

# Potential-Based Advice for Stochastic Policy Learning

Baicen Xiao<sup>1</sup>, Bhaskar Ramasubramanian<sup>1</sup>, Andrew Clark<sup>2</sup>,  
Hannaneh Hajishirzi<sup>1</sup>, Linda Bushnell<sup>1</sup>, and Radha Poovendran<sup>1</sup>

**Abstract**—This paper augments the reward received by a reinforcement learning agent with potential functions in order to help the agent learn (possibly stochastic) optimal policies. We show that a potential-based reward shaping scheme is able to preserve optimality of stochastic policies, and demonstrate that the ability of an agent to learn an optimal policy is not affected when this scheme is augmented to soft Q-learning. We propose a method to impart potential-based advice schemes to policy gradient algorithms. An algorithm that considers an advantage actor-critic architecture augmented with this scheme is proposed, and we give guarantees on its convergence. Finally, we evaluate our approach on a puddle-jump grid world with indistinguishable states, and the continuous state and action mountain car environment from classical control. Our results indicate that these schemes allow the agent to learn a stochastic optimal policy faster and obtain a higher average reward.

## I. INTRODUCTION

Reinforcement learning (RL) is a framework that allows an agent to complete tasks in an environment, even when a model of the environment is not known. The agent ‘learns’ to complete a task by maximizing its expected long-term reward, where the reward signal is supplied by the environment. RL algorithms have been successfully implemented in many fields, including robotics [1], [2], and games [3], [4]. However, it remains difficult for an RL agent to master new tasks in unseen environments. This is especially true when the reward given by the environment is sparse/ significantly delayed.

It may be possible to guide an RL agent towards more promising solutions faster, if it is equipped with some form of *prior knowledge* about the environment. This can be encoded by modifying the reward signal received by the agent during training. However, the modification must be carried out in a principled manner, since providing an additional reward at each step might distract the agent from the true goal [5]. Potential-based reward shaping (PBRS) is one such method that augments the reward in an environment specified by a Markov Decision Process (MDP) with a term that is a difference of *potentials* [6]. This method is attractive since it easily allows for the recovery of optimal policies, while enabling the agent to learn these policies faster.

Potential functions are typically functions of states. This could be a limitation, since in some cases, such a function may not be able to encode all information available in the environment. To allow for imparting more information to the agent, a potential-based advice (PBA) scheme was proposed in [7]. The potential functions in PBA include both states and actions as their arguments.

To the best of our knowledge, PBRS and PBA schemes in the literature [6], [7], [8] assume that an optimal policy is deterministic. This will not always be the case, since an optimal policy might be a stochastic policy. This is especially true when there are states in the environment that are partially observable or indistinguishable from each other. Moreover, the aforementioned papers limit their focus to discrete state and action spaces.

In this paper, we study the addition of PBRS and PBA schemes to the reward, in settings where: *i*) the optimal policy will be stochastic, and *ii*) state and action spaces may be continuous. We additionally provide guarantees on the convergence of an advantage actor-critic architecture that is augmented with a PBA scheme. We make the following contributions:

- We prove that the ability of an agent to learn an optimal stochastic policy remains unaffected when augmenting PBRS to soft Q-learning.
- We propose a technique for adapting PBA in policy-based methods, in order to use these schemes in environments with continuous state and action spaces.
- We present an Algorithm, **AC-PBA**, describing an advantage actor-critic architecture augmented with PBA, and provide guarantees on its convergence.
- We evaluate our approach on two experimental domains: a discrete-state, discrete-action *Puddle-jump Gridworld* that has indistinguishable states, and a continuous-state, continuous-action *Mountain Car*.

The remainder of this paper is organized as follows: Section II presents related work in reward shaping. Required preliminaries to RL, PBRS and PBA is presented in Section III. Section IV presents our results on using PBRS for stochastic policy learning. We present a method to augment PBA to policy gradient frameworks and an algorithm detailing this in Section V. Experiments validating our approach are reported in Section VI, and we conclude the paper in Section VII.

## II. RELATED WORK

Shaping or augmenting the reward received by an RL agent in order to enable it to learn optimal policies faster

<sup>1</sup>University of Washington, Seattle, WA 98195, USA. {bcxiao, bhaskarr, hannaneh, lb2, rp3}@uw.edu

<sup>2</sup>Worcester Polytechnic Institute, Worcester, MA 01609, USA. aclark@wpi.edu

is an active area of research. Reward modification via human feedback was used in [9], [10] to interactively shape an agent’s response so that it learned a desired behavior. However, frequent human supervision is usually costly and may not be possible in every situation. A curiosity-based RL algorithm for sparse reward environments was presented in [11], where an intrinsic reward signal characterized the prediction error of the agent as a curiosity reward. The reward received by the agent was augmented with a function that represented the number of times the agent had visited a state in [12].

Entropy regularization as a way to encourage exploration of policies during the early stages of learning was studied in [13] and [14]. This was used to lead a policy towards states with a high reward in [15] and [16].

Static potential-based functions were shown to preserve the optimality of deterministic policies in [6]. This property was extended to dynamic potential-based functions in [8]. The authors of [17] showed that when an agent learned a policy using Q-learning, applying PBRS at each training step was equivalent to initializing the Q-function with the potentials. They studied value-based methods, but restricted their focus to learning deterministic policies. The authors of [18] demonstrated a method to transform a reward function into a potential-based function during training. The potential function in PBA was obtained using an ‘experience filter’ in [19].

The use of PBRS in model-based RL was studied in [20], and for episodic RL in [21]. PBRS was extended to planning in partially observable domains in [22]. However, these papers only considered the finite-horizon, discounted cost setting in this paper.

In control theoretic settings, RL algorithms have been used to establish guarantees on convergence to an optimal controller for the Linear Quadratic Regulator, when a model of the underlying system was not known in [23], [24]. A survey of using RL for control is presented in [25]. OpenAI Gym [26] enables the solving of several problems in classical control using RL algorithms.

### III. PRELIMINARIES

#### A. Reinforcement Learning

An MDP [27] is a tuple  $(S, A, \mathbb{T}, \rho_0, R)$ .  $S$  is the set of states,  $A$  the set of actions,  $\mathbb{T} : S \times A \times S \rightarrow [0, 1]$  encodes  $\mathbb{P}(s_{t+1}|s_t, a_t)$ , the probability of transition to  $s_{t+1}$ , given current state  $s_t$  and action  $a_t$ .  $\rho_0$  is a probability distribution over the initial states.  $R : S \times A \rightarrow \mathbb{R}$  denotes the reward that the agent receives when transitioning from  $s_t$  while taking action  $a_t$ . In this paper,  $R < \infty$ .

The goal for an RL agent [28] is to learn a *policy*  $\pi$ , in order to maximize  $J := \mathbb{E}_{\tau \sim \pi}[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t)]$ . Here,  $\gamma$  is a discounting factor, and the expectation is taken over the trajectory  $\tau = (s_0, a_0, r_0, s_1, \dots)$  induced by policy  $\pi$ . If  $\pi : S \rightarrow A$ , the policy is *deterministic*. On the other hand, a randomized policy returns a probability distribution over the set of actions, and is denoted  $\pi : S \times A \rightarrow [0, 1]$ .

