\text{TV}}} + \frac{2}{1-\gamma^m} \|Q - \tilde{Q}\|_d \end{aligned} \quad (174)$$

Using [Lemma D.1](#), we finally obtain,

$$\begin{aligned} \sqrt{\mathbb{E}[\|Q - \bar{Q}\|_d^2]} &\leq \frac{\|\beta_0 - \beta_*\|_2}{\sqrt{2\alpha K(1-\gamma^m)}} + \frac{1+\gamma^m}{1-\gamma^m} \|\hat{Q}^* - Q\|_d + \sqrt{\frac{\alpha}{2(1-\gamma^m)}} \left( \frac{1}{1-\gamma} + 2B \right) \\ &\quad + \left( B + \frac{1}{1-\gamma} \right) \sqrt{\frac{2\gamma^m}{1-\gamma^m}} \sqrt{\|d_m - d\|_{\text{TV}}} \\ &\quad + \left( \frac{2}{1-\gamma^m} \right) \frac{1-\gamma^m}{1-\gamma} \left\| \mathbb{E} \left[ \sum_{k=0}^{\infty} \gamma^{km} \|\hat{b}_{km} - b_{km}\|_{\text{TV}} \middle| Z_0 = \cdot \right] \right\|_d \end{aligned} \quad (175)$$

$$\begin{aligned}
 &\leq \frac{\|\beta_0 - \beta_*\|_2}{\sqrt{2\alpha K(1-\gamma^m)}} + \frac{1+\gamma^m}{1-\gamma^m} \min_{f \in \mathcal{F}_\phi^B} \|f - \mathcal{Q}\|_d + \sqrt{\frac{\alpha}{2(1-\gamma^m)}} \left( \frac{1}{1-\gamma} + 2B \right) \\
 &\quad + \left( B + \frac{1}{1-\gamma} \right) \sqrt{\frac{2\gamma^m}{1-\gamma^m}} \sqrt{\|d_m - d\|_{\text{TV}}} \\
 &\quad + \frac{2}{1-\gamma} \left\| \mathbb{E} \left[ \sum_{k=0}^{\infty} \gamma^{km} \|\hat{b}_{km} - b_{km}\|_{\text{TV}} \middle| Z_0 = \cdot \right] \right\|_d.
 \end{aligned} \tag{176}$$

By setting  $\alpha = \frac{1}{\sqrt{K}}$  and upper bounding  $\|\beta_0 - \beta_*\|$  by  $2B$ , we get,

$$\begin{aligned}
 \sqrt{\mathbb{E} [\|Q - \bar{Q}\|_d^2]} &\leq \sqrt{\frac{4B^2 + \left(\frac{1}{1-\gamma} + 2B\right)^2}{2\sqrt{K}(1-\gamma^m)}} + \frac{1+\gamma^m}{1-\gamma^m} \min_{f \in \mathcal{F}_\phi^B} \|f - \mathcal{Q}\|_d \\
 &\quad + \left( B + \frac{1}{1-\gamma} \right) \sqrt{\frac{2\gamma^m}{1-\gamma^m}} \sqrt{\|d_m - d\|_{\text{TV}}} \\
 &\quad + \frac{2}{1-\gamma} \left\| \mathbb{E} \left[ \sum_{k=0}^{\infty} \gamma^{km} \|\hat{b}_{km} - b_{km}\|_{\text{TV}} \middle| Z_0 = \cdot \right] \right\|_d.
 \end{aligned} \tag{177}$$

This concludes the proof.  $\square$

## E. Proof of the Finite-Time Bound for the Natural Actor-Critic

Let us first give the performance difference lemma for POMDP proved by Cayci et al. (2024). Note that this proof is completely agnostic about the critic used to compute  $\pi_1, \pi_2 \in \Pi_{\mathcal{M}}$  and is thus applicable both to the asymmetric setting and the symmetric setting.

**Lemma E.1** (Performance difference (Cayci et al., 2024)). For any two agent-state policies  $\pi_1, \pi_2 \in \Pi_{\mathcal{M}}$ ,

$$V^{\pi_2}(z_0) - V^{\pi_1}(z_0) \leq \frac{1}{1-\gamma} \mathbb{E}^{d^{\pi_2}} [A^{\pi_1}(Z, A) | Z_0 = z_0] + \frac{2}{1-\gamma} \varepsilon_{\text{inf}}^{\pi_2}(z_0), \tag{178}$$

where,

$$\varepsilon_{\text{inf}}^{\pi_2}(z_0) = \mathbb{E}^{\pi_2} \left[ \sum_{k=0}^{\infty} \gamma^k \|\hat{b}_k - b_k\|_{\text{TV}} \middle| Z_0 = z_0 \right]. \tag{179}$$

*Proof.* The proof is similar to the one of Cayci et al. (2024). First, let us decompose the performance difference in the following terms,

$$V^{\pi_2}(z_0) - V^{\pi_1}(z_0) = \mathbb{E}^{\pi_2} \left[ \sum_{t=0}^{\infty} \gamma^t R_t \middle| Z_0 = z_0 \right] - V^{\pi_1}(z_0) \tag{180}$$

$$= \mathbb{E}^{\pi_2} \left[ \sum_{t=0}^{\infty} \gamma^t (R_t - V^{\pi_1}(Z_t) + V^{\pi_1}(Z_t)) \middle| Z_0 = z_0 \right] - V^{\pi_1}(z_0) \tag{181}$$

$$= \mathbb{E}^{\pi_2} \left[ \sum_{t=0}^{\infty} \gamma^t (R_t - V^{\pi_1}(Z_t) + \gamma V^{\pi_1}(Z_{t+1})) \middle| Z_0 = z_0 \right] \tag{182}$$

$$\begin{aligned}
 &= \mathbb{E}^{\pi_2} \left[ \sum_{t=0}^{\infty} \gamma^t (R_t + \gamma \mathcal{V}^{\pi_1}(S_{t+1}, Z_{t+1}) - V^{\pi_1}(Z_t)) \middle| Z_0 = z_0 \right] \\
 &\quad + \mathbb{E}^{\pi_2} \left[ \sum_{t=0}^{\infty} \gamma^t (\gamma V^{\pi_1}(Z_{t+1}) - \gamma \mathcal{V}^{\pi_1}(S_{t+1}, Z_{t+1})) \middle| Z_0 = z_0 \right]
 \end{aligned} \tag{183}$$

$$\begin{aligned}
 &= \mathbb{E}^{\pi_2} \left[ \sum_{t=0}^{\infty} \gamma^t (R_t + \gamma \mathcal{V}^{\pi_1}(S_{t+1}, Z_{t+1}) - V^{\pi_1}(Z_t)) \middle| Z_0 = z_0 \right] \\
 &\quad + \mathbb{E}^{\pi_2} \left[ \sum_{t=0}^{\infty} \gamma^{t+1} (V^{\pi_1}(Z_{t+1}) - \mathcal{V}^{\pi_1}(S_{t+1}, Z_{t+1})) \middle| Z_0 = z_0 \right].
 \end{aligned} \tag{184}$$

