

Natural actor critic algorithms<sup>1</sup>Shalabh Bhatnagar<sup>a</sup>, Richard S. Sutton<sup>b</sup>, Mohammad Ghavamzadeh<sup>c</sup>, Mark Lee<sup>b</sup><sup>a</sup> Department of Computer Science and Automation, Indian Institute of Science, Bangalore 560 012, India<sup>b</sup> The RLAI Laboratory, Department of Computing Science, University of Alberta, Edmonton, Alberta, Canada T6G 2E8<sup>c</sup> INRIA Lille - Nord Europe, Team Sequel, France

## article info

## Article history:

Received 23 May 2007

Received in revised form

1 January 2009

Accepted 3 July 2009

Available online xxxx

## Keywords:

Actor critic reinforcement learning algorithms

Policy-gradient methods

Approximate dynamic programming

Function approximation

Two-timescale stochastic approximation

Temporal difference learning

Natural gradient

## abstract

We present four new reinforcement learning algorithms based on actor critic, natural-gradient and function-approximation ideas, and we provide their convergence proofs. Actor critic reinforcement learning methods are online approximations to policy iteration in which the value-function parameters are estimated using temporal difference learning and the policy parameters are updated by stochastic gradient descent. Methods based on policy gradients in this way are of special interest because of their compatibility with function-approximation methods, which are needed to handle large or infinite state spaces. The use of temporal difference learning in this way is of special interest because in many applications it dramatically reduces the variance of the gradient estimates. The use of the natural gradient is of interest because it can produce better conditioned parameterizations and has been shown to further reduce variance in some cases. Our results extend prior two-timescale convergence results for actor critic methods by Konda and Tsitsiklis by using temporal difference learning in the actor and by incorporating natural gradients. Our results extend prior empirical studies of natural actor critic methods by Peters, Vijayakumar and Schaal by providing the first convergence proofs and the first fully incremental algorithms.

© 2009 Elsevier Ltd. All rights reserved.

## 1. Introduction

Many problems of scientific and economic importance are optimal sequential decision problems and as such can be formulated as Markov decision processes (MDPs) (Bertsekas & Tsitsiklis, 1996; Rust, 1996; White, 1993). In some cases, MDPs can be solved analytically, and in many cases they can be solved iteratively by dynamic programming or linear programming. However, in other cases these methods cannot be applied either because the state space is too large, a system model is available only as a simulator, or no system model is available. It is in these cases that the techniques and algorithms of reinforcement learning (RL) may be helpful.

Reinforcement learning (Bertsekas & Tsitsiklis, 1996; Sutton & Barto, 1998) can be viewed as a broad class of sample-based methods for solving MDPs. In place of a model, these methods use sample trajectories of the system and the controller interacting, such as could be obtained from a simulation. It is not unusual in

practical applications for such a simulator to be available when an explicit transition-probability model of the sort suitable for use by dynamic or linear programming is not (Crites & Barto, 1998; Tesauro, 1995). Reinforcement learning methods can also be used with no model at all, by obtaining sample trajectories by direct interaction with the system (Kohl & Stone, 2004; Ng et al., 2004).

One of the biggest challenges to solve MDPs with conventional methods is handling large state (and action) spaces. This is sometimes known as the "curse of dimensionality" because of the tendency of the size of a state space to grow exponentially with the number of its dimensions. The computational effort required to solve an MDP thus increases exponentially with the dimension and cardinality of the state space. A natural and venerable way of addressing the curse is to approximate the value function and policy parametrically with a number of parameters much smaller than the size of the state space (Bellman & Dreyfus, 1959). However a straightforward application of such function-approximation methods to dynamic programming has not proved effective on large problems. Some work with RL and function approximation has also run into problems of convergence and instability (Baird, 1995; Boyan & Moore, 1995), but about a decade ago it was established that if trajectories were sampled according to their distribution under the target policy (the on-policy distribution) then convergence could be assured for linear feature-based function approximators (Sutton, 1996; Tadic, 2001; Tsitsiklis & Van Roy, 1997).

<sup>1</sup> This paper was not presented at any IFAC meeting. This paper was recommended for publication in revised form by Associate Editor Giuseppe De Nicolao under the direction of Editor Ian R. Petersen.

Corresponding author. Tel.: +91 80 2293 2987; fax: +91 80 2360 2911.

E-mail addresses: [shalabh@csa.iisc.ernet.in](mailto:shalabh@csa.iisc.ernet.in) (S. Bhatnagar), [sutton@cs.ualberta.ca](mailto:sutton@cs.ualberta.ca) (R.S. Sutton), [mohammad.ghavamzadeh@inria.fr](mailto:mohammad.ghavamzadeh@inria.fr) (M. Ghavamzadeh), [mlee@cs.ualberta.ca](mailto:mlee@cs.ualberta.ca) (M. Lee).

Reinforcement learning's most impressive successes have in fact been on problems with extremely large state spaces that could not have been solved without function approximation (Crites & Barto, 1998; Ng et al., 2004; Tesauro, 1995). The ability of sample-based methods to use function approximation effectively is one of the most important reasons for interest in RL within the engineering disciplines.

Policy-gradient methods are some of the simplest RL algorithms and provide both a good illustration of RL and a foundation for the actor–critic methods that are the primary focus of this paper. In policy-gradient methods, the policy is taken to be an arbitrary differentiable function of a parameter vector  $\theta \in \mathbb{R}^d$ . Given some performance measure  $J: \mathbb{R}^d \rightarrow \mathbb{R}$ , we would like to update the policy parameter in the direction of the gradient:

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} J(\theta) \quad (1)$$

The gradient is not directly available of course, but sample trajectories can be used to construct unbiased estimators of it, estimators that can be used in a stochastic approximation of the actual gradient. This is the basic idea behind all policy-gradient methods (Aleksandrov, Sysoyev, & Shemeneva, 1968; Baxter & Bartlett, 2001; Bhatnagar, 2005, 2007; Ghavamzadeh & Mahadevan, 2003; Ghavamzadeh & Engel, 2007a,b; Glynn, 1990; Konda & Tsitsiklis, 2003; Marbach & Tsitsiklis, 2001; Peters & Schaal, 2008; Sutton, McAllester, Singh, & Mansour, 2000; Williams, 1992). Theoretical analysis and empirical evaluations have highlighted a major shortcoming of these algorithms, namely, the high variance of their gradient estimates, and thus the slow convergence and sample inefficiency.

One possible solution to this problem, proposed by Kakade (2002) and then refined and extended by Bagnell and Schneider (2003) and by Peters, Vijayakumar, and Schaal (2003), is based on the idea of *natural* gradients previously developed for supervised learning by Amari (1998). In the application to RL, the policy gradient in (1) is replaced with a natural version. This is motivated by the intuition that a change in the policy parameterization should not influence the result of the policy update. In terms of the policy update rule (1), the move to natural gradient amounts to linearly transforming the gradient using the inverse Fisher information matrix of the policy. In empirical evaluations, natural policy gradient has sometimes been shown to outperform conventional policy-gradient methods (Bagnell & Schneider, 2003; Kakade, 2002; Peters et al., 2003; Richter, Aberdeen, & Yu, 2007). Moreover, the use of natural gradients can lead to simpler, and in some cases, more computationally efficient algorithms. Three of the four algorithms we introduce in this paper incorporate natural gradients.

In this paper we focus on a sub-class of policy-gradient methods known as actor–critic algorithms. These methods can be thought of as reinforcement learning analogs of dynamic programming's policy iteration method. Actor–critic methods are based on the simultaneous online estimation of the parameters of two structures, called the *actor* and the *critic*. The actor corresponds to a conventional action–selection policy, mapping states to actions in a probabilistic manner. The critic corresponds to a conventional state-value function, mapping states to expected cumulative future reward. Thus, the critic addresses a problem of prediction, whereas the actor is concerned with control. These problems are separable, but are solved simultaneously to find an optimal policy. A variety of methods can be used to solve the prediction problem, but the ones that have proved most effective are those based on some form of temporal difference (TD) learning (Sutton, 1988), in which estimates are updated on the basis of other estimates. Such “bootstrapping methods” (Sutton & Barto, 1998) can be viewed as a way of accelerating learning by trading bias for variance.

Actor–critic methods were among the earliest to be investigated in reinforcement learning (Barto, Sutton, & Anderson, 1983;

Sutton, 1984). They were largely supplanted in the 1990s by methods that estimate action-value functions (mappings from states and actions to the subsequent expected return) that are then used directly to select actions without constructing an explicit policy structure. The action-value approach was initially appealing because of its simplicity, but theoretical complications arose when it was combined with function approximation: these methods do not converge in the normal sense, but rather may “chatter” in the neighborhood of a good solution (Gordon, 1995). These complications lead to renewed interest in policy-gradient methods. Policy-gradient methods without bootstrapping can easily be proved convergent, but can suffer from high variance resulting in slow convergence as mentioned above, motivating their combination with bootstrapping temporal difference methods as in actor–critic algorithms.

In this paper we introduce four novel actor–critic algorithms along these lines. For all four methods we prove convergence of the parameters of the policy and state-value function to a small neighborhood of the set of local maxima of the average reward when the TD error inherent in the function approximation is small. Our results are an extension of our prior work (Bhatnagar, Sutton, Ghavamzadeh, & Lee, 2008), and of prior work on the convergence of two-timescale stochastic approximation recursions (Abdulla & Bhatnagar, 2007; Bhatnagar & Kumar, 2004; Konda & Borkar, 1999; Konda & Tsitsiklis, 2003). That work had previously shown convergence to a locally optimal policy for several non-bootstrapping algorithms with or without function approximation. Convergence of general two-timescale stochastic approximation algorithms has been shown under some assumptions in Borkar (1997). Konda and Tsitsiklis (2003) have shown convergence for an actor–critic algorithm that uses bootstrapping in the critic, but our results are the first to prove convergence when the actor is bootstrapping as well. Our results also extend prior two-timescale results by incorporating natural gradients. Our results and algorithms differ in a number of other, smaller ways from those of Konda and Tsitsiklis; we detail these in Section 6 after the analysis has been presented.

Two other aspects of the theoretical results presented here should be mentioned at the outset. First, one of the issues that arises in policy-gradient methods is the selection of a baseline reward level. In contrast to previous work, we show that, in an actor–critic setting when compatible features are used, the baseline that minimizes the estimator variance for any given policy is in fact the state-value function. Second, for the case of a fixed policy we use a recent result by Borkar and Meyn (2000) to provide an alternative, simpler proof of convergence (cf. Tsitsiklis & Van Roy, 1997; Tsitsiklis & Van Roy, 1999) in the Euclidean norm of TD recursions.

In this paper we do not explicitly consider the treatment of eligibility traces ( $\gamma > 0$  in TD(  $\gamma$  ) (Sutton, 1988)), which have been shown to improve performance in cases of function approximation or partial observability, but we believe the extension of all of our results to general  $\gamma$  would be straightforward. Less clear is how or whether our results could be extended to least-squares TD methods (Boyan, 1999; Bradtke & Barto, 1996; Farahmand, Ghavamzadeh, Szepesvári, & Mannor, 2009; Lagoudakis & Parr, 2003). It is not clear how to satisfactorily incorporate these methods in a context in which the policy is changing. Our proof techniques do not immediately extend to this case and we leave it for future work. We do consider the use of approximate advantages as in the works of Baird (1993) and of Peters and Schaal (2008). Because of space limitations, we do not present empirical results obtained from our algorithms in this paper but these can be seen in Section 8 of our technical report (Bhatnagar, Sutton, Ghavamzadeh, & Lee, 2009).

