



**HAL**  
open science

## Off-Policy Actor-Critic

Thomas Degris, Martha White, Richard S. Sutton

► **To cite this version:**

Thomas Degris, Martha White, Richard S. Sutton. Off-Policy Actor-Critic. International Conference on Machine Learning, Jun 2012, Edinburgh, United Kingdom. hal-00764021

**HAL Id: hal-00764021**

**<https://inria.hal.science/hal-00764021v1>**

Submitted on 12 Dec 2012

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

---

# Off-Policy Actor-Critic

---

**Thomas Degris**

Flowers Team, INRIA, Talence, ENSTA-ParisTech, Paris, France

THOMAS.DEGRIS@INRIA.FR

**Martha White**

**Richard S. Sutton**

RLAI Laboratory, Department of Computing Science, University of Alberta, Edmonton, Canada

WHITEM@CS.UALBERTA.CA

SUTTON@CS.UALBERTA.CA

## Abstract

This paper presents the first actor-critic algorithm for off-policy reinforcement learning. Our algorithm is online and incremental, and its per-time-step complexity scales linearly with the number of learned weights. Previous work on actor-critic algorithms is limited to the on-policy setting and does not take advantage of the recent advances in off-policy gradient temporal-difference learning. Off-policy techniques, such as Greedy-GQ, enable a target policy to be learned while following and obtaining data from another (behavior) policy. For many problems, however, actor-critic methods are more practical than action value methods (like Greedy-GQ) because they explicitly represent the policy; consequently, the policy can be stochastic and utilize a large action space. In this paper, we illustrate how to practically combine the generality and learning potential of off-policy learning with the flexibility in action selection given by actor-critic methods. We derive an incremental, linear time and space complexity algorithm that includes eligibility traces, prove convergence under assumptions similar to previous off-policy algorithms, and empirically show better or comparable performance to existing algorithms on standard reinforcement-learning benchmark problems.

methods with convergence guarantees have been restricted to the *on-policy* setting, in which the agent learns only about the policy it is executing.

In an *off-policy* setting, on the other hand, an agent learns about a policy or policies different from the one it is executing. Off-policy methods have a wider range of applications and learning possibilities. Unlike on-policy methods, off-policy methods are able to, for example, learn about an optimal policy while executing an exploratory policy (Sutton & Barto, 1998), learn from demonstration (Smart & Kaelbling, 2002), and learn multiple tasks in parallel from a single sensorimotor interaction with an environment (Sutton et al., 2011). Because of this generality, off-policy methods are of great interest in many application domains.

The most well known off-policy method is Q-learning (Watkins & Dayan, 1992). However, while Q-Learning is guaranteed to converge to the optimal policy for the tabular (non-approximate) case, it may diverge when using linear function approximation (Baird, 1995). Least-squares methods such as LSTD (Bradtke & Barto, 1996) and LSPI (Lagoudakis & Parr, 2003) can be used off-policy and are sound with linear function approximation, but are computationally expensive; their complexity scales quadratically with the number of features and weights. Recently, these problems have been addressed by the new family of gradient-TD (Temporal Difference) methods (e.g., Sutton et al., 2009), such as Greedy-GQ (Maei et al., 2010), which are of linear complexity and convergent under off-policy training with function approximation.

All action-value methods, including gradient-TD methods such as Greedy-GQ, suffer from three important limitations. First, their target policies are deterministic, whereas many problems have stochastic optimal policies, such as in adversarial settings or in partially observable Markov decision processes. Second, finding the greedy action with respect to the action-

The reinforcement learning framework is a general temporal learning formalism that has, over the last few decades, seen a marked growth in algorithms and applications. Until recently, however, practical online

value function becomes problematic for larger action spaces. Finally, a small change in the action-value function can cause large changes in the policy, which creates difficulties for convergence proofs and for some real-time applications.

The standard way of avoiding the limitations of action-value methods is to use policy-gradient algorithms (Sutton et al., 2000) such as actor-critic methods (e.g., Bhatnagar et al., 2009). For example, the natural actor-critic, an on-policy policy-gradient algorithm, has been successful for learning in continuous action spaces in several robotics applications (Peters & Schaal, 2008).

The first and main contribution of this paper is to introduce the first actor-critic method that can be applied off-policy, which we call Off-PAC, for Off-Policy Actor-Critic. Off-PAC has two learners: the actor and the critic. The actor updates the policy weights. The critic learns an off-policy estimate of the value function for the current actor policy, different from the (fixed) behavior policy. This estimate is then used by the actor to update the policy. For the critic, in this paper we consider a version of Off-PAC that uses GTD( $\lambda$ ) (Maei, 2011), a gradient-TD method with eligibility traces for learning state-value functions. We define a new objective for our policy weights and derive a valid backward-view update using eligibility traces. The time and space complexity of Off-PAC is linear in the number of learned weights.

The second contribution of this paper is an off-policy policy-gradient theorem and a convergence proof for Off-PAC when  $\lambda = 0$ , under assumptions similar to previous off-policy gradient-TD proofs.

Our third contribution is an empirical comparison of Q( $\lambda$ ), Greedy-GQ, Off-PAC, and a soft-max version of Greedy-GQ that we call Softmax-GQ, on three benchmark problems in an off-policy setting. To the best of our knowledge, this paper is the first to provide an empirical evaluation of gradient-TD methods for off-policy control (the closest known prior work is the work of Delp (2011)). We show that Off-PAC outperforms other algorithms on these problems.

## 1. Notation and Problem Setting

In this paper, we consider Markov decision processes with a discrete state space  $\mathcal{S}$ , a discrete action space  $\mathcal{A}$ , a distribution  $P : \mathcal{S} \times \mathcal{S} \times \mathcal{A} \rightarrow [0, 1]$ , where  $P(s'|s, a)$  is the probability of transitioning into state  $s'$  from state  $s$  after taking action  $a$ , and an expected reward function  $\mathcal{R} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$  that provides an expected reward for taking action  $a$  in state  $s$  and transitioning

into  $s'$ . We observe a stream of data, which includes states  $s_t \in \mathcal{S}$ , actions  $a_t \in \mathcal{A}$ , and rewards  $r_t \in \mathbb{R}$  for  $t = 1, 2, \dots$  with actions selected from a fixed behavior policy,  $b(a|s) \in (0, 1]$ .

Given a termination condition  $\gamma : \mathcal{S} \rightarrow [0, 1]$  (Sutton et al., 2011), we define the value function for  $\pi : \mathcal{S} \times \mathcal{A} \rightarrow (0, 1]$  to be:

$$V^{\pi, \gamma}(s) = \mathbb{E}[r_{t+1} + \dots + r_{t+T} | s_t = s] \quad \forall s \in \mathcal{S} \quad (1)$$

where policy  $\pi$  is followed from time step  $t$  and terminates at time  $t + T$  according to  $\gamma$ . We assume termination always occurs in a finite number of steps.

The action-value function,  $Q^{\pi, \gamma}(s, a)$ , is defined as:

$$Q^{\pi, \gamma}(s, a) = \sum_{s' \in \mathcal{S}} P(s'|s, a) [\mathcal{R}(s, a, s') + \gamma(s') V^{\pi, \gamma}(s')] \quad (2)$$

for all  $a \in \mathcal{A}$  and for all  $s \in \mathcal{S}$ . Note that  $V^{\pi, \gamma}(s) = \sum_{a \in \mathcal{A}} \pi(a|s) Q^{\pi, \gamma}(s, a)$ , for all  $s \in \mathcal{S}$ .

The policy  $\pi_{\mathbf{u}} : \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$  is an arbitrary, differentiable function of a weight vector,  $\mathbf{u} \in \mathbb{R}^{N_{\mathbf{u}}}$ ,  $N_{\mathbf{u}} \in \mathbb{N}$ , with  $\pi_{\mathbf{u}}(a|s) > 0$  for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}$ . Our aim is to choose  $\mathbf{u}$  so as to maximize the following scalar objective function:

$$J_{\gamma}(\mathbf{u}) = \sum_{s \in \mathcal{S}} d^b(s) V^{\pi_{\mathbf{u}}, \gamma}(s) \quad (3)$$

where  $d^b(s) = \lim_{t \rightarrow \infty} P(s_t = s | s_0, b)$  is the limiting distribution of states under  $b$  and  $P(s_t = s | s_0, b)$  is the probability that  $s_t = s$  when starting in  $s_0$  and executing  $b$ . The objective function is weighted by  $d^b$  because, in the off-policy setting, data is obtained according to this behavior distribution. For simplicity of notation, we will write  $\pi$  and implicitly mean  $\pi_{\mathbf{u}}$ .

## 2. The Off-PAC Algorithm

In this section, we present the Off-PAC algorithm in three steps. First, we explain the basic theoretical ideas underlying the gradient-TD methods used in the critic. Second, we present our off-policy version of the policy-gradient theorem. Finally, we derive the forward view of the actor and convert it to a backward view to produce a complete mechanistic algorithm using eligibility traces.

### 2.1. The Critic: Policy Evaluation

Evaluating a policy  $\pi$  consists of learning its value function,  $V^{\pi, \gamma}(s)$ , as defined in Equation 1. Since it is often impractical to explicitly represent every

state  $s$ , we learn a linear approximation of  $V^{\pi, \gamma}(s)$ :  $\hat{V}(s) = \mathbf{v}^\top \mathbf{x}_s$  where  $\mathbf{x}_s \in \mathbb{R}^{N_v}$ ,  $N_v \in \mathbb{N}$ , is the feature vector of the state  $s$ , and  $\mathbf{v} \in \mathbb{R}^{N_v}$  is another weight vector.

Gradient-TD methods (Sutton et al., 2009) incrementally learn the weights,  $\mathbf{v}$ , in an off-policy setting, with a guarantee of stability and a linear per-time-step complexity. These methods minimize the  $\lambda$ -weighted mean-squared projected Bellman error:

$$\text{MSPBE}(\mathbf{v}) = \|\hat{V} - \Pi T_\pi^{\lambda, \gamma} \hat{V}\|_D^2$$

where  $\hat{V} = X\mathbf{v}$ ;  $X$  is the matrix whose rows are all  $\mathbf{x}_s$ ;  $\lambda$  is the decay of the eligibility trace;  $D$  is a matrix with  $d^b(s)$  on its diagonal;  $\Pi$  is a projection operator that projects a value function to the nearest representable value function given the function approximator; and  $T_\pi^{\lambda, \gamma}$  is the  $\lambda$ -weighted Bellman operator for the target policy  $\pi$  with termination probability  $\gamma$  (e.g., see Maei & Sutton, 2010). For a linear representation,  $\Pi = X(X^\top D X)^{-1} X^\top D$ .

In this paper, we consider the version of Off-PAC that updates its critic weights by the GTD( $\lambda$ ) algorithm introduced by Maei (2011).

## 2.2. Off-policy Policy-gradient Theorem

Like other policy gradient algorithms, Off-PAC updates the weights approximately in proportion to the gradient of the objective:

$$\mathbf{u}_{t+1} - \mathbf{u}_t \approx \alpha_{u,t} \nabla_{\mathbf{u}} J_\gamma(\mathbf{u}_t) \quad (4)$$

where  $\alpha_{u,t} \in \mathbb{R}$  is a positive step-size parameter. Starting from Equation 3, the gradient can be written:

$$\begin{aligned} \nabla_{\mathbf{u}} J_\gamma(\mathbf{u}) &= \nabla_{\mathbf{u}} \left[ \sum_{s \in \mathcal{S}} d^b(s) \sum_{a \in \mathcal{A}} \pi(a|s) Q^{\pi, \gamma}(s, a) \right] \\ &= \sum_{s \in \mathcal{S}} d^b(s) \sum_{a \in \mathcal{A}} [\nabla_{\mathbf{u}} \pi(a|s) Q^{\pi, \gamma}(s, a) \\ &\quad + \pi(a|s) \nabla_{\mathbf{u}} Q^{\pi, \gamma}(s, a)] \end{aligned}$$

The final term in this equation,  $\nabla_{\mathbf{u}} Q^{\pi, \gamma}(s, a)$ , is difficult to estimate in an incremental off-policy setting. The first approximation involved in the theory of Off-PAC is to omit this term. That is, we work with an approximation to the gradient, which we denote  $\mathbf{g}(\mathbf{u}) \in \mathbb{R}^{N_u}$ , defined by

$$\nabla_{\mathbf{u}} J_\gamma(\mathbf{u}) \approx \mathbf{g}(\mathbf{u}) = \sum_{s \in \mathcal{S}} d^b(s) \sum_{a \in \mathcal{A}} \nabla_{\mathbf{u}} \pi(a|s) Q^{\pi, \gamma}(s, a) \quad (5)$$

The two theorems below provide justification for this approximation.