The value of a state-action pair  $(s, a)$  following policy  $\pi$  is represented by the *Q-function*, written  $Q^\pi(s, a) = \mathbb{E}_{\tau \sim \pi}[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) | s_0 = s, a_0 = a]$ . The Q-function allows us to calculate the state value  $V^\pi(s) = \mathbb{E}_{a \sim \pi}[Q^\pi(s, a)]$ . The advantage of a particular action  $a$ , over other actions at a state  $s$  is defined by  $A^\pi(s, a) := Q^\pi(s, a) - V^\pi(s)$ .

#### B. Value-based and Policy-based Methods

The RL problem has two general solution techniques. *Value-based* methods determine an optimal policy by maintaining a set of reward estimates when following a particular policy. At each state, an action that achieves the highest (expected) reward is taken. Typical value-based methods to learn greedy (deterministic) policies include Q-learning and Sarsa-learning [28]. Recently, the authors of [29] proposed *soft Q-learning*, which is a value-based method that is able to learn stochastic policies.

In comparison, *policy-based* methods directly search over the policy space [28]. Starting from an initial policy, specified by a set of parameters, these methods compute the expected reward for this policy, and update the parameter set according to certain rules to improve the policy. Policy gradient [30] is one way to achieve policy improvement. This method repeatedly computes (an estimate of) the gradient of the expected reward with respect to the policy parameters. Policy-based approaches usually exhibit better convergence properties, and can be used in continuous action spaces [31]. They can also be used to learn stochastic policies. REINFORCE and actor-critic are examples of policy gradient algorithms [28].

#### C. PBRS and PBA

Reward shaping methods augment the environment reward  $R$  with an additional reward  $F \in \mathbb{R}$ ,  $F < \infty$ . This changes the structure of the original MDP  $M = (S, A, \mathbb{T}, \rho_0, R)$  to  $M' = (S, A, \mathbb{T}, \rho_0, R + F)$ . The goal is to choose  $F$  so that an optimal policy for  $M'$ ,  $\pi_{M'}^*$ , is also optimal for the original MDP  $M$ . *Potential-based reward shaping* (PBRS) schemes were shown to be able to preserve the optimality of deterministic policies in [6].

In PBRS, the function  $F$  is defined as a difference of *potentials*,  $\phi(\cdot)$ . Specifically,  $F(s_t, a_t, s_{t+1}) := \gamma \phi(s_{t+1}) - \phi(s_t)$ . Then, the Q-function,  $Q_M^*(s, a)$ , of the optimal greedy policy for  $M$  and the optimal Q-function  $Q_{M'}^*(s, a)$  for  $M'$  are related by:  $Q_{M'}^*(s, a) = Q_M^*(s, a) - \phi(s)$ . Therefore, the optimal greedy policy is not changed [6], [8], since:

$$\begin{aligned} \pi_{M'}^*(s) &\in \operatorname{argmax}_{a \in A} Q_{M'}^*(s, a) \\ &= \operatorname{argmax}_{a \in A} (Q_M^*(s, a) - \phi(s)) = \operatorname{argmax}_{a \in A} Q_M^*(s, a). \end{aligned}$$

The authors of [7] augmented  $\phi(s)$  to include action  $a$  as an argument. They termed this *potential-based advice* (PBA). There are two forms— *look-ahead PBA* and *look-back PBA*— respectively defined by:

$$F(s_t, a_t, s_{t+1}, a_{t+1}) = \gamma \phi(s_{t+1}, a_{t+1}) - \phi(s_t, a_t) \quad (1)$$

$$F(s_t, a_t, s_{t-1}, a_{t-1}) = \phi(s_t, a_t) - \gamma^{-1} \phi(s_{t-1}, a_{t-1}). \quad (2)$$

For the look-ahead PBA scheme, the state-action value function for  $M$  following policy  $\pi$  is given by:

$$Q_M^\pi(s, a) = Q_{M'}^\pi(s, a) + \phi(s, a). \quad (3)$$

The optimal greedy policy for  $M$  can be recovered from the optimal state-action value function for  $M'$  from:

$$\pi_M^*(s_t) \in \operatorname{argmax}_{a \in A} (Q_{M'}^*(s_t, a) + \phi(s_t, a)). \quad (4)$$

The optimal greedy policy for  $M$  using look-back PBA can be recovered similarly.

#### IV. PBRs FOR STOCHASTIC POLICY LEARNING

The existing literature on PBRs has focused on augmenting value-based methods to learn optimal deterministic policies. In this section, we first show that PBRs preserves optimality, when the optimal policy is stochastic. Then, we show that the *learnability* will not be changed when using PBRs in soft Q-learning.

*Proposition 1:* Assume that the optimal policy is stochastic. Then, with  $F := \gamma\phi(s_{t+1}) - \phi(s_t)$ , PBRs preserves the optimality of stochastic policies.

*Proof:* The goal in the original MDP  $M$  was to find a policy  $\pi$  in order to maximize:

$$\pi_M^* = \operatorname{argmax}_{\pi} \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \right]. \quad (5)$$

In PBRs, the goal is to determine a policy so that:

$$\begin{aligned} \pi_{M'}^* &= \operatorname{argmax}_{\pi} \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t (R(s_t, a_t) + F(s_t, a_t, s_{t+1}, a_{t+1})) \right] \\ &= \operatorname{argmax}_{\pi} \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t (R(s_t, a_t) + \gamma\phi(s_{t+1}) - \phi(s_t)) \right] \\ &= \operatorname{argmax}_{\pi} \left[ \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \right] - \mathbb{E}_{\tau \sim \pi} [\phi(s_0)] \right] \\ &= \operatorname{argmax}_{\pi} \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t R(s_t, a_t) \right] - \int_s \rho_0(s) \phi(s) ds. \end{aligned} \quad (6)$$

The last term in Equation (6) is constant, and doesn't affect the identity of the maximizing policy of (5). ■

Next, we examine the effect on learnability when using PBRs with soft Q-learning. Soft Q-learning is a value-based method for stochastic policy learning that was proposed in [29]. Different from Equation (5), the goal is to maximize both, the accumulated reward, and the policy entropy at each visited state:

$$\pi_{\text{soft}}^* = \operatorname{argmax}_{\pi} \mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t (R(s_t, a_t) + \alpha \mathcal{H}(\pi(\cdot|s_t))) \right]. \quad (7)$$

The entropy term  $\mathcal{H}(\pi(\cdot|s_t))$  encourages exploration of the state space, and the parameter  $\alpha$  is a trade-off between exploitation and exploration.