The rest of the paper is organized as follows. In Section 2 we present our RL framework and provide an overview of policy-gradient methods. In Section 3 we discuss policy-gradient methods with function approximation and present some preliminary results. We show here in particular that the minimum variance baseline for the action-value function corresponds to the state-value function and obtain a form of bias in gradient estimates that results from the use of function approximation. Our four actor-critic algorithms and their convergence analysis are presented in Sections 4 and 5, respectively. In Section 6 we discuss the relationship of our algorithms to the actor-critic algorithm of Konda and Tsitsiklis (2003) and to the natural actor-critic algorithm of Peters et al. (2003). Section 7 contains concluding remarks.

## 2. The policy-gradient framework

We consider the standard reinforcement learning framework (e.g., see Bertsekas & Tsitsiklis, 1996; Sutton & Barto, 1998) in which a learning agent interacts with a stochastic environment. The overall model we consider is that of a discrete time Markov decision process (MDP) with finite numbers of states and actions, and bounded rewards. We allow  $S$  and  $A$  to respectively denote the state and action spaces of this MDP. For simplicity, we assume that  $S$  is the set  $S = \{1, \dots, n\}$ . We denote by  $s_t$ ,  $a_t$ , and  $r_t$ , the state, action, and reward at time  $t$ , respectively. We assume that reward is stochastic, real-valued and uniformly bounded. For simplicity and ease of notation, we assume that all actions in  $A$  are feasible in each state. The state transition probabilities for the environment will be characterized by  $P(s'; a; s) = \Pr\{s_{t+1} = s' | s_t = s, a_t = a\}$ ,  $s' \in S$ ,  $s \in S$ ,  $a \in A$ . Further, the single-stage expected reward when action  $a$  is taken in state  $s$  will be denoted  $R(s; a) = \mathbb{E}r_{t+1} | s_t = s, a_t = a$ .

An admissible policy  $\pi$  is a decision rule that is described by a sequence of functions  $\pi = \{\pi_0, \pi_1, \dots\}$  such that each  $\pi_t : S \times \{0, \dots, t\} \rightarrow A$ , with action  $\pi_t(s)$  taken in state  $s$  at instant  $t \geq 0$ . A stationary policy is a time invariant decision rule, i.e., one for which  $\pi_t = \pi$ ,  $\forall t \geq 0$ , for some  $\pi : S \rightarrow A$ . Most often, one refers to the function  $\pi$  itself as the stationary policy. A stationary randomized policy  $\pi$  that we refer to as simply a randomized policy is specified via a probability distribution  $\pi(s; a)$  over  $A$ , for  $s \in S$ . Under the long-run average reward setting considered in this paper, it can be shown that a stationary optimal policy exists e.g., (see Puterman, 1994). Note that any stationary policy is trivially a randomized policy as well. We motivate the following discussion from the viewpoint of randomized policies as we consider a parameterized class of these in this paper. From now on, for simplicity, we shall refer to a randomized policy as a policy. For a given policy, the sequence of states produced by the MDP is a Markov chain. Throughout the paper we assume

**(A1)** Under any policy  $\pi$ , the Markov chain resulting from the MDP is irreducible and aperiodic.

Let  $d(s)$  denote the stationary probability of the Markov chain being in state  $s \in S$ , and let  $d = \{d(s) : s \in S\}$ . Our aim is to find a policy  $\pi$  that maximizes the long-run average reward  $J(\pi)$  given by

$$J(\pi) = \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \sum_{t=0}^{T-1} r_{t+1} \quad (2)$$

$$= \sum_{s \in S} d(s) \sum_{a \in A} \pi(s; a) R(s; a)$$

The limit in (2) is well defined by (A1). Let  $\pi^{opt}$  denote an optimal policy  $\pi^{opt} = \arg \max_{\pi} J(\pi)$ . Further, we shall denote by  $Q(s; a)$ ,

the expected differential reward associated with a state-action pair  $(s; a)$ , given policy  $\pi$ , that is defined by  $Q(s; a) = R(s; a) - V(s)$ ,

$$Q(s; a) = R(s; a) - V(s) = \mathbb{E}r_{t+1} | s_t = s, a_t = a - V(s)$$

Likewise, we denote by  $V(s)$ , the expected differential reward associated with a state  $s$  when actions are selected according to policy  $\pi$ . Here  $V(s) = \sum_{a \in A} \pi(s; a) Q(s; a)$ . The Poisson equation under policy  $\pi$  is given by

$$J(\pi) - V(s) = \sum_{a \in A} \pi(s; a) Q(s; a) \quad (3)$$

$s \in S$  (Puterman, 1994). In policy-gradient methods, we define a class of parameterized randomized policies  $\pi_\theta(s; a)$ ,  $\theta \in \mathbb{R}^{d_\theta}$ , estimate the gradient of the average reward with respect to the policy parameters  $\theta$  from the observed states, actions, and rewards, and then improve the policy by adjusting its parameters in the direction of an estimate of the gradient of  $J$  with respect to  $\theta$ . Because in this setting a policy  $\pi$  is represented by its parameters  $\theta$ ,  $J$  can be viewed as a function of  $\theta$  and by abuse of notation, we let  $J(\theta)$  denote the long-run average reward when the parameter is  $\theta$ . In what follows, we shall interchangeably use  $J(\theta)$  or  $J(\pi)$  to denote the long-run average reward when the policy  $\pi$  or its associated parameter  $\theta$  are to be emphasized. We also drop  $\pi$  from  $\pi(s; a)$  and simply denote this quantity as  $\pi(s; a)$ . The optimum parameter can now be obtained as  $\theta^{opt} = \arg \max_{\theta} J(\theta)$ . The following assumption is a standard requirement in policy-gradient methods.

**(A2)** For any state-action pair  $(s; a)$ ,  $Q(s; a)$  is continuously differentiable in the parameter  $\theta$ .

Previous works (Baxter & Bartlett, 2001; Konda & Tsitsiklis, 2003; Marbach & Tsitsiklis, 2001; Sutton et al., 2000) have shown that the gradient of the average reward for parameterized policies that satisfy (A1) and (A2) is given by<sup>1</sup>

$$\nabla_{\theta} J(\theta) = \sum_{s \in S} d(s) \sum_{a \in A} \pi(s; a) \nabla_{\theta} Q(s; a) \quad (4)$$

For the case of Markov processes with a parameterized infinitesimal generator, a similar expression was obtained by Cao and Chen (1997). Observe that if  $b(s)$  is any given function of  $s$  (also called a baseline), then

$$\nabla_{\theta} J(\theta) = \sum_{s \in S} d(s) \sum_{a \in A} \pi(s; a) \nabla_{\theta} Q(s; a) = \sum_{s \in S} d(s) \sum_{a \in A} \pi(s; a) \nabla_{\theta} (Q(s; a) - b(s)) \quad (5)$$

and thus, for any baseline  $b(s)$ , the gradient of the average reward can be written as

$$\nabla_{\theta} J(\theta) = \sum_{s \in S} d(s) \sum_{a \in A} \pi(s; a) \nabla_{\theta} Q(s; a) - b(s) \nabla_{\theta} J(\theta) \quad (6)$$

The baseline  $b(s)$  can be chosen in a way that the variance of the gradient estimates  $\nabla_{\theta} J(\theta)$  is minimized (Greensmith, Bartlett, & Baxter, 2004).

The natural gradient, denoted  $\nabla^{\natural} J(\theta)$ , can be calculated by linearly transforming the regular gradient,  $\nabla_{\theta} J(\theta)$ , using the inverse Fisher information matrix of the policy:  $\nabla^{\natural} J(\theta) = G(\theta)^{-1} \nabla_{\theta} J(\theta)$ .

<sup>1</sup> In the rest of the paper we use the notation  $\nabla$  to denote  $\nabla_{\theta}$  the gradient with respect to the policy parameters.

The Fisher information matrix  $G_{\pi}$  can be seen (Bagnell & Schneider, 2003; Kakade, 2002; Peters et al., 2003) to be

$$G_{\pi} = \mathbb{E}_{s \sim d} \left[ \mathbb{E}_{a \sim \pi(s)} \left[ \frac{\nabla \log \pi(s, a)}{\pi(s, a)} \frac{\nabla \log \pi(s, a)}{\pi(s, a)} \right] \right] \quad (7)$$

Matrix  $G_{\pi}$  plays an important role in the algorithms that use natural gradients (Kakade, 2002; Peters & Schaal, 2008). Here  $\mathbb{E}_{s \sim d}$  denotes the expectation under the conditional joint distribution where states are first selected according to distribution  $d$ , and then given that a state  $s$  is selected, actions are selected according to distribution  $\pi(s, \cdot)$ . The Fisher information matrix is clearly positive definite (Kakade, 2002).

A well-studied example of parameterized randomized policies, which we use in our experiments, is the Gibbs (or Boltzmann) distribution having the form

$$\pi(s, a) = \frac{e^{\phi(s, a)} \pi_0(s, a)}{\sum_{a'} e^{\phi(s, a')} \pi_0(s, a')}; \quad \phi(s, a) = \sum_{i=1}^{d_1} \phi_i(s, a) \quad (8)$$

where each  $\phi_i$  is a  $d_1$ -dimensional feature vector for the state-action pair  $(s, a)$ .

### 3. Policy gradient with function approximation

Now consider the case in which the action-value function for a fixed policy  $\pi$ ,  $Q_{\pi}$ , is approximated by a learned function approximator. If the approximation is sufficiently good, we might hope to use it in place of  $Q_{\pi}$  in Eqs. (4) and (6), and still point roughly in the direction of the true gradient. Sutton et al. (2000) showed that if the approximation  $\hat{Q}_w$  with parameter  $w \in \mathbb{R}^{d_1}$  is compatible, i.e.,  $\nabla_w \hat{Q}_w(s, a) \propto \nabla \log \pi(s, a)$ , and minimizes the mean-squared error

$$\mathbb{E}_{s \sim d} \left[ \mathbb{E}_{a \sim \pi(s, \cdot)} \left[ \left( Q_{\pi}(s, a) - \hat{Q}_w(s, a) \right)^2 \right] \right] \quad (9)$$

for parameter value  $w^*$ , then we can replace  $Q_{\pi}$  with  $\hat{Q}_{w^*}$  in Eqs. (4) and (6). We work with linear approximation  $\hat{Q}_w(s, a) = \phi(s, a)^T w$  in which the  $\phi$ 's are compatible features defined according to  $\phi(s, a) \propto \nabla \log \pi(s, a)$ . Convergence of a temporal difference critic under a linear approximation when trajectories are sampled according to their distribution under the target policy has been established earlier (Sutton, 1996; Tadic, 2001; Tsitsiklis & Van Roy, 1997). Note that compatible features are well defined under (A2). As an example, the compatible features for the Gibbs policy in Eq. (8) are  $\phi_i(s, a) = \phi_i(s, a) / \pi_0(s, a)$ . The Fisher information matrix of Eq. (7) can be written using the compatible features as

$$G_{\pi} = \mathbb{E}_{s \sim d} \left[ \mathbb{E}_{a \sim \pi(s, \cdot)} \left[ \left( \sum_{i=1}^{d_1} \phi_i(s, a) \frac{\nabla \log \pi(s, a)}{\pi(s, a)} \right) \left( \sum_{j=1}^{d_1} \phi_j(s, a) \frac{\nabla \log \pi(s, a)}{\pi(s, a)} \right)^T \right] \right] \quad (10)$$

Suppose  $\mathbb{E}_{s \sim d} \left[ \mathbb{E}_{a \sim \pi(s, \cdot)} \left[ \left( Q_{\pi}(s, a) - \hat{Q}_w(s, a) \right)^2 \right] \right]$  denotes the mean-squared error

$$\mathbb{E}_{s \sim d} \left[ \mathbb{E}_{a \sim \pi(s, \cdot)} \left[ \left( Q_{\pi}(s, a) - \hat{Q}_w(s, a) \right)^2 \right] \right] \quad (11)$$

of our compatible linear parameterized approximation,  $\hat{Q}_{w^*}$ , and an arbitrary baseline  $b(s)$ . Let  $w^* = \arg \min_w \mathbb{E}_{s \sim d} \left[ \mathbb{E}_{a \sim \pi(s, \cdot)} \left[ \left( Q_{\pi}(s, a) - \hat{Q}_w(s, a) \right)^2 \right] \right]$  denote the optimal parameter. Lemma 1 shows that the value of  $w^*$  does not depend on the given baseline  $b(s)$ ; as a result the mean-squared error problems of Eqs. (9) and (11) have the same solutions. Next in Lemma 2, we show that if the parameter is set to be equal to  $w^*$ , then the resulting mean-squared error  $\mathbb{E}_{s \sim d} \left[ \mathbb{E}_{a \sim \pi(s, \cdot)} \left[ \left( Q_{\pi}(s, a) - \hat{Q}_{w^*}(s, a) \right)^2 \right] \right]$  (now treated as a function of the baseline  $b(s)$ ) is further minimized when  $b(s) = V_{\pi}(s)$  (see also Chapter 11 of Meyn, 2007). In other words, the variance in the action-value-function estimator

is minimized if the baseline is chosen to be the value function itself.<sup>2</sup> The proofs of Lemmas 1 and 2 can be found in Bhatnagar et al. (2009, 2008).