Let us focus on bounding the first term in equation (184). We have, for any  $T > 0$ ,

$$\left| \sum_{t=0}^T \gamma^t (R_t + \gamma \mathcal{V}^{\pi_1}(S_{t+1}, Z_{t+1}) - V^{\pi_1}(Z_t)) \right| \leq \frac{2}{(1-\gamma)^2} < \infty. \tag{185}$$

By Lebesgue's dominated convergence, we have,

$$\begin{aligned}
 &\mathbb{E}^{\pi_2} \left[ \sum_{t=0}^{\infty} \gamma^t (R_t + \gamma \mathcal{V}^{\pi_1}(S_{t+1}, Z_{t+1}) - V^{\pi_1}(Z_t)) \middle| Z_0 = z_0 \right] \\
 &= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}^{\pi_2} [R_t + \gamma \mathcal{V}^{\pi_1}(S_{t+1}, Z_{t+1}) - V^{\pi_1}(Z_t) | Z_0 = z_0].
 \end{aligned} \tag{186}$$

Then, by the law of total expectation, we have at any timestep  $t \geq 0$ ,

$$\begin{aligned}
 &\mathbb{E}^{\pi_2} [R_t + \gamma \mathcal{V}^{\pi_1}(S_{t+1}, Z_{t+1}) - V^{\pi_1}(Z_t) | Z_0 = z_0] \\
 &= \mathbb{E} [\mathbb{E}^{\pi_2} [R_t + \gamma \mathcal{V}^{\pi_1}(S_{t+1}, Z_{t+1}) | H_t, Z_t] - V^{\pi_1}(Z_t) | Z_0 = z_0].
 \end{aligned} \tag{187}$$

And, we have,

$$\begin{aligned}
 &\mathbb{E}^{\pi_2} [R_t + \gamma \mathcal{V}^{\pi_1}(S_{t+1}, Z_{t+1}) | H_t = h_t, Z_t = z_t] \\
 &= \sum_{s_t, a_t} b_t(s_t | h_t) \pi_2(a_t | z_t) \mathcal{Q}^{\pi_1}(s_t, z_t, a_t)
 \end{aligned} \tag{188}$$

$$= \sum_{a_t} \pi_2(a_t | z_t) Q^{\pi_1}(z_t, a_t) + \sum_{s_t, a_t} b_t(s_t | h_t) \pi_2(a_t | z_t) \mathcal{Q}^{\pi_1}(s_t, z_t, a_t) - \sum_{a_t} \pi_2(a_t | z_t) Q^{\pi_1}(z_t, a_t) \tag{189}$$

$$\begin{aligned}
 &= \sum_{a_t} \pi_2(a_t | z_t) Q^{\pi_1}(z_t, a_t) + \sum_{s_t, a_t} b_t(s_t | h_t) \pi_2(a_t | z_t) \mathcal{Q}^{\pi_1}(s_t, z_t, a_t) \\
 &\quad - \sum_{s_t, a_t} \hat{b}_t(s_t | z_t) \pi_2(a_t | z_t) \mathcal{Q}^{\pi_1}(s_t, z_t, a_t)
 \end{aligned} \tag{190}$$

$$= \sum_{a_t} \pi_2(a_t | z_t) Q^{\pi_1}(z_t, a_t) + \sum_{s_t, a_t} (b_t(s_t | h_t) - \hat{b}_t(s_t | z_t)) \pi_2(a_t | z_t) \mathcal{Q}^{\pi_1}(s_t, z_t, a_t). \tag{191}$$

By noting that  $\sup_{s, z} |\sum_a \pi_2(a | z) \mathcal{Q}^{\pi_1}(s, z, a)| \leq \sup_{s, z, a} |\mathcal{Q}^{\pi_1}(s, z, a)| \leq \frac{1}{1-\gamma}$ , we obtain,

$$\begin{aligned}
 &\mathbb{E}^{\pi_2} [R_t + \gamma \mathcal{V}^{\pi_1}(S_{t+1}, Z_{t+1}) | H_t = h_t, Z_t = z_t] \\
 &\leq \sum_{a_t} \pi_2(a_t | z_t) Q^{\pi_1}(z_t, a_t) + \frac{1}{1-\gamma} \|b_t(\cdot | h_t) - \hat{b}_t(\cdot | z_t)\|_{\text{TV}}.
 \end{aligned} \tag{192}$$

Finally, the expectation at time  $t \geq 0$  can be written as,

$$\begin{aligned}
 &\mathbb{E}^{\pi_2} [R_t + \gamma \mathcal{V}^{\pi_1}(S_{t+1}, Z_{t+1}) - V^{\pi_1}(Z_t) | Z_0 = z_0] \\
 &= \mathbb{E} [\mathbb{E}^{\pi_2} [R_t + \gamma \mathcal{V}^{\pi_1}(S_{t+1}, Z_{t+1}) | H_t, Z_t] - V^{\pi_1}(Z_t) | Z_0 = z_0]
 \end{aligned} \tag{193}$$

$$\leq \mathbb{E}^{\pi_2} \left[ Q^{\pi_1}(Z_t, A_t) + \frac{1}{1-\gamma} \|b_t(\cdot | H_t) - \hat{b}_t(\cdot | Z_t)\|_{\text{TV}} - V^{\pi_1}(Z_t) \middle| Z_0 = z_0 \right] \tag{194}$$

$$= \mathbb{E}^{\pi_2} \left[ A^{\pi_1}(Z_t, A_t) - \frac{1}{1-\gamma} \|b_t(\cdot | H_t) - \hat{b}_t(\cdot | Z_t)\|_{\text{TV}} \middle| Z_0 = z_0 \right] \tag{195}$$



Now, by using Lebesgue's dominated theorem in the reverse direction, we have,

$$\begin{aligned} & \mathbb{E}^{\pi_2} \left[ \sum_{t=0}^{\infty} \gamma^t (R_t + \gamma \mathcal{V}^{\pi_1}(S_{t+1}, Z_{t+1}) - V^{\pi_1}(Z_t)) \middle| Z_0 = z_0 \right] \\ & \leq \mathbb{E}^{\pi_2} \left[ \sum_{t=0}^{\infty} \gamma^t A^{\pi_1}(Z_t, A_t) \middle| Z_0 = z_0 \right] + \frac{1}{1-\gamma} \mathbb{E}^{\pi_2} \left[ \sum_{t=0}^{\infty} \gamma^t \|\hat{b}_t - b_t\|_{\text{TV}} \middle| Z_0 = z_0 \right] \end{aligned} \quad (196)$$

$$= \mathbb{E}^{\pi_2} \left[ \sum_{t=0}^{\infty} \gamma^t A^{\pi_1}(Z_t, A_t) \middle| Z_0 = z_0 \right] + \frac{1}{1-\gamma} \varepsilon_{\inf}^{\pi_2}(z_0) \quad (197)$$