Before stating our result, we summarize the soft Q-learning update procedure. From [29], the optimal value-function,  $V_{\text{soft}}^*(s_t)$ , is given by:

$$V_{\text{soft}}^*(s_t) = \alpha \log \int_A \exp\left(\frac{1}{\alpha} Q_{\text{soft}}^*(s_t, a)\right) da. \quad (8)$$

The optimal soft Q-function is determined by solving the soft Bellman equation:

$$Q_{\text{soft}}^*(s_t, a_t) = r_t + \gamma \mathbb{E}_{s_{t+1}} [V_{\text{soft}}^*(s_{t+1})]. \quad (9)$$

The optimal policy can be obtained from Equation (9) as:

$$\pi_{\text{soft}}^*(a_t|s_t) = \exp\left(\frac{1}{\alpha} (Q_{\text{soft}}^*(s_t, a_t) - V_{\text{soft}}^*(s_t))\right), \quad (10)$$

In the rest of this Section, we assume both, states and actions are discrete and no function approximator is used. We also omit subscripts for  $Q_{\text{soft}}$  and  $V_{\text{soft}}$ , and set  $\alpha = 1$  for simplicity. From Equation (9), and as in Q-learning, soft Q-learning updates the soft Q-function by minimizing the soft Bellman error:

$$\delta Q_k(s_k, a_k) = r(s_k, a_k) + \gamma V_k(s_{k+1}) - Q_k(s_k, a_k), \quad (11)$$

where  $V_k(s_{t+1}) = \log \sum_{a \in A} \exp(Q_k(s_{t+1}, a))$ . During training,  $\pi_k(a_t|s_t) = \exp(Q_k(s_t, a_t) - V_k(s_t))$ . With  $\lambda$  denoting the learning rate, the Q-function update is given by:

$$Q_{k+1}(s_k, a_k) = Q_k(s_k, a_k) + \lambda \delta Q_k(s_k, a_k). \quad (12)$$

The main result of this section shows that the ability of an agent to learn an optimal policy is unaffected when using soft Q-learning augmented with PBRs. We define a notion of *learnability*, and use this to establish our claim.

During training, an agent encounters a sequence of states, actions, and rewards that serves as ‘raw-data’ which is fed to the RL algorithm. Let  $L$  and  $L'$  denote two RL agents. Let  $\mathcal{D}_k = (s_k, a_k, r_k, s_{k+1})$  and  $\mathcal{D}'_k = (s'_k, a'_k, r'_k, s'_{k+1})$  denote the experience tuple at learning step  $k$  from a trajectory used by  $L$  and  $L'$ , respectively.

*Definition 1 (Learnability):* Denote the accumulated difference in the Q-functions of  $L$  and  $L'$  after learning for  $k$  steps by  $\Delta Q_k(s, a)$  and  $\Delta Q'_k(s, a)$ , respectively. Then, given identical sample experiences, (that is,  $\mathcal{D}_{k'} = \mathcal{D}'_{k'}$ ,  $\forall k' \leq k$ ),  $L$  and  $L'$  are said to have the same learnability if  $\Delta Q_{k'}(s, a) = \Delta Q'_{k'}(s, a) \forall k' \leq k \forall s \forall a$ .

*Proposition 2:* Soft Q-learning, with initial soft Q-values  $Q(s, a) = Q_0(s, a)$  and augmented with PBRs where state potential is  $\phi(s)$ , has the same learnability as soft Q-learning without PBRs but with its soft Q-values initialized to  $Q(s, a) = Q_0(s, a) + \phi(s)$ .

*Proof:* Consider an agent  $L$  that uses a PBRs scheme during learning and an agent  $L'$  that does not use PBRs, but has its soft Q-values initialized as  $Q'_0(s, a) := Q_0(s, a) + \phi(s)$ , where  $Q_0(s, a)$  is the initial Q-value of  $L$ . We further assume that  $L$  and  $L'$  adopt the same learning rate. From Definition 1, to show that  $L$  and  $L'$  have the same learnability, we need to show that the soft Bellman errors  $\delta Q_k(s_t, a_t)$  and  $\delta Q'_k(s_k, a_k)$  are equal at each training step  $k$ , given the same experience sets  $\mathcal{D}_k$  and  $\mathcal{D}'_k$ . From Equation (11), the soft Bellman errors for  $L$  and  $L'$  can be respectively written as:

$$\begin{aligned} \delta Q_k(s_k, a_k) &= r(s_k, a_k) + \gamma \phi(s_{k+1}) - \phi(s_k) + \\ &\quad \gamma V_k(s_{k+1}) - Q_k(s_k, a_k) \\ \delta Q'_k(s'_k, a'_k) &= r(s'_k, a'_k) + \gamma V'_k(s'_{k+1}) - Q'_k(s'_k, a'_k). \end{aligned}$$

Since  $\mathcal{D}_{k'} = \mathcal{D}'_{k'}$  for each  $k' \leq k$ , comparing  $\delta Q'_k(s_k, a_k)$  and  $\delta Q_k(s'_k, a'_k)$  is reduced to comparing  $\delta Q'_k(s_k, a_k)$  and  $\delta Q_k(s_k, a_k)$ . We show this by induction.

At training step  $k = 0$  there is no update. Thus,  $\delta Q_0(s_0, a_0) = \delta Q'_0(s_0, a_0)$ . Assume that the Bellman errors are identical up to a step  $k = K$ . That is,  $\delta Q_k(s_k, a_k) = \delta Q'_k(s_k, a_k) \forall k \leq K$ . Then, the accumulated errors for the two agents until this step are also identical. That is,  $\Delta Q_K(s, a) = \Delta Q'_K(s, a) \forall s \forall a$ . Consider training step  $k = K + 1$ . The state values at this step are:  $V_K(s_{K+1}) = \log \sum_{a \in A} \exp [Q_0(s_{K+1}, a) + \Delta Q_K(s_{K+1}, a)]$  and  $V'_K(s_{K+1}) = \log \sum_{a \in A} \exp [Q_0(s_{K+1}, a) + \phi(s_{K+1}) + \Delta Q'_K(s_{K+1}, a)]$  respectively. The Bellman errors at  $k = K + 1$  are:

$$\begin{aligned} \delta Q_{K+1}(s_K, a_K) &= r(s_K, a_K) + \gamma \phi(s_{K+1}) - \phi(s_K) \\ &\quad + \gamma V_K(s_{K+1}) - Q_K(s_K, a_K) \\ &= r(s_K, a_K) + \gamma \phi(s_{K+1}) - \phi(s_K) + \gamma V_K(s_{K+1}) \\ &\quad - Q_0(s_K, a_K) - \Delta Q_K(s_K, a_K) \end{aligned}$$

$$\begin{aligned} \delta Q'_{K+1}(s_K, a_K) &= r(s_K, a_K) + \gamma V'_K(s_{K+1}) - Q'_K(s_K, a_K) \\ &= r(s_K, a_K) + \gamma V'_K(s_{K+1}) - Q_0(s_K, a_K) - \phi(s_K) - \Delta Q'_K(s_K, a_K) \\ &= \delta Q_{K+1}(s_K, a_K) - \gamma \phi(s_{K+1}) + \gamma (V'_K(s_{K+1}) - V_K(s_{K+1})) \\ &= \delta Q_{K+1}(s_K, a_K) - \gamma \phi(s_{K+1}) + \gamma \phi(s_{K+1}) \\ &= \delta Q_{K+1}(s_K, a_K). \end{aligned}$$

It follows that  $\Delta Q_{K+1}(s, a) = \Delta Q'_{K+1}(s, a) \forall s \forall a$ .  $\blacksquare$

*Remark 1:* If the Q-function is represented by a function approximator (as is typical for continuous action spaces), then Proposition 2 may not hold. This is because the Q-function in this scenario is updated using gradient descent, instead of Equation (12). Gradient descent is sensitive to initialization. Thus, different initial values will result in different updates of the Q-function.