**Theorem 1** (Policy Improvement). *Given any policy parameter  $\mathbf{u}$ , let*

$$\mathbf{u}' = \mathbf{u} + \alpha \mathbf{g}(\mathbf{u})$$

*Then there exists an  $\epsilon > 0$  such that, for all positive  $\alpha < \epsilon$ ,*

$$J_\gamma(\mathbf{u}') \geq J_\gamma(\mathbf{u})$$

*Further, if  $\pi$  has a tabular representation (i.e., separate weights for each state), then  $V^{\pi_{\mathbf{u}'}, \gamma}(s) \geq V^{\pi_{\mathbf{u}}, \gamma}(s)$  for all  $s \in \mathcal{S}$ .*

(Proof in Appendix).

In the conventional on-policy theory of policy-gradient methods, the policy-gradient theorem (Marbach & Tsitsiklis, 1998; Sutton et al., 2000) establishes the relationship between the gradient of the objective function and the expected action values. In our notation, that theorem essentially says that our approximation is exact, that  $\mathbf{g}(\mathbf{u}) = \nabla_{\mathbf{u}} J_\gamma(\mathbf{u})$ . Although, we can not show this in the off-policy case, we can establish a relationship between the solutions found using the true and approximate gradient:

**Theorem 2** (Off-Policy Policy-Gradient Theorem). *Given  $\mathcal{U} \subset \mathbb{R}^{N_u}$  a non-empty, compact set, let*

$$\begin{aligned} \tilde{\mathcal{Z}} &= \{\mathbf{u} \in \mathcal{U} \mid \mathbf{g}(\mathbf{u}) = 0\} \\ \mathcal{Z} &= \{\mathbf{u} \in \mathcal{U} \mid \nabla_{\mathbf{u}} J_\gamma(\mathbf{u}) = 0\} \end{aligned}$$

*where  $\mathcal{Z}$  is the true set of local maxima and  $\tilde{\mathcal{Z}}$  the set of local maxima obtained from using the approximate gradient,  $\mathbf{g}(\mathbf{u})$ . If the value function can be represented by our function class, then  $\mathcal{Z} \subset \tilde{\mathcal{Z}}$ . Moreover, if we use a tabular representation for  $\pi$ , then  $\mathcal{Z} = \tilde{\mathcal{Z}}$ .*

(Proof in Appendix).

The proof of Theorem 2, showing that  $\mathcal{Z} = \tilde{\mathcal{Z}}$ , requires tabular  $\pi$  to avoid update overlap: updates to a single parameter influence the action probabilities for only one state. Consequently, both parts of the gradient (one part with the gradient of the policy function and the other with the gradient of the action-value function) locally greedily change the action probabilities for only that one state. Extrapolating from this result, in practice, more generally a local representation for  $\pi$  will likely suffice, where parameter updates influence only a small number of states. Similarly, in the non-tabular case, the claim will likely hold if  $\gamma$  is small (the return is myopic), again because changes to the policy mostly affect the action-value function locally.

Fortunately, from an optimization perspective, for all  $\mathbf{u} \in \tilde{\mathcal{Z}} \setminus \mathcal{Z}$ ,  $J_\gamma(\mathbf{u}) < \min_{\mathbf{u}' \in \mathcal{Z}} J_\gamma(\mathbf{u}')$ , in other words,

$\mathcal{Z}$  represents all the largest local maxima in  $\tilde{\mathcal{Z}}$  with respect to the objective,  $J_\gamma$ . Local optimization techniques, like random restarts, should help ensure that we converge to larger maxima and so to  $\mathbf{u} \in \mathcal{Z}$ . Even with the true gradient, these approaches would be incorporated into learning because our objective,  $J_\gamma$ , is non-convex.

### 2.3. The Actor: Incremental Update Algorithm with Eligibility Traces

We now derive an incremental update algorithm using observations sampled from the behavior policy. First, we rewrite Equation 5 as an expectation:

$$\begin{aligned} \mathbf{g}(\mathbf{u}) &= \mathbb{E} \left[ \sum_{a \in \mathcal{A}} \nabla_{\mathbf{u}} \pi(a|s) Q^{\pi, \gamma}(s, a) \Big| s \sim d^b \right] \\ &= \mathbb{E} \left[ \sum_{a \in \mathcal{A}} b(a|s) \frac{\pi(a|s)}{b(a|s)} \frac{\nabla_{\mathbf{u}} \pi(a|s)}{\pi(a|s)} Q^{\pi, \gamma}(s, a) \Big| s \sim d^b \right] \\ &= \mathbb{E} [\rho(s, a) \psi(s, a) Q^{\pi, \gamma}(s, a) | s \sim d^b, a \sim b(\cdot|s)] \\ &= \mathbb{E}_b [\rho(s_t, a_t) \psi(s_t, a_t) Q^{\pi, \gamma}(s_t, a_t)] \end{aligned}$$

where  $\rho(s, a) = \frac{\pi(a|s)}{b(a|s)}$ ,  $\psi(s, a) = \frac{\nabla_{\mathbf{u}} \pi(a|s)}{\pi(a|s)}$ , and we introduce the new notation  $\mathbb{E}_b[\cdot]$  to denote the expectation implicitly conditional on all the random variables (indexed by time step) being drawn from their limiting stationary distribution under the behavior policy. A standard result (e.g., see Sutton et al., 2000) is that an arbitrary function of state can be introduced into these equations as a baseline without changing the expected value. We use the approximate state-value function provided by the critic,  $\hat{V}$ , in this way:

$$\mathbf{g}(\mathbf{u}) = \mathbb{E}_b \left[ \rho(s_t, a_t) \psi(s_t, a_t) \left( Q^{\pi, \gamma}(s_t, a_t) - \hat{V}(s_t) \right) \right]$$

The next step is to replace the action value,  $Q^{\pi, \gamma}(s_t, a_t)$ , by the off-policy  $\lambda$ -return. Because these are not exactly equal, this step introduces a further approximation:

$$\mathbf{g}(\mathbf{u}) \approx \widehat{\mathbf{g}}(\mathbf{u}) = \mathbb{E}_b \left[ \rho(s_t, a_t) \psi(s_t, a_t) \left( R_t^\lambda - \hat{V}(s_t) \right) \right]$$

where the off-policy  $\lambda$ -return is defined by:

$$\begin{aligned} R_t^\lambda &= r_{t+1} + (1 - \lambda) \gamma (s_{t+1}) \hat{V}(s_{t+1}) \\ &\quad + \lambda \gamma (s_{t+1}) \rho(s_{t+1}, a_{t+1}) R_{t+1}^\lambda \end{aligned}$$

Finally, based on this equation, we can write the forward view of Off-PAC:

$$\mathbf{u}_{t+1} - \mathbf{u}_t = \alpha_{u,t} \rho(s_t, a_t) \psi(s_t, a_t) \left( R_t^\lambda - \hat{V}(s_t) \right)$$

The forward view is useful for understanding and analyzing algorithms, but for a mechanistic implementation it must be converted to a backward view that

---

#### Algorithm 1 The Off-PAC algorithm

---

Initialize the vectors  $\mathbf{e}_v$ ,  $\mathbf{e}_u$ , and  $\mathbf{w}$  to zero  
 Initialize the vectors  $\mathbf{v}$  and  $\mathbf{u}$  arbitrarily  
 Initialize the state  $s$   
 For each step:  
   Choose an action,  $a$ , according to  $b(\cdot|s)$   
   Observe resultant reward,  $r$ , and next state,  $s'$   
    $\delta \leftarrow r + \gamma (s') \mathbf{v}^\top \mathbf{x}_{s'} - \mathbf{v}^\top \mathbf{x}_s$   
    $\rho \leftarrow \pi_{\mathbf{u}}(a|s) / b(a|s)$   
   Update the critic (GTD( $\lambda$ ) algorithm):  
      $\mathbf{e}_v \leftarrow \rho (\mathbf{x}_s + \gamma(s) \lambda \mathbf{e}_v)$   
      $\mathbf{v} \leftarrow \mathbf{v} + \alpha_v [\delta \mathbf{e}_v - \gamma(s') (1 - \lambda) (\mathbf{w}^\top \mathbf{e}_v) \mathbf{x}_s]$   
      $\mathbf{w} \leftarrow \mathbf{w} + \alpha_w [\delta \mathbf{e}_v - (\mathbf{w}^\top \mathbf{x}_s) \mathbf{x}_s]$   
   Update the actor:  
      $\mathbf{e}_u \leftarrow \rho \left[ \frac{\nabla_{\mathbf{u}} \pi_{\mathbf{u}}(a|s)}{\pi_{\mathbf{u}}(a|s)} + \gamma(s) \lambda \mathbf{e}_u \right]$   
      $\mathbf{u} \leftarrow \mathbf{u} + \alpha_u \delta \mathbf{e}_u$   
    $s \leftarrow s'$

---

does not involve the  $\lambda$ -return. The key step, proved in the appendix, is the observation that

$$\mathbb{E}_b \left[ \rho(s_t, a_t) \psi(s_t, a_t) \left( R_t^\lambda - \hat{V}(s_t) \right) \right] = \mathbb{E}_b [\delta_t \mathbf{e}_t] \quad (6)$$

where  $\delta_t = r_{t+1} + \gamma (s_{t+1}) \hat{V}(s_{t+1}) - \hat{V}(s_t)$  is the conventional temporal difference error, and  $\mathbf{e}_t \in \mathbb{R}^{N_u}$  is the eligibility trace of  $\psi$ , updated by:

$$\mathbf{e}_t = \rho(s_t, a_t) (\psi(s_t, a_t) + \lambda \mathbf{e}_{t-1})$$

Finally, combining the three previous equations, the backward view of the actor update can be written simply as:

$$\mathbf{u}_{t+1} - \mathbf{u}_t = \alpha_{u,t} \delta_t \mathbf{e}_t$$

The complete Off-PAC algorithm is given above as Algorithm 1. Note that although the algorithm is written in terms of states  $s$  and  $s'$ , it really only ever needs access to the corresponding feature vectors,  $\mathbf{x}_s$  and  $\mathbf{x}_{s'}$ , and to the behavior policy probabilities,  $b(\cdot|s)$ , for the current state. All of these are typically available in large-scale applications with function approximation. Also note that Off-PAC is fully incremental and has per-time step computation and memory complexity that is linear in the number of weights,  $N_u + N_v$ .

With discrete actions, a common policy distribution is the Gibbs distribution, which uses a linear combination of features  $\pi(a|s) = \frac{e^{\mathbf{u}^\top \phi_{s,a}}}{\sum_b e^{\mathbf{u}^\top \phi_{s,b}}}$  where  $\phi_{s,a}$  are state-action features for state  $s$ , action  $a$ , and where  $\psi(s, a) = \frac{\nabla_{\mathbf{u}} \pi(a|s)}{\pi(a|s)} = \phi_{s,a} - \sum_b \pi(b|s) \phi_{s,b}$ . The state-action features,  $\phi_{s,a}$ , are potentially unrelated to the feature vectors  $\mathbf{x}_s$  used in the critic.

### 3. Convergence Analysis

Our algorithm has the same recursive stochastic form as the off-policy value-function algorithms

$$\mathbf{u}_{t+1} = \mathbf{u}_t + \alpha_t(h(\mathbf{u}_t, \mathbf{v}_t) + M_{t+1})$$

where  $h : \mathbb{R}^N \rightarrow \mathbb{R}^N$  is a differentiable function and  $\{M_t\}_{t \geq 0}$  is a noise sequence. Following previous off-policy gradient proofs (Maei, 2011), we study the behavior of the ordinary differential equation

$$\dot{\mathbf{u}}(t) = \mathbf{u}(h(\mathbf{u}(t), \mathbf{v}))$$

The two updates (for the actor and for the critic) are not independent on each time step; we analyze two separate ODEs using a two timescale analysis (Borkar, 2008). The actor update is analyzed given fixed critic parameters, and vice versa, iteratively (until convergence). We make the following assumptions.

