

Derivatives of Logarithmic Stationary Distributions for Policy Gradient Reinforcement Learning

Tetsuro Morimura

tetsuro@jp.ibm.com

IBM Research – Tokyo, Yamato, Kanagawa 242-8502, Japan

Eiji Uchibe

uchibe@oist.jp

Okinawa Institute of Science and Technology, Uruma, Okinawa 904-2234, Japan

Junichiro Yoshimoto

jun-y@oist.jp

*Okinawa Institute of Science and Technology, Uruma, Okinawa 904-2234, Japan,
and Nara Institute of Science and Technology, Ikoma, Nara 630-0192, Japan*

Jan Peters

jan.peters@tuebingen.mpg.de

Max Planck Institute for Biological Cybernetics, 72076, Tübingen, Germany

Kenji Doya

doya@oist.jp

*Okinawa Institute of Science and Technology, Uruma, Okinawa 904-2234, Japan;
Nara Institute of Science and Technology, Ikoma, Nara 630-0192, Japan; and ATR
Computational Neuroscience Laboratories, Soraku, Kyoto 619-0288, Japan*

Most conventional policy gradient reinforcement learning (PGRL) algorithms neglect (or do not explicitly make use of) a term in the average reward gradient with respect to the policy parameter. That term involves the derivative of the stationary state distribution that corresponds to the sensitivity of its distribution to changes in the policy parameter. Although the bias introduced by this omission can be reduced by setting the forgetting rate γ for the value functions close to 1, these algorithms do not permit γ to be set exactly at $\gamma = 1$. In this article, we propose a method for estimating the log stationary state distribution derivative (LSD) as a useful form of the derivative of the stationary state distribution through backward Markov chain formulation and a temporal difference learning framework. A new policy gradient (PG) framework with an LSD is also proposed, in which the average reward gradient can be estimated by setting $\gamma = 0$, so it becomes unnecessary to learn the value functions. We also test the performance of the proposed algorithms using simple

benchmark tasks and show that these can improve the performances of existing PG methods.

1 Introduction

Policy gradient reinforcement learning (PGRL) is a popular family of algorithms in reinforcement learning (RL). It improves a policy parameter to maximize the average reward (also called the expected return) by using the average reward gradients with respect to the policy parameter, which are called the policy gradients (PGs) (Gullapalli, 1990; Williams, 1992; Kimura & Kobayashi, 1998; Baird & Moore, 1999; Sutton, McAllester, Singh, & Mansour, 2000; Baxter & Bartlett, 2001; Konda & Tsitsiklis, 2003; Peters & Schaal, 2006). However, most conventional PG algorithms for infinite-horizon problems neglect (or do not explicitly make use of) the term associated with the derivative of the stationary (state) distribution in the PGs, with the exception of Ng, Parr, and Koller (2000), since to date there is not an efficient algorithm to estimate this derivative. This derivative is an indicator of how sensitive the stationary distribution is to changes in the policy parameter. While the biases introduced by this omission can be reduced by using a forgetting (or discounting)¹ rate “ γ ” for the value functions close to 1, that tends to increase the variance of the PG estimates, and for $\gamma = 1$, the variance can become infinite, which violates the conditions of these algorithms. This trade-off makes it difficult to find an appropriate γ in practice. Furthermore, while the solution to discounted reinforcement learning is well defined if the optimal control solution can be perfectly represented by the policy, this is no longer true in the case where function approximation is employed. For approximations of the policies, the solution will always be determined by the start-state distributions and thus is in general an ill-defined problem. Average reward RL, on the other hand, is a well-posed problem, as it depends only on the stationary distribution.

Here, we propose a new PG framework for estimating the derivative of the logarithmic stationary state distribution (log stationary distribution derivative, or LSD) as an alternative and useful form of the derivative of the stationary distribution for estimating the PG.² It is our main result and contribution of this article that a method for estimating the LSD is derived

¹Note that the parameter γ has two different meanings: discounting and forgetting. γ is sometimes interpreted as a discounting rate to define the objective function. On the other hand, the role of γ can be regarded as the forgetting rate to enforce a horizon change for the approach of Baxter and Bartlett (2001), where the objective function is the average reward. That is, while the discounting rate is seen as part of the problem, the forgetting rate is part of algorithm. Since we focus on the average reward as the infinite-horizon problem, we use the term *forgetting rate* for γ in this article.