Now, let us focus on bounding the second term in equation (184). We have, for any  $T > 0$ ,

$$\left| \sum_{t=0}^T \gamma^{t+1} (V^{\pi_1}(Z_{t+1}) - \mathcal{V}^{\pi_1}(S_{t+1}, Z_{t+1})) \right| \leq \frac{2}{(1-\gamma)^2} < \infty. \quad (198)$$

Using Lebesgue dominated convergence theorem, we can write,

$$\begin{aligned} & \mathbb{E}^{\pi_2} \left[ \sum_{t=0}^{\infty} \gamma^{t+1} (V^{\pi_1}(Z_{t+1}) - \mathcal{V}^{\pi_1}(S_{t+1}, Z_{t+1})) \middle| Z_0 = z_0 \right] \\ & = \sum_{t=0}^{\infty} \gamma^{t+1} \mathbb{E}^{\pi_2} [V^{\pi_1}(Z_{t+1}) - \mathcal{V}^{\pi_1}(S_{t+1}, Z_{t+1}) | Z_0 = z_0]. \end{aligned} \quad (199)$$

By the law of total expectation, we have at any timestep  $t \geq 0$ ,

$$\begin{aligned} & \mathbb{E}^{\pi_2} [V^{\pi_1}(Z_{t+1}) - \mathcal{V}^{\pi_1}(S_{t+1}, Z_{t+1}) | Z_0 = z_0] \\ & = \mathbb{E} [V^{\pi_1}(Z_{t+1}) - \mathbb{E}^{\pi_2} [\mathcal{V}^{\pi_1}(S_{t+1}, Z_{t+1}) | H_{t+1}, Z_{t+1}] | Z_0 = z_0]. \end{aligned} \quad (200)$$

And, we have,

$$\begin{aligned} & \mathbb{E}^{\pi_2} [\mathcal{V}^{\pi_1}(S_{t+1}, z_{t+1}) | H_{t+1} = h_{t+1}, Z_{t+1} = z_{t+1},] \\ & = \sum_{s_{t+1}} b_{t+1}(s_{t+1} | h_{t+1}) \mathcal{V}^{\pi_1}(s_{t+1}, z_{t+1}) \end{aligned} \quad (201)$$

$$= V^{\pi_1}(z_{t+1}) + \sum_{s_{t+1}} b_{t+1}(s_{t+1} | h_{t+1}) \mathcal{V}^{\pi_1}(s_{t+1}, z_{t+1}) - V^{\pi_1}(z_{t+1}) \quad (202)$$

$$= V^{\pi_1}(z_{t+1}) + \sum_{s_{t+1}} b_{t+1}(s_{t+1} | h_{t+1}) \mathcal{V}^{\pi_1}(s_{t+1}, z_{t+1}) - \sum_{s_{t+1}} \hat{b}_{t+1}(s_{t+1} | z_{t+1}) \mathcal{V}^{\pi_1}(s_{t+1}, z_{t+1}) \quad (203)$$

$$= V^{\pi_1}(z_{t+1}) + \sum_{s_{t+1}} (b_{t+1}(s_{t+1} | h_{t+1}) - \hat{b}_{t+1}(s_{t+1} | z_{t+1})) \mathcal{V}^{\pi_1}(s_{t+1}, z_{t+1}). \quad (204)$$

From there, by noting that  $\sup_{s,z} |\mathcal{V}^{\pi_1}(s, z)| \leq \frac{1}{1-\gamma}$ , we obtain,

$$\begin{aligned} & \mathbb{E}^{\pi_2} [\mathcal{V}^{\pi_1}(S_{t+1}, z_{t+1}) | H_{t+1} = h_{t+1}, Z_{t+1} = z_{t+1},] \\ & \geq V^{\pi_1}(z_{t+1}) - \frac{1}{1-\gamma} \left\| b_{t+1}(\cdot | h_{t+1}) - \hat{b}_{t+1}(\cdot | z_{t+1}) \right\|_{\text{TV}}. \end{aligned} \quad (205)$$

Finally, the expectation at time  $t \geq 0$  can be written as,

$$\begin{aligned} & \mathbb{E}^{\pi_2} [V^{\pi_1}(Z_{t+1}) - \mathcal{V}^{\pi_1}(S_{t+1}, Z_{t+1}) | Z_0 = z_0] \\ & = \mathbb{E} [V^{\pi_1}(Z_{t+1}) - \mathbb{E}^{\pi_2} [\mathcal{V}^{\pi_1}(S_{t+1}, Z_{t+1}) | H_{t+1}, Z_{t+1}] | Z_0 = z_0] \end{aligned} \quad (206)$$

$$\leq \mathbb{E} \left[ V^{\pi_1}(Z_{t+1}) - V^{\pi_1}(Z_{t+1}) + \frac{1}{1-\gamma} \left\| b_{t+1}(\cdot | H_{t+1}) - \hat{b}_{t+1}(\cdot | Z_{t+1}) \right\|_{\text{TV}} \middle| Z_0 = z_0 \right] \quad (207)$$

$$\leq \mathbb{E} \left[ \frac{1}{1-\gamma} \left\| b_{t+1}(\cdot|H_{t+1}) - \hat{b}_{t+1}(\cdot|Z_{t+1}) \right\|_{\text{TV}} \middle| Z_0 = z_0 \right]. \quad (208)$$

Now, by using Lebesgue's dominated theorem in the reverse direction, we have,

$$\begin{aligned} & \mathbb{E}^{\pi_2} \left[ \sum_{t=0}^{\infty} \gamma^{t+1} (V^{\pi_1}(Z_{t+1}) - \mathcal{V}^{\pi_1}(S_{t+1}, Z_{t+1})) \middle| Z_0 = z_0 \right] \\ & \leq \frac{1}{1-\gamma} \mathbb{E}^{\pi_2} \left[ \sum_{t=0}^{\infty} \gamma^{t+1} \left\| b_{t+1}(\cdot|H_{t+1}) - \hat{b}_{t+1}(\cdot|Z_{t+1}) \right\|_{\text{TV}} \middle| Z_0 = z_0 \right] \end{aligned} \quad (209)$$

$$= \frac{1}{1-\gamma} \mathbb{E}^{\pi_2} \left[ \sum_{t=0}^{\infty} \gamma^t \left\| b_t(\cdot|H_t) - \hat{b}_t(\cdot|Z_t) \right\|_{\text{TV}} - \left\| b_0(\cdot|H_0) - \hat{b}_0(\cdot|Z_0) \right\|_{\text{TV}} \middle| Z_0 = z_0 \right] \quad (210)$$