**Lemma 1.** The optimum weight parameter  $w^*$  for any given policy  $\pi$  satisfies<sup>3</sup>

$$w^* = \mathbb{E}_{s \sim d} \left[ \mathbb{E}_{a \sim \pi(s, \cdot)} \left[ \frac{\nabla \log \pi(s, a)}{\pi(s, a)} \right] \right]$$

**Lemma 2.** For any given policy  $\pi$ , the minimum variance baseline  $b^*(s)$  in the action-value-function estimator corresponds to the state-value function  $V_{\pi}(s)$ .

From Lemma 1,  $w^*$  is a least-squared optimal parametric representation for the action-value function  $Q_{\pi}(s, a)$ . On the other hand, from Lemma 2, the same is also a least-squared optimal parametric representation for the advantage function  $A_{\pi}(s, a) = Q_{\pi}(s, a) - V_{\pi}(s)$ . The mean-squared error (11) is seen to be minimized w.r.t. the baseline  $b(s)$  for  $b^*(s) = V_{\pi}(s)$ , thereby making it more meaningful to consider  $w^*$  to be the least-squared optimal parametric representation for the advantage function rather than the action-value function itself.

The temporal difference (TD) error  $\delta_t$  is a random quantity that is defined according to

$$\delta_t = r_{t+1} - \bar{Q}_{t+1} - \bar{Q}_{s_t, a_t} - \bar{Q}_{s_t}; \quad (12)$$

where  $\bar{Q}_{s_t}$  is an unbiased estimate of the differential reward in states  $s_t$ ,  $i \in \{1, \dots, d_1\}$ . Likewise,  $\bar{Q}_{t+1}$  is an unbiased estimate of the average reward. Thus, in particular, these estimates satisfy  $\mathbb{E}[\bar{Q}_{s_t} | s_t] = V_{\pi}(s_t)$  and  $\mathbb{E}[\bar{Q}_{t+1} | s_t] = J_{\pi}$ , for any  $t \geq 0$ , respectively. We assume here that actions are chosen according to policy  $\pi$ . The next lemma is also a simple result that shows that  $\delta_t$  is an unbiased estimate of the advantage function  $A_{\pi}$ , see Bhatnagar et al. (2009, 2008) for a proof. A proof of this lemma is also available in Peters et al. (2003) and Peters and Schaal (2008).

**Lemma 3.** Under given policy  $\pi$  with actions chosen according to it, we have

$$\mathbb{E}[\delta_t | s_t, a_t] = A_{\pi}(s_t, a_t)$$

By setting the baseline  $b(s)$  equal to the value function  $V_{\pi}(s)$ , Eq. (6) can be written as

$$\nabla J_{\pi} = \mathbb{E}_{s \sim d} \left[ \mathbb{E}_{a \sim \pi(s, \cdot)} \left[ \frac{\nabla \log \pi(s, a)}{\pi(s, a)} \delta_t \right] \right]$$

From Lemma 3,  $\delta_t$  is an unbiased estimate of the advantage function  $A_{\pi}(s, a)$ . Thus,  $\bar{Q}_{t+1} - \bar{Q}_{s_t, a_t} - \bar{Q}_{s_t}$  is an unbiased estimate of  $\nabla J_{\pi}$ . However, calculating  $\delta_t$  requires having estimates,  $\bar{Q}_{t+1}$ , of the average reward and the value function. While an average reward estimate is simple enough to obtain given the single-stage reward function, the same is not necessarily true for the value function. We use function approximation for the value functions as well. Suppose  $f_s$  is a  $d_2$ -dimensional feature vector for state  $s$  (for some  $d_2 \geq 1$ ). We denote  $f_s = [f_{s,1}, \dots, f_{s,d_2}]^T$ . One may then approximate  $V_{\pi}(s)$  with  $v^T f_s$ , where  $v$  is a  $d_2$ -dimensional

<sup>2</sup> It is important to note that Lemma 2 is not about the minimum variance baseline for gradient estimation. It is about the minimum variance baseline of the action-value-function estimator.

<sup>3</sup> This lemma is similar to Theorem 1 by Kakade (2002), except that we consider baseline  $b(s)$  which again can be seen as additional basis functions in the sense of Peters et al. (2003) and Peters and Schaal (2008).



weight vector which can be tuned (for a fixed policy  $\pi$ ) using a TD algorithm. In our algorithms, we then use

$$e_t = r_{tC1} - \hat{V}_{tC1} = V_t - \hat{V}_{tC1} \quad (13)$$

as an estimate for the TD error, where  $V_t$  corresponds to the value function parameter at time  $t$ . From now on, unless explicitly mentioned, we shall consider  $e_t$  to be defined according to (13). Let  $\hat{V}_{\infty}$  denote the quantity

$$\hat{V}_{\infty} = \sum_{a \in A} d(s; a) R(s; a) + \gamma \sum_{s' \in S} P(s; a; s') \hat{V}_{\infty} \quad (14)$$

where  $\hat{V}_{\infty}$  is an estimate of the differential value function  $V_{\infty}$  that is obtained upon convergence of a TD recursion (above) viz.,  $\lim_{t \rightarrow \infty} V_t = V_{\infty}$  with probability one. Also, let  $e_t$  denote the associated quantity

$$e_t = r_{tC1} - \hat{V}_{tC1} = V_t - \hat{V}_{tC1} \quad (15)$$

Here  $r_{tC1}$  and  $\hat{V}_{tC1}$  are the same as before. Then  $e_t$  corresponds to a stationary estimate of the TD error (with function approximation) under policy  $\pi$ . We have the following analog of Theorem 1 of Sutton et al. (2000).

**Lemma 4.**

$$\mathbb{E} \|e_t\| \leq \sum_{s \in S} d(s) \|r - \gamma \sum_{s' \in S} P(s; a; s') \hat{V}_{\infty}\| \quad (16)$$

**Proof.** A simple calculation shows that

$$\mathbb{E} \|e_t\| \leq \sum_{s \in S} d(s) \|r - \gamma \sum_{s' \in S} P(s; a; s') \hat{V}_{\infty}\| \quad (16)$$

Now from (14),

$$\hat{V}_{\infty} = \sum_{a \in A} d(s; a) R(s; a) + \gamma \sum_{s' \in S} P(s; a; s') \hat{V}_{\infty} \quad (17)$$

From (16) and the above, we get

$$\sum_{s \in S} d(s) \|r - \gamma \sum_{s' \in S} P(s; a; s') \hat{V}_{\infty}\| \leq \sum_{s \in S} d(s) \|r - \gamma \sum_{s' \in S} P(s; a; s') \hat{V}_{\infty}\| \quad (17)$$

Now observe that  $d(s)$  correspond to the stationary probabilities that satisfy,  $\sum_{s \in S} d(s) = 1$ ,

$$\sum_{s \in S} d(s) = \sum_{s \in S} d(s) \sum_{a \in A} p(s; a) \sum_{s' \in S} d(s') = \sum_{s \in S} d(s) \sum_{a \in A} p(s; a) \sum_{s' \in S} d(s') = 1 \quad (18)$$

where  $p(s; a) = \sum_{s' \in S} d(s') P(s; a; s')$  are the transition probabilities of the resulting Markov chain under policy  $\pi$ . Hence,

$$\sum_{s \in S} d(s) \|r - \gamma \sum_{s' \in S} P(s; a; s') \hat{V}_{\infty}\| \leq \sum_{s \in S} d(s) \|r - \gamma \sum_{s' \in S} P(s; a; s') \hat{V}_{\infty}\| \quad (19)$$

The claim now follows from (17).

Note that according to Theorem 1 of Sutton et al. (2000),  $\mathbb{E} \|e_t\| \leq \sum_{s \in S} d(s) \|r - \gamma \sum_{s' \in S} P(s; a; s') \hat{V}_{\infty}\|$ , provided  $e_t$  is defined according to (12). For the case with function approximation that we study, from Lemma 4, the quantity  $\sum_{s \in S} d(s) \|r - \gamma \sum_{s' \in S} P(s; a; s') \hat{V}_{\infty}\|$  may be viewed as the error or bias in the estimate of the gradient of average reward that results from the use of function approximation. It is interesting to observe that this does not depend on the differential reward  $V_{\infty}$  that is obtained as a solution to (3). We also have

**Corollary 1.**  $\sum_{s \in S} d(s) \|r - \gamma \sum_{s' \in S} P(s; a; s') \hat{V}_{\infty}\| \leq 0$ :

**Proof.** This follows directly from the definition of  $\hat{V}_{\infty}$  in (14), the definition of  $J_{\infty}$  in (2), and an analogous equation as (19) with  $V_{\infty}$  in place of  $r - \gamma \sum_{s' \in S} P(s; a; s') \hat{V}_{\infty}$ .

#### 4. Actor-critic algorithms

We present four new actor-critic algorithms in this section. They update the policy parameters along the direction of the average reward gradient. While estimates of the regular gradient are used for this purpose in Algorithm 1, natural gradient estimates are used in Algorithms 2–4. Let  $\hat{V}_{\infty}$  denote the parameterized approximation to the differential value function in state  $s$ . One can also denote the same as  $\hat{V}_{\infty} = \mathcal{G}V_{\infty}$ , where  $\mathcal{G}$  is an  $n \times d_2$  matrix whose  $k$ th column ( $k = 1, \dots, d_2$ ) is  $f_k(s) = \sum_{a \in A} \pi(a|s) \nabla_a V_{\infty}(s)$ . We make the following assumption as in Tsitsiklis and Van Roy (1999) (see also Tsitsiklis & Van Roy, 1997).

**(A3)** The basis functions  $f_k(s)$  ( $k = 1, \dots, d_2$ ) are linearly independent. In particular,  $d_2 = n$  and  $\mathcal{G}$  has full rank. Also, for every  $v \in \mathbb{R}^{d_2}$ ,  $\mathcal{G}v \in \mathcal{C}$ , where  $\mathcal{C}$  is the  $n$ -dimensional vector with all entries equal to one.

Let  $\eta_t$  and  $\bar{\eta}_t$  be two step-size schedules that satisfy (20). Further, let  $\eta_t$  defined by  $\eta_t = c \bar{\eta}_t$  for some  $c > 0$  be the step-size schedule for the average reward recursions in our algorithms.

From the last condition in (20),  $\eta_t \rightarrow 0$  faster than  $\bar{\eta}_t$ . Thus the critic is a faster recursion than the actor.