²While the log stationary distribution derivative with respect to the policy parameter is sometimes referred to as the *likelihood ratio gradient* or *score function*, we call it the LSD in this article.

through backward Markov chain formulation and a temporal difference learning method. Then the learning agent estimates the LSD instead of estimating the value functions in this PG framework. In addition, the use of LSD estimation will open other possibilities for RL. In particular, it allows implementing state-of-the-art natural gradient learning for RL (Morimura, Uchibe, Yoshimoto, & Doya, 2008, in press), which was reported to be especially effective in the randomly synthesized large-scale MDPs. The Fisher information matrix including the LSD can be used as the Riemannian metric, which defines the transformation of an ordinary gradient to a natural gradient.

This article is an extended version of an earlier technical report (Morimura, Uchibe, Yoshimoto, & Doya, 2007), with new results, and is organized as follows. In section 2, we review the conventional PGRL methods and describe a motivation to estimate LSD. In section 3, we propose an $\mathcal{L}\text{SLSD}(\lambda)$ algorithm for the estimation of LSD by a *Least Squares* temporal difference method based on the backward Markov chain formulation. In section 4, the $\mathcal{L}\text{SLSD}(\lambda)$ -PG algorithm is instantly derived as a novel PG algorithm utilizing $\mathcal{L}\text{SLSD}(\lambda)$. We also propose a baseline function for $\mathcal{L}\text{SLSD}(\lambda)$ -PG, which decreases the variance of the PG estimate. To verify the performances of the proposed algorithms, numerical results for simple Markov decision processes (MDPs) are shown in section 5. In section 6, we review existing (stationary) state distribution derivative estimating and average reward PG methods. In section 7, we give a summary and discussion. We also suggest other significant possibilities brought by the realization of the LSD estimation.

2 Policy Gradient Reinforcement Learning

We briefly review conventional PGRL methods and present the motivation to estimate LSD. A discrete-time MDP with finite sets of states $s \in \mathcal{S}$ and actions $a \in \mathcal{A}$ is defined by a state transition probability $p(s_{+1} | s, a) \equiv \Pr(s_{+1} | s, a)$ and a (bounded) reward function $r_{+1} = r(s, a, s_{+1})$, where s_{+1} is the state followed by the action a at the state s and r_{+1} is the observed immediate reward at s_{+1} (Bertsekas, 1995; Sutton & Barto, 1998). The state s_{+k} and the action a_{+k} denote a state and an action after k time steps from the state s and the action a , respectively, and backward for $-k$. The decision-making rule follows a stochastic policy $\pi(s, a; \theta) \equiv \Pr(a | s, \theta)$, parameterized by $\theta \in \mathcal{R}^d$. We assume the policy $\pi(s, a; \theta)$ is always differentiable with respect to θ . We also posit the following assumption:

Assumption 1. The Markov chain $M(\theta) = \{\mathcal{S}, \mathcal{A}, p, \pi, \theta\}$ is ergodic (irreducible and aperiodic) for all policy parameters θ . Then there exists a unique stationary state distribution $d_{M(\theta)}(s) = \lim_{k \rightarrow \infty} \Pr(s_{+k} = s | M(\theta)) > 0$, which is independent of the initial state and satisfies the recursion

$$d_{M(\theta)}(s) = \sum_{s_{-1} \in \mathcal{S}} \sum_{a_{-1} \in \mathcal{A}} p_{M(\theta)}(s, a_{-1} | s_{-1}) d_{M(\theta)}(s_{-1}), \quad (2.1)$$

where $p_{M(\theta)}(s, a_{-1} | s_{-1}) \equiv \pi(s_{-1}, a_{-1}; \theta)p(s | s_{-1}, a_{-1})$.