$$\begin{aligned} & = \frac{1}{1-\gamma} \mathbb{E}^{\pi_2} \left[ \sum_{t=0}^{\infty} \gamma^t \left\| b_t(\cdot|H_t) - \hat{b}_t(\cdot|Z_t) \right\|_{\text{TV}} \middle| Z_0 = z_0 \right] \\ & \quad - \mathbb{E}^{\pi_2} \left[ \left\| b_0(\cdot|H_0) - \hat{b}_0(\cdot|Z_0) \right\|_{\text{TV}} \middle| Z_0 = z_0 \right] \end{aligned} \quad (211)$$

$$= \frac{1}{1-\gamma} \varepsilon_{\text{inf}}^{\pi_2}(z_0) - \mathbb{E}^{\pi_2} \left[ \left\| b_0(\cdot|H_0) - \hat{b}_0(\cdot|Z_0) \right\|_{\text{TV}} \middle| Z_0 = z_0 \right] \quad (212)$$

$$\leq \frac{1}{1-\gamma} \varepsilon_{\text{inf}}^{\pi_2}(z_0). \quad (213)$$

Finally, by substituting the upper bound (197) on the first term and the upper bound (213) on the second term into equation (184), we obtain,

$$V^{\pi_2}(z_0) - V^{\pi_1}(z_0) \leq \mathbb{E}^{\pi_2} \left[ \sum_{t=0}^{\infty} \gamma^t A^{\pi_1}(Z_t, A_t) \middle| Z_0 = z_0 \right] + \frac{2}{1-\gamma} \varepsilon_{\text{inf}}^{\pi_2}(z_0) \quad (214)$$

$$= \frac{1}{1-\gamma} \mathbb{E}^{d^{\pi_2}} [A^{\pi_1}(Z, A) | Z_0 = z_0] + \frac{2}{1-\gamma} \varepsilon_{\text{inf}}^{\pi_2}(z_0). \quad (215)$$

This concludes the proof.  $\square$

Using Lemma E.1, we can prove Theorem 5, that is recalled below. The proof from Cayci et al. (2024) is generalized to the asymmetric setting.

**Theorem 5** (Finite-time bound for asymmetric and symmetric natural actor-critic algorithm). For any agent-state process  $\mathcal{M} = (\mathcal{Z}, U)$ , we have for Algorithm 2 with  $\alpha = \frac{1}{\sqrt{K}}$ ,  $\zeta = \frac{B\sqrt{1-\gamma}}{\sqrt{2N}}$ ,  $\eta = \frac{1}{\sqrt{T}}$  and arbitrary  $B > 0$ ,

$$(1-\gamma) \min_{0 \leq t < T} \mathbb{E} [J(\pi^*) - J(\pi_t)] \leq \varepsilon_{\text{nac}} + 2\varepsilon_{\text{inf}} + \bar{C}_{\infty} \left( \varepsilon_{\text{actor}} + 2\varepsilon_{\text{grad}} + 2\sqrt{6} \frac{1}{T} \sum_{t=0}^{T-1} \varepsilon_{\text{critic}}^{\pi_t} \right), \quad (41)$$

where the different terms may differ for asymmetric and symmetric critics,

$$\varepsilon_{\text{nac}} = \frac{B^2 + 2\log|\mathcal{A}|}{2\sqrt{T}} \quad (42)$$

$$\varepsilon_{\text{actor}} = \sqrt{\frac{(2-\gamma)B}{(1-\gamma)\sqrt{N}}} \quad (43)$$

$$\varepsilon_{\text{inf,asym}} = 0 \quad (44)$$

$$\varepsilon_{\text{inf,sym}} = \mathbb{E}^{\pi^*} \left[ \sum_{k=0}^{\infty} \gamma^k \left\| \hat{b}_k - b_k \right\|_{\text{TV}} \right] \quad (45)$$

$$\varepsilon_{\text{grad,asym}} = \sup_{0 \leq t < T} \sqrt{\min_w \mathcal{L}_t(w)} \quad (46)$$

$$\varepsilon_{\text{grad,sym}} = \sup_{0 \leq t < T} \sqrt{\min_w L_t(w)} \quad (47)$$

and  $\varepsilon_{\text{critic}}^{\pi_t}$  is given in [Theorem 3](#) and [Theorem 4](#).

*Proof.* The proof is based on a Lyapounov drift result using the following Lyapounov function,

$$\Lambda(\pi) = \sum_{z \in \mathcal{Z}} d^{\pi^*}(z) \text{KL}(\pi^*(\cdot|z) \parallel \pi(\cdot|z)). \quad (216)$$

The Lyapounov drift is given by,

$$\Lambda(\pi_{t+1}) - \Lambda(\pi_t) = \sum_{z \in \mathcal{Z}} d^{\pi^*}(z) \sum_{a \in \mathcal{A}} \pi^*(a|z) \log \frac{\pi_t(a|z)}{\pi_{t+1}(a|z)} \quad (217)$$

$$= \sum_{z,a} d^{\pi^*}(z,a) \log \frac{\pi_t(a|z)}{\pi_{t+1}(a|z)}. \quad (218)$$

Since  $\sup_{z,a} \|\psi(z,a)\|_2 \leq 1$ , we have that  $\log \pi_\theta(a|z)$  is 1-smooth ([Agarwal et al., 2021](#)), which implies,

$$\log \pi_{\theta_2}(a|z) \leq \log \pi_{\theta_1}(a|z) + \langle \nabla_\theta \log \pi_{\theta_1}(a|z), \theta_2 - \theta_1 \rangle + \frac{1}{2} \|\theta_2 - \theta_1\|_2^2. \quad (219)$$

By selecting  $\theta_2 = \theta_t$  and  $\theta_1 = \theta_{t+1}$  and noting that  $\theta_{t+1} - \theta_t = \eta \bar{w}_t = \eta \frac{1}{N} \sum_{n=0}^{N-1} w_{t,n}$  we obtain,

$$\log \frac{\pi_t(a|z)}{\pi_{t+1}(a|z)} \leq \frac{\eta^2}{2} \|\bar{w}_t\|_2^2 - \eta \langle \nabla_\theta \log \pi_t(a|z), \bar{w}_t \rangle. \quad (220)$$