- (A1) The policy viewed as a function of  $\mathbf{u}$ ,  $\pi_{(\cdot)}(a|s) : \mathbb{R}^{N_u} \rightarrow (0, 1]$ , is continuously differentiable,  $\forall s \in \mathcal{S}, a \in \mathcal{A}$ .
- (A2) The update on  $\mathbf{u}_t$  includes a projection operator,  $\Gamma : \mathbb{R}^{N_u} \rightarrow \mathbb{R}^{N_u}$ , that projects any  $\mathbf{u}$  to a compact set  $\mathcal{U} = \{\mathbf{u} \mid q_i(\mathbf{u}) \leq 0, i = 1, \dots, s\} \subset \mathbb{R}^{N_u}$ , where  $q_i(\cdot) : \mathbb{R}^{N_u} \rightarrow \mathbb{R}$  are continuously differentiable functions specifying the constraints of the compact region. For  $\mathbf{u}$  on the boundary of  $\mathcal{U}$ , the gradients of the active  $q_i$  are linearly independent. Assume the compact region is large enough to contain at least one (local) maximum of  $J_\gamma$ .
- (A3) The behavior policy has a minimum positive value  $b_{\min} \in (0, 1]$ :  $b(a|s) \geq b_{\min} \forall s \in \mathcal{S}, a \in \mathcal{A}$
- (A4) The sequence  $(\mathbf{x}_t, \mathbf{x}_{t+1}, r_{t+1})_{t \geq 0}$  is i.i.d. and has uniformly bounded second moments.
- (A5) For every  $\mathbf{u} \in \mathcal{U}$  (the compact region to which  $\mathbf{u}$  is projected),  $V^{\pi, \gamma} : \mathcal{S} \rightarrow \mathbb{R}$  is bounded.

**Remark 1:** It is difficult to prove the boundedness of the iterates without the projection operator. Since we have a bounded function (with range  $(0, 1]$ ), we could instead assume that the gradient goes to zero exponentially as  $\mathbf{u} \rightarrow \infty$ , ensuring boundedness. Previous work, however, has illustrated that the stochasticity in practice makes convergence to an unstable equilibrium unlikely (Pemantle, 1990); therefore, we avoid restrictions on the policy function and do not include the projection in our algorithm

Finally, we have the following (standard) assumptions on features and step-sizes.

- (P1)  $\|\mathbf{x}_t\|_\infty < \infty, \forall t$ , where  $\mathbf{x}_t \in \mathbb{R}^{N_v}$

(P2) Matrices  $C = E[\mathbf{x}_t \mathbf{x}_t^\top]$ ,  $A = E[\mathbf{x}_t(\mathbf{x}_t - \gamma \mathbf{x}_{t+1})^\top]$  are non-singular and uniformly bounded.  $A$ ,  $C$  and  $E[r_{t+1} \mathbf{x}_t]$  are well-defined because the distribution of  $(\mathbf{x}_t, \mathbf{x}_{t+1}, r_{t+1})$  does not depend on  $t$ .

(S1)  $\alpha_{v,t}, \alpha_{w,t}, \alpha_{u,t} > 0, \forall t$  are deterministic such that  $\sum_t \alpha_{v,t} = \sum_t \alpha_{w,t} = \sum_t \alpha_{u,t} = \infty$  and  $\sum_t \alpha_{v,t}^2 < \infty, \sum_t \alpha_{w,t}^2 < \infty$  and  $\sum_t \alpha_{u,t}^2 < \infty$  with  $\frac{\alpha_{u,t}}{\alpha_{v,t}} \rightarrow 0$ .

(S2) Define  $H(A) \doteq (A + A^\top)/2$  and let  $\lambda_{\min}(C^{-1}H(A))$  be the minimum eigenvalue of the matrix  $C^{-1}H(A)$ <sup>1</sup>. Then  $\alpha_{w,t} = \eta \alpha_{v,t}$  for some  $\eta > \max(0, -\lambda_{\min}(C^{-1}H(A)))$ .

**Remark 2:** The assumption  $\alpha_{u,t}/\alpha_{v,t} \rightarrow 0$  in (S1) states that the actor step-sizes go to zero at a faster rate than the value function step-sizes: the actor update moves on a slower timescale than the critic update (which changes more from its larger step sizes). This timescale is desirable because we effectively want a converged value function estimate for the current policy weights,  $\mathbf{u}_t$ . Examples of suitable step sizes are  $\alpha_{v,t} = \frac{1}{t}$ ,  $\alpha_{u,t} = \frac{1}{1+t \log t}$  or  $\alpha_{v,t} = \frac{1}{t^{2/3}}$ ,  $\alpha_{u,t} = \frac{1}{t}$ . (with  $\alpha_{w,t} = \eta \alpha_{v,t}$  for  $\eta$  satisfying (S2)).

The above assumptions are actually quite unrestrictive. Most algorithms inherently assume bounded features with bounded value functions for all policies; unbounded values trivially result in unbounded value function weights. Common policy distributions are smooth, making  $\pi(a|s)$  continuously differentiable in  $\mathbf{u}$ . The least practical assumption is that the tuples  $(\mathbf{x}_t, \mathbf{x}_{t+1}, r_{t+1})$  are i.i.d., in other words, Martingale noise instead of Markov noise. For Markov noise, our proof as well as the proofs for GTD( $\lambda$ ) and GQ( $\lambda$ ), require Borkar's (2008) two-timescale theory to be extended to Markov noise (which is outside the scope of this paper). Finally, the proof for Theorem 3 assumes  $\lambda = 0$ , but should extend to  $\lambda > 0$  similarly to GTD( $\lambda$ ) (see Maei, 2011, Section 7.4, for convergence remarks).

We give a proof sketch of the following convergence theorem, with the full proof in the appendix.

**Theorem 3** (Convergence of Off-PAC). *Let  $\lambda = 0$  and consider the Off-PAC iterations with GTD(0)<sup>2</sup> for the critic. Assume that (A1)-(A5), (P1)-(P2) and (S1)-(S2) hold. Then the policy weights,  $\mathbf{u}_t$ , converge to  $\hat{\mathcal{Z}} = \{u \in \mathcal{U} \mid \widehat{\mathbf{g}}(\mathbf{u}) = 0\}$  and the value function weights,  $\mathbf{v}_t$ , converge to the corresponding TD-solution with probability one.*

**Proof Sketch:** We follow a similar outline to the two timescale analysis for on-policy policy gradient

<sup>1</sup>Minimum exists as all eigenvalues real-valued (Lemma 4)

<sup>2</sup>GTD(0) is GTD( $\lambda$ ) with  $\lambda = 0$ , not the different algorithm called GTD(0) by Sutton, Szepesvari & Maei (2008)

actor-critic (Bhatnagar et al., 2009) and for nonlinear GTD (Maei et al., 2009). We analyze the dynamics for our two weights,  $\mathbf{u}_t$  and  $\mathbf{z}_t^\top = (\mathbf{w}_t^\top \mathbf{v}_t^\top)$ , based on our update rules. The proof involves satisfying seven requirements from Borkar (2008, p. 64) to ensure convergence to an asymptotically stable equilibrium. ■

## 4. Empirical Results

This section compares the performance of Off-PAC to three other off-policy algorithms with linear memory and computational complexity: 1)  $Q(\lambda)$  (called Q-Learning when  $\lambda = 0$ ), 2) Greedy-GQ (GQ( $\lambda$ ) with a greedy target policy), and 3) Softmax-GQ (GQ( $\lambda$ ) with a Softmax target policy). The policy in Off-PAC is a Gibbs distribution as defined in section 2.3.

We used three benchmarks: mountain car, a pendulum problem and a continuous grid world. These problems all have a discrete action space and a continuous state space, for which we use function approximation. The behavior policy is a uniform distribution over all the possible actions in the problem for each time step. Note that  $Q(\lambda)$  may not be stable in this setting (Baird, 1995), unlike all the other algorithms.

The goal of the mountain car problem (see Sutton & Barto, 1998) is to drive an underpowered car to the top of a hill. The state of the system is composed of the current position of the car (in  $[-1.2, 0.6]$ ) and its velocity (in  $[-.07, .07]$ ). The car was initialized with a position of -0.5 and a velocity of 0. Actions are a throttle of  $\{-1, 0, 1\}$ . The reward at each time step is  $-1$ . An episode ends when the car reaches the top of the hill on the right or after 5,000 time steps.

The second problem is a pendulum problem (Doya, 2000). The state of the system consists of the angle (in radians) and the angular velocity (in  $[-78.54, 78.54]$ ) of the pendulum. Actions, the torque applied to the base, are  $\{-2, 0, 2\}$ . The reward is the cosine of the angle of the pendulum with respect to its fixed base. The pendulum is initialized with an angle and an angular velocity of 0 (i.e., stopped in a horizontal position). An episode ends after 5,000 time steps.

For the pendulum problem, it is unlikely that the behavior policy will explore the optimal region where the pendulum is maintained in a vertical position. Consequently, this experiment illustrates which algorithms make best use of limited behavior samples.

The last problem is a continuous grid-world. The state is a 2-dimensional position in  $[0, 1]^2$ . The actions are the pairs  $\{(0.0, 0.0), (-.05, 0.0), (.05, 0.0), (0.0, -.05), (0.0, .05)\}$ , representing moves in both dimensions. Uniform noise in  $[-.025, .025]$  is added

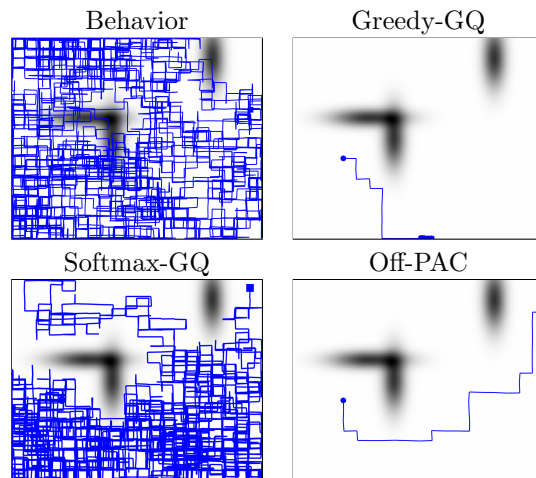
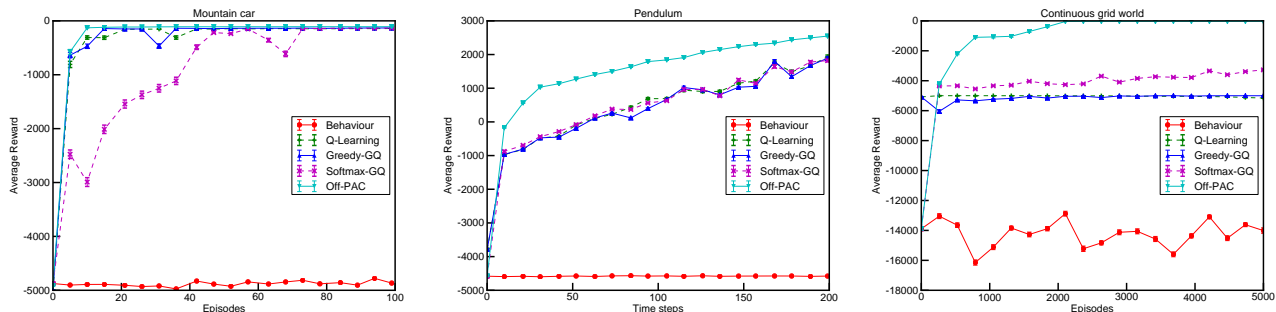


Figure 1. Example of one trajectory for each algorithm in the continuous 2D grid world environment after 5,000 learning episodes from the behavior policy. Off-PAC is the only algorithm that learned to reach the goal reliably.

to each action component. The reward at each time step for arriving in a position  $(p_x, p_y)$  is defined as:  $-1 + -2(\mathcal{N}(p_x, .3, .1) \cdot \mathcal{N}(p_y, .6, .03) + \mathcal{N}(p_x, .4, .03) \cdot \mathcal{N}(p_y, .5, .1) + \mathcal{N}(p_x, .8, .03) \cdot \mathcal{N}(p_y, .9, .1))$  where  $\mathcal{N}(p, \mu, \sigma) = e^{-\frac{(p-\mu)^2}{2\sigma^2}} / \sigma\sqrt{2\pi}$ . The start position is  $(0.2, 0.4)$  and the goal position is  $(1.0, 1.0)$ . An episode ends when the goal is reached, that is when the distance from the current position to the goal is less than 0.1 (using the L1-norm), or after 5,000 time steps. Figure 1 shows a representation of the problem.