## V. PBA FOR STOCHASTIC POLICY LEARNING

Although PBRS can preserve the optimality of policies in several settings, it suffers from the drawback of being unable to encode richer information, such as desired relations between states and actions. The authors of [7] proposed *potential-based advice* (PBA), a scheme that augments the potential function by including actions as an argument together with states. In this section, we show that while using PBA, recovering the optimal policy can be difficult if the optimal policy is stochastic. Then, we propose a novel way to impart prior information in order to learn a stochastic policy with PBA.

### A. Stochastic policy learning with PBA

Assume that we can compute  $Q_M^*(s, a)$ , the optimal value for state-action pair  $(s, a)$  in MDP  $M$ . The optimal stochastic policy for  $M$  is  $\pi_M^* = \arg \max_{\pi} \mathbb{E}_{\pi} [Q_M^*(s, a)]$ . From Equation (3), the optimal stochastic policy for the modified MDP  $M'$  that has its reward augmented with PBA is given by  $\pi_{M'}^* = \arg \max_{\pi} \mathbb{E}_{\pi} [Q_M^*(s, a) - \phi(s, a)]$ . Without loss of generality,  $\pi_M^* \neq \pi_{M'}^*$ . If the optimal policy is deterministic, then the policy for  $M$  can be recovered

easily from that for  $M'$  using Equation (4). However, when it is stochastic, we need to average over trajectories in the MDP, which makes it difficult to recover the optimal policy for  $M$  from that of  $M'$ .

In the sequel, we will propose a novel way to take advantage of PBA in the policy gradient framework in order to directly learn a stochastic policy.

### B. Imparting PBA in policy gradient

Let  $J_M(\theta)$  denote the value of a parameterized policy  $\pi_{\theta}$  in MDP  $M$ . That is,  $J_M(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} [\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t)]$ . Following the policy gradient theorem [28], and defining  $G(s_t, a_t) := \sum_{i=t}^{\infty} \gamma^{i-t} r_i$ , the gradient of  $J(\theta)$  with respect to the parameter  $\theta$  is given by:

$$\nabla_{\theta} J_M(\theta) = \mathbb{E}_{\tau \sim \pi_{\theta}} [G(s_t, a_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)]. \quad (13)$$

Then,  $\mathbb{E}_{\tau \sim \pi_{\theta}} [G(s_t, a_t)] = Q^{\pi_{\theta}}(s_t, a_t)$ .

REINFORCE [28] is a policy gradient method that uses Monte Carlo simulation to learn  $\theta$ , where the parameter update is performed only at the end of an episode (a trajectory of length  $T$ ). If we apply a look-ahead PBA scheme as in Equation (1) along with REINFORCE, then the total return from time  $t$  is given by:

$$\begin{aligned} G^a(s_t, a_t) &= \sum_{i=t}^{i=T} \gamma^{i-t} r_i + \gamma^{T-t} \phi(s_T, a_T) - \phi(s_t, a_t) \\ &= G(s_t, a_t) + \gamma^{T-t} \phi(s_T, a_T) - \phi(s_t, a_t). \end{aligned} \quad (14)$$

Notice that if  $G^a(s_t, a_t)$  is used in Equation (13) instead of  $G(s_t, a_t)$ , then the policy gradient is biased. One way to resolve the problem is to add the difference  $-\gamma^{T-t} \phi(s_T, a_T) + \phi(s_t, a_t)$  to  $G^a(s_t, a_t)$ . However, this makes the learning process identical to the original REINFORCE and PBA is not used. While using PBA in a policy gradient setup, it is important to add the term  $\phi(s, a)$  so that the policy gradient is unbiased, and also leverage the advantage that PBA offers during learning.

To apply PBA in policy gradient, we turn to temporal difference (TD) methods. TD methods update estimates of the accumulated return based in part on other learned estimates, before the end of an episode. A popular TD-based policy gradient method is the actor-critic framework [28]. In this setup, after performing action  $a_t$  at step  $t$ , the accumulated return  $G(s_t, a_t)$  is estimated by  $Q_M(s_t, a_t)$  which, in turn, is estimated by  $r_t + \gamma V_M(s_{t+1})$ . It should be noted that the estimates are unbiased.

When the reward is augmented with look-ahead PBA, the accumulated return is changed to  $Q_{M'}(s_t, a_t)$ , which is estimated by  $r_t + \gamma \phi(s_{t+1}, a_{t+1}) - \phi(s_t, a_t) + \gamma V_{M'}(s_{t+1})$ . From Equation (3), at steady state,  $Q_M(s_t, a_t) - Q_{M'}(s_t, a_t) = \phi(s_t, a_t)$ . Intuitively, to keep policy gradient unbiased when augmented with look-ahead PBA, we can add  $\phi(s_t, a_t)$  at each training step. In other words, we can use  $r_t + \gamma \phi(s_{t+1}, a_{t+1}) + \gamma V_{M'}(s_{t+1})$  as the estimated return. It should be noted that before the policy reaches steady state, adding  $\phi(s_t, a_t)$  at each time step will not cancel out the effect of PBA. This is unlike in REINFORCE,

where the addition of this term negates the effect of using PBA. In the advantage actor-critic, an advantage term is used instead of the Q-function in order to reduce the variance of the estimated policy gradient. In this case also, the potential term  $\phi(s_t, a_t)$  can be added in order to keep the policy gradient unbiased.

---

**Algorithm AC-PBA** : Actor-critic augmented with PBA

---

**Input:** Differentiable policy function  $\pi_\theta(a|s)$   
Differentiable value function  $V^\omega(s)$   
Potential-based advice  $\phi(s, a)$   
Maximum episode  $T_{max}$

**Initialization:**

policy parameter  $\theta$ , value parameter  $\omega$ , learning rate  $\alpha^\theta$  and  $\alpha^\omega$ , discount factor  $\gamma$ , episode counter  $T \leftarrow 0$

**repeat**

initialize state  $s_0$ ,  $t \leftarrow 0$

**repeat**

Sample action  $a_t \sim \pi_\theta(\cdot|s_t)$

Take action  $a_t$ , observe reward  $r_t$ , next state  $s_{t+1}$

$$R = \begin{cases} 0, & \text{if } s_{t+1} \text{ is a terminal state,} \\ V^\omega(s_{t+1}), & \text{otherwise.} \end{cases}$$

**if** use look-ahead advice **then**

$$\delta_t = r_t + \gamma\phi(s_{t+1}, a_{t+1}) - \phi(s_t, a_t) + \gamma R - V^\omega(s_t)$$

Update  $\theta \leftarrow \theta + \alpha^\theta (\delta_t + \phi(s_t, a_t)) \nabla_\theta \log \pi_\theta(a_t|s_t)$