Now, we separately bound the Lyapounov drift for the asymmetric and symmetric settings. In the following, some notations are overloaded across both setting when their meaning is clear from context. For the asymmetric setting, we have,

$$\Lambda(\pi_{t+1}) - \Lambda(\pi_t) = \sum_{z,a} d^{\pi^*}(z,a) \log \frac{\pi_t(a|z)}{\pi_{t+1}(a|z)} \quad (221)$$

$$\leq \frac{\eta^2}{2} \|\bar{w}_t\|_2^2 - \eta \sum_{z,a} d^{\pi^*}(z,a) \langle \nabla_\theta \log \pi_t(a|z), \bar{w}_t \rangle \quad (222)$$

$$= \frac{\eta^2}{2} B^2 - \eta \sum_{s,z,a} d^{\pi^*}(s,z,a) \mathcal{A}^{\pi_t}(s,z,a) - \eta \sum_{s,z,a} d^{\pi^*}(s,z,a) (\langle \nabla_\theta \log \pi_t(a|z), \bar{w}_t \rangle - \mathcal{A}^{\pi_t}(s,z,a)) \quad (223)$$

$$\leq \frac{\eta^2}{2} B^2 - \eta \sum_{s,z,a} d^{\pi^*}(s,z,a) \mathcal{A}^{\pi_t}(s,z,a) + \eta \sum_{z,a} d^{\pi^*}(s,z,a) \sqrt{(\langle \nabla_\theta \log \pi_t(a|z), \bar{w}_t \rangle - \mathcal{A}^{\pi_t}(s,z,a))^2}. \quad (224)$$

For the symmetric setting, we observe instead,

$$\Lambda(\pi_{t+1}) - \Lambda(\pi_t) = \sum_{z,a} d^{\pi^*}(z,a) \log \frac{\pi_t(a|z)}{\pi_{t+1}(a|z)} \quad (225)$$

$$\leq \frac{\eta^2}{2} B^2 - \eta \sum_{z,a} d^{\pi^*}(z,a) \mathcal{A}^{\pi_t}(z,a) + \eta \sum_{z,a} d^{\pi^*}(z,a) \sqrt{(\langle \nabla_\theta \log \pi_t(a|z), \bar{w}_t \rangle - \mathcal{A}^{\pi_t}(z,a))^2}. \quad (226)$$

Now, let  $\mathfrak{H}_t$  denote the sigma field of all samples used in the computation of  $\pi_t$  (which excludes the samples used for computing  $\bar{w}_t$ ), along with all the samples used in the computation of  $\bar{Q}^{\pi_t}$ . We define the ideal and approximate loss functions, both in the asymmetric and the symmetric setting,

$$\mathcal{L}_t(w) = \mathbb{E} \left[ (\langle \nabla_\theta \log \pi_t(A|Z), w \rangle - \mathcal{A}^{\pi_t}(S,Z,A))^2 \middle| \mathfrak{H}_t \right] \quad (227)$$

$$\bar{\mathcal{L}}_t(w) = \mathbb{E} \left[ \left( \langle \nabla_{\theta} \log \pi_t(A|Z), w \rangle - \bar{\mathcal{A}}^{\pi_t}(S, Z, A) \right)^2 \middle| \mathfrak{H}_t \right] \quad (228)$$

$$L_t(w) = \mathbb{E} \left[ \left( \langle \nabla_{\theta} \log \pi_t(A|Z), w \rangle - A^{\pi_t}(Z, A) \right)^2 \middle| \mathfrak{H}_t \right] \quad (229)$$

$$\bar{L}_t(w) = \mathbb{E} \left[ \left( \langle \nabla_{\theta} \log \pi_t(A|Z), w \rangle - \bar{A}^{\pi_t}(Z, A) \right)^2 \middle| \mathfrak{H}_t \right]. \quad (230)$$

Because  $\mathbb{E} \left[ \|\mathcal{V}^{\pi_t} - \bar{\mathcal{V}}^{\pi_t}\|_{d^{\pi_t}}^2 \middle| \mathfrak{H}_t \right] \leq \mathbb{E} \left[ \|\bar{\mathcal{Q}}^{\pi_t} - \mathcal{Q}^{\pi_t}\|_{d^{\pi_t}}^2 \middle| \mathfrak{H}_t \right]$ , the error between the asymmetric advantage  $\mathcal{A}$  and its approximation  $\bar{\mathcal{A}}$  is upper bounded by,

$$\sqrt{\mathbb{E} \left[ \left( \bar{\mathcal{A}}^{\pi_t}(S, Z, A) - \mathcal{A}^{\pi_t}(S, Z, A) \right)^2 \middle| \mathfrak{H}_t \right]} = \sqrt{\mathbb{E} \left[ \|\bar{\mathcal{A}}^{\pi_t} - \mathcal{A}^{\pi_t}\|_{d^{\pi_t}}^2 \middle| \mathfrak{H}_t \right]} \quad (231)$$

$$= \sqrt{\mathbb{E} \left[ \|\bar{\mathcal{Q}}^{\pi_t} - \bar{\mathcal{V}}^{\pi_t} - \mathcal{Q}^{\pi_t} + \mathcal{V}^{\pi_t}\|_{d^{\pi_t}}^2 \middle| \mathfrak{H}_t \right]} \quad (232)$$

$$= \sqrt{\mathbb{E} \left[ \|\bar{\mathcal{Q}}^{\pi_t} - \mathcal{Q}^{\pi_t} + \mathcal{V}^{\pi_t} - \bar{\mathcal{V}}^{\pi_t}\|_{d^{\pi_t}}^2 \middle| \mathfrak{H}_t \right]} \quad (233)$$

$$\leq \sqrt{\mathbb{E} \left[ \|\bar{\mathcal{Q}}^{\pi_t} - \mathcal{Q}^{\pi_t}\|_{d^{\pi_t}}^2 + \|\mathcal{V}^{\pi_t} - \bar{\mathcal{V}}^{\pi_t}\|_{d^{\pi_t}}^2 \middle| \mathfrak{H}_t \right]} \quad (234)$$

$$\leq \sqrt{\mathbb{E} \left[ \|\bar{\mathcal{Q}}^{\pi_t} - \mathcal{Q}^{\pi_t}\|_{d^{\pi_t}}^2 \middle| \mathfrak{H}_t \right]} + \sqrt{\mathbb{E} \left[ \|\mathcal{V}^{\pi_t} - \bar{\mathcal{V}}^{\pi_t}\|_{d^{\pi_t}}^2 \middle| \mathfrak{H}_t \right]} \quad (235)$$

$$\leq 2\varepsilon_{\text{critic,asym}}^{\pi_t}, \quad (236)$$

where  $\varepsilon_{\text{critic,asym}}^{\pi_t} = \varepsilon_{\text{td,asym}}^{\pi_t} + \varepsilon_{\text{app,asym}}^{\pi_t} + \varepsilon_{\text{shift,asym}}^{\pi_t}$  is given by the upper bound (29) in Theorem 3. Similarly, the error between the symmetric advantage  $A$  and its approximation  $\bar{A}$  is upper bounded by,