The feature vectors  $\mathbf{x}_s$  were binary vectors constructed according to the standard tile-coding technique (Sutton & Barto, 1998). For all problems, we used ten tilings, each of roughly  $10 \times 10$  over the joint space of the two state variables, then hashed to a vector of dimension  $10^6$ . An addition feature was added that was always 1. State-action features,  $\psi_{s,a}$ , were also  $10^6 + 1$  dimensional vectors constructed by also hashing the actions. We used a constant  $\gamma = 0.99$ . All the weight vectors were initialized to 0. We performed a parameter sweep to select the following parameters: 1) the step size  $\alpha_v$  for  $Q(\lambda)$ , 2) the step-sizes  $\alpha_v$  and  $\alpha_w$  for the two vectors in Greedy-GQ, 3)  $\alpha_v$ ,  $\alpha_w$  and the temperature  $\tau$  of the target policy distribution for Softmax-GQ and 4) the step sizes  $\alpha_v$ ,  $\alpha_w$  and  $\alpha_u$  for Off-PAC. For the step sizes, the sweep was done over the following values:  $\{10^{-4}, 5 \cdot 10^{-4}, 10^{-3}, \dots, .5, 1.\}$  divided by  $10+1=11$ , that is the number of tilings plus 1. To compare TD methods to gradient-TD methods, we also used  $\alpha_w = 0$ . The temperature parameter,  $\tau$ , was chosen from  $\{.01, .05, .1, .5, 1, 5, 10, 50, 100\}$  and  $\lambda$  from  $\{0, .2, .4, .6, .8, .99\}$ . We ran thirty runs



		Mountain car					Pendulum					Continuous grid world				
		$\alpha_w$	$\alpha_u, \tau$	$\alpha_v$	$\lambda$	Reward	$\alpha_w$	$\alpha_u, \tau$	$\alpha_v$	$\lambda$	Reward	$\alpha_w$	$\alpha_u, \tau$	$\alpha_v$	$\lambda$	Reward
Behaviour:	final	na	na	na	na	$-4822 \pm 6$	na	na	na	na	$-4582 \pm 0$	na	na	na	na	$-13814 \pm 127$
	overall	na	na	na	na	$-4880 \pm 2$	na	na	na	na	$-4580 \pm 3$	na	na	na	na	$-14237 \pm 33$
Q( $\lambda$ ):	final	na	na	.1	.6	$-143 \pm 4$	na	na	.5	.99	$1802 \pm 35$	na	na	.0001	0	$-5138 \pm 4$
	overall	na	na	.1	0	$-442 \pm 4$	na	na	.5	.99	$376 \pm 15$	na	na	.0001	0	$-5034 \pm 2$
Greedy-GQ:	final	.0001	na	.1	.4	$-131.9 \pm 4$	0	na	.5	.4	$1782 \pm 31$	0.05	na	1.0	.2	$-5002 \pm 2$
	overall	.0001	na	.1	.2	$-434 \pm 4$	.0001	na	.01	.4	$785 \pm 11$	0	na	.0001	0	$-5034 \pm 2$
Softmax-GQ:	final	.0005	.1	.1	.4	$-133.4 \pm 4$	0	.1	.5	.4	$1789 \pm 32$	.1	50	.5	.6	$-3332 \pm 20$
	overall	.0001	.1	.05	.2	$-470 \pm 7$	.0001	.05	.005	.6	$620 \pm 11$	.1	50	.5	.6	$-4450 \pm 11$
Off-PAC:	final	.0001	1.0	.05	0	$-108.6 \pm 0.4$	.005	.5	.5	0	$2521 \pm 17$	0	.001	.1	.4	$-37 \pm 0.1$
	overall	.001	1.0	.5	0	$-356 \pm 4$	0	.5	.5	0	$1432 \pm 10$	0	.001	.005	.6	$-1003 \pm 6$

Figure 2. Performance of Off-PAC compared to the performance of Q( $\lambda$ ), Greedy-GQ, and Softmax-GQ when learning off-policy from a random behavior policy. Final performance selected the parameters for the best performance for the last 10% of the run, whereas the overall performance was over all the runs. The plots on the top show the learning curve for the best parameters for the final performance. Off-PAC had always the best performance and was the only algorithm able to learn to reach the goal reliably in the continuous grid world. Performance is indicated with the standard error.

with each setting of the parameters.

For each parameter combination, the learning algorithm updates a target policy online from the data generated by the behavior policy. For all the problems, the target policy was evaluated at 20 points in time during the run by running it 5 times on another instance of the problem. The target policy was *not* updated during evaluation, ensuring that it was learned only with data from the behavior policy.

Figure 2 shows results on three problems. Softmax-GQ and Off-PAC improved their policy compared to the behavior policy on all problems, while the improvements for Q( $\lambda$ ) and Greedy-GQ is limited on the continuous grid world. Off-PAC performed best on all problems. On the continuous grid world, Off-PAC was the only algorithm able to learn a policy that reliably found the goal after 5,000 episodes (see Figure 1). On all problems, Off-PAC had the lowest standard error.

### 5. Discussion

Off-PAC, like other two-timescale update algorithms, can be sensitive to parameter choices, particularly the step-sizes. Off-PAC has four parameters:  $\lambda$  and the

three step sizes,  $\alpha_v$  and  $\alpha_w$  for the critic and  $\alpha_u$  for the actor. In practice, the following procedure can be used to set these parameters. The value of  $\lambda$ , as with other algorithms, will depend on the problem and it is often better to start with low values (less than .4). A common heuristic is to set  $\alpha_v$  to 0.1 divided by the norm of the feature vector,  $\mathbf{x}_s$ , while keeping the value of  $\alpha_w$  low. Once GTD( $\lambda$ ) is stable learning the value function with  $\alpha_u = 0$ ,  $\alpha_u$  can be increased so that the policy of the actor can be improved. This corroborates the requirements in the proof, where the step-sizes should be chosen so that the slow update (the actor) is not changing as quickly as the fast inner update to the value function weights (the critic).

As mentioned by Borkar (2008, p. 75), another scheme that works well in practice is to use the restrictions on the step-sizes in the proof and to also subsample updates for the slow update. Subsampling updates means only updating every  $\{tN, t \geq 0\}$ , for some  $N > 1$ : the actor is fixed in-between  $tN$  and  $(t+1)N$  while the critic is being updated. This further slows the actor update and enables an improved value function estimate for the current policy,  $\pi$ .

In this work, we did not explore incremental *natural*



actor-critic methods (Bhatnagar et al., 2009), which use the *natural gradient* as opposed to the conventional gradient. The extension to off-policy natural actor-critic should be straightforward, involving only a small modification to the update and analysis of this new dynamical system (which will have similar properties to the original update).

Finally, as pointed out by Precup et al. (2006), off-policy updates can be more noisy compared to on-policy learning. The results in this paper suggest that Off-PAC is more robust to such noise because it has lower variance than the action-value based methods. Consequently, we think Off-PAC is a promising direction for extending off-policy learning to a more general setting such as continuous action spaces.

## 6. Conclusion

This paper proposed a new algorithm for learning control off-policy, called Off-PAC (Off-Policy Actor-Critic). We proved that Off-PAC converges in a standard off-policy setting. We provided one of the first empirical evaluations of off-policy control with the new gradient-TD methods and showed that Off-PAC has the best final performance on three benchmark problems and consistently has the lowest standard error. Overall, Off-PAC is a significant step toward robust off-policy control.

## 7. Acknowledgments

This work was supported by MPrime, the Alberta Innovates Centre for Machine Learning, the Glenrose Rehabilitation Hospital Foundation, Alberta Innovates—Technology Futures, NSERC and the ANR MACSi project. Computational time was provided by Westgrid and the Mésocentre de Calcul Intensif Aquitain.

**Appendix:** See <http://arXiv.org/abs/1205.4839>

## References

- Baird, L. (1995). Residual algorithms: Reinforcement learning with function approximation. In *Proceedings of the Twelfth International Conference on Machine Learning*, pp. 30–37. Morgan Kaufmann.
- Bhatnagar, S., Sutton, R. S., Ghavamzadeh, M., Lee, M. (2009). Natural actor-critic algorithms. *Automatica* 45(11):2471–2482.
- Borkar, V. S. (2008). *Stochastic approximation: A dynamical systems viewpoint*. Cambridge Univ Press.
- Bradtke, S. J., Barto, A. G. (1996). Linear least-squares algorithms for temporal difference learning. *Machine Learning* 22:33–57.
- Delp, M. (2010). Experiments in off-policy reinforcement learning with the GQ( $\lambda$ ) algorithm. Masters thesis, University of Alberta.
- Doya, K. (2000). Reinforcement learning in continuous time and space. *Neural computation* 12:219–245.
- Lagoudakis, M., Parr, R. (2003). Least squares policy iteration. *Journal of Machine Learning Research* 4:1107–1149.
- Maei, H. R., Sutton, R. S. (2010). GQ( $\lambda$ ): A general gradient algorithm for temporal-difference prediction learning with eligibility traces. In *Proceedings of the Third Conf. on Artificial General Intelligence*.
- Maei, H. R. (2011). *Gradient Temporal-Difference Learning Algorithms*. PhD thesis, University of Alberta.
- Maei, H. R., Szepesvári, C., Bhatnagar, S., Precup, D., Silver, D., Sutton, R. S. (2009). Convergent temporal-difference learning with arbitrary smooth function approximation. *Advances in Neural Information Processing Systems* 22:1204–1212.
- Maei, H. R., Szepesvári, C., Bhatnagar, S., Sutton, R. S. (2010). Toward off-policy learning control with function approximation. *Proceedings of the 27th International Conference on Machine Learning*.
- Marbach, P., Tsitsiklis, J. N. (1998). Simulation-based optimization of Markov reward processes. Technical report LIDS-P-2411.
- Pemantle, R. (1990). Nonconvergence to unstable points in urn models and stochastic approximations. *The Annals of Probability* 18(2):698–712.
- Peters, J., Schaal, S. (2008). Natural actor-critic. *Neurocomputing* 71(7):1180–1190.
- Precup, D., Sutton, R.S., Paduraru, C., Koop, A., Singh, S. (2006). Off-policy learning with recognizers. *Neural Information Processing Systems* 18.
- Smart, W.D., Pack Kaelbling, L. (2002). Effective reinforcement learning for mobile robots. In *Proceedings of International Conference on Robotics and Automation*, volume 4, pp. 3404–3410.
- Sutton, R. S., Barto, A. G. (1998). *Reinforcement Learning: An Introduction*. MIT Press.
- Sutton, R. S., McAllester, D., Singh, S., Mansour, Y. (2000). Policy gradient methods for reinforcement learning with function approximation. *Advances in Neural Information Processing Systems* 12.
- Sutton, R. S., Szepesvári, Cs., Maei, H. R. (2008). A convergent  $O(n)$  algorithm for off-policy temporal-difference learning with linear function approximation. In *Advances in Neural Information Processing Systems* 21, pp. 1609–1616.
- Sutton, R. S., Maei, H. R., Precup, D., Bhatnagar, S., Silver, D., Szepesvári, Cs., Wiewiora, E. (2009). Fast gradient-descent methods for temporal-difference learning with linear function approximation. In *Proceedings of the 26th Annual International Conference on Machine Learning*, pp. 993–1000.
- Sutton, R. S., Modayil, J., Delp, M., Degris, T., Pilarski, P. M., and Precup, D. (2011). Horde: A scalable real-time architecture for learning knowledge from unsupervised sensorimotor interaction. In *Proceedings of the 10th International Conference on Autonomous Agents and Multiagent Systems*.
- Watkins, C. J. C. H., Dayan, P. (1992). Q-learning. *Machine Learning* 8(3):279–292.

## A. Appendix of Off-Policy Actor-Critic

### A.1. Policy Improvement and Policy Gradient Theorems

**Theorem 1** [Off-Policy Policy Improvement Theorem]

Given any policy parameter  $\mathbf{u}$ , let

$$\mathbf{u}' = \mathbf{u} + \alpha \mathbf{g}(\mathbf{u})$$

Then there exists an  $\epsilon > 0$  such that, for all positive  $\alpha < \epsilon$ ,

$$J_\gamma(\mathbf{u}') \geq J_\gamma(\mathbf{u})$$

Further, if  $\pi$  has a tabular representation (i.e., separate weights for each state), then  $V^{\pi_{\mathbf{u}'}, \gamma}(s) \geq V^{\pi_{\mathbf{u}}, \gamma}(s)$  for all  $s \in \mathcal{S}$ .