We now present our actor-critic algorithms. For the actor updates in our algorithms, we use a projection operator  $O: \mathbb{R}^{d_1} \rightarrow \mathbb{R}^{d_1}$  that projects any  $x \in \mathbb{R}^{d_1}$  to a compact set  $\mathcal{C} \subset \mathbb{R}^{d_1}$  defined by  $\mathcal{C} = \{x \in \mathbb{R}^{d_1} : 0 \leq x_i \leq 1, i = 1, \dots, d_1\}$ , where  $q_i, i = 1, \dots, d_1$  are real-valued, continuously differentiable functions on  $\mathbb{R}^{d_1}$  that represent the constraints specifying the (above) compact region. Here for each  $x$  on the boundary of  $\mathcal{C}$ , the gradients of the active constraints are considered to be linearly independent. This is the setting considered for projection-based algorithms in Chapter 5 of Kushner and Clark (1978). For any  $x \in \mathbb{R}^{d_1}$ ,  $Ox \in \mathcal{C}$  and in particular for  $x \in \mathcal{C}$ ,  $Ox = x$  itself. As explained in Chapter 2 of Kushner and Clark (1978), any compact hyperrectangle in  $\mathbb{R}^{d_1}$  is a special case of  $\mathcal{C}$  (above). The projection method is an often used technique to ensure boundedness of iterates in stochastic approximation algorithms, see for instance, Abounadi, Bertsekas, and Borkar (2001), where it has been used in the context of a stochastic shortest path Q-learning algorithm. Some discussion on this is also available in Tsitsiklis (1994). The other approach (that is also usually taken, which we do not follow) is to simply assume that the iterates (see below) (23), (29), (34) and (39) are bounded without the projection, and then show their convergence under this assumption. In our experiments, however (see Bhatnagar et al., 2009), we do not project the iterates to a constraint region as they are seen to remain bounded (without projection). In Remark 2 (that follows Theorem 1), we explain the difficulties in proving boundedness of iterates in the absence of the projection operator  $O$ .

**Algorithm 1** (Regular-Gradient Actor–Critic).

$$J_{tC1} \leftarrow J_t + \alpha_t (J_t - J_{tC1}); \quad (21)$$

$$V_{tC1} \leftarrow V_t + \alpha_t (V_t - V_{tC1}); \quad (22)$$

$$\pi_{tC1} \leftarrow \pi_t + \alpha_t (C_t - V_{tC1}); \quad (23)$$

with  $\pi_t$  as in (13). This is the only actor–critic algorithm presented in the paper that is based on the regular-gradient estimate. It stores two parameter vectors  $\pi$  and  $V$ . Its per time-step computational cost is linear in the number of policy and value-function parameters.

The next algorithm is based on the natural-gradient estimate  $\hat{J}_t = J_t / D G_t^{-1} \pi_t$  in place of the regular-gradient estimate in Algorithm 1. We derive a procedure below for recursively estimating  $G_t^{-1}$  on a faster timescale. The above estimation is done on a faster scale so that convergence of the associated iterates is achieved prior to a  $\pi$ -update. Suppose  $G_t^{-1}$  denote the  $t$ th estimate of  $G_t^{-1}$ . Our procedure is obtained in a similar manner as the method described on pp. 147–152 of Widrow and Stearns (1985). The latter approach however considers the estimates as being obtained via a “fading memory” condition in which the most recent observation is given the highest weight. The weights themselves decrease geometrically over past observations. On the other hand, unlike Widrow & Stearns, 1985, we consider stationary averages that depend on parameter  $\pi$ , that in turn gets updated along the “slower timescale”. This constitutes a natural setting for our algorithm. We show in Lemma 6 that  $G_t^{-1} \rightarrow G_t^{-1}$  as  $t \rightarrow \infty$  with probability one. This is required for proving convergence of our algorithm. On the other hand, showing the same for the corresponding estimates in Widrow and Stearns (1985) does not seem possible because  $G_t \neq G_t$  there.

We consider  $G_t, t \geq 0$ , defined as (the sample averages)

$$G_t \leftarrow \frac{1}{tC1} \sum_{i=0}^{t-1} s_i a_i s_i a_i^T.$$

Thus, one may obtain recursively

$$G_t \leftarrow (1 - \frac{1}{tC1}) G_{t-1} + \frac{1}{tC1} s_t a_t s_t a_t^T. \quad (24)$$

More generally, one may consider the recursion

$$G_t \leftarrow (1 - \alpha_t) G_{t-1} + \alpha_t s_t a_t s_t a_t^T, \quad (25)$$

where the step-size  $\alpha_t$  is as before. This would correspond to a case of weighted averages (with the weights corresponding to the step-sizes  $\alpha_t$ ). However, through a stochastic approximation argument, one can see that (25) would asymptotically converge to  $G_t$  almost surely, if  $\pi$  is held fixed. In fact, with an appropriate choice of  $\alpha_t$ , one can obtain faster convergence of iterates in (25) over those in (24). Using Sherman–Morrison matrix inversion lemma, one obtains

$$G_t^{-1} \leftarrow \frac{1}{1 - \alpha_t} G_{t-1}^{-1} + \frac{\alpha_t G_{t-1}^{-1} s_t a_t s_t a_t^T G_{t-1}^{-1}}{1 - \alpha_t s_t a_t s_t a_t^T G_{t-1}^{-1}}; \quad (26)$$

The following assumption is on the matrices  $G_t, G_t^{-1}$ .

**(A4)** We have  $\sup_{t \geq 0} \lambda_{\min}(G_t) \geq k, \sup_{t \geq 0} \lambda_{\max}(G_t^{-1}) < 1/k$ .

This assumption will be used in proving the convergence of our Algorithms 2 and 4. It is similar to a corresponding requirement in the case of certain Hessian matrices in the Newton-based simulation optimization schemes in Bhatnagar (2005, 2007). A sufficient condition for both the requirements in (A4) is that (cf. pp. 35 of Bertsekas, 1999) for some scalars  $c_1, c_2 > 0$ ,

$$c_1 k \leq \lambda_{\min}(G_t) \leq c_2 k, \quad \forall t \geq 0.$$

for all  $s \in S, a \in A, z \in \mathbb{R}^{d_1}$  and  $\lambda$ . It is then easy to see that  $\lambda_{\min}(G_t) \geq k \lambda_{\min}(G_0) \geq k \lambda_{\min}(G_0)$ , for all  $t \geq 0$ , and the eigenvalues of  $G_t$  lie between  $\lambda_{\min}(G_0)$  and  $\lambda_{\max}(G_0)$ . Here  $\lambda_{\min}(G_0) \geq \min_{s,a} s a^T G_0^{-1} s a$  and  $\lambda_{\max}(G_0) \leq \max_{s,a} s a^T G_0^{-1} s a$ . Also,  $\lambda_{\min}(G_0), \lambda_{\max}(G_0) > 0$ . Hence, the procedure (below) does not get stuck at a nonstationary point. Under the above sufficient condition, (A4) follows from Propositions A.9 and A.15 of Bertsekas (1999).

Our second algorithm stores matrix  $G^{-1}$  and two parameter vectors  $\pi$  and  $V$ . Its per time-step computational cost is linear in the number of value-function parameters and quadratic in the number of policy parameters.

**Algorithm 2** (Natural-Gradient Actor–Critic with Fisher Information Matrix).

$$J_{tC1} \leftarrow J_t + \alpha_t (J_t - J_{tC1}); \quad (27)$$

$$V_{tC1} \leftarrow V_t + \alpha_t (V_t - V_{tC1}); \quad (28)$$

$$\pi_{tC1} \leftarrow \pi_t + \alpha_t G_t^{-1} (C_t - V_{tC1}); \quad (29)$$

with  $\pi_t$  as in (13). Also, the estimate of the inverse Fisher information matrix updated according to Eq. (26). Like Widrow and Stearns (1985), we let  $G_0^{-1} \leftarrow kI$ , where  $I$  is a  $d_1 \times d_1$ -dimensional identity matrix and  $k > 0$ . Thus  $G_0^{-1}$  and hence also  $G_0$  are positive definite and symmetric matrices. From (25),  $G_t^{-1}, t \geq 1$  can be seen to be positive definite and symmetric because these are convex combinations of positive definite and symmetric matrices. Hence,  $G_t^{-1}, t \geq 1$  are positive definite and symmetric matrices as well.

As we mentioned in Section 3, it is better to think of the compatible approximation  $W^{\pi, \pi}$  as an approximation of the advantage function rather than of the action-value function. In our next algorithm, we tune the weight parameters  $W$  in such a way as to minimize an estimate of the least-squared error  $E \|W - D E_s \pi\|^2$ . Note that the gradient of  $E \|W - D E_s \pi\|^2$  w.r.t.  $W$  is

We use the following estimate of  $\nabla_W E \|W - D E_s \pi\|^2$ .

$$\nabla_W E \|W - D E_s \pi\|^2 \approx \frac{2}{tC1} \sum_{i=0}^{t-1} (s_i a_i - D E_s \pi) (s_i a_i - D E_s \pi)^T. \quad (30)$$

Hence, we update advantage parameters  $W$  along with value-function parameters  $V$  in the critic update of this algorithm as

$$W_{tC1} \leftarrow W_t + \alpha_t \nabla_W E \|W_t - D E_s \pi_t\|^2.$$

The factor 2 on the RHS of (30) does not play a role because of the diminishing step-size sequence  $\alpha_t, t \geq 0$  and so has been dropped in the above recursion. We maximize the long-run average reward  $J$  along the slower timescale and use the natural-gradient estimate for this purpose. Like Peters et al. (2003) and Peters and Schaal (2008), the natural-gradient estimate that we use in the actor update of Algorithm 3 is  $\hat{J}_t = J_t / D W_{tC1}$ . This algorithm stores three parameter vectors,  $V, W$ , and  $\pi$ . Its per time-step computational cost is linear in the number of value-function parameters and quadratic in the number of policy parameters.

**Algorithm 3** (Natural-Gradient Actor–Critic with Advantage Parameters).

$$J_{tC1} \leftarrow J_t + \alpha_t (J_t - J_{tC1}); \quad (31)$$

$$V_{tC1} \leftarrow V_t + \alpha_t (V_t - V_{tC1}); \quad (32)$$

$$W_{tC1} \leftarrow W_t + \alpha_t \nabla_W E \|W_t - D E_s \pi_t\|^2; \quad (33)$$

$$\pi_{tC1} \leftarrow \pi_t + \alpha_t W_{tC1}^{-1} \hat{J}_t; \quad (34)$$

with  $\pi_t$  as in (13). Although the estimates of  $G_t^{-1}$  are not explicitly computed and used in Algorithm 3, the convergence analysis of this algorithm in the next section shows that the overall scheme still moves in the direction of the natural gradient of average reward.

In Algorithm 4, however, we explicitly estimate  $G_t^{-1}$  (as in Algorithm 2), and use it in the critic update for  $W$ . The overall scheme is again seen to follow the direction of the natural gradient of average reward. Here, we let

$$\hat{G}_t^{-1} = \frac{W/D}{2G_t^{-1}} \cdot \frac{s_t a_t}{s_t a_t} W \quad (35)$$

be the estimate of the natural gradient of the least-squared error  $E_t(W)$ . This also simplifies the critic update for  $W$ . Further, we remove the factor 2 from the natural-gradient estimate (35) because of diminishing  $s_t, t \rightarrow 0$ , as before. Algorithm 4 stores a matrix  $G_t^{-1}$  and three parameter vectors,  $V, W$ , and  $C$ . Its per time-step computational cost is linear in the number of value-function parameters and quadratic in the number of policy parameters.

**Algorithm 4** (Natural-Gradient Actor-Critic with Advantage Parameters and Fisher Information Matrix).

$$\hat{G}_{tC1} \leftarrow \frac{1}{t} \hat{G}_t C \quad (36)$$

$$V_{tC1} \leftarrow V_t C \quad (37)$$

$$W_{tC1} \leftarrow \frac{1}{t} \hat{G}_t W_t C \quad (38)$$

$$C_{tC1} \leftarrow \frac{1}{t} C \quad (39)$$

with  $\hat{G}_t$  as in (13) and where the estimate of the inverse of the Fisher information matrix is updated according to (26). As with Algorithm 2, we let  $G_0^{-1} \leftarrow kI$  with  $k > 0$ .