$$\sqrt{\mathbb{E} \left[ \left( \bar{A}^{\pi_t}(Z, A) - A^{\pi_t}(Z, A) \right)^2 \middle| \mathfrak{H}_t \right]} \leq 2\varepsilon_{\text{critic,sym}}^{\pi_t}, \quad (237)$$

where  $\varepsilon_{\text{critic,sym}}^{\pi_t} = \varepsilon_{\text{td,sym}}^{\pi_t} + \varepsilon_{\text{app,sym}}^{\pi_t} + \varepsilon_{\text{shift,sym}}^{\pi_t} + \varepsilon_{\text{alias,sym}}^{\pi_t}$  is given by the upper bound (35) in Theorem 4. By using the inequality  $(x + y)^2 \leq 2x^2 + 2y^2$ ,

$$\bar{\mathcal{L}}_t(w) = \mathbb{E} \left[ \left( \langle \nabla_{\theta} \log \pi_t(A|Z), w \rangle - \bar{\mathcal{A}}^{\pi_t}(S, Z, A) \right)^2 \middle| \mathfrak{H}_t \right] \quad (238)$$

$$= \mathbb{E} \left[ \left( \langle \nabla_{\theta} \log \pi_t(A|Z), w \rangle - \mathcal{A}^{\pi_t}(S, Z, A) + \mathcal{A}^{\pi_t}(S, Z, A) - \bar{\mathcal{A}}^{\pi_t}(S, Z, A) \right)^2 \middle| \mathfrak{H}_t \right] \quad (239)$$

$$\leq 2\mathbb{E} \left[ \left( \langle \nabla_{\theta} \log \pi_t(A|Z), w \rangle - \mathcal{A}^{\pi_t}(S, Z, A) \right)^2 \middle| \mathfrak{H}_t \right] + 2\mathbb{E} \left[ \left( \mathcal{A}^{\pi_t}(S, Z, A) - \bar{\mathcal{A}}^{\pi_t}(S, Z, A) \right)^2 \middle| \mathfrak{H}_t \right] \quad (240)$$

$$\leq 2\mathcal{L}_t(w) + 2(2\varepsilon_{\text{critic,asym}}^{\pi_t})^2. \quad (241)$$

Similarly, we obtain in the symmetric case,

$$\bar{L}_t(w) \leq 2L_t(w) + 2(2\varepsilon_{\text{critic,sym}}^{\pi_t})^2. \quad (242)$$

Starting from the ideal objective and following a similar technique, we also obtain,

$$\mathcal{L}_t(w) \leq 2\bar{\mathcal{L}}_t(w) + 2(2\varepsilon_{\text{critic,asym}}^{\pi_t})^2 \quad (243)$$

$$L_t(w) \leq 2\bar{L}_t(w) + 2(2\varepsilon_{\text{critic,sym}}^{\pi_t})^2. \quad (244)$$

By using Theorem 14.8 in (Shalev-Shwartz & Ben-David, 2014) with step size  $\zeta = \frac{B\sqrt{1-\gamma}}{\sqrt{2N}}$ , we obtain for the average iterate  $\bar{w}_t$  under the asymmetric loss and symmetric loss, respectively,

$$\bar{\mathcal{L}}_t(\bar{w}_t) \leq \varepsilon_{\text{actor}}^2 + \min_{\|w\|_2 \leq B} \bar{\mathcal{L}}_t(w) \quad (245)$$

$$\bar{L}_t(\bar{w}_t) \leq \varepsilon_{\text{actor}}^2 + \min_{\|w\|_2 \leq B} \bar{L}_t(w), \quad (246)$$

where  $\varepsilon_{\text{actor}}^2 = \frac{(2-\gamma)B}{2(1-\gamma)\sqrt{N}}$ . On expectation, for the ideal asymmetric objective  $\mathcal{L}_t$ , we obtain,

$$\mathbb{E}[\mathcal{L}_t(\bar{w}_t)] \leq 2\mathbb{E}[\bar{\mathcal{L}}_t(\bar{w}_t)] + 2(2\varepsilon_{\text{critic,asym}}^{\pi_t})^2 \quad (247)$$

$$\leq 2\varepsilon_{\text{actor}}^2 + 2 \min_{\|w\|_2 \leq B} \bar{\mathcal{L}}_t(w) + 2(2\varepsilon_{\text{critic,asym}}^{\pi_t})^2 \quad (248)$$

$$\leq 2\varepsilon_{\text{actor}}^2 + 2 \left( 2 \min_{\|w\|_2 \leq B} \mathcal{L}_t(w) + 2(2\varepsilon_{\text{critic,asym}}^{\pi_t})^2 \right) + 2(2\varepsilon_{\text{critic,asym}}^{\pi_t})^2 \quad (249)$$

$$= 2\varepsilon_{\text{actor}}^2 + 4 \min_{\|w\|_2 \leq B} \mathcal{L}_t(w) + 6(2\varepsilon_{\text{critic,asym}}^{\pi_t})^2 \quad (250)$$

$$= 2\varepsilon_{\text{actor}}^2 + 4 \left( \varepsilon_{\text{grad,asym}}^{\pi_t} \right)^2 + 6(2\varepsilon_{\text{critic,asym}}^{\pi_t})^2, \quad (251)$$

where we define the actor gradient function approximation error as,

$$\left( \varepsilon_{\text{grad,asym}}^{\pi_t} \right)^2 = \min_{\|w\|_2 \leq B} \mathcal{L}_t(w). \quad (252)$$

Similarly, we obtain on expectation for the ideal symmetric objective  $L_t$ ,

$$\mathbb{E}[L_t(\bar{w}_t)] \leq 2\varepsilon_{\text{actor}}^2 + 4 \left( \varepsilon_{\text{grad,sym}}^{\pi_t} \right)^2 + 6(2\varepsilon_{\text{critic,sym}}^{\pi_t})^2, \quad (253)$$

where we define the actor gradient function approximation error as,

$$\left( \varepsilon_{\text{grad,sym}}^{\pi_t} \right)^2 = \min_{\|w\|_2 \leq B} L_t(w). \quad (254)$$

Now, let us go back to the asymmetric and symmetric Lyapounov drift functions of equation (224) and (226). First, we assume that there exists  $\bar{C}_\infty < \infty$  such that  $\sup_{t \geq 0} \mathbb{E}[C_t] \leq \bar{C}_\infty$  with,

$$C_t = \sup_{s, z, a} \left| \frac{d^{\pi^*}(s, z) \pi^*(a|z)}{d^{\pi_{\theta_t}}(s, z) \pi_{\theta_t}(a|z)} \right|. \quad (255)$$

Second, we leverage the performance difference lemma to bound the advantage. For the asymmetric setting, the performance difference lemma for MDP (Kakade & Langford, 2002) holds because of the Markovianity of  $(S_t, Z_t)$ ,

$$(1 - \gamma) \left( V^{\pi^*}(s_0, z_0) - V^{\pi_t}(s_0, z_0) \right) = \mathbb{E}^{d^{\pi^*}}[\mathcal{A}^{\pi_t}(S, Z, A) | S_0 = s_0, Z_0 = z_0]. \quad (256)$$