**else**

$$\delta_t = r_t + \phi(s_t, a_t) - \gamma^{-1} \phi(s_{t-1}, a_{t-1}) + \gamma R - V^\omega(s_t)$$

Update  $\theta \leftarrow \theta + \alpha^\theta \delta_t \nabla_\theta \log \pi_\theta(a_t|s_t)$

**end if**

Update  $\omega \leftarrow \omega - \alpha^\omega \delta_t \nabla_\omega V^\omega(s_t)$

**until**  $s_{t+1}$  is a terminal state

$T \leftarrow T + 1$

**until**  $T > T_{max}$

---

A procedure for augmenting the advantage actor-critic with PBA is presented in Algorithm AC-PBA.  $\alpha^\theta$  and  $\alpha^\omega$  denote learning rates for the actor and critic respectively. When applying look-ahead PBA, at training step  $t$ , parameter  $\omega$  of the critic  $V^\omega(s)$  is updated as follows:

$$\delta_t^a = r_t + \gamma\phi(s_{t+1}, a_{t+1}) - \phi(s_t, a_t) + \gamma V^\omega(s_{t+1}) - V^\omega(s_t)$$

$$\omega = \omega - \alpha^\omega \delta_t^a \nabla_\omega V^\omega(s_t),$$

where  $\delta_t^a$  is the estimation error of the state value after receiving new reward  $[r_t + \gamma\phi(s_{t+1}, a_{t+1}) - \phi(s_t, a_t)]$  at step  $t$ . To ensure an unbiased estimate of the policy gradient, the potential term  $\phi(s_t, a_t)$  is added while updating  $\theta$  as:

$$\theta = \theta + \alpha^\theta (\delta_t^a + \phi(s_t, a_t)) \nabla_\theta \log \pi_\theta(a_t|s_t).$$

A similar method can be used when learning with look-back PBA. In this case, the critic and the policy parameter are updated as follows:

$$\delta_t^b = r_t + \phi(s_t, a_t) - \gamma^{-1} \phi(s_{t-1}, a_{t-1}) + \gamma V^\omega(s_{t+1}) - V^\omega(s_t)$$

$$\omega = \omega - \alpha^\omega \delta_t^b \nabla_\omega V^\omega(s_t),$$

$$\theta = \theta + \alpha (\delta_t^b + \gamma^{-1} \mathbb{E}[\phi(s_{t-1}, a_{t-1})|s_t]) \nabla_\theta \log \pi_\theta(a_t|s_t) \quad (15)$$

In fact, the potential term need not be added to ensure an unbiased estimate in this case. Then, the policy parameter update becomes:

$$\theta = \theta + \alpha \delta_t^b \nabla_\theta \log \pi_\theta(a_t|s_t), \quad (16)$$

which is exactly the policy update of the advantage actor-critic. This is formally stated in Proposition 3

*Proposition 3:* When the actor-critic is augmented with look-back PBA, Equations (15) and (16) are equal in the sense of expectation. That is

$$\mathbb{E}_{(s_t, a_t) \sim \rho^{\pi_\theta}} [(\delta_t^b + \gamma^{-1} \mathbb{E}[\phi(s_{t-1}, a_{t-1})|s_t]) \nabla_\theta \log \pi_\theta(a_t|s_t)]$$

$$= \mathbb{E}_{(s_t, a_t) \sim \rho^{\pi_\theta}} [\delta_t^b \nabla_\theta \log \pi_\theta(a_t|s_t)], \quad (17)$$

where  $\rho^{\pi_\theta}$  is the distribution induced by the policy  $\pi_\theta$ .

*Proof:* It is equivalent to show that:

$$\mathbb{E}_{(s_t, a_t) \sim \rho^{\pi_\theta}} [\mathbb{E}[\phi(s_{t-1}, a_{t-1})|s_t] \nabla_\theta \log \pi_\theta(a_t|s_t)] = 0. \quad (18)$$

The inner expectation  $\mathbb{E}[\phi(s_{t-1}, a_{t-1})|s_t]$  is a function of  $s_t$ , policy  $\pi_\theta$ , and transition probability  $\mathbb{T}$ . Denoting this expectation by  $f(s_t, \pi_\theta, \mathbb{T})$ , we obtain:

$$\mathbb{E}_{(s_t, a_t) \sim \rho^{\pi_\theta}} [f(s_t, \pi_\theta, \mathbb{T}) \nabla_\theta \log \pi_\theta(a_t|s_t)]$$

$$= \mathbb{E}_{s_t \sim \rho^{\pi_\theta}} [\mathbb{E}_{a_t \sim \pi_\theta} [f(s_t, \pi_\theta, \mathbb{T}) \nabla_\theta \log \pi_\theta(a_t|s_t)]]$$

$$= \mathbb{E}_{s_t \sim \rho^{\pi_\theta}} \left[ \int_{\mathcal{A}} \pi_\theta(a_t|s_t) f(s_t, \pi_\theta, \mathbb{T}) \frac{\nabla_\theta \pi_\theta(a_t|s_t)}{\pi_\theta(a_t|s_t)} da \right]$$

$$= \mathbb{E}_{s_t \sim \rho^{\pi_\theta}} \left[ f(s_t, \pi_\theta, \mathbb{T}) \nabla_\theta \int_{\mathcal{A}} \pi_\theta(a_t|s_t) da \right] = 0. \quad (19)$$

The last equality follows from the fact that the integral evaluates to 1, and its gradient is 0. ■

The main result of this paper presents guarantees on the convergence of Algorithm AC-PBA using the theory of ‘two time-scale stochastic analysis’ [32]. Assume that:

- **A1:** The value function  $V^\omega(s)$  belongs to a linear family. That is,  $V^\omega = \Phi\omega$ , where  $\Phi \in \mathbb{R}^{|S| \times k}$ ,  $k < S$  is a known full-rank feature matrix, and  $\omega \in \Omega \subseteq \mathbb{R}^k$ .
- **A2:** For the set of policies  $\{\pi_\theta, \theta \in \Theta \subseteq \mathbb{R}^d\}$ , there exists a constant  $C_\Theta$  such that  $\|\nabla_\theta \log \pi_\theta\|_2 \leq C_\Theta$ .
- **A3:** Learning rates of the actor and critic satisfy:  $\sum_t \alpha_t^\theta = \sum_t \alpha_t^\omega = \infty$ ,  $\sum_t [(\alpha_t^\theta)^2 + (\alpha_t^\omega)^2] < \infty$ ,  $\lim_{t \rightarrow \infty} \frac{\alpha_t^\theta}{\alpha_t^\omega} = 0$ .

For any probability measure  $\mu$  on a finite set  $\mathcal{M}$ , the  $\ell_2$ -norm of  $f$  with respect to  $\mu$  is given by  $\|f\|_\mu := [\int_{\mathcal{M}} |f(x)|^2 d\mu(x)]^{\frac{1}{2}}$ . Theorem 1 gives a bound on the error introduced as a result of approximating the value function  $V_{M'}$  with  $V_{M'}^\omega$  as in assumption **A1**. This error term is small if the family  $\Omega$  is rich. In fact, if the critic is updated in batches, a tighter bound can be achieved, as shown in Proposition 1 of [33]. Extending the result to the case of online updates is a subject of future work.