## 5. Convergence analysis

We now present the convergence analysis of our algorithms. The analysis mainly follows the ordinary differential equation (ODE) approach (Benveniste, Metivier, & Priouret, 1990; Kushner & Clark, 1978; Kushner & Yin, 1997). Note that the problem we consider is a maximization and not a minimization problem. For the purpose of analysis, we consider an associated problem with costs defined as negative rewards and our aim is to minimize the associated long-run average cost. The negative of the minimum cost thus obtained then corresponds to the maximum reward in the original problem. This is useful in pushing through certain stability arguments and showing convergence of iterates. Our algorithms use function approximation and aim at finding the local maxima of the average rewards. All our convergence results are in the Euclidean norm. Further, for any matrix  $A$ , we define its norm as the induced matrix norm  $\|A\|_1 = \max_{j \in \{1, \dots, n\}} \sum_{i=1}^n |A_{ij}|$ .

### 5.1. Convergence analysis for Algorithm 1

We require Assumptions (A1)–(A3) here. As explained above, one may view  $r_{tC1}$  as the cost incurred at instant  $t$  in a transformed problem. Because of the above, a change occurs only in the actor recursion (23) due to this transformation, and it becomes

$$C_{tC1} \leftarrow C_t + \frac{1}{t} (r_{tC1} - C_t) \quad (40)$$

Recursions for the average reward (21), TD error (13), and critic (22), being fixed point recursions (see Tsitsikis & Van Roy, 1999), are left unchanged. For any given policy (along the faster timescale), average reward (21), TD error (13), and critic (22) recursions correspond to the TD recursions in Tsitsikis and Van Roy (1999) with  $D = 0$ .

Let  $D$  denote the diagonal matrix with elements  $d_i = s_i / \sum_{j=1}^n s_j$  along its diagonal. Let  $P$  be the probability matrix  $P \leftarrow \sum_{j=1}^n s_j s_j^T / \sum_{j=1}^n s_j$  for the Markov chain under policy  $\pi$  and  $R$  be the corresponding column vector of average rewards whose  $i$ th element is  $\sum_{j=1}^n s_j a_j / \sum_{j=1}^n s_j$ . Also, let  $T \in \mathbb{R}^n \rightarrow \mathbb{R}^n$  be the operator given by  $T(J) \leftarrow R + \gamma J - P J$ . The proof of

convergence of TD by Tsitsikis and Van Roy (1999) is based on a result by Benveniste et al. (1990). We provide in Lemma 5 an alternative simpler proof of convergence under the same assumptions used by Tsitsikis and Van Roy (1999), using a recently developed result by Borkar and Meyn (2000). We consider  $D = 0$  to suit our algorithm. The proof however carries through quite easily for  $D > 0$  as well. We have

**Lemma 5.** For any given  $\pi$  and  $f, g$  as in (21) and (22), respectively, we have  $\hat{J}_t \rightarrow J$  and  $\hat{V}_t \rightarrow V$  with probability one, where

$$J = \sum_{s \in \mathcal{S}} d_s \sum_{a \in \mathcal{A}} \pi(a|s) \frac{s \cdot a}{R \cdot s + a} \quad (41)$$

is the average reward under  $\pi$  and  $V$  is obtained as the unique solution to

$$D \delta V = D \delta T \delta V \quad (42)$$

**Proof.** The proof is based on verifying the Assumptions (A1)–(A2) of Borkar and Meyn (2000). First consider the average reward recursion (21). The ODE describing the asymptotic behavior of this recursion corresponds to

$$\dot{P} = \sum_{s \in \mathcal{S}} C_s \sum_{a \in \mathcal{A}} d_s \frac{s \cdot a}{R \cdot s + a} \quad (43)$$

Let  $f$  denote the RHS of (43). Then  $f$  is Lipschitz continuous in  $P$ . Let  $f_1 = \lim_{r \rightarrow 1} \frac{f \cdot r}{r}$ . The function  $f_1$  exists and is simply  $f_1 = \sum_{s \in \mathcal{S}} C_s \sum_{a \in \mathcal{A}} d_s \frac{s \cdot a}{R \cdot s}$ . The origin is clearly an asymptotically stable equilibrium for the ODE  $\dot{P} = f_1 P$ .

Now consider recursions for TD error (13) and critic (22). Consider the following ODE (in vector-matrix notation) associated with them.

$$\dot{V} = D \delta T \delta V \quad (44)$$

Let  $g^1 \cdot V$  denote the RHS of (44). Then  $g^1 \cdot V$  is also Lipschitz continuous in  $V$ . Further, for  $g_1^1 = \lim_{r \rightarrow 1} \frac{g^1 \cdot r V}{r}$ , it can be seen that  $g_1^1 \cdot V$  exists and equals  $g_1^1 \cdot V = D \delta T \delta P \cdot I / \delta V$ , where  $I$  is the identity matrix. Consider now the system

$$\dot{V} = g_1^1 \cdot V \quad (45)$$

Note that the matrix  $P$  has a simple eigenvalue of one and its remaining eigenvalues have real parts that are less than one. Thus  $P \cdot I$  will have one eigenvalue of zero and other eigenvalues with negative real parts. Also, corresponding to the eigenvalue zero, the matrix  $P \cdot I$  has a left eigenvector  $d^T$  and a right eigenvector  $e \leftarrow \frac{1}{\sum_{j=1}^n s_j} \mathbf{1}$  (the  $n$ -dimensional unit vector), respectively. Thus, in principle, the set of asymptotically stable fixed points of (45) would correspond to the set  $\{V \in \mathbb{R}^n \mid d^T V = 0 \text{ and } V \in \mathbb{R}^n\}$ . Now by the second part of Assumption (A3),  $\delta V \in \mathbb{R}^n$  for every  $V \in \mathbb{R}^n$ . Thus the only asymptotically stable equilibrium for (45) is the origin.

Next, define  $N^1 \cdot t, M^1 \cdot t, t \geq 0$ , according to  $N^1 \cdot t \leftarrow r_{tC1}$ ,  $E[r_{tC1} \mid F_1, t] \leftarrow M^1 \cdot t$ ,  $E[t_{tC1} \mid F_1, t] \leftarrow T$ , respectively, where  $F_1 \cdot t \leftarrow V_t, \hat{J}_t, M^1 \cdot t, N^1 \cdot t, t \geq 0$ . It is easy to see that

$$E[k N^1 \cdot t \mid F_1, t] \leq C_1 \cdot k \hat{J}_t^2 \leq C_1 \cdot k V_t^2$$

$$E[k M^1 \cdot t \mid F_1, t] \leq C_2 \cdot k V_t^2 \leq C_2 \cdot k \hat{J}_t^2$$

$t \geq 0$ , for some  $C_1, C_2 < \infty$ . In fact, quantities  $N^1 \cdot t$  can be directly seen to be uniformly bounded almost surely. Thus Assumptions (A1) and (A2) of Borkar and Meyn (2000) can be seen to be satisfied in the case of the average reward (21), TD error (13), and critic (22) recursions. From Theorem 2.1 of Borkar and Meyn (2000), average reward, TD error, and critic iterates are uniformly bounded with

probability one. Now note that (43) has  $J_*$  defined as in (41) as its unique globally asymptotically stable equilibrium.

Next, suppose that  $v \in \mathcal{V}$  is a solution to the system

$$\dot{v} = D\delta v + D^* \delta v. \quad (46)$$

We show that  $v_*$  is the unique globally asymptotically stable equilibrium of the ODE (44) with the function  $W_*$  defined by  $W_* = D^* M_* v_* / M_* v_*$  with  $M_* v_* \in \mathcal{S}^> D.T. \delta v / \delta v$  serving as an associated strict Lyapunov function. Thus

$$\frac{dW_* v}{dt} \leq D^* W_* v \leq D^* M_* v \leq \mathcal{S}^> D.P \quad I/8 M_* v;$$

In lieu of (A3), for any  $r \in \mathbb{R}^{d_2}$ ,  $\delta r$  is a nonconstant vector (i.e., one that is not of the form  $e$  for  $\theta \in \mathbb{D}$ ). Thus,

$$r^> \mathcal{S}^> D.P \quad I/8 r < 0; \quad \delta r \in \mathbb{D};$$

( $\mathbb{D}$  being the vector in  $\mathbb{R}^{d_2}$  with all entries 0), i.e., the matrix  $\mathcal{S}^> D.P \quad I/8$  is negative definite (see also the proof of Lemma 7, p. 1803 of Tsitsikis and Van Roy (1999) for a similar conclusion). Now any  $\hat{v} \in \mathcal{V} \subset \mathcal{V}_*$  with  $\hat{v} \in \mathbb{R}$ ,  $\hat{v} \in \mathbb{D}$  and  $\hat{v}$  such that  $\delta \hat{v} \in e$  will also be a solution to the linear system of Eqs. (46). However, again by Assumption (A3),  $\delta \hat{v} \in e$  for any  $\hat{v} \in \mathbb{R}^{d_2}$ . Thus any  $\hat{v}$  as above will not be a solution and the only solution is  $\hat{v} \in \mathcal{V}_*$  which is therefore unique. Thus,  $\frac{dW_* v}{dt} < 0$  on the set  $\mathcal{F} \subset \mathbb{R}^{d_2} \setminus \mathcal{V}_*$  and  $\frac{dW_* v}{dt} = 0$  on the set  $\mathcal{F} \subset \mathcal{V}_*$ . Thus for (44),  $v_*$  is the unique globally asymptotically stable equilibrium. Assumptions (A1)–(A2) of Borkar and Meyn (2000) are now verified and the claim follows from their Theorem 2.2, p. 450 of Borkar and Meyn (2000).

**Remark 1.** Note that (A3) has also been used in the analysis of average cost TD learning by Tsitsikis and Van Roy (1999) (cf. Assumption 2, p. 1800). We also require this assumption as our TD recursions are exactly the same as those in Tsitsikis and Van Roy (1999). On the other hand, in a recent paper Borkar (2008) develops a variant of TD learning with function approximation that is based on the relative value iteration scheme. For such a scheme, one would not require the later part of (A3) (i.e., Assumption 2(b) of Tsitsikis & Van Roy, 1999).

Consider now recursion (40) along the slower timescale corresponding to  $t$ . Let  $v_t$  be a vector field on  $\mathcal{C}$ . Define another vector field

$$\hat{v}_t = v_t + \lim_{\theta \rightarrow 0} \frac{\partial y_{\mathcal{C}}(v_t, y)}{\partial y} y;$$

In case the above limit is not unique, we let  $\hat{v}_t = v_t + y$  be the set of all possible limit points (see p. 191 of Kushner & Clark, 1978). Consider now the ODE (in lieu of Lemma 4, see Bhatnagar et al., 2009)

$$\dot{P} = D \hat{v}_t + P J_t + e; \quad (47)$$

where  $e = \sum_{s \in \mathcal{S}} d(s, r) \psi(s, r) v^> f_s$ . Consider also an associated ODE:

$$\dot{P} = D \hat{v}_t + P J_t; \quad (48)$$

Let  $Z$  denote the set of asymptotically stable equilibria of (48), i.e., the local minima of  $J$ , and let  $Z^*$  be the  $\epsilon$ -nbd of  $Z$ , i.e.,  $Z^* = \{x \in \mathcal{C} \mid \exists z \in Z, \|x - z\| < \epsilon\}$ . Further, let  $\mathcal{Y}$  denote the set of asymptotically stable equilibria of (47). We have

**Theorem 1.** Under Assumptions A1–A3, given  $\epsilon > 0$ ,  $\eta > 0$  such that for  $t, t_0$  obtained using Algorithm 1, if  $\sup_{t \in [t_0, t_0 + \epsilon]} \|k_t - k^*\| < \eta$ , then  $t \in Z^*$  as  $t \rightarrow \infty$ , with probability one.