We note that  $\mathbb{E}[V^{\pi^*}(S_0, Z_0) - V^{\pi_t}(S_0, Z_0)] = \mathbb{E}[J(\pi^*) - J(\pi_t)]$ , such that,

$$-\mathbb{E}^{d^{\pi^*}}[\mathcal{A}^{\pi_t}(S, Z, A)] = -(1 - \gamma)(J(\pi^*) - J(\pi_t)). \quad (257)$$

$$= -(1 - \gamma)(J(\pi^*) - J(\pi_t)) + \varepsilon_{\text{inf,asym}}, \quad (258)$$

where  $\varepsilon_{\text{inf,asym}} = 0$ . For the symmetric setting, using Lemma E.1 with  $\pi_2 = \pi^*$  and  $\pi_1 = \pi_t$ , we note that,

$$(1 - \gamma) \left( V^{\pi^*}(z_0) - V^{\pi_t}(z_0) \right) \leq \mathbb{E}^{d^{\pi^*}}[A^{\pi_t}(Z, A) | Z_0 = z_0] + 2\varepsilon_{\text{inf}}^{\pi^*}(z_0), \quad (259)$$

which implies,

$$-\mathbb{E}^{d^{\pi^*}}[A^{\pi_t}(Z, A) | Z_0 = z_0] \leq -(1 - \gamma) \left( V^{\pi^*}(z_0) - V^{\pi_t}(z_0) \right) + 2\varepsilon_{\text{inf}}^{\pi^*}(z_0). \quad (260)$$

We note that  $\mathbb{E}[V^{\pi^*}(Z_0) - V^{\pi_t}(Z_0)] = \mathbb{E}[J(\pi^*) - J(\pi_t)]$  and we denote  $\mathbb{E}[\varepsilon_{\text{inf}}^{\pi^*}(Z_0)]$  with  $\varepsilon_{\text{inf,sym}}$ , so that,

$$\varepsilon_{\text{inf,sym}} = \mathbb{E}^{\pi^*} \left[ \sum_{k=0}^{\infty} \gamma^k \left\| \hat{b}_k - b_k \right\|_{\text{TV}} \middle| Z_0 = z_0 \right] \quad (261)$$

$$= \mathbb{E}^{\pi^*} \left[ \sum_{k=0}^{\infty} \gamma^k \left\| \hat{b}_k - b_k \right\|_{\text{TV}} \right]. \quad (262)$$

By rearranging, we have,

$$-\mathbb{E}^{d^{\pi^*}} [A^{\pi_t}(Z, A)] \leq -(1 - \gamma) \mathbb{E} [J(\pi^*) - J(\pi_t)] + 2\varepsilon_{\text{inf}, \text{sym}}. \quad (263)$$

Note that  $\sum_{s, z, a} d^{\pi^*}(s, z, a) f(s, z, a) = \sum_{s, z, a} \frac{d^{\pi^*}(s, z, a)}{d^{\pi_t}(s, z, a)} d^{\pi_t}(s, z, a) f(s, z, a) \leq C_t \sum_{s, z, a} d^{\pi_t}(s, z, a) f(s, z, a)$  for positive  $f$ . Taking expectation over the asymmetric Lyapounov drift of equation (224), we obtain using equation (255),

$$\begin{aligned} \mathbb{E} [\Lambda(\pi_{t+1}) - \Lambda(\pi_t)] &\leq \frac{\eta^2}{2} B^2 - \eta \sum_{z, a} d^{\pi^*}(z, a) A^{\pi_t}(z, a) \\ &\quad + \eta \sum_{s, z, a} d^{\pi^*}(s, z, a) \sqrt{(\langle \nabla_{\theta} \log \pi_t(a|z), \bar{w}_t \rangle - \mathcal{A}^{\pi_t}(s, z, a))^2} \end{aligned} \quad (264)$$

$$\begin{aligned} &\leq \frac{\eta^2}{2} B^2 - \eta(1 - \gamma) \mathbb{E} [J(\pi^*) - J(\pi_t)] + 2\eta \varepsilon_{\text{inf}, \text{asym}} \\ &\quad + \eta \bar{C}_{\infty} \sqrt{2\varepsilon_{\text{actor}}^2 + 4 \left( \varepsilon_{\text{grad}, \text{asym}}^{\pi_t} \right)^2 + 6(2\varepsilon_{\text{critic}, \text{asym}}^{\pi_t})^2} \end{aligned} \quad (265)$$

$$\begin{aligned} &\leq \frac{\eta^2}{2} B^2 - \eta(1 - \gamma) \mathbb{E} [J(\pi^*) - J(\pi_t)] + 2\eta \varepsilon_{\text{inf}, \text{asym}} \\ &\quad + \eta \bar{C}_{\infty} \left( \sqrt{2}\varepsilon_{\text{actor}} + 2\varepsilon_{\text{grad}, \text{asym}}^{\pi_t} + 2\sqrt{6}\varepsilon_{\text{critic}, \text{asym}}^{\pi_t} \right). \end{aligned} \quad (266)$$

Similarly, taking expectation over the symmetric drift of equation (226), we obtain a similar expression,

$$\begin{aligned} \mathbb{E} [\Lambda(\pi_{t+1}) - \Lambda(\pi_t)] &\leq \frac{\eta^2}{2} B^2 - \eta \sum_{z, a} d^{\pi^*}(z, a) A^{\pi_t}(z, a) \\ &\quad + \eta \sum_{z, a} d^{\pi^*}(z, a) \sqrt{(\langle \nabla_{\theta} \log \pi_t(a|z), \bar{w}_t \rangle - A^{\pi_t}(z, a))^2} \end{aligned} \quad (267)$$

$$\begin{aligned} &\leq \frac{\eta^2}{2} B^2 - \eta(1 - \gamma) \mathbb{E} [J(\pi^*) - J(\pi_t)] + 2\eta \varepsilon_{\text{inf}, \text{sym}} \\ &\quad + \eta \bar{C}_{\infty} \left( \sqrt{2}\varepsilon_{\text{actor}} + 2\varepsilon_{\text{grad}, \text{sym}}^{\pi_t} + 2\sqrt{6}\varepsilon_{\text{critic}, \text{sym}}^{\pi_t} \right). \end{aligned} \quad (268)$$

Given the similarity of equation (266) and equation (268), in the following we denote the upper bounds using  $\varepsilon_{\text{inf}}$ ,  $\varepsilon_{\text{grad}}^{\pi_t}$  and  $\varepsilon_{\text{critic}}^{\pi_t}$ , irrespectively of the setting (i.e., asymmetric or symmetric).