*Theorem 1:* Let  $\mathcal{E}(\theta) := \left\| V_{M'}^{\omega(\theta)}(s) - V_{M'}^{\pi_\theta}(s) \right\|_{\rho^{\pi_\theta}}$ . Then, for any limit point  $(\theta^*, \omega^*) := \lim_{T_{max} \rightarrow \infty} (\theta_{T_{max}}, \omega_{T_{max}})$  of Algorithm AC-PBA,  $\|\nabla_\theta \mathcal{J}_M(\theta^*)\|_2 \leq C \mathcal{E}(\theta^*)$ .

*Proof:* We consider only look-ahead PBA. The proof for look-back PBA follows similarly. Define  $F := F(s, a, s', a')$ . From assumption **A3**, the actor is updated

at a slower rate than the critic. This allows us to fix the actor to study the asymptotic behavior of the critic [34]. The update dynamics of the critic can be represented by:

$$\dot{\omega} = \mathbb{E}_{\rho^{\pi_\theta}} [\delta_\omega \nabla_\omega V_{M'}^\omega(s)], \quad (20)$$

where  $\delta_\omega = r(s, a) + \gamma \phi(s', a') - \phi(s, a) + \gamma V^\omega(s') - V^\omega(s)$  if look-ahead PBA is applied. When the critic is approximated by a linear function (assumption **A1**),  $\omega$  will converge to  $\omega(\theta)$ , an asymptotically stable equilibrium of Equation (20). The update of the actor is then:

$$\dot{\theta} = \mathbb{E}_{\rho^{\pi_\theta}} [\nabla_\theta \log \pi_\theta(a|s)(r(s, a) + F + \gamma V_{M'}^{\omega(\theta)}(s') + \phi(s, a))]. \quad (21)$$

Let  $\Theta_s$  denote the set of asymptotic stable equilibria in Equation (21). Any  $\theta \in \Theta_s$  will satisfy  $\dot{\theta} = 0$  in Equation (21). Then,  $\{(\theta_t, \omega_t)\}_{t>0}$  will converge to  $\{(\theta, \omega(\theta)) : \theta \in \Theta_s\}$ .

Now, consider the evaluation of  $\pi_\theta$ ,  $\theta \in \Theta_s$ , in the original MDP  $M$ . We obtain the following equations:

$$\begin{aligned} \nabla_\theta J_M(\theta) &= \mathbb{E}_{\rho^{\pi_\theta}} [\nabla_\theta \log \pi_\theta(a|s) Q_M^{\pi_\theta}(s, a)] \\ &= \mathbb{E}_{\rho^{\pi_\theta}} [\nabla_\theta \log \pi_\theta(a|s) (Q_{M'}^{\pi_\theta}(s, a) + \phi(s, a))] \\ &= \mathbb{E}_{\rho^{\pi_\theta}} [\nabla_\theta \log \pi_\theta(a|s) (r(s, a) + F + \gamma V_{M'}^{\pi_\theta}(s') + \phi(s, a))]. \end{aligned} \quad (22)$$

Subtracting Equation (21) from Equation (22), and applying the Cauchy-Schwarz inequality to the result yields:

$$\begin{aligned} \nabla_\theta J_M(\theta) &= \gamma \mathbb{E}_{\rho^{\pi_\theta}} [\nabla_\theta \log \pi_\theta(a|s) (V_{M'}^{\omega(\theta)}(s') - V_{M'}^{\pi_\theta}(s))] \\ \therefore \|\nabla_\theta J_M(\theta)\|_2 &\leq \gamma \|\nabla_\theta \log \pi_\theta(a|s)\|_{\rho^{\pi_\theta}} \left\| V_{M'}^{\omega(\theta)}(s) - V_{M'}^{\pi_\theta}(s) \right\|_{\rho^{\pi_\theta}}. \end{aligned}$$

The result follows by applying assumption **A2**. ■

*Remark 2:* Look-back PBA could result in better performance compared to look-ahead PBA since look-back PBA does not involve estimating a future action.

## VI. EXPERIMENTS

Our experiments seek to compare the performance of an actor-critic architecture augmented with PBA and with PBRs with the ‘vanilla’ advantage actor-critic (A2C). We consider two setups. The first is a *Puddle-Jump Gridworld* [35], where the state and action spaces are discrete. The second environment we study is a continuous state and action space *mountain car* [26].

In each experiment, we compare the rewards received by the agent when it uses the following schemes: *i*): ‘vanilla’ (A2C); *ii*): A2C augmented with PBRs; *iii*): A2C with look-ahead PBA; *iv*): A2C with look-back PBA.

### A. Puddle-Jump Gridworld

Figure 1 depicts the *Puddle-jump gridworld* environment as a 10x10 grid. The state space is  $s = (x, y)$  denoting the position of the agent in the grid, where  $x, y \in \{0, 1, \dots, 9\}$ . The goal of the agent is to navigate from the start state  $S = (0, 0)$  to the goal  $G = (9, 9)$ . At each step, the agent can choose from actions in the set  $A = \{up, down, left, right, jump\}$ . There is a *puddle* along row 2 which the agent should jump over. Further, the states (9, 8) and (8, 9) (blue squares in Figure 1) are

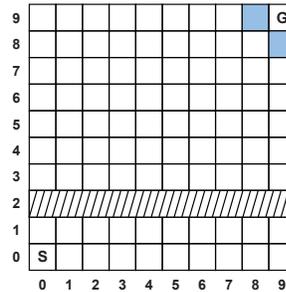


Fig. 1: Schematic of the puddle-jump gridworld. The state of the agent is its position  $(x, y)$ . The shaded row (row 2) represents the puddle the agent should jump over. The two blue grids denote states that are indistinguishable to the agent. The agent can choose an action from the set  $\{up, down, left, right, jump\}$  at each step.

indistinguishable to the agent. As a result, any optimal policy for the agent is a stochastic policy.

If the *jump* action is chosen in rows 3 or 1, the agent will land on the other side of the puddle with probability  $p_j$ , and remain in the same state otherwise. This action chosen in other rows will keep the agent in its current state. Any action that will move the agent off the grid will keep its state unchanged. The agent receives a reward of  $-0.05$  for each action, and  $+1000$  for reaching  $G$ .

When using PBRs, we set  $\phi^{PBRs}(s) := u_0$  for states in rows 0 and 1, and  $\phi^{PBRs}(s) := u_1$  for all other states. We need  $u_1 > u_0$  to encourage the agent to jump over the puddle. Unlike in PBRs, PBA can provide the agent with more information about the actions it can take. We set  $\phi^{PBA}(s, a)$  to a ‘large’ value if action  $a$  at state  $s$  results in the agent moving closer to the goal according to the  $\ell_1$  norm,  $(|G-x| + |G-y|)$ . We additionally stipulate that  $\frac{1}{|A|} \sum_{a \in A} \phi^{PBA}(s, a) = \phi^{PBRs}(s)$ . That is, the state potential of PBA is the same as the state potential of PBRs under a uniform distribution over the actions. This is to ensure a fair comparison between PBRs and PBA.