**Proof.** Let  $F_{2,t}/D$  denote the sequence of  $\epsilon$ -fields generated by  $r_t, r_0 = 0$ . We have

$$t \in \mathcal{D} \cap \mathcal{O}_t \quad t \in \mathcal{E} \cap \mathcal{S}_{at} \quad F_{2,t}/D \quad t \in \mathcal{D} \cap \mathcal{O}_t \quad t \in \mathcal{E} \cap \mathcal{S}_{at}$$

where  $t \in \mathcal{D} \cap \mathcal{O}_t \quad t \in \mathcal{E} \cap \mathcal{S}_{at} \quad F_{2,t}/D$  and  $t \in \mathcal{D} \cap \mathcal{O}_t \quad t \in \mathcal{E} \cap \mathcal{S}_{at} \quad F_{2,t}/D$ . Here,  $t$  is the policy corresponding to  $t$ . Because the critic converges along the faster timescale, from Lemma 5, it follows that  $t \in \mathcal{D} \cap \mathcal{O}_t \quad t \in \mathcal{E} \cap \mathcal{S}_{at}$ . Now let  $M^2_t/D$  be the quantities  $t$  can be seen to be uniformly bounded because from the proof in Lemma 5,  $f_{tC1g}$  and  $f_{tV1g}$  are bounded sequences. It is now easy to see (Bhatnagar & Kumar, 2004) using (20) that  $fM^2_t/g$  is a convergent martingale sequence.

Thus, for any  $T > 0$ , with  $n_T = \min_{r \in \mathcal{R}} \inf_{t \in [0, T]} n_j = \min_{r \in \mathcal{R}} \inf_{t \in [0, T]} n_j$ , we have that  $n_T = \min_{r \in \mathcal{R}} \inf_{t \in [0, T]} n_j = \min_{r \in \mathcal{R}} \inf_{t \in [0, T]} n_j$ .

Next, it can be seen using similar arguments as before (see proof of Lemma 4) that  $\mathcal{E} \cap \mathcal{S}_{at} \quad t \in \mathcal{D} \cap \mathcal{O}_t \quad h^1_t$ , where  $\mathcal{E} \cap \mathcal{S}_{at} \quad t \in \mathcal{D} \cap \mathcal{O}_t \quad h^1_t$ . We now show that  $h^1$  is Lipschitz continuous. Here  $v^t$  corresponds to the weight vector to which the critic update converges along the faster timescale when the corresponding policy is  $t$  (see Lemma 5). A simple calculation shows that for  $s \in \mathcal{S}, a \in \mathcal{A}, r \in \mathcal{R}$ ,  $t \in \mathcal{T}$ ,  $s/a/r$  exists and is bounded. Further, from (18), it can be seen that  $d^t(s, s) \in \mathcal{S}$  are continuously differentiable in  $t$  and have bounded derivatives. Also,  $J_t$  is continuously differentiable as well and has bounded derivative as can also be seen from (41). Further,  $v^t$  can be seen to be continuously differentiable with bounded derivatives. Thus,  $h^1_t$  is a Lipschitz continuous function and the ODE (47) is well posed. Using Hirsch lemma (Theorem 1, p. 339 of Hirsch, 1989; see also Lemma 6 of Bhatnagar et al., 2009), it is easy to see that  $t \in \mathcal{D} \cap \mathcal{O}_t \quad h^1_t$  w.p. 1. Now as  $\sup_{k \in \mathcal{K}} \|k_t - k^*\| \rightarrow 0$ , the trajectories of (47) converge to those of (48) uniformly on compacts for the same initial condition in both. The claim follows from the Hirsch lemma (Theorem 1, p. 339 of Hirsch, 1989; also given as Lemma 6 in Bhatnagar et al. (2009)). See Theorem 2 of Bhatnagar et al. (2009) for detailed arguments.

**Remark 2.** From Theorem 1, it follows that if the error term  $\sum_{s \in \mathcal{S}} d(s, r) \psi(s, r) v^> f_s$  is small, the algorithm will converge almost surely to a small neighborhood of a local minimum of  $J$ . (For the original problem, this corresponds to a small neighborhood of a local maximum of  $J$ .) Note that, in principle, the stochastic approximation scheme may get trapped in an unstable equilibrium. Pemantle (1990), with noise assumed to be sufficiently 'omnidirectional' in addition, showed that convergence to unstable fixed points will not occur; see also Brandiere (1998) for conditions on avoidance of unstable equilibria that lie in certain compact connected chain recurrent sets. However, in most cases (even without extra noise conditions) due to the inherent randomness, stochastic approximation algorithms converge to stable equilibria.

We discuss now the difficulties involved in proving boundedness of iterates when projection  $\mathcal{O}_P$  is not used in (40). Suppose we rewrite  $h^1_P$  as  $h^1_t$ .  $\sum_{s \in \mathcal{S}} d(s, r) \psi(s, r) v^> f_s$ . Note here that we write  $\mathcal{E} \cap \mathcal{S}_{at}$  in place of  $\mathcal{E} \cap \mathcal{S}_{at}$  in order to show explicit dependence of  $\mathcal{E} \cap \mathcal{S}_{at}$  on  $t$ . Then defining  $h^1_t = \sum_{s \in \mathcal{S}} d(s, r) \psi(s, r) v^> f_s$ , one obtains  $h^1_t = \sum_{s \in \mathcal{S}} d(s, r) \psi(s, r) v^> f_s$ . It is not clear whether the limit above exists because of the complex dependence of  $d$  and  $v$  on  $t$ . Note that  $v$  is obtained as a solution to a linear system of equations (see Lemma 5) with the matrix  $D$  therein also depending on  $t$ . Assumption (A1') on p. 454 in Borkar and Meyn (2000) considers the case where the above limits may not exist. However, it requires that for



$r \in R$  and  $t \in T$ , for some  $R, T > 0$ , the trajectories  $\phi_t$  of the ODE  $\dot{p}_t \in \frac{h^1(r, t)}{r}$  should lie within a ball of radius  $1/2$  around the origin. This can be shown provided the origin is a unique asymptotically stable attractor for the above ODEs for all  $r \in R$ . Again, it is not clear if this is the case here. Next, note that the methods described by Abounadi et al. (2001) and Tsitsiklis (1994) for stability of iterates are for different classes of algorithms, largely of the Q-learning type, and are not directly applicable in our setting.

Finally, we discuss the use of the stochastic Lyapunov function method (Kushner & Yin, 1997) for stability of iterates in (40). The prime requirement here is that there exists a real-valued nonnegative function  $W_t$  that satisfies

$$E[W_t | \mathcal{F}_t] \leq W_t + K_t \quad (48)$$

for all  $t \geq 0$  and  $W_t \leq g_t$ , where  $K_t \geq 0$  is continuous on  $Q$ . Then by Thm. 4.1, pp. 80–81 of Kushner and Yin (1997), the stability and convergence of iterates would follow. Hence consider the recursion (40). By Taylor's expansion for "small"  $t$  assuming a smooth  $W_t$ , one gets

$$E[W_t | \mathcal{F}_t] \leq W_t + t \nabla W_t \cdot \nabla \phi_t + \frac{t^2}{2} \nabla^2 W_t : \nabla^2 \phi_t + O(t^3) \quad (49)$$

It appears difficult to obtain such a  $W_t$  here. On the other hand, if we use the look up table representation (viz.,  $d_2 \in \mathbb{R}^n$  in Assumption (A3) or that  $\phi_t$  is as in (12)), then from Lemma 4 above, as also Theorem 1 of Sutton et al. (2000), one would get  $E[W_t | \mathcal{F}_t] \leq W_t + t \nabla W_t \cdot \nabla \phi_t + O(t^3)$ . Then  $W_t \leq g_t$  would serve as a Lyapunov function and the iterates (40) (without the projection) will be bounded and almost surely convergent, in lieu of Theorem 4.1 of Kushner and Yin (1997). It is only because of the use of function approximation in the iterates that a Lyapunov function is hard to obtain. However, in our experiments, we do not use projection but still observe that the iterates remain bounded and convergence is achieved.

Note that if function approximation is not used,  $J_t$  also serves as a Lyapunov function for the ODE associated with (40) without the projection. When function approximation is used (as with our case), the above problem of finding a suitable Lyapunov function (now) for the associated ODE also carries over and it is difficult to suitably characterize the set of stable attractors.

Remark 1 and many of the arguments in the analysis of Alg. 1 are valid for the analysis of the other algorithms. We skip the details in such cases to avoid repetition.

## 5.2. Convergence analysis for Algorithm 2

The analysis in Lemma 5 of the recursions for the average reward (27), TD error (13), and critic (28) proceeds in the same manner as for Algorithm 1. We thus concentrate on showing convergence of the recursion for the inverse of the Fisher information matrix (26) and the actor recursion (29). We assume (A1)–(A4) for our analysis here. We have

**Lemma 6.** For any given parameter  $\gamma, G_t^{-1}, t \geq 1$ , in (26) satisfy  $G_t^{-1} \leq G_t^{-1} + \gamma$  as  $t \geq 1$  with probability one.

**Proof.** It is easy to see from recursion (25) that  $G_t \leq G_t + \gamma$  as  $t \geq 1$  with probability one, for any given  $\gamma$  held fixed. Now for fixed  $\gamma$ , we have,

$$\begin{aligned} k G_t^{-1} &\leq G_t^{-1} k \leq D k G_t^{-1} G_t \leq G_t / G_t^{-1} k \\ \sup_{t \geq 1} k G_t^{-1} k &\leq \sup_{t \geq 1} k G_t^{-1} k \leq k G_t \leq G_t k \\ &\leq 0 \text{ as } t \rightarrow \infty; \end{aligned}$$

by (A4). In the above,  $I$  denotes the  $d_2 \times d_2$  identity matrix. The inequality above follows from the property on induced matrix norms (see Proposition A.12 of Bertsekas & Tsitsiklis, 1989). The claim follows.

As with Algorithm 1, we consider again the transformed problem with rewards replaced with costs (see above). This transformation, however, only affects the actor recursion (29) that now becomes

$$\dot{c}_t = D c_t - \dot{G}_t^{-1} \nabla c_t \cdot \nabla \phi_t \quad (50)$$

We have

**Theorem 2.** Under Assumptions (A1)–(A4), given  $\gamma > 0, \eta > 0$ , such that for  $t, t_0$  obtained using Algorithm 2, if  $\sup_{t \geq t_0} k e^{-t} k < \gamma$ , then  $\phi_t \rightarrow \phi^*$  as  $t \rightarrow \infty$ , with probability one.

**Proof.** As with the proof of Theorem 1, let  $F_3(t) = \nabla c_t / \nabla \phi_t$ . Note that

$$\dot{c}_t = D c_t - \dot{G}_t^{-1} \nabla c_t \cdot \nabla \phi_t \leq F_3(t) \leq \nabla c_t / \nabla \phi_t$$

where  $F_3(t) = D c_t - \dot{G}_t^{-1} \nabla c_t \cdot \nabla \phi_t \leq \nabla c_t / \nabla \phi_t$ . In lieu of Lemmas 5 and 6,  $\nabla c_t / \nabla \phi_t \leq 0.1$ . As before, the critic recursion (28) converges faster for a given policy  $\phi_t$ , corresponding to an actor update  $\phi_t$ , and converges to  $v^*$ . For  $t \geq 1$ , let  $M^3(t) = \nabla c_t / \nabla \phi_t$ . The quantities  $\nabla c_t$  and  $G_t^{-1}$  are uniformly bounded from Lemmas 5 and 6, and from (A4) respectively. Now using (20), it can be seen (Bhatnagar & Kumar, 2004) that  $M^3(t)$  is a convergent martingale sequence. Hence,  $\nabla c_t / \nabla \phi_t \rightarrow 0$  a.s. as  $t \rightarrow \infty$ , with  $n_T$  as before (see proof of Theorem 1). As before, also note that

$$\begin{aligned} E[G_t^{-1} \nabla c_t \cdot \nabla \phi_t] &\leq E[G_t^{-1} \nabla c_t \cdot \nabla \phi_t] \leq \sum_{s \in S} d^s \cdot \nabla c_t / \nabla \phi_t \\ &\leq \sum_{s \in S} \nabla c_t / \nabla \phi_t \leq \sum_{s \in S} P(s; a_t^0 / V^{t_0}) \leq \gamma \end{aligned}$$

Consider now the ODE associated with (50) (in lieu of Lemma 4)

$$\dot{c}_t = D c_t - \dot{G}_t^{-1} \nabla c_t \cdot \nabla \phi_t \quad (51)$$

Consider also the associated ODE

$$\dot{c}_t = D c_t - \nabla c_t \cdot \nabla \phi_t \quad (52)$$

As with Theorem 1, from the Hirsch lemma,  $\phi_t \rightarrow \phi^*$  as  $t \rightarrow \infty$  w.p. 1. Now as  $\sup_{t \geq 1} k e^{-t} k \leq 0$ , the trajectories of (51) converge to those of (52) uniformly on compact sets for the same initial condition. See the proof of Theorem 3 of Bhatnagar et al. (2009) for details. The claim follows.