*Proof.* Notice first that for any  $(s, a)$ , the gradient  $\nabla_{\mathbf{u}} \pi(a|s)$  is the direction to increase the probability of action  $a$  according to function  $\pi(\cdot|s)$ . For an appropriate step size  $\alpha_{u,t}$  (so that the update to  $\pi_{\mathbf{u}'}$  increases the objective with the action-value function  $Q^{\pi_{\mathbf{u}}, \gamma}$ , fixed as the old action-value function), we can guarantee that

$$\begin{aligned} J_\gamma(\mathbf{u}) &= \sum_{s \in \mathcal{S}} d^b(s) \sum_{a \in \mathcal{A}} \pi_{\mathbf{u}}(a|s) Q^{\pi_{\mathbf{u}}, \gamma}(s, a) \\ &\leq \sum_{s \in \mathcal{S}} d^b(s) \sum_{a \in \mathcal{A}} \pi_{\mathbf{u}'}(a|s) Q^{\pi_{\mathbf{u}}, \gamma}(s, a) \end{aligned}$$

Now we can proceed similarly to the Policy Improvement theorem proof provided by Sutton and Barto (1998) by extending the right-hand side using the definition of  $Q^{\pi, \gamma}(s, a)$  (equation 2):

$$\begin{aligned} J_\gamma(\mathbf{u}_t) &\leq \sum_{s \in \mathcal{S}} d^b(s) \sum_{a \in \mathcal{A}} \pi_{\mathbf{u}'}(a|s) \mathbb{E}[r_{t+1} + \gamma_{t+1} V^{\pi_{\mathbf{u}}, \gamma}(s_{t+1}) | \pi_{\mathbf{u}'}, \gamma] \\ &\leq \sum_{s \in \mathcal{S}} d^b(s) \sum_{a \in \mathcal{A}} \pi_{\mathbf{u}'}(a|s) \mathbb{E}[r_{t+1} + \gamma_{t+1} r_{t+2} + \gamma_{t+2} V^{\pi_{\mathbf{u}}, \gamma}(s_{t+2}) | \pi_{\mathbf{u}'}, \gamma] \\ &\vdots \\ &\leq \sum_{s \in \mathcal{S}} d^b(s) \sum_{a \in \mathcal{A}} \pi_{\mathbf{u}'}(a|s) Q^{\pi_{\mathbf{u}'}, \gamma}(s, a) \\ &= J_\gamma(\mathbf{u}') \end{aligned}$$

The second part of the Theorem has similar proof to the above. With a tabular representation for  $\pi$ , we know that the gradient satisfies:

$$\sum_{a \in \mathcal{A}} \pi_{\mathbf{u}}(a|s) Q^{\pi_{\mathbf{u}}, \gamma}(s, a) \leq \sum_{a \in \mathcal{A}} \pi_{\mathbf{u}'}(a|s) Q^{\pi_{\mathbf{u}}, \gamma}(s, a)$$

because the probabilities can be updated independently for each state with separate weights for each state.

Now for any  $s \in \mathcal{S}$ :

$$\begin{aligned} V^{\pi_{\mathbf{u}}, \gamma}(s) &= \sum_{a \in \mathcal{A}} \pi_{\mathbf{u}}(a|s) Q^{\pi_{\mathbf{u}}, \gamma}(s, a) \\ &\leq \sum_{a \in \mathcal{A}} \pi_{\mathbf{u}'}(a|s) Q^{\pi_{\mathbf{u}}, \gamma}(s, a) \\ &\leq \sum_{a \in \mathcal{A}} \pi_{\mathbf{u}'}(a|s) \mathbb{E}[r_{t+1} + \gamma_{t+1} V^{\pi_{\mathbf{u}}, \gamma}(s_{t+1}) | \pi_{\mathbf{u}'}, \gamma] \\ &\leq \sum_{a \in \mathcal{A}} \pi_{\mathbf{u}'}(a|s) \mathbb{E}[r_{t+1} + \gamma_{t+1} r_{t+2} + \gamma_{t+2} V^{\pi_{\mathbf{u}}, \gamma}(s_{t+2}) | \pi_{\mathbf{u}'}, \gamma] \\ &\vdots \\ &\leq \sum_{a \in \mathcal{A}} \pi_{\mathbf{u}'}(a|s) Q^{\pi_{\mathbf{u}'}, \gamma}(s, a) \\ &= V^{\pi_{\mathbf{u}'}, \gamma}(s) \end{aligned}$$

□

**Theorem 2** [Off-Policy Policy Gradient Theorem]

Let  $\mathcal{Z} = \{\mathbf{u} \in \mathcal{U} \mid \nabla_{\mathbf{u}} J_{\gamma}(\mathbf{u}) = 0\}$  and  $\tilde{\mathcal{Z}} = \{\mathbf{u} \in \mathcal{U} \mid \mathbf{g}(\mathbf{u}) = 0\}$ , which are both non-empty by Assumption (A2). If the value function can be represented by our function class, then

$$\mathcal{Z} \subset \tilde{\mathcal{Z}}$$

Moreover, if we use a tabular representation for  $\pi$ , then

$$\mathcal{Z} = \tilde{\mathcal{Z}}$$

*Proof.* This theorem follows from our policy improvement theorem.

Assume there exists  $\mathbf{u}^* \in \mathcal{Z}$  such that  $\mathbf{u}^* \notin \tilde{\mathcal{Z}}$ . Then  $\nabla_{\mathbf{u}^*} J_{\gamma}(\mathbf{u}) = 0$  but  $\mathbf{g}(\mathbf{u}^*) \neq 0$ . By the policy improvement theorem (Theorem 1), we know that  $J_{\gamma}(\mathbf{u}^* + \alpha_{u,t} \mathbf{g}(\mathbf{u}^*)) > J_{\gamma}(\mathbf{u}^*)$ , for some positive  $\alpha_{u,t}$ . However, this is a contradiction, as the true gradient is zero. Therefore, such an  $\mathbf{u}^*$  cannot exist.

For the second part of the theorem, we have a tabular representation, in other words, each weight corresponds to exactly one state. Without loss of generality, assume each state  $s$  is represented with  $m \in \mathbb{N}$  weights, indexed by let  $i_{s,1} \dots i_{s,m}$  in the vector  $\mathbf{u}$ . Therefore, for any state,  $s$

$$\sum_{s' \in \mathcal{S}} d^b(s') \sum_{a \in \mathcal{A}} \frac{\partial}{\partial \mathbf{u}_{i_{s,j}}} \pi_{\mathbf{u}}(a|s') Q^{\pi_{\mathbf{u}}, \gamma}(s', a) = d^b(s) \sum_{a \in \mathcal{A}} \frac{\partial}{\partial \mathbf{u}_{i_{s,j}}} \pi_{\mathbf{u}}(a|s) Q^{\pi_{\mathbf{u}}, \gamma}(s, a) \doteq \mathbf{g}_1(\mathbf{u}_{i_{s,j}})$$

Assume there exists  $s \in \mathcal{S}$  such that  $\mathbf{g}_1(\mathbf{u}_{i_{s,j}}) = 0 \forall j$  but there exists  $1 \leq k \leq m$  for  $\mathbf{g}_2(\mathbf{u}_{i_{s,k}}) \doteq \sum_{s' \in \mathcal{S}} d^b(s') \sum_{a \in \mathcal{A}} \pi_{\mathbf{u}}(a|s') \frac{\partial}{\partial \mathbf{u}_{i_{s,k}}} Q^{\pi_{\mathbf{u}}, \gamma}(s', a)$  such that  $\mathbf{g}_2(\mathbf{u}_{i_{s,k}}) \neq 0$ .  $\frac{\partial}{\partial \mathbf{u}_{i_{s,k}}} Q^{\pi_{\mathbf{u}}, \gamma}(s', a)$  can only increase the value of  $Q^{\pi_{\mathbf{u}}, \gamma}(s, a)$  locally (i.e., shift the probabilities of the actions to increase return), because it cannot change the value in other states ( $\mathbf{u}_{i_s}$  is only used for state  $s$  and the remaining weights are fixed when this partial derivative is computed). Therefore, since  $\mathbf{g}_2(\mathbf{u}_{i_{s,k}}) \neq 0$ , we must be able to increase the value of state  $s$  by changing the probabilities of the actions in state  $s$

$$\implies \sum_{j=1}^m \sum_{a \in \mathcal{A}} \frac{\partial}{\partial \mathbf{u}_{i_{s,j}}} \pi_{\mathbf{u}}(a|s) Q^{\pi_{\mathbf{u}}, \gamma}(s, a) \neq 0$$

which is a contradiction (since we assumed  $\mathbf{g}_1(\mathbf{u}_{i_{s,j}}) = 0 \forall j$ ).

Therefore, in the tabular case, whenever  $\sum_s d^b(s) \sum_a \nabla_{\mathbf{u}} \pi_{\mathbf{u}}(a|s) Q^{\pi_{\mathbf{u}}, \gamma}(s, a) = 0$ , then  $\sum_s d^b(s) \sum_a \pi_{\mathbf{u}}(a|s) \nabla_{\mathbf{u}} Q^{\pi_{\mathbf{u}}, \gamma}(s, a) = 0$ , implying that  $\tilde{\mathcal{Z}} \subset \mathcal{Z}$ . Since we already know that  $\mathcal{Z} \subset \tilde{\mathcal{Z}}$ , then we can conclude that for a tabular representation for  $\pi$ ,  $\mathcal{Z} = \tilde{\mathcal{Z}}$ .  $\square$

## A.2. Forward/Backward view analysis

In this section, we prove the key relationship between the forward and backward views:

$$\mathbb{E}_b \left[ \rho(s_t, a_t) \psi(s_t, a_t) \left( R_t^\lambda - \hat{V}(s_t) \right) \right] = \mathbb{E}_b [\delta_t \mathbf{e}_t] \quad (6)$$

where, in these expectations, and in all the expectations in this section, the random variables (indexed by time step) are from their stationary distributions under the behavior policy. We assume that the behavior policy is stationary and that the Markov chain is aperiodic and irreducible (i.e., that we have reached the limiting distribution,  $d^b$ , over  $s \in \mathcal{S}$ ). Note that under these definitions:

$$\mathbb{E}_b [X_t] = \mathbb{E}_b [X_{t+k}] \quad (7)$$

for all integer  $k$  and for all random variables  $X_t$  and  $X_{t+k}$  that are simple temporal shifts of each other. To simplify the notation in this section, we define  $\rho_t = \rho(s_t, a_t)$ ,  $\psi_t = \psi(s_t, a_t)$ ,  $\gamma_t = \gamma(s_t)$ , and  $\delta_t^\lambda = R_t^\lambda - \hat{V}(s_t)$ .

*Proof.* First we note that  $\delta_t^\lambda$ , which might be called the forward-view TD error, can be written recursively:

$$\begin{aligned}
 \delta_t^\lambda &= R_t^\lambda - \hat{V}(s_t) \\
 &= r_{t+1} + (1 - \lambda)\gamma_{t+1}\hat{V}(s_{t+1}) + \lambda\gamma_{t+1}\rho_{t+1}R_{t+1}^\lambda - \hat{V}(s_t) \\
 &= r_{t+1} + \gamma_{t+1}\hat{V}(s_{t+1}) - \lambda\gamma_{t+1}\hat{V}(s_{t+1}) + \lambda\gamma_{t+1}\rho_{t+1}R_{t+1}^\lambda - \hat{V}(s_t) \\
 &= r_{t+1} + \gamma_{t+1}\hat{V}(s_{t+1}) - \hat{V}(s_t) + \lambda\gamma_{t+1}\left(\rho_{t+1}R_{t+1}^\lambda - \hat{V}(s_{t+1})\right) \\
 &= \delta_t + \lambda\gamma_{t+1}\left(\rho_{t+1}R_{t+1}^\lambda - \rho_{t+1}\hat{V}(s_{t+1}) - (1 - \rho_{t+1})\hat{V}(s_{t+1})\right) \\
 &= \delta_t + \lambda\gamma_{t+1}\left(\rho_{t+1}\delta_{t+1}^\lambda - (1 - \rho_{t+1})\hat{V}(s_{t+1})\right)
 \end{aligned} \tag{8}$$

where  $\delta_t = r_{t+1} + \gamma_{t+1}\hat{V}(s_{t+1}) - \hat{V}(s_t)$  is the conventional one-step TD error.