The goal of PGRL is to find a policy parameter θ^* that maximizes the time average of immediate rewards called the *average reward* or *expected return*:

$$\eta(\theta) \equiv \lim_{K \rightarrow \infty} \frac{1}{K} \mathbb{E}_{M(\theta)} \left\{ \sum_{k=1}^K r_{+k} | s \right\},$$

where $\mathbb{E}_{M(\theta)}$ denotes the expectation over the Markov chain $M(\theta)$. Under assumption 1, the average reward is independent of the initial state s and can be shown to be equal to (Bertsekas, 1995),

$$\begin{aligned} \eta(\theta) &= \mathbb{E}_{M(\theta)} \{ r(s, a, s_{+1}) \} \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \sum_{s_{+1} \in \mathcal{S}} d_{M(\theta)}(s) \pi(s, a; \theta) p(s_{+1} | s, a) r(s, a, s_{+1}) \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_{M(\theta)}(s) \pi(s, a; \theta) \bar{r}(s, a), \end{aligned} \tag{2.2}$$

where $\bar{r}(s, a) = \sum_{s_{+1} \in \mathcal{S}} p(s_{+1} | s, a) r(s, a, s_{+1})$ does not depend on the policy parameter. The policy gradient RL algorithms update the policy parameters θ in the direction of the gradient of the average reward $\eta(\theta)$ with respect to θ so that

$$\nabla_{\theta} \eta(\theta) \equiv \left[\frac{\partial \eta(\theta)}{\partial \theta_1}, \dots, \frac{\partial \eta(\theta)}{\partial \theta_d} \right]^{\top},$$

where \top denotes the transpose. This derivative is often referred to as the policy gradient (PG). Using equation 2.2, the PG is directly determined to be

$$\begin{aligned} \nabla_{\theta} \eta(\theta) &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_{M(\theta)}(s) \pi(s, a; \theta) (\nabla_{\theta} \ln d_{M(\theta)}(s) \\ &\quad + \nabla_{\theta} \ln \pi(s, a; \theta)) \bar{r}(s, a). \end{aligned} \tag{2.3}$$

However, since the derivation of the gradient of the log stationary state distribution $\nabla_{\theta} \ln d_{M(\theta)}(s)$ is nontrivial, the conventional PG algorithms (Baxter & Bartlett, 2001; Kimura & Kobayashi, 1998) utilize an alternative representation of the PG as (we give this derivation in appendix A)

$$\begin{aligned} \nabla_{\theta} \eta(\theta) &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_{M(\theta)}(s) \pi(s, a; \theta) \nabla_{\theta} \ln \pi(s, a; \theta) Q_{\gamma}^{\pi}(s, a) \\ &\quad + (1 - \gamma) \sum_{s \in \mathcal{S}} d_{M(\theta)}(s) \nabla_{\theta} \ln d_{M(\theta)}(s) V_{\gamma}^{\pi}(s), \end{aligned} \tag{2.4}$$

where

$$Q_\gamma^\pi(s, a) \equiv \lim_{K \rightarrow \infty} \mathbb{E}_{M(\theta)} \left\{ \sum_{k=1}^K \gamma^{k-1} r_{+k} \mid s, a \right\},$$

$$V_\gamma^\pi(s) \equiv \lim_{K \rightarrow \infty} \mathbb{E}_{M(\theta)} \left\{ \sum_{k=1}^K \gamma^{k-1} r_{+k} \mid s \right\},$$

are the cumulative rewards defined with forgetting rate $\gamma \in [0, 1)$, known as the state-action and state value functions, respectively (Sutton & Barto, 1998). Since the contribution of the second term in equation 2.4 shrinks as γ approaches 1 (Baxter & Bartlett, 2001), the conventional algorithms omitted the second term and approximated the PG as a biased estimate by taking $\gamma \approx 1$: the PG estimate was composed of only the first term in equation 2.4. Although the bias introduced by this omission shrinks as γ approaches 1, the variance of the estimate becomes larger (Baxter & Bartlett, 2001; Baxter, Bartlett, & Weaver, 2001). In addition, these algorithms prohibit γ from being set exactly at 1, though the bias disappears in the limit of $\gamma \rightarrow 1$.