In our experiment, we set the discount factor  $\gamma = 1$ . Since the dimensions of the state and action spaces is not large, we do not use a function approximator for the policy  $\pi$ . A parameter  $\theta_{s,a}$  is associated to each state-action pair, and the policy is computed as:  $\pi_\theta(a|s) = \frac{\exp(\theta_{s,a})}{\sum_{a \in A} \exp(\theta_{s,a})}$ . We fix  $\alpha^\omega = 0.001$ , and  $\alpha^\theta = 0.2$  for all cases.

From Figure 2, we observe that the look-back PBA scheme performs the best, in that the agent converges to the goal in **five times** fewer episodes (25 vs. 125 episodes) than A2C without advice. When A2C is augmented with PBRs, convergence to the goal is slightly faster than without any reward shaping. When augmented with look-ahead PBA, in the first few episodes, the reward increases faster than in the case of A2C augmented with PBRs. However, this slows down after the early training stages and the policy converges to the goal in about the same number of episodes as a policy trained without advice. A reason for this could be that

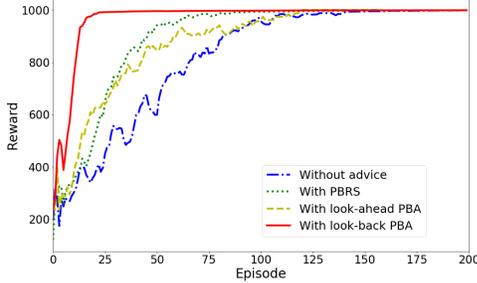


Fig. 2: Average rewards in puddle-jump gridworld when jump success probability  $p_j = 0.2$ . The baseline is the advantage actor-critic without advice.

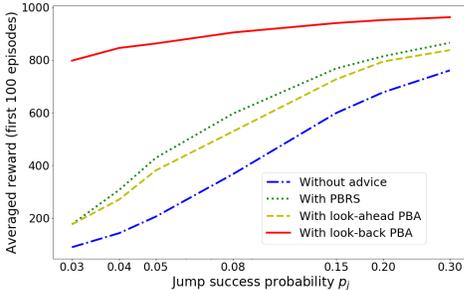


Fig. 3: Average reward for the first 100 episodes with respect to the jump success probability  $p_j$ .

during later stages of training, a look-ahead PBA scheme might advise an agent with ‘bad’ actions, leading to bad policies, thereby impeding the progress of learning. For example, an action  $a_t$  might be a good choice at state  $s_t$ , but the look-ahead PBA scheme might indicate that  $a_t$  is bad, due to a poor estimate of the future action  $a_{t+1}$ .

A smaller jump success probability  $p_j$  is an indication that it is more difficult for the agent to reach the goal state  $G$ . Figure 3 shows that look-back PBA results in the highest reward for a more difficult task (lower  $p_j$ ), when compared with the other reward shaping schemes.

### B. Continuous Mountain Car

In the mountain car (MC) environment, an under powered car in a valley has to drive up a steep hill to reach the goal. In order to achieve this, the car should learn how to accumulate momentum. A schematic for this environment is shown in Figure 4. This MC environment has continuous state and action spaces. The state  $s = (p, v)$  denotes position  $p \in [-1.2, 0.6]$  and velocity  $v \in [-0.07, 0.07]$ . The action  $a \in [-1, +1]$ . The continuous action space makes it difficult to use classic value-based methods, such as Q-learning and Sarsa-learning. The reward provided by the environment depends on the action and whether the car reaches the goal. Specifically, once the car reaches the goal it receives +100, and before

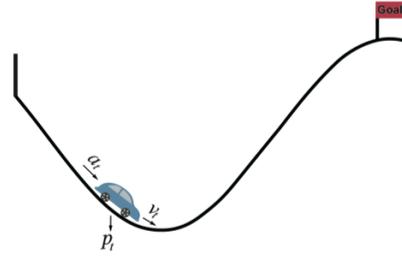


Fig. 4: Schematic of the mountain-car environment. The agent’s state is represented by its position  $p_t$  (along the  $x$ -coordinate) and velocity  $v_t$ . The action  $a_t$  is a force applied to the car. The goal is marked as a flag.

that, the reward at time  $t$  is  $-|a_t|^2$ . This reward structure therefore discourages the waste of energy. This acts as a barrier for learning, because there appears to be a sub-optimal solution where the agent remains at the bottom of the valley. Moreover, the reward for reaching the goal is significantly delayed, which makes it difficult for the conventional actor-critic algorithm to learn a good policy.

One choice of a potential function while using PBRs in this environment is  $\phi^{PBRs}(s_t) := p_t + 2$ , where the offset is so that the potential is positive. An interpretation of this scheme is: ‘state value is larger when the car is horizontally closer to the goal.’ The PBA scheme we use for this environment encourages the accumulation of momentum by the car– the direction of the action is encouraged to be the same as the current direction of the car’s velocity. In the meanwhile, we discourage inaction. Mathematically, the potential advice function has a larger value if  $a_t \neq 0$ . We let  $\phi^{PBA}(s_t, a_t) = 1$ , if  $a_t v_t > 0$ , and  $\phi^{PBA}(s_t, a_t) = 0$ , otherwise.

In our experiments, we set  $\gamma = 0.99$ . To deal with the continuous state space, we use a neural network (NN) as a function approximator. The policy distribution  $\pi_\theta(a|s)$  is approximated by a normal distribution, the mean and variance of which are the outputs of the NN. The value function is also represented by an NN. We set  $a^\theta = 1 \times 10^{-5}$  and  $a^\omega = 5.6 \times 10^{-4}$ , and use Adam [36] to update the NN parameters. The results we report are averaged over 10 different environment seeds.

No advice	PBRs	Look-ahead PBA	Look-back PBA
10%	20%	40%	100%

TABLE I: Percentage of trials where policy converges correctly in continuous mountain car problem.

Our experiments indicate that the policy makes the agent converge to one of two points: the goal, or remain stationary at the bottom of the valley. The percentage of solutions that converge to the goal is shown in Table I. From Figure 5 and Table I, when learning with the vanilla A2C, the agent is able to reach the goal only in 10% of the trials (out of 10 trials), and was stuck at the sub-optimal solution for the remaining trials. With

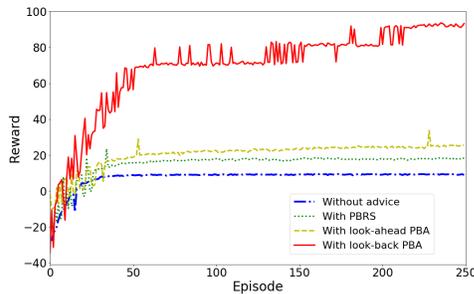


Fig. 5: Average rewards for continuous mountain car problem (averaged over 10 different environment random seeds). The baseline is the A2C without advice.