## 5.3. Convergence analysis for Algorithm 3

As stated previously, the main idea in this algorithm is to minimize the least-squares error in estimating the advantage function via function approximation. The analysis of the average reward (31), TD error (13), and critic (32) recursions proceeds in the same manner as before (cf. Lemma 5). We thus concentrate on recursion (33) and the actor recursion (34). We require Assumptions (A1)–(A3) here. In the transformed problem (with costs in place of rewards), recursion (33) can be rewritten as

$$\dot{w}_t = D w_t - \nabla w_t \cdot \nabla \phi_t \quad (53)$$

with the actor recursion (34) the same as before. Note that (53) moves on a faster timescale as compared to the actor recursion. Hence, on the timescale of the former recursion, one may consider the parameter  $\phi_t$  to be fixed. We have the following result:

**Lemma 7.** Under a given parameter  $\gamma$ , the  $w_t$  in (53) satisfy  $w_t \leq G_t^{-1} \nabla w_t \cdot \nabla \phi_t$  as  $t \geq 1$  with probability one, where  $\phi_t$  is the policy corresponding to  $\gamma$ .

**Proof.** Consider the following ODE associated with (53) for given

$$\dot{W} = D_{s_t} E_{s_t} \left( \frac{g^2}{r} \right) - \frac{1}{r} \left( \frac{g^2}{r} \right) W_t \quad (54)$$

Note that  $g^2/r$  is Lipschitz continuous in  $W$ . Now let  $g^2/r \in D_{s_t} E_{s_t} \left( \frac{g^2}{r} \right)$ . The function  $g^2/r$  exists and can be seen to satisfy  $g^2/r \in D_{s_t} E_{s_t} \left( \frac{g^2}{r} \right)$ . For the ODE  $\dot{W} = D_{s_t} E_{s_t} \left( \frac{g^2}{r} \right) - \frac{1}{r} \left( \frac{g^2}{r} \right) W_t$ , the origin is an asymptotically stable equilibrium (because  $G_{s_t}/r$  is positive definite). Define now  $f_{M^4}(t)$  as  $M^4(t) \in D_{s_t} E_{s_t} \left( \frac{g^2}{r} \right) - \frac{1}{r} \left( \frac{g^2}{r} \right) W_t$ , where  $\hat{G}^2(t) \in D_{s_t} E_{s_t} \left( \frac{g^2}{r} \right) - \frac{1}{r} \left( \frac{g^2}{r} \right) W_t$  and  $F_4(t) \in D_{s_t} E_{s_t} \left( \frac{g^2}{r} \right) - \frac{1}{r} \left( \frac{g^2}{r} \right) W_t$ . It is easy to see that there exists a constant  $C_0 < 1$  such that

$$E \left[ \|M^4(t)\| \right] \leq C_0 \left( \|W_t\| + \|F_4(t)\| \right)$$

for all  $t \geq 0$ . For the ODE (54), it can be easily verified that  $W \in D_{s_t} E_{s_t} \left( \frac{g^2}{r} \right) - \frac{1}{r} \left( \frac{g^2}{r} \right) W_t$  is an asymptotically stable equilibrium (see Lemma 8 of Bhatnagar et al., 2009). Now from Thm. 2.2 of Borkar and Meyn (2000), (53) converges with probability one to  $W$ .

We now consider the actor recursion (34), which is the slower recursion.

**Theorem 3.** Under Assumptions A1–A3, given  $\gamma > 0$ ,  $\rho > 0$  such that for  $t, t' \geq 0$ , obtained using Algorithm 3, if  $\sup_t \|k_t\| \leq \rho$ , then  $\|k_t\| \leq \rho$  as  $t \rightarrow \infty$ , with probability one.

**Proof.** Note that the recursion (34) can be written as

$$k_{t+1} = \gamma \left( \frac{1}{r} \left( \frac{g^2}{r} \right) - \frac{1}{r} \left( \frac{g^2}{r} \right) W_t \right) + (1-\gamma) k_t$$

where  $\gamma \in (0, 1]$  by Lemma 7. The rest can be shown in a similar manner as Theorem 2.

#### 5.4. Convergence analysis for Algorithm 4

As with Algorithm 2, we require Assumptions (A1)–(A4). The result in Lemma 6 continues to hold here and we get for fixed  $\gamma$ ,  $G_t \rightarrow G$  as  $t \rightarrow \infty$  with probability one. The recursions for average reward (36), TD error (13), and critic (37) are the same as before and have been analyzed earlier (cf. Lemma 5). We now concentrate on recursion (38) and the actor recursion (39). Under the transformed problem (with costs in place of rewards), recursion (38) can be rewritten as

$$W_{t+1} = \gamma \left( \frac{1}{r} \left( \frac{g^2}{r} \right) - \frac{1}{r} \left( \frac{g^2}{r} \right) W_t \right) + (1-\gamma) W_t \quad (55)$$

with the actor recursion the same as before. An exactly similar result as Lemma 7 holds in this case as well (below); see Lemma 9 of Bhatnagar et al. (2009) for a proof.

**Lemma 8.** Under a given parameter  $\gamma$ , the  $W_t, t \geq 1$ , defined by (55) satisfy  $W_t \rightarrow G$  as  $t \rightarrow \infty$  with probability one, with  $\pi$  being the policy corresponding to  $G$ .

We finally consider the actor recursion (39) and have the following result whose proof follows as in Theorems 2 and 3.

**Theorem 4.** Under Assumptions A1–A4, given  $\gamma > 0$ ,  $\rho > 0$  such that for  $t, t' \geq 0$ , obtained using Algorithm 4, if  $\sup_t \|k_t\| \leq \rho$ , then  $\|k_t\| \leq \rho$  as  $t \rightarrow \infty$ , with probability one.

## 6. Relation to the previous algorithms

As we mentioned in Section 1, the actor critic algorithms presented in this paper extend prior actor critic methods, especially those of Konda and Tsitsiklis (2003) and of Peters et al. (2003). In this section, we discuss these relationships further.

**Actor–Critic Algorithm of Konda and Tsitsiklis (2003):** Contrary to Algorithms 2–4, this algorithm does not use estimates of the natural gradient in its actor's update. It is somewhat similar to our Algorithm 1, but with some key differences. (1) Konda's algorithm uses the Markov process of state–action pairs and thus its critic update is based on an action-value function. Algorithm 1 uses the state process and therefore its critic update is based on a value function. (2) Whereas Algorithm 1 uses TD error in both critic and actor recursions, Konda's algorithm uses TD error only in its critic update. The actor recursion in Konda's algorithm uses a  $Q$ -value estimate instead. Because the TD error is an unbiased estimate of the advantage function (Lemma 3), the actor recursion in Algorithm 1 uses estimates of advantages instead of  $Q$ -values, which may result in lower variances. (3) The convergence analysis of Konda's algorithm is based on the martingale approach and aims at bounding error terms and directly showing convergence. Convergence to a local optimum is shown when TD(1) critic is used. For the case when  $\gamma < 1$ , they show that given  $\epsilon > 0$ , there exists  $\delta$  close enough to one such that when a TD  $\gamma$ -critic is used, one gets  $\liminf_t J_t \geq J^* - \epsilon$  with probability one. Unlike Konda and Tsitsiklis, we primarily use the ordinary differential equation (ODE) approach for our convergence analysis. Even though we also use martingale arguments in our analysis, these are restricted to showing that the noise terms asymptotically diminish and the resulting scheme can be viewed as a Euler-discretization of the associated ODE.

**Natural Actor–Critic Algorithm of Peters et al. (2003):** Algorithms 2–4 extend this algorithm by being fully incremental and by providing convergence proofs. Peters's algorithm uses a least-squares TD method in its critic's update, whereas our algorithms are all fully incremental. It is not entirely clear how to satisfactorily incorporate least-squares TD methods in a context in which the policy is changing. Our proof techniques do not immediately extend to this case. However, we use estimates of the advantage function in Algorithms 3 and 4 as in Peters's algorithm.

## 7. Conclusions and future work

We have introduced and analyzed four actor–critic reinforcement learning algorithms utilizing linear function approximation. All the algorithms are based on existing ideas such as temporal difference learning, natural policy gradients, and two-timescale convergence analysis, but we combine them in new ways. The main contribution of this paper is the proof of convergence of the four algorithms to a local maximum in the space of policy and value-function parameters. Our four algorithms are the first actor–critic algorithms to be shown convergent that utilize both function approximation and bootstrapping, a combination which seems essential to large-scale applications of reinforcement learning.

Our Algorithms 2–4 are explorations of the use of natural gradients within an actor–critic policy-gradient architecture. The way we use natural gradients is distinctive in that it is totally incremental: the policy is changed on every time step yet we never reset the gradient computation as is done in the algorithm of Peters and Schaal (2008). Algorithm 3 is perhaps the most interesting of the three natural-gradient algorithms. It never explicitly stores an estimate of the inverse of the Fisher information matrix, and as a result, it requires less computation. In our empirical experiments (Bhatnagar et al., 2009), we found it easier to find

good parameter settings for Algorithm 3 than for the other natural-gradient algorithms, and perhaps because of this, it converged more rapidly than them and than Konda's and Tsitsiklis' algorithm. These empirical observations should be taken only as suggestive; more experiments to properly assess the relative performance of these algorithms must be carried out.

The most important potential extension of our results would be to characterize the quality of the converged solution. It may be possible to bound the performance loss due to bootstrapping and approximation error in a way similar to how it was bounded by Tsitsiklis and Van Roy (1997). Because of the use of function approximation, our convergence analysis would carry through for the case of continuously valued state-action spaces as well. It would be interesting to study empirical evaluations of our algorithms in such settings in order to evaluate their applicability in such scenarios. There are a number of other ways in which our results are limited and suggest future work. (1) There is the issue of rate of convergence. Ideally one would like analytic results but, short of that, it would be useful to conduct a thorough empirical study, varying parameters and schedules in a more extensive and sophisticated way than what we have done in Bhatnagar et al. (2009). (2) The algorithms could be extended to incorporate eligibility traces and least-squares methods. As discussed earlier, the former seems straightforward whereas the latter seems to require more fundamental extensions. (3) A thorough study of the sensitivity of our algorithms to the various system parameters and settings is needed. (4) A study of the choice of the basis functions for the critic to obtain a good estimate of the policy gradient needs to be done. (5) Application of these ideas and algorithms to a real-world problem is needed to assess their ultimate utility.