In this article, we propose an estimation approach for the log stationary distribution derivative (LSD), $\nabla_\theta \ln d_{M(\theta)}(s)$. The realization of the LSD estimation instantly makes an alternative PG approach feasible, which uses equation 2.3 for computing the PG estimate with the LSD estimate. An important feature is that this approach does not need to learn value function Q_γ^π or V_γ^π , and therefore the resulting algorithms are free from the bias-variance trade-off in the choice of the forgetting rate γ .

3 Estimation of the Log Stationary Distribution Derivative

In this section, we propose an LSD estimation algorithm based on least squares, $\mathcal{L}SLS(\lambda)$. For this purpose, we formulate the backwardness of the ergodic Markov chain $M(\theta)$ and show that LSD can be estimated on the basis of the temporal difference learning (Sutton, 1988; Bradtke & Barto, 1996; Boyan, 2002).

3.1 Properties of Forward and Backward Markov Chains. According to Bayes's theorem, a backward probability q of a previous state-action pair (s_{-1}, a_{-1}) leading to the current state s is given by

$$q(s_{-1}, a_{-1} \mid s) = \frac{p(s \mid s_{-1}, a_{-1}) \Pr(s_{-1}, a_{-1})}{\sum_{s_{-1} \in \mathcal{S}} \sum_{a_{-1} \in \mathcal{A}} p(s \mid s_{-1}, a_{-1}) \Pr(s_{-1}, a_{-1})}.$$

The posterior $q(s_{-1}, a_{-1} \mid s)$ depends on the prior distribution $\Pr(s_{-1}, a_{-1})$. When the prior distribution follows a stationary distribution under a fixed policy π , that is, $\Pr(s_{-1}, a_{-1}) = \pi(s_{-1}, a_{-1}; \theta) d_{M(\theta)}(s_{-1})$, the posterior is called

the stationary backward probability denoted by $q_{B(\theta)}(s_{-1}, a_{-1} | s)$ and satisfies

$$\begin{aligned}
 q_{B(\theta)}(s_{-1}, a_{-1} | s) &= \frac{p(s | s_{-1}, a_{-1})\pi(s_{-1}, a_{-1}; \theta)d_{M(\theta)}(s_{-1})}{d_{M(\theta)}(s)} \\
 &= \frac{p_{M(\theta)}(s, a_{-1} | s_{-1})d_{M(\theta)}(s_{-1})}{d_{M(\theta)}(s)}. \tag{3.1}
 \end{aligned}$$

If a Markov chain follows $q_{B(\theta)}(s_{-1}, a_{-1} | s)$, we call it the backward Markov chain $B(\theta)$ associated with $M(\theta)$ following $p_{M(\theta)}(s, a_{-1} | s_{-1})$. Both Markov chains— $M(\theta)$ and $B(\theta)$ —are closely related, as described in the following two propositions:

Proposition 1. *Let a Markov chain $M(\theta)$ characterized by a transition probability $p_{M(\theta)}(s | s_{-1}) \equiv \sum_{a_{-1} \in \mathcal{A}} p_{M(\theta)}(s, a_{-1} | s_{-1})$, which is irreducible and ergodic. Then the backward Markov chain $B(\theta)$ characterized by the backward (stationary) transition probability $q_{B(\theta)}(s_{-1} | s) \equiv \sum_{a_{-1} \in \mathcal{A}} q_{B(\theta)}(s_{-1}, a_{-1} | s)$ with respect to $p_{M(\theta)}$ is also ergodic and has the same unique stationary distribution as $M(\theta)$:*

$$d_{M(\theta)}(s) = d_{B(\theta)}(s), \quad \forall s \in \mathcal{S}, \tag{3.2}$$

where $d_{M(\theta)}(s)$ and $d_{B(\theta)}(s)$ are the stationary distributions of $M(\theta)$ and $B(\theta)$, respectively.