Second, we note that the following expectation is zero:

$$\begin{aligned}
 &\mathbb{E}_b \left[ \rho_t \psi_t \gamma_{t+1} (1 - \rho_{t+1}) \hat{V}(s_{t+1}) \right] \\
 &= \sum_s d^b(s) \sum_a b(a|s) \rho(s, a) \psi(s, a) \sum_{s'} P(s'|s, a) \gamma(s') \left( 1 - \sum_{a'} b(a'|s') \rho(s', a') \right) \hat{V}(s') \\
 &= \sum_s d^b(s) \sum_a b(a|s) \rho(s, a) \psi(s, a) \sum_{s'} P(s'|s, a) \gamma(s') \left( 1 - \sum_{a'} b(a'|s') \frac{\pi(a'|s')}{b(a'|s')} \right) \hat{V}(s') \\
 &= \sum_s d^b(s) \sum_a b(a|s) \rho(s, a) \psi(s, a) \sum_{s'} P(s'|s, a) \gamma(s') \left( 1 - \sum_{a'} \pi(a'|s') \right) \hat{V}(s') \\
 &= 0
 \end{aligned} \tag{9}$$

We are now ready to prove Equation 6 simply by repeated unrolling and rewriting of the right-hand side, using Equations 8, 9, and 7 in sequence, until the pattern becomes clear:

$$\begin{aligned}
 &\mathbb{E}_b \left[ \rho_t \psi_t \left( R_t^\lambda - \hat{V}(s_t) \right) \right] \\
 &= \mathbb{E}_b \left[ \rho_t \psi_t \left( \delta_t + \lambda\gamma_{t+1} \left( \rho_{t+1} \delta_{t+1}^\lambda - (1 - \rho_{t+1}) \hat{V}(s_{t+1}) \right) \right) \right] \tag{using (8)} \\
 &= \mathbb{E}_b [\rho_t \psi_t \delta_t] + \mathbb{E}_b [\rho_t \psi_t \lambda\gamma_{t+1} \rho_{t+1} \delta_{t+1}^\lambda] - \mathbb{E}_b [\rho_t \psi_t \lambda\gamma_{t+1} (1 - \rho_{t+1}) \hat{V}(s_{t+1})] \\
 &= \mathbb{E}_b [\rho_t \psi_t \delta_t] + \mathbb{E}_b [\rho_t \psi_t \lambda\gamma_{t+1} \rho_{t+1} \delta_{t+1}^\lambda] \tag{using (9)} \\
 &= \mathbb{E}_b [\rho_t \psi_t \delta_t] + \mathbb{E}_b [\rho_{t-1} \psi_{t-1} \lambda\gamma_t \rho_t \delta_t^\lambda] \tag{using (7)} \\
 &= \mathbb{E}_b [\rho_t \psi_t \delta_t] + \mathbb{E}_b \left[ \rho_{t-1} \psi_{t-1} \lambda\gamma_t \rho_t \left( \delta_t + \lambda\gamma_{t+1} \left( \rho_{t+1} \delta_{t+1}^\lambda - (1 - \rho_{t+1}) \hat{V}(s_{t+1}) \right) \right) \right] \tag{using (8)} \\
 &= \mathbb{E}_b [\rho_t \psi_t \delta_t] + \mathbb{E}_b [\rho_{t-1} \psi_{t-1} \lambda\gamma_t \rho_t \delta_t] + \mathbb{E}_b [\rho_{t-1} \psi_{t-1} \lambda\gamma_t \rho_t \lambda\gamma_{t+1} \rho_{t+1} \delta_{t+1}^\lambda] \tag{using (9)} \\
 &= \mathbb{E}_b [\rho_t \delta_t (\psi_t + \lambda\gamma_t \rho_{t-1} \psi_{t-1})] + \mathbb{E}_b [\lambda^2 \rho_{t-2} \psi_{t-2} \gamma_{t-1} \rho_{t-1} \gamma_t \rho_t \delta_t^\lambda] \tag{using (7)} \\
 &= \mathbb{E}_b [\rho_t \delta_t (\psi_t + \lambda\gamma_t \rho_{t-1} \psi_{t-1})] + \mathbb{E}_b \left[ \lambda^2 \rho_{t-2} \psi_{t-2} \gamma_{t-1} \rho_{t-1} \gamma_t \rho_t \left( \delta_t + \lambda\gamma_{t+1} \left( \rho_{t+1} \delta_{t+1}^\lambda - (1 - \rho_{t+1}) \hat{V}(s_{t+1}) \right) \right) \right] \\
 &= \mathbb{E}_b [\rho_t \delta_t (\psi_t + \lambda\gamma_t \rho_{t-1} \psi_{t-1})] + \mathbb{E}_b [\lambda^2 \rho_{t-2} \psi_{t-2} \gamma_{t-1} \rho_{t-1} \gamma_t \rho_t \delta_t] + \mathbb{E}_b [\lambda^2 \rho_{t-2} \psi_{t-2} \gamma_{t-1} \rho_{t-1} \gamma_t \rho_t \lambda\gamma_{t+1} \rho_{t+1} \delta_{t+1}^\lambda] \\
 &= \mathbb{E}_b [\rho_t \delta_t (\psi_t + \lambda\gamma_t \rho_{t-1} (\psi_{t-1} + \lambda\gamma_{t-1} \rho_{t-2} \psi_{t-2}))] + \mathbb{E}_b [\lambda^3 \rho_{t-3} \psi_{t-3} \gamma_{t-2} \rho_{t-2} \gamma_{t-1} \rho_{t-1} \gamma_t \rho_t \delta_t^\lambda] \\
 &\vdots \\
 &= \mathbb{E}_b [\rho_t \delta_t (\psi_t + \lambda\gamma_t \rho_{t-1} (\psi_{t-1} + \lambda\gamma_{t-1} \rho_{t-2} (\psi_{t-2} + \lambda\gamma_{t-2} \rho_{t-3} \dots)))] \\
 &= \mathbb{E}_b [\delta_t \mathbf{e}_t]
 \end{aligned}$$

where  $\mathbf{e}_t = \rho_t (\psi_t + \lambda\gamma_t \mathbf{e}_{t-1})$ . □

### A.3. Convergence Proofs

Our algorithm has the same recursive stochastic form that the two-timescale off-policy value-function algorithms have:

$$\begin{aligned} u_{t+1} &= u_t + \alpha_t(h(u_t, z_t) + M_{t+1}) \\ z_{t+1} &= z_t + \alpha_t(f(u_t, z_t) + N_{t+1}) \end{aligned}$$

where  $x \in \mathbb{R}^d$ ,  $h : \mathbb{R}^d \rightarrow \mathbb{R}^d$  is a differentiable functions,  $\{\alpha_t\}_{k \geq 0}$  is a positive step-size sequence and  $\{M_t\}_{k \geq 0}$  is a noise sequence. Again, following the GTD( $\lambda$ ) and GQ( $\lambda$ ) proofs, we study the behavior of the ordinary differential equation

$$\dot{u}(t) = h(u(t), z)$$

Since we have two updates, one for the actor and one for the critic, and those time updates are not linearly separable, we have to do a two timescale analysis (Borkar, 2008). In order to satisfy the conditions for the two-timescale analysis, we will need the following assumptions on our objective, the features and the step-sizes. Note that it is difficult to prove the boundedness of the iterates without the projection operator we describe below, though the projection was not necessary during experiments.

- (A1) The policy function,  $\pi_{(\cdot)}(a|s) : \mathbb{R}^{N_u} \rightarrow [0, 1]$ , is continuously differentiable in  $\mathbf{u}$ ,  $\forall s \in \mathcal{S}, a \in \mathcal{A}$ .
- (A2) The update on  $\mathbf{u}_t$  includes a projection operator,  $\Gamma : \mathbb{R}^{N_u} \rightarrow \mathbb{R}^{N_u}$  that projects any  $\mathbf{u}$  to a compact set  $\mathcal{U} = \{\mathbf{u} \mid q_i(\mathbf{u}) \leq 0, i = 1, \dots, s\} \subset \mathbb{R}^{N_u}$ , where  $q_i(\cdot) : \mathbb{R}^{N_u} \rightarrow \mathbb{R}$  are continuously differentiable functions specifying the constraints of the compact region. For each  $\mathbf{u}$  on the boundary of  $\mathcal{U}$ , the gradients of the active  $q_i$  are considered to be linearly independent. Assume that the compact region,  $\mathcal{U}$ , is large enough to contain at least one local maximum of  $J_\gamma$ .
- (A3) The behavior policy has a minimum positive weight for all actions in every state, in other words,  $b(a|s) \geq b_{\min} \forall s \in \mathcal{S}, a \in \mathcal{A}$ , for some  $b_{\min} \in (0, 1]$ .
- (A4) The sequence  $(\mathbf{x}_t, \mathbf{x}_{t+1}, r_{t+1})_{t \geq 0}$  is i.i.d. and has uniformly bounded second moments.
- (A5) For every  $\mathbf{u} \in \mathcal{U}$  (the compact region to which  $\mathbf{u}$  is projected),  $V^{\pi_{\mathbf{u}}, \gamma} : \mathcal{S} \rightarrow \mathbb{R}$  is bounded.
- (P1)  $\|\mathbf{x}_t\|_\infty < \infty, \forall t$ , where  $\mathbf{x}_t \in \mathbb{R}^{N_v}$
- (P2) The matrices  $C = E[\mathbf{x}_t \mathbf{x}_t^\top]$  and  $A = E[\mathbf{x}_t(\mathbf{x}_t - \gamma \mathbf{x}_{t+1})^\top]$  are non-singular and uniformly bounded.  $A, C$  and  $E[r_{t+1} \mathbf{x}_t]$  are well-defined because the distribution of  $(\mathbf{x}_t, \mathbf{x}_{t+1}, r_{t+1})$  does not depend on  $t$ .
- (S1)  $\alpha_{v,t}, \alpha_{w,t}, \alpha_{u,t} > 0, \forall t$  are deterministic such that  $\sum_t \alpha_{v,t} = \sum_t \alpha_{w,t} = \sum_t \alpha_{u,t} = \infty$  and  $\sum_t \alpha_{v,t}^2 < \infty, \sum_t \alpha_{w,t}^2 < \infty$  and  $\sum_t \alpha_{u,t}^2 < \infty$  with  $\frac{\alpha_{u,t}}{\alpha_{v,t}} \rightarrow 0$ .
- (S2) Define  $H(A) \doteq (A + A^\top)/2$  and let  $\chi_{\min}(C^{-1}H(A))$  be the minimum eigenvalue of the matrix  $C^{-1}H(A)$ . Then  $\alpha_{w,t} = \eta \alpha_{v,t}$  for some  $\eta > \max\{0, -\chi_{\min}(C^{-1}H(A))\}$ .

**Theorem 3** (Convergence of Off-PAC) Let  $\lambda = 0$  and consider the Off-PAC iterations for the critic (GTD( $\lambda$ ), i.e., TDC with importance sampling correction) and the actor (for weights  $\mathbf{u}_t$ ). Assume that (A1)-(A5), (P1)-(P2) and (S1)-(S2) hold. Then the policy weights,  $\mathbf{u}_t$ , converge to  $\hat{\mathcal{Z}} = \{u \in \mathcal{U} \mid \widehat{\mathbf{g}}(\mathbf{u}) = 0\}$  and the value function weights,  $\mathbf{v}_t$ , converge to the corresponding TD-solution with probability one.

*Proof.* We follow a similar outline to that of the two timescale analysis proof for TDC (Sutton et al., 2009). We will analyze the dynamics for our two weights,  $\mathbf{u}_t$ , and  $\mathbf{z}_t^\top = (\mathbf{w}_t^\top \mathbf{v}_t^\top)$ , based on our update rules. We will take  $\mathbf{u}_t$  as the slow timescale update and  $\mathbf{z}_t$  as the fast inner update.