## References

- Abdulla, M. S., & Bhatnagar, S. (2007). Reinforcement learning based algorithms for average cost Markov decision processes. *Discrete Event Dynamic Systems: Theory and Applications*, 17(1), 23–52.
- Abounadi, J., Bertsekas, D., & Borkar, V. S. (2001). Learning algorithms for Markov decision processes. *SIAM Journal on Control and Optimization*, 40, 681–698.
- Aleksandrov, V., Sysoyev, V., & Shemeneva, V. (1968). Stochastic optimization. *Engineering Cybernetics*, 5, 11–16.
- Amari, S. (1998). Natural gradient works efficiently in learning. *Neural Computation*, 10(2), 251–276.
- Baird, L. (1993). Advantage updating. *Technical Report WL-TR-93-1146*, Wright Laboratory, OH.
- Baird, L. (1995). Residual algorithms: Reinforcement learning with function approximation. In *Proceedings of the twelfth international conference on machine learning* (pp. 30–37).
- Bagnell, J., & Schneider, J. (2003). Covariant policy search. In *Proceedings of international joint conference on artificial intelligence* (pp. 1019–1024).
- Barto, A., Sutton, R. S., & Anderson, C. (1983). Neuron-like elements that can solve difficult learning control problems. *IEEE Transactions on Systems, Man and Cybernetics*, 13, 835–846.
- Baxter, J., & Bartlett, P. (2001). Infinite-horizon policy-gradient estimation. *Journal of Artificial Intelligence Research*, 15, 319–350.
- Bellman, R. E., & Dreyfus, S. E. (1959). Functional approximations and dynamic programming. *Mathematical Tables and Other Aids to Computation*, 13, 247–251.
- Benveniste, A., Metivier, M., & Priouret, P. (1990). *Adaptive algorithms and stochastic approximations*. Berlin: Springer.
- Bertsekas, D. (1999). *Nonlinear programming*. Belmont, MA: Athena Scientific.
- Bertsekas, D., & Tsitsiklis, J. (1989). *Parallel and distributed computation*. New Jersey: Prentice Hall.
- Bertsekas, D., & Tsitsiklis, J. (1996). *Neuro-dynamic programming*. Belmont, MA: Athena Scientific.
- Bhatnagar, S., & Kumar, S. (2004). A simultaneous perturbation stochastic approximation based actor-critic algorithm for Markov decision processes. *IEEE Transactions on Automatic Control*, 49(4), 592–598.
- Bhatnagar, S. (2005). Adaptive multivariate three-timescale stochastic approximation algorithms for simulation based optimization. *ACM Transactions on Modeling and Computer Simulation*, 15(1), 74–107.
- Bhatnagar, S. (2007). Adaptive Newton-based multivariate smoothed functional algorithms for simulation optimization. *ACM Transactions on Modeling and Computer Simulation*, 18(1), 2:1–2:35.
- Bhatnagar, S., Sutton, R. S., Ghavamzadeh, M., & Lee, M. (2009). Natural actor-critic algorithms. Technical Report, *Department of Computing Science, University of Alberta, Canada* [<http://www.cs.ualberta.ca/research/techreports/2009/TR09-10.php>].
- Bhatnagar, S., Sutton, R. S., Ghavamzadeh, M., & Lee, M. (2008). Incremental natural actor-critic algorithms. *Advances in Neural Information Processing Systems*, 20, 105–112.
- Borkar, V. S. (1997). Stochastic approximation with two timescales. *Systems and Control Letters*, 29, 291–294.
- Borkar, V. S. (2008). Reinforcement learning—a bridge between numerical methods and Monte-Carlo. *Preprint*.
- Borkar, V. S., & Meyn, S. (2000). The O.D.E. method for convergence of stochastic approximation and reinforcement learning. *SIAM Journal on Control and Optimization*, 38(2), 447–469.
- Boyan, J. (1999). Least-squares temporal difference learning. In *Proceedings of the sixteenth international conference on machine learning* (pp. 49–56).
- Boyan, J., & Moore, A. (1995). Generalization in reinforcement learning: Safely approximating the value function. *Advances in Neural Information Processing Systems*, 7, 369–376.
- Brandiere, O. (1998). Some pathological traps for stochastic approximation. *SIAM Journal on Control and Optimization*, 36, 1293–1314.
- Bradtke, S., & Barto, A. (1996). Linear least-squares algorithms for temporal difference learning. *Machine Learning*, 22, 33–57.
- Cao, X., & Chen, H. (1997). Perturbation realization, potentials and sensitivity analysis of Markov processes. *IEEE Transactions on Automatic Control*, 42, 1382–1393.
- Crites, R., & Barto, A. (1998). Elevator group control using multiple reinforcement learning agents. *Machine Learning*, 33, 235–262.
- Farahmand, A. M., Ghavamzadeh, M., Szepesvári, Cs., & Mannor, S. (2009). Regularized policy iteration. *Advances in Neural Information Processing Systems*, 21, 441–448.
- Ghavamzadeh, M., & Mahadevan, S. (2003). Hierarchical policy gradient algorithms. In *Proceedings of the twentieth international conference on machine learning* (pp. 226–233).
- Ghavamzadeh, M., & Engel, Y. (2007a). Bayesian policy gradient algorithms. *Advances in Neural Information Processing Systems*, 19, 457–464.
- Ghavamzadeh, M., & Engel, Y. (2007b). Bayesian actor-critic algorithms. In *Proceedings of the twenty-fourth international conference on machine learning* (pp. 297–304).
- Glynn, P. (1990). Likelihood ratio gradient estimation for stochastic systems. *Communications of the ACM*, 33, 75–84.
- Gordon, G. (1995). Stable function approximation in dynamic programming. In *Proceedings of the twelfth international conference on machine learning* (pp. 261–268).
- Greensmith, E., Bartlett, P., & Baxter, J. (2004). Variance reduction techniques for gradient estimates in reinforcement learning. *Journal of Machine Learning Research*, 5, 1471–1530.
- Hirsch, M. (1989). Convergent activation dynamics in continuous time networks. *Neural Networks*, 2, 331–349.
- Kakade, S. (2002). A natural policy gradient. *Advances in Neural Information Processing Systems*, 14.
- Kohl, N., & Stone, P. (2004). Policy gradient reinforcement learning for fast quadrupedal locomotion. In *Proceedings of the IEEE international conference on robotics and automation* (pp. 2619–2624).
- Konda, V., & Borkar, V. S. (1999). Actor-critic like learning algorithms for Markov decision processes. *SIAM Journal on Control and Optimization*, 38(1), 94–123.
- Konda, V., & Tsitsiklis, J. (2003). On actor-critic algorithms. *SIAM Journal on Control and Optimization*, 42(4), 1143–1166.
- Kushner, H., & Clark, D. (1978). *Stochastic approximation methods for constrained and unconstrained systems*. New York: Springer Verlag.
- Kushner, H., & Yin, G. (1997). *Stochastic approximation algorithms and applications*. New York: Springer Verlag.
- Lagoudakis, M., & Parr, R. (2003). Least-squares policy iteration. *Journal of Machine Learning Research*, 4, 1107–1149.
- Marbach, P., & Tsitsiklis, J. (2001). Simulation-based optimization of Markov reward processes. *IEEE Transactions on Automatic Control*, 46, 191–209.
- Meyn, S. (2007). *Control techniques for complex networks*. Cambridge, UK: Cambridge Univ. Press.
- Ng, A., Coates, A., Diel, M., Ganapathi, V., Schulte, J., & Tse, B. et al. (2004). Inverted autonomous helicopter flight via reinforcement learning. In *International symposium on experimental robotics*.
- Pemantle, R. (1990). Nonconvergence to unstable points in urn models and stochastic approximations. *Annals of Probability*, 18, 698–712.
- Peters, J., Vijayakumar, S., & Schaal, S. (2003). Reinforcement learning for humanoid robotics. In *Proceedings of the third IEEE-RAS international conference on humanoid robots*.
- Peters, J., & Schaal, S. (2008). Natural actor-critic. *Neurocomputing*, 71(7–9), 1180–1190.
- Puterman, M. (1994). *Markov decision processes: Discrete stochastic dynamic programming*. New York: John Wiley.
- Richter, S., Aberdeen, D., & Yu, J. (2007). Natural actor-critic for road traffic optimization. *Advances in Neural Information Processing Systems*, 19, 1169–1176.
- Rust, J. (1996). Numerical dynamic programming in economics. In *Handbook of computational economics* (pp. 614–722). Amsterdam: Elsevier.
- Sutton, R. S. (1984). Temporal credit assignment in reinforcement learning. *Doctoral dissertation*, Amherst: University of Massachusetts.

- Sutton, R. S. (1988). Learning to predict by the method of temporal differences. *Machine Learning*, 3, 9–44.
- Sutton, R. S. (1996). Generalization in reinforcement learning: Successful examples using sparse coarse coding. *Advances in Neural Information Processing Systems*, 8, 1038–1044.
- 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, 1057–1063.
- Sutton, R. S., & Barto, A. (1998). *Reinforcement learning: An introduction*. Cambridge, MA: MIT Press.
- Tadic, V. (2001). On the convergence of temporal difference learning with linear function approximation. *Machine Learning*, 42(3), 241–267.
- Tesauro, G. (1995). Temporal difference learning and TD-Gammon. *Communications of the ACM*, 38, 58–68.
- Tsitsiklis, J. (1994). Asynchronous stochastic approximation and Q-learning. *Machine Learning*, 16, 185–202.
- Tsitsiklis, J., & Van Roy, B. (1997). An analysis of temporal-difference learning with function approximation. *IEEE Transactions on Automatic Control*, 42(5), 674–690.
- Tsitsiklis, J., & Van Roy, B. (1999). Average cost temporal-difference learning. *Automatica*, 35, 1799–1808.
- White, D. (1993). A survey of applications of Markov decision processes. *Journal of the Operational Research Society*, 44, 1073–1096.
- Widrow, B., & Stearns, S. (1985). *Adaptive signal processing*. Englewood Cliffs, NJ: Prentice-Hall.
- Williams, R. (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8, 229–256.

**Shalabh Bhatnagar** received his Ph.D. in Electrical Engineering from the Indian Institute of Science, Bangalore, in 1997. He was a postdoctoral research associate from 1997 to 2000 at the Institute for Systems Research, University of Maryland, College Park and from 2000 to 2001 at the Free University, Amsterdam, Netherlands. He is currently an Associate Professor in the Department of Computer Science and Automation, Indian Institute of Science, Bangalore. He has held visiting positions at the Indian Institute of Technology, Delhi and the RLAI research laboratory, University of Alberta, Canada.

His research interests are in the areas of stochastic control and simulation-based stochastic optimization with applications specifically in communication and wireless networks.

**Richard S. Sutton** is a professor and iCORE chair in the department of computing science at the University of Alberta. He is a fellow of the Association for the Advancement of Artificial Intelligence and co-author of the textbook *Reinforcement Learning: An Introduction*. Before joining the University of Alberta in 2003, he worked in industry at AT&T and GTE Labs, and in academia at the University of Massachusetts. He received a Ph.D. in computer science from the University of Massachusetts in 1984 and a BA in psychology from Stanford University in 1978. Rich's research interests center on the learning problems facing a decision-maker interacting with its environment, which he sees as central to artificial intelligence. He is also interested in animal learning psychology, in connectionist networks, and generally in systems that continually improve their representations and models of the world.

**Mohammad Ghavamzadeh** received a Ph.D. degree in Computer Science from the University of Massachusetts Amherst in 2005. From 2005 to 2008, he was a postdoctoral fellow at the Department of Computing Science at the University of Alberta, Canada. He has been a researcher at INRIA Lille - Nord Europe, Team Sequel in France since November 2008. The main objective of his research is to investigate the principles of scalable decision-making under uncertainty. In the last four years, Mohammad's research has been mostly focused on using recent advances in statistical machine learning, especially Bayesian reasoning and kernel methods, to develop more scalable reinforcement learning algorithms.

**Mark Lee** is a professional web developer and programmer, and the co-author of the book "C++ Programming for the Absolute Beginner". He received a B.Sc. degree in computer science from the University of Alberta in 2005.