By summing all Laypounov drifts, we obtain,

$$\begin{aligned} \mathbb{E} [\Lambda(\pi_T) - \Lambda(\pi_0)] &\leq T \frac{\eta^2}{2} B^2 - \eta(1 - \gamma) \sum_{t=0}^{T-1} \mathbb{E} [J(\pi^*) - J(\pi_t)] + 2\eta T \varepsilon_{\text{inf}} \\ &\quad + \eta \sum_{t=0}^{T-1} \bar{C}_{\infty} \left( \sqrt{2}\varepsilon_{\text{actor}} + 2\varepsilon_{\text{grad}}^{\pi_t} + 2\sqrt{6}\varepsilon_{\text{critic}}^{\pi_t} \right) \end{aligned} \quad (269)$$

$$\begin{aligned} &\leq T \frac{\eta^2}{2} B^2 - \eta(1 - \gamma) \sum_{t=0}^{T-1} \mathbb{E} [J(\pi^*) - J(\pi_t)] + 2\eta T \varepsilon_{\text{inf}} \\ &\quad + \eta \bar{C}_{\infty} \left( \sqrt{2}T \varepsilon_{\text{actor}} + 2 \sum_{t=0}^{T-1} \varepsilon_{\text{grad}}^{\pi_t} + 2\sqrt{6} \sum_{t=0}^{T-1} \varepsilon_{\text{critic}}^{\pi_t} \right). \end{aligned} \quad (270)$$

Since  $\pi_0$  is initialized at the uniform policy with  $\theta_0 := 0$ , we have,

$$\Lambda(\pi_0) = \sum_{z \in \mathcal{Z}} d^{\pi^*}(z) \text{KL}(\pi^*(\cdot|z) \parallel \pi_0(\cdot|z)) \quad (271)$$

$$= \sum_{z \in \mathcal{Z}} d^{\pi^*}(z) \left( \sum_{a \in \mathcal{A}} \pi^*(a|z) \log \pi^*(a|z) - \sum_{a \in \mathcal{A}} \pi^*(a|z) \log \pi_0(a|z) \right) \quad (272)$$

$$= \sum_{z \in \mathcal{Z}} d^{\pi^*}(z) \left( \sum_{a \in \mathcal{A}} \pi^*(a|z) \log \pi^*(a|z) - \sum_{a \in \mathcal{A}} \pi^*(a|z) \log \frac{1}{|\mathcal{A}|} \right) \quad (273)$$

$$= \sum_{z \in \mathcal{Z}} d^{\pi^*}(z) \left( \sum_{a \in \mathcal{A}} \pi^*(a|z) \log \pi^*(a|z) + \log |\mathcal{A}| \right) \quad (274)$$

$$= \sum_{z \in \mathcal{Z}} d^{\pi^*}(z) (\log |\mathcal{A}| - H(\pi^*(\cdot|z))) \quad (275)$$

$$\leq \sum_{z \in \mathcal{Z}} d^{\pi^*}(z) \log |\mathcal{A}| \quad (276)$$

$$\leq \log |\mathcal{A}|, \quad (277)$$

where  $H$  denotes the Shannon entropy. Rearranging and dividing by  $\eta T$ , we obtain after neglecting  $\mathcal{L}(\pi_T) > 0$ ,

$$(1 - \gamma) \frac{1}{T} \sum_{t=0}^{T-1} \mathbb{E}[J(\pi^*) - J(\pi_t)] \leq \frac{\log |\mathcal{A}|}{\eta T} + \frac{\eta}{2} B^2 + 2\varepsilon_{\text{inf}} + \bar{C}_{\infty} \left( \sqrt{2}\varepsilon_{\text{actor}} + 2\frac{1}{T} \sum_{t=0}^{T-1} \varepsilon_{\text{grad}}^{\pi_t} + 2\sqrt{6}\frac{1}{T} \sum_{t=0}^{T-1} \varepsilon_{\text{critic}}^{\pi_t} \right). \quad (278)$$

It can also be noted that  $\min_{0 \leq t < T} [x_t] \leq \frac{1}{T} \sum_{t=0}^T x_t$ , which implies that,

$$(1 - \gamma) \min_{0 \leq t < T} \mathbb{E}[J(\pi^*) - J(\pi_t)] \leq \frac{\log |\mathcal{A}|}{\eta T} + \frac{\eta}{2} B^2 + 2\varepsilon_{\text{inf}} + \bar{C}_{\infty} \left( \sqrt{2}\varepsilon_{\text{actor}} + 2\frac{1}{T} \sum_{t=0}^{T-1} \varepsilon_{\text{grad}}^{\pi_t} + 2\sqrt{6}\frac{1}{T} \sum_{t=0}^{T-1} \varepsilon_{\text{critic}}^{\pi_t} \right). \quad (279)$$

Let us define the worse actor gradient function approximation error,

$$\varepsilon_{\text{grad}} = \sup_{0 \leq t < T} \varepsilon_{\text{grad}}^{\pi_t} \quad (280)$$

$$= \sup_{0 \leq t < T} \sqrt{\min_{\|w\|_2 \leq B} L_t(w)}, \quad (281)$$

and let us note that,

$$\frac{1}{T} \sum_{t=0}^{T-1} \varepsilon_{\text{grad}}^{\pi_t} \leq \varepsilon_{\text{grad}}. \quad (282)$$

By setting  $\eta = \frac{1}{\sqrt{T}}$ , we obtain,

$$(1 - \gamma) \min_{0 \leq t < T} \mathbb{E}[J(\pi^*) - J(\pi_t)] \leq \frac{\log |\mathcal{A}|}{\sqrt{T}} + \frac{B^2}{2\sqrt{T}} + 2\varepsilon_{\text{inf}} + \bar{C}_{\infty} \left( \sqrt{2}\varepsilon_{\text{actor}} + 2\frac{1}{T} \sum_{t=0}^{T-1} \varepsilon_{\text{grad}}^{\pi_t} + 2\sqrt{6}\frac{1}{T} \sum_{t=0}^{T-1} \varepsilon_{\text{critic}}^{\pi_t} \right) \quad (283)$$

$$= \frac{B^2 + 2\log |\mathcal{A}|}{2\sqrt{T}} + 2\mathbb{E}^{\pi^*} \left[ \sum_{k=0}^{\infty} \gamma^k \left\| \hat{b}_k - b_k \right\|_{\text{TV}} \right] + \bar{C}_{\infty} \left( \sqrt{\frac{(2 - \gamma)B}{(1 - \gamma)\sqrt{N}}} + 2\varepsilon_{\text{grad}} + 2\sqrt{6}\frac{1}{T} \sum_{t=0}^{T-1} \varepsilon_{\text{critic}}^{\pi_t} \right). \quad (284)$$

This concludes the proof.  $\square$