First, we need to rewrite our updates for  $\mathbf{v}$ ,  $\mathbf{w}$  and  $\mathbf{u}$ , amenable to a two timescale analysis:

$$\begin{aligned}\mathbf{v}_{t+1} &= \mathbf{v}_t + \alpha_{v,t}\rho_t[\delta_t\mathbf{x}_t - \gamma\mathbf{x}_t^\top\mathbf{w}\mathbf{x}_t] \\ \mathbf{w}_{t+1} &= \mathbf{w}_t + \alpha_{w,t}\eta[\rho_t\delta_t\mathbf{x}_t - \mathbf{x}_t^\top\mathbf{w}\mathbf{x}_t] \\ \mathbf{z}_{t+1} &= \mathbf{z}_t + \alpha_{v,t}\rho_t[G_{\mathbf{u},t+1}\mathbf{z}_t + q_{\mathbf{u},t+1}]\end{aligned}\tag{10}$$

$$\mathbf{u}_{t+1} = \Gamma\left(\mathbf{u}_t + \alpha_{\mathbf{u},t}\delta_t\frac{\nabla_{\mathbf{u}_t}\pi_t(a_t|s_t)}{b(a_t|s_t)}\right)\tag{11}$$

where  $\rho_t = \rho(s_t, a_t)$ ,  $\delta_t = r_{t+1} + \gamma(\hat{V}(s_{t+1}) - \hat{V}(s_t))$ ,  $\eta = \alpha_{w,t}/\alpha_{v,t}$ ,  $q_{\mathbf{u},t+1}^\top = (\eta\rho_t r_{t+1}\mathbf{x}_t^\top, \rho_t r_{t+1}\mathbf{x}_t^\top)$ , and

$$G_{\mathbf{u},t+1} = \begin{pmatrix} -\eta\mathbf{x}_t\mathbf{x}_t^\top & \eta\rho_t(\mathbf{u}_t)\mathbf{x}_t(\gamma\mathbf{x}_{t+1} - \mathbf{x}_t)^\top \\ -\gamma\rho_t(\mathbf{u}_t)\mathbf{x}_{t+1}\mathbf{x}_t^\top & \rho_t(\mathbf{u}_t)\mathbf{x}_t(\gamma\mathbf{x}_{t+1} - \mathbf{x}_t)^\top \end{pmatrix}.$$

Note that  $G_{\mathbf{u}} = E[G_{\mathbf{u},t}|\mathbf{u}]$  and  $q_{\mathbf{u}} = E[q_{\mathbf{u},t}|\mathbf{u}]$  are well defined because we assumed that the process  $(\mathbf{x}_t, \mathbf{x}_{t+1}, r_{t+1})_{t \geq 0}$  is i.i.d.,  $0 < \rho_t \leq b_{\min}^{-1}$ , and we have fixed  $\mathbf{u}_t$ . Now we can define  $h$  and  $f$ :

$$\begin{aligned}h(\mathbf{z}_t, \mathbf{u}_t) &= G_{\mathbf{u}_t}\mathbf{z}_t + q_{\mathbf{u}_t} \\ f(\mathbf{z}_t, \mathbf{u}_t) &= E\left[\delta_t\frac{\nabla_{\mathbf{u}_t}\pi_t(a_t|s_t)}{b(a_t|s_t)}|\mathbf{z}_t, \mathbf{u}_t\right] \\ M_{t+1} &= (G_{\mathbf{u}_t,t+1} - G_{\mathbf{u}_t})\mathbf{z}_t + q_{\mathbf{u}_t,t+1} - q_{\mathbf{u}_t} \\ N_{t+1} &= \delta_t\frac{\nabla_{\mathbf{u}_t}\pi_t(a_t|s_t)}{b(a_t|s_t)} - f(\mathbf{z}_t, \mathbf{u}_t)\end{aligned}$$

We have to satisfy the following conditions from Borkar (2008, p. p64):

- (B1)  $h : \mathbb{R}^{N_{\mathbf{u}}+2N_{\mathbf{v}}} \rightarrow \mathbb{R}^{2N_{\mathbf{v}}}$  and  $f : \mathbb{R}^{N_{\mathbf{u}}+2N_{\mathbf{v}}} \rightarrow \mathbb{R}^{N_{\mathbf{u}}}$  are Lipschitz.
- (B2)  $\alpha_{v,t}$ ,  $\alpha_{u,t} \forall t$  are deterministic and  $\sum_t \alpha_{v,t} = \sum_t \alpha_{u,t} = \infty$ ,  $\sum_t \alpha_{v,t}^2 < \infty$ ,  $\sum_t \alpha_{u,t}^2 < \infty$ ,  $\frac{\alpha_{u,t}}{\alpha_{v,t}} \rightarrow 0$  (i.e., the system in Equation 11 moves on a slower timescale than Equation 10).
- (B3) The sequences  $\{M_t\}_{k \geq 0}$  and  $\{N_t\}_{k \geq 0}$  are Martingale difference sequences w.r.t. the increasing  $\sigma$ -fields,  $\mathcal{F}_t \doteq \sigma(\mathbf{z}_m, \mathbf{u}_m, M_m, N_m, m \leq t)$  (i.e.,  $E[M_{t+1}|\mathcal{F}_t] = 0$ )
- (B4) For some constant  $K > 0$ ,  $E[\|M_{t+1}\|^2|\mathcal{F}_t] \leq K(1+\|x_t\|^2+\|y_t\|^2)$  and  $E[\|N_{t+1}\|^2|\mathcal{F}_t] \leq K(1+\|x_t\|^2+\|y_t\|^2)$  holds for any  $k \geq 0$ .
- (B5) The ODE  $\dot{\mathbf{z}}(t) = h(\mathbf{z}(t), \mathbf{u})$  has a globally asymptotically stable equilibrium  $\chi(\mathbf{u})$  where  $\chi : \mathbb{R}^{N_{\mathbf{u}}} \rightarrow \mathbb{R}^{N_{\mathbf{v}}}$  is a Lipschitz map.
- (B6) The ODE  $\dot{\mathbf{u}}(t) = f(\chi(\mathbf{u}(t)), \mathbf{u}(t))$  has a globally asymptotically stable equilibrium,  $\mathbf{u}^*$ .
- (B7)  $\sup_t(\|\mathbf{z}_t\| + \|\mathbf{u}_t\|) < \infty$ , a.s.

An asymptotically stable equilibrium for a dynamical system is an attracting point for which small perturbations still cause convergence back to that point. If we can verify these conditions, then we can use Theorem 2 by Borkar (2008) that states that  $(\mathbf{z}_t, \mathbf{u}_t) \rightarrow (\chi(\mathbf{u}^*), \mathbf{u}^*)$  a.s. Note that the previous actor-critic proofs transformed the update to the negative update, assuming they were minimizing costs,  $-R$ , rather than maximizing and so converging to a (local) minimum. This is unnecessary because we simply need to prove we have a stable equilibrium, whether a maximum or minimum; therefore, we keep the update as in the algorithm and assume a (local) maximum.

First note that because we have a bounded function,  $\pi_{(\cdot)}(s, a) : U \rightarrow (0, 1]$ , we can more simply satisfy some of the properties from Borkar (2008). Mainly, we know our policy function is Lipschitz (because it is bounded and continuously differentiable), so we know the gradient is bounded, in other words, there exists  $B_{\nabla_{\mathbf{u}}} \in \mathbb{R}$  such that  $\|\nabla_{\mathbf{u}}\pi(a|s)\| \leq B_{\nabla_{\mathbf{u}}}$ .

**For requirement (B1)**,  $h$  is clearly Lipschitz because it is linear in  $\mathbf{z}$  and  $\rho_t(\mathbf{u})$  is continuously differentiable and bounded ( $\rho_t(\mathbf{u}) \leq b_{\min}^{-1}$ ).  $f$  is Lipschitz because it is linear in  $\mathbf{v}$  and  $\nabla_{\mathbf{u}}\pi(a|s)$  is bounded and continuously differentiable (making  $J_\gamma$  with a fixed  $\hat{Q}^{\pi,\gamma}$  continuously differentiable with a bounded derivative).

**Requirement (B2)** is satisfied by our assumptions.

**Requirement (B3)** is satisfied by the construction of  $M_t$  and  $N_t$ .

**For requirement (B4)**, we can first notice that  $M_t$  satisfies the requirement because  $r_{t+1}, \mathbf{x}_t$  and  $\mathbf{x}_{t+1}$  have uniformly bounded second moments (which is the justification used in the TDC proof (Sutton et al., 2009) and because  $0 < \rho_t \leq b_{\min}^{-1}$ ).

$$\begin{aligned} E[|M_{t+1}|^2|\mathcal{F}_t] &= E[|(G_{\mathbf{u}_t,t} - G_{\mathbf{u}_t})\mathbf{z}_t + (q_{\mathbf{u}_t,t} - q_{\mathbf{u}_t})|^2|\mathcal{F}_t] \\ &\leq E[|(G_{\mathbf{u}_t,t} - G_{\mathbf{u}_t})\mathbf{z}_t|^2 + |(q_{\mathbf{u}_t,t} - q_{\mathbf{u}_t})|^2|\mathcal{F}_t] \\ &\leq E[|c_1\mathbf{z}_t|^2 + c_2|\mathcal{F}_t] \\ &\leq K(\|\mathbf{z}_t\|^2 + 1) \leq K(\|\mathbf{z}_t\|^2 + \|\mathbf{u}_t\|^2 + 1) \end{aligned}$$

where the second inequality is by the Cauchy Schwartz inequality,  $(G_{\mathbf{u}_t,t} - G_{\mathbf{u}_t})\mathbf{z}_t \leq c_1\|\mathbf{z}_t\|$  and  $\|q_{\mathbf{u}_t,t} - q_{\mathbf{u}_t}\|^2 \leq c_2$  (because  $r_{t+1}, \mathbf{x}_t$  and  $\mathbf{x}_{t+1}$  have uniformly bounded second moments), with  $c_1, c_2 \in \mathbb{R}_+$ . When then simply set  $K = \max(c_1, c_2)$ .

For  $N_t$ , since the iterates are bounded as we show below for requirement (B7) (giving  $\sup_t \|\mathbf{u}_t\| < B_{\mathbf{u}}$  and  $\sup_t \|\mathbf{z}_t\| < B_{\mathbf{z}}$  for some  $B_{\mathbf{u}}, B_{\mathbf{z}} \in \mathbb{R}$ ), we see that

$$\begin{aligned} E[|N_{t+1}|^2|\mathcal{F}_t] &\leq E \left[ \left\| \delta_t \frac{\nabla_{\mathbf{u}_t}\pi_t(a_t|s_t)}{b(a_t|s_t)} \right\|^2 + \left\| E \left[ \delta_t \frac{\nabla_{\mathbf{u}_t}\pi_t(a_t|s_t)}{b(a_t|s_t)} \mid \mathbf{z}_t, \mathbf{u}_t \right] \right\|^2 \mid \mathcal{F}_t \right] \\ &\leq E \left[ \left\| \delta_t \frac{\nabla_{\mathbf{u}_t}\pi_t(a_t|s_t)}{b(a_t|s_t)} \right\|^2 + E \left[ \left\| \delta_t \frac{\nabla_{\mathbf{u}_t}\pi_t(a_t|s_t)}{b(a_t|s_t)} \right\|^2 \mid \mathbf{z}_t, \mathbf{u}_t \right] \mid \mathcal{F}_t \right] \\ &\leq 2E \left[ \left| \frac{\delta_t}{b(a_t|s_t)} \right|^2 \|\nabla_{\mathbf{u}_t}\pi_t(a_t|s_t)\|^2 \mid \mathcal{F}_t \right] \\ &\leq \frac{2}{b_{\min}^2} E [|\delta_t|^2 B_{\nabla_{\mathbf{u}}}^2 \mid \mathcal{F}_t] \\ &\leq K(\|\mathbf{v}_t\|^2 + 1) \leq K(\|\mathbf{z}_t\|^2 + \|\mathbf{u}_t\|^2 + 1) \end{aligned}$$

for some  $K \in \mathbb{R}$  because  $E[|\delta|^2|\mathcal{F}_t] \leq c_1(1 + \|\mathbf{v}_t\|)$  for some  $c_1 \in \mathbb{R}$  because  $r_{t+1}, \mathbf{x}_t$  and  $\mathbf{x}_{t+1}$  have uniformly bounded second moments and since  $\|\nabla_{\mathbf{u}}\pi(a|s)\| \leq B_{\nabla_{\mathbf{u}}} \forall s \in \mathcal{S}, a \in \mathcal{A}$  (as stated above because  $\pi(a|s)$  is Lipschitz continuous).