PBRs, the agent could converge correctly in only 20% of the trials. This is because the agent might have to take an action that moves it away from the goal in order to accumulate momentum. However, the potential function  $\phi^{PBRs}(\cdot)$  discourages such actions. In comparison, the average reward when using look-ahead PBA is slightly higher, but the agent is able to reach the goal in only 40% of the trials. Similar to the gridworld setup, look-back PBA performs the best, where the agent is able to reach the goal in 100% of the trials.

## VII. CONCLUSION

This paper presented a framework for augmenting the reward received by an RL agent with PBRs and with PBA. Different from prior work, we demonstrated that our approach can be used in environments with continuous states and actions, and when the optimal policy is stochastic. We presented guarantees on the convergence of an algorithm that augments an A2C architecture with these schemes. Our experiments indicated that these schemes allowed the agent to achieve higher average rewards, and learn an optimal policy faster. Future work will focus on establishing tighter bounds for Theorem 1, and extending our approach to the average reward case.

## REFERENCES

- [1] R. Hafner and M. Riedmiller, “Reinforcement learning in feedback control,” *Machine Learning*, vol. 84, pp. 137–169, 2011.
- [2] T. P. Lillicrap *et al.*, “Continuous control with deep reinforcement learning,” in *International Conference on Learning and Representations*, 2016.
- [3] V. Mnih *et al.*, “Human-level control through deep reinforcement learning,” *Nature*, vol. 518, no. 7540, 2015.
- [4] D. Silver *et al.*, “Mastering the game of Go with deep neural networks and tree search,” *Nature*, vol. 529, no. 7587, 2016.
- [5] J. Randløv and P. Alstrøm, “Learning to drive a bicycle using reinforcement learning and shaping,” in *International Conference on Machine Learning*, 1998.
- [6] A. Y. Ng, D. Harada, and S. Russell, “Policy invariance under reward transformations: Theory and application to reward shaping,” in *International Conference on Machine Learning*, 1999.
- [7] E. Wiewiora, G. W. Cottrell, and C. Elkan, “Principled methods for advising reinforcement learning agents,” in *International Conference on Machine Learning*, 2003, pp. 792–799.
- [8] S. M. Devlin and D. Kudenko, “Dynamic potential-based reward shaping,” in *Autonomous Agents and Multiagent Systems*, 2012, pp. 433–440.

- [9] A. L. Thomaz and C. Breazeal, “Reinforcement learning with human teachers: Evidence of feedback and guidance with implications for learning performance,” in *AAAI*, 2006, pp. 1000–1005.
- [10] W. B. Knox and P. Stone, “Combining manual feedback with subsequent MDP reward signals for reinforcement learning,” in *Autonomous Agents and Multiagent Systems*, 2010, pp. 5–12.
- [11] D. Pathak, P. Agrawal, A. A. Efros, and T. Darrell, “Curiosity-driven exploration by self-supervised prediction,” in *International Conference on Machine Learning*, 2017.
- [12] H. Tang *et al.*, “# Exploration: A study of count-based exploration for deep reinforcement learning,” in *Advances in Neural Information Processing Systems*, 2017.
- [13] R. J. Williams and J. Peng, “Function optimization using connectionist reinforcement learning algorithms,” *Connection Science*, vol. 3, no. 3, pp. 241–268, 1991.
- [14] V. Mnih *et al.*, “Asynchronous methods for deep reinforcement learning,” in *International Conference on Machine Learning*, 2016.
- [15] S. Levine and V. Koltun, “Guided policy search,” in *International Conference on Machine Learning*, 2013, pp. 1–9.
- [16] S. Levine, C. Finn, T. Darrell, and P. Abbeel, “End-to-end training of deep visuomotor policies,” *The Journal of Machine Learning Research*, vol. 17, no. 1, pp. 1334–1373, 2016.
- [17] E. Wiewiora, “Potential-based shaping and Q-value initialization are equivalent,” *Journal of Artificial Intelligence Research*, pp. 205–208, 2003.
- [18] A. Harutyunyan, S. Devlin, P. Vrancx, and A. Nowé, “Expressing arbitrary reward functions as potential-based advice,” in *AAAI*, 2015, pp. 2652–2658.
- [19] M. Li, T. Brys, and D. Kudenko, “Introspective reinforcement learning and learning from demonstration,” in *Autonomous Agents and MultiAgent Systems*, 2018, pp. 1992–1994.
- [20] J. Asmuth, M. L. Littman, and R. Zinkov, “Potential-based shaping in model-based RL,” in *AAAI*, 2008, pp. 604–609.
- [21] M. Grześ, “Reward shaping in episodic reinforcement learning,” in *Autonomous Agents and MultiAgent Systems*, 2017, pp. 565–573.
- [22] A. Eck, L.-K. Soh, S. Devlin, and D. Kudenko, “Potential-based reward shaping for finite horizon online POMDP planning,” *Autonomous Agents and Multi-Agent Systems*, vol. 30, no. 3, 2016.
- [23] S. J. Bradtke, “RL applied to linear quadratic regulation,” in *Advances in Neural Information Processing Systems*, 1993.
- [24] M. Fazal, R. Ge, S. Kakade, and M. Mesbahi, “Global convergence of policy gradient methods for the linear quadratic regulator,” in *International Conference on Machine Learning*, 2018.
- [25] L. Busoniu, T. de Bruin, D. Tolić, J. Kober, and I. Palunko, “Reinforcement learning for control: Performance, stability, and deep approximators,” *Annual Reviews in Control*, 2018.
- [26] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba, “Openai Gym,” *arXiv:1606.01540*, 2016.
- [27] M. L. Puterman, *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, 2014.
- [28] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. MIT press, 2018.
- [29] T. Haarnoja, H. Tang, P. Abbeel, and S. Levine, “Reinforcement Learning with Deep Energy-Based Policies,” in *International Conference on Machine Learning*, 2017, pp. 1352–1361.
- [30] R. S. Sutton, D. A. McAllester, S. P. Singh, and Y. Mansour, “Policy gradient methods for reinforcement learning with function approximation,” in *Advances in Neural Information Processing Systems*, 2000, pp. 1057–1063.
- [31] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine, “Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor,” *arXiv:1801.01290*, 2018.
- [32] V. Borkar and S. Meyn, “The ODE method for convergence of stochastic approximation and reinforcement learning,” *SIAM Journal on Control and Optimization*, vol. 38, no. 2, 2000.
- [33] Z. Yang, K. Zhang, M. Hong, and T. Başar, “A finite sample analysis of the actor-critic algorithm,” in *IEEE Conference on Decision and Control (CDC)*, 2018, pp. 2759–2764.
- [34] S. Bhatnagar, R. S. Sutton, M. Ghavamzadeh, and M. Lee, “Natural actor-critic algorithms,” *Automatica*, vol. 45, no. 11, 2009.
- [35] O. Marom and B. Rosman, “Belief reward shaping in reinforcement learning,” in *AAAI*, 2018, pp. 3762–3769.
- [36] D. P. Kingma and J. Ba., “Adam: A method for stochastic optimization,” *arXiv:1412.6980*, 2014.