**For requirement (B5)**, we know that every policy,  $\pi$ , has a corresponding bounded  $V^{\pi,\gamma}$  (by assumption). We need to show that for each  $\mathbf{u}$ , there is a globally asymptotically stable equilibrium of the system,  $h(\mathbf{z}(t), \mathbf{u})$  (which has yet to be shown for weighted importance sampling TDC, i.e., GTD( $\lambda = 0$ )). To do so, we use the Hartman-Grobman Theorem, that requires us to show that  $G$  has all negative eigenvalues. For readability, we show this in a separate lemma (Lemma 4 below). Using Lemma 4, we know that there exists a function  $\chi : \mathbb{R}^{N_{\mathbf{u}}} \rightarrow \mathbb{R}^{N_{\mathbf{v}}}$  such that  $\chi(\mathbf{u}) = (\mathbf{v}_{\mathbf{u}}^\top \mathbf{w}_{\mathbf{u}}^\top)^\top$ , where  $\mathbf{v}_{\mathbf{u}}$  is the unique TD-solution value-function weights for policy  $\pi$  and  $\mathbf{w}_{\mathbf{u}}$  is the corresponding expectation estimate. This function,  $\chi$ , is continuously differentiable with bounded gradient (by Lemma 5 below) and is therefore a Lipschitz map.

**For requirement (B6)**, we need to prove that our update  $f(\chi(\cdot), \cdot)$  has an asymptotically stable equilibrium. This requirement can be relaxed to a local rather than global asymptotically stable equilibrium, because we simply need convergence. Our objective function,  $J_\gamma$ , is not concave because our policy function,  $\pi(a|s)$  may not be concave in  $\mathbf{u}$ . Instead, we need to prove that all (local) equilibria are asymptotically stable.

We define a vector field operator,  $\hat{\Gamma} : \mathcal{C}(\mathbb{R}^{N_{\mathbf{u}}}) \rightarrow \mathcal{C}(\mathbb{R}^{N_{\mathbf{u}}})$  that projects any gradients leading outside the compact

region,  $\mathcal{U}$ , back into  $\mathcal{U}$ :

$$\hat{\Gamma}(g(y)) = \lim_{h \rightarrow 0} \frac{\Gamma(y + hg(y)) - y}{h}$$

By our forward-backward view analysis and from the same arguments following from Lemma 3 by Bhatnagar et al. (2009), we know that the ODE  $\dot{\mathbf{u}}(t) = f(\chi(\mathbf{u}(t)), \mathbf{u}(t))$  is  $\mathbf{g}(\mathbf{u})$ . Given that we have satisfied requirements 1-5 and given our step-size conditions, using standard arguments (c.f. Lemma 6 in Bhatnagar et al., 2009), we can deduce that  $\mathbf{u}_t$  converges almost surely to the set of asymptotically stable fixed points,  $\tilde{\mathcal{Z}}$ , of  $\dot{\mathbf{u}} = \hat{\Gamma}\mathbf{g}(\mathbf{u})$ .

**For requirement (B7)**, we know that  $\mathbf{u}_t$  is bounded because it is always projected to  $\mathcal{U}$ . Since  $\mathbf{u}$  stays in  $\mathcal{U}$ , we know that  $\mathbf{v}$  stays bounded (by assumption, otherwise  $V^{\pi, \gamma}$  would not be bounded) and correspondingly  $\mathbf{w}(\mathbf{v})$  must stay bounded, by the same argument as by Sutton et al. (2009). Therefore, we have that  $\sup_t \|\mathbf{u}_t\| < B_{\mathbf{u}}$  and that  $\sup_t \|\mathbf{z}_t\| < B_{\rho}$  for some  $B_{\mathbf{u}}, B_{\mathbf{z}} \in \mathbb{R}$ . □

**Lemma 4.** *Under assumptions (A1)-(A5), (P1)-(P2) and (S1)-(S2), for any fixed set of actor weights,  $\mathbf{u} \in \mathcal{U}$ , the GTD( $\lambda = 0$ ) update for the critic weights,  $\mathbf{v}_t$ , converge to the TD solution with probability one.*

*Proof.* Recall that

$$G_{\mathbf{u}, t+1} = \begin{pmatrix} -\eta \mathbf{x}_t \mathbf{x}_t^\top & \eta \rho_t(\mathbf{u}) \mathbf{x}_t (\gamma \mathbf{x}_{t+1} - \mathbf{x}_t)^\top \\ -\gamma \rho_t(\mathbf{u}) \mathbf{x}_{t+1} \mathbf{x}_t^\top & \rho_t(\mathbf{u}) \mathbf{x}_t (\gamma \mathbf{x}_{t+1} - \mathbf{x}_t)^\top \end{pmatrix}.$$

and  $G_{\mathbf{u}} = \mathbb{E}[G_{\mathbf{u}, t}]$ , meaning

$$G_{\mathbf{u}} = \begin{pmatrix} -\eta C & -\eta A_{\rho}(\mathbf{u}) \\ -F_{\rho}(\mathbf{u})^\top & -A_{\rho}(\mathbf{u}) \end{pmatrix}.$$

where  $F_{\rho}(\mathbf{u}) = \gamma \mathbb{E}[\rho_t(\mathbf{u}) \mathbf{x}_{t+1} \mathbf{x}_t^\top]$ , with  $C_{\rho}(\mathbf{u}) = A_{\rho}(\mathbf{u}) - F_{\rho}(\mathbf{u})$ . For the remainder of the proof, we will simply write  $A_{\rho}$  and  $C_{\rho}$ , because it is clear that we have a fixed  $\mathbf{u} \in \mathcal{U}$ .

Because GTD( $\lambda$ ) is solely for value function approximation, the feature vector,  $\mathbf{x}$ , is only dependent on the state:

$$\begin{aligned} \mathbb{E}[\rho_t \mathbf{x}_t \mathbf{x}_t^\top] &= \sum_{s_t, a_t} d(s_t) b(a_t | s_t) \rho_t \mathbf{x}(s_t) \mathbf{x}_t^\top \\ &= \sum_{s_t, a_t} d(s_t) \pi(a_t | s_t) \mathbf{x}(s_t) \mathbf{x}_t^\top \\ &= \sum_{s_t} d(s_t) \mathbf{x}(s_t) \mathbf{x}_t^\top \left( \sum_{a_t} \pi(a_t | s_t) \right) \\ &= \sum_{s_t} d(s_t) \mathbf{x}(s_t) \mathbf{x}_t^\top = \mathbb{E}[\mathbf{x}_t \mathbf{x}_t^\top] \end{aligned}$$

because  $\sum_{a_t} \pi(a_t | s_t) = 1$ . A similar argument shows that  $\mathbb{E}[\rho_t \mathbf{x}_{t+1} \mathbf{x}_t^\top] = \mathbb{E}[\mathbf{x}_{t+1} \mathbf{x}_t^\top]$ . Therefore, we get that  $F_{\rho}(\mathbf{u}) = \gamma \mathbb{E}[\mathbf{x} \mathbf{x}_t^\top]$  and  $A_{\rho}(\mathbf{u}) = \mathbb{E}[\mathbf{x}_t (\gamma \mathbf{x}_{t+1} - \mathbf{x}_t)^\top]$ . The expected value of the update,  $G$ , therefore, is that same as for TDC, which has been shown to converge under our assumptions (see Maei, 2011). □

**Lemma 5.** *Under assumptions (A1)-(A5), (P1)-(P2) and (S1)-(S2), let  $\chi : \mathcal{U} \rightarrow \mathcal{V}$  be the map from policy weights to corresponding value function,  $V^{\pi, \gamma}$ , obtained from using GTD( $\lambda = 0$ ) (proven to exist by Lemma 4). Then  $\chi$  is continuously differentiable with a bounded gradient for all  $\mathbf{u} \in \mathcal{U}$ .*

*Proof.* To show that  $\chi$  is continuous, we use the Weierstrass definition ( $\delta - \epsilon$  definition). Because  $\chi(\mathbf{u}) = -G(\mathbf{u})^{-1} q(\mathbf{u}) = \mathbf{z}_{\mathbf{u}}$ , which is a complicated function of  $\mathbf{u}$ , we can luckily break it up and prove continuity about parts of it. Recall that 1) the inverse of a continuous function is continuous at every point that represents a non-singular matrix and 2) the multiplication of two continuous functions is continuous. Since  $G(\mathbf{u})$  is always nonsingular, we simply need to proof that  $a(\mathbf{u}) \rightarrow G(\mathbf{u})$  and  $b(\mathbf{u}) \rightarrow q(\mathbf{u})$  are continuous.  $G(\mathbf{u})$  is composed of



several block matrices, including  $C$ ,  $F_\rho(\mathbf{u})$  and  $A_\rho(\mathbf{u})$ . We will start by showing that  $\mathbf{u} \rightarrow F_\rho(\mathbf{u})$  is continuous, where  $F_\rho(\mathbf{u}) = -\mathbb{E}[\eta\rho_t(\mathbf{u})\mathbf{x}_{t+1}\mathbf{x}_t^\top|b]$ . The remaining entries are similar.

Take any  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}$ , and  $\mathbf{u} \in \mathcal{U}$ . We know that  $\pi(a|s) : \mathcal{U} \rightarrow [0, 1]$  is continuous for all  $\mathbf{u} \in \mathcal{U}$  (by assumption). Let  $\epsilon_1 = \frac{\epsilon}{\gamma|\mathcal{A}|\mathbb{E}[\mathbf{x}_{t+1}\mathbf{x}_t^\top|b]}$  (well-defined because  $\mathbb{E}[\mathbf{x}_{t+1}\mathbf{x}_t^\top|b]$  is nonsingular). Then we know there exists a  $\delta > 0$  such that for any  $\mathbf{u}_2 \in \mathcal{U}$  with  $\|\mathbf{u}_1 - \mathbf{u}_2\| < \delta$ , then  $\|\pi_{\mathbf{u}_1}(a_t|s_t) - \pi_{\mathbf{u}_2}(a_t|s_t)\| < \epsilon_1$ . Now

$$\begin{aligned}
 \|F_\rho(\mathbf{u}_1) - F_\rho(\mathbf{u}_2)\| &= \gamma \|\mathbb{E}[\rho_t(\mathbf{u}_1)\mathbf{x}_{t+1}\mathbf{x}_t^\top] - \mathbb{E}[\rho_t(\mathbf{u}_2)\mathbf{x}_{t+1}\mathbf{x}_t^\top]\| \\
 &= \gamma \left\| \sum_{s_t, a_t} d^b(s_t)b(a_t|s_t) \frac{\pi_{\mathbf{u}_1}(a_t|s_t)}{b(a_t|s_t)} \mathbf{x}_{t+1}\mathbf{x}_t^\top - \sum_{s_t, a_t} d^b(s_t)b(a_t|s_t) \frac{\pi_{\mathbf{u}_2}(a_t|s_t)}{b(a_t|s_t)} \mathbf{x}_{t+1}\mathbf{x}_t^\top \right\| \\
 &= \gamma \left\| \sum_{s_t, a_t} d^b(s_t)[\pi_{\mathbf{u}_1}(a_t|s_t) - \pi_{\mathbf{u}_2}(a_t|s_t)] \mathbf{x}_{t+1}\mathbf{x}_t^\top \right\| \\
 &< \gamma \sum_{s_t, a_t} d^b(s_t) \|\pi_{\mathbf{u}_1}(a_t|s_t) - \pi_{\mathbf{u}_2}(a_t|s_t)\| \|\mathbf{x}_{t+1}\mathbf{x}_t^\top\| \\
 &< \gamma \epsilon_1 \sum_{s_t, a_t} d^b(s_t) \|\mathbf{x}_{t+1}\mathbf{x}_t^\top\| \\
 &= \gamma \epsilon_1 |\mathcal{A}| \mathbb{E}[\mathbf{x}_{t+1}\mathbf{x}_t^\top|b] = \epsilon
 \end{aligned}$$

Therefore,  $\mathbf{u} \rightarrow F_\rho(\mathbf{u})$  is continuous. This same process can be done for  $A_\rho(\mathbf{u})$  and  $\mathbb{E}[\rho_t(\mathbf{u})r_t\mathbf{x}_t|b]$  in  $q(\mathbf{u})$ .

Since  $\mathbf{u} \rightarrow G$  and  $\mathbf{u} \rightarrow q$  are continuous for all  $\mathbf{u}$ , we know that  $\chi(\mathbf{u}) = -G(\mathbf{u})^{-1}q(\mathbf{u})$  is continuous.

The above can also be accomplished to show that  $\nabla_{\mathbf{u}}\chi$  is continuous, simply by replacing  $\pi$  with  $\nabla_{\mathbf{u}}\pi$  above. Finally, because our policy function is Lipschitz (because it is bounded and continuously differentiable), we know that it has a bounded gradient. As a result, the gradient of  $\chi$  is bounded (since we have nonsingular and bounded expectation matrices), which would again follow from a similar analysis as above.  $\square$