Proof. By multiplying both sides of equation 3.1 by $d_{M(\theta)}(s)$ and summing over all possible $a_{-1} \in \mathcal{A}$, we obtain a “detailed balance equations” (MacKay, 2003):

$$q_{B(\theta)}(s_{-1} | s)d_{M(\theta)}(s) = p_{M(\theta)}(s | s_{-1})d_{M(\theta)}(s_{-1}), \quad \forall s_{-1} \in \mathcal{S}, \forall s \in \mathcal{S}. \tag{3.3}$$

Then

$$\sum_{s \in \mathcal{S}} q_{B(\theta)}(s_{-1} | s)d_{M(\theta)}(s) = d_{M(\theta)}(s_{-1})$$

holds by summing both sides of equation 3.3 over all possible $s \in \mathcal{S}$, indicating that (i) $B(\theta)$ has the same stationary distribution as $M(\theta)$. By assumption 1, (i) directly proves that (ii) $B(\theta)$ is irreducible. Equation 3.3 is reformulated by the matrix notation: both transition probabilities $p_{M(\theta)}(s | s_{-1})$ and

$q_{B(\theta)}(s_{-1} | s)$ are assembled into $P_{M(\theta)}$ and $Q_{B(\theta)}$, respectively,³ and the stationary distribution into d_θ ⁴

$$Q_{B(\theta)} = \text{diag}(d_\theta)^{-1} P_{M(\theta)}^\top \text{diag}(d_\theta).$$

We can easily see that the diagonal components of $(P_{M(\theta)})^n$ are equal to those of $(Q_{B(\theta)})^n$ for any natural number n . This implies that (iii) $B(\theta)$ has the same aperiodic property as $M(\theta)$. Proposition 1, equation 3.2, is directly proven by (i), (ii), and (iii) (Schinazi, 1999).

Proposition 2. *Let the distribution of s_{-K} follow $d_{M(\theta)}(s)$, and let $f(s_k, a_k)$ be an arbitrary function of a state-action pair at time k . Then*

$$\begin{aligned} \mathbb{E}_{B(\theta)} \left\{ \sum_{k=1}^K f(s_{-k}, a_{-k}) | s \right\} &= \mathbb{E}_{M(\theta)} \left\{ \sum_{k=1}^K f(s_{-k}, a_{-k}) | s, d_{M(\theta)}(s_{-K}) \right\} \\ &= \mathbb{E}_{M(\theta)} \left\{ \sum_{k=0}^{K-1} f(s_{+k}, a_{+k}) | s_{+K}, d_{M(\theta)}(s) \right\}, \end{aligned} \tag{3.4}$$

where $\mathbb{E}_{B(\theta)}$ and $\mathbb{E}_{M(\theta)}$ are the expectations over the backward and forward Markov chains, $B(\theta)$ and $M(\theta)$, respectively, and $\mathbb{E}\{\cdot | d_{M(\theta)}(s)\} \equiv \mathbb{E}\{\cdot | \Pr(s) = d_{M(\theta)}(s)\}$. Equation 3.4 holds even at the limit $K \rightarrow \infty$.

Proof. By utilizing the Markov property and substituting equation 3.1, we have the following relationship:

$$\begin{aligned} q_{B(\theta)}(s_{-1}, a_{-1}, \dots, s_{-K}, a_{-K} | s) &= q_{B(\theta)}(s_{-1}, a_{-1} | s) \cdots q_{B(\theta)}(s_{-K}, a_{-K} | s_{-K+1}) \\ &= \frac{p_{M(\theta)}(s, a_{-1} | s_{-1}) \cdots p_{M(\theta)}(s_{-K+1}, a_{-K} | s_{-K}) d_{M(\theta)}(s_{-K})}{d_{M(\theta)}(s)} \\ &\propto p_{M(\theta)}(s, a_{-1} | s_{-1}) \cdots p_{M(\theta)}(s_{-K+1}, a_{-K} | s_{-K}) d_{M(\theta)}(s_{-K}). \end{aligned}$$

This directly implies the proposition in the case of finite K . Since the following equations are derived from proposition 1, the proposition in the limit case $K \rightarrow \infty$ is also instantly proven:

$$\begin{aligned} \lim_{K \rightarrow \infty} \mathbb{E}_{B(\theta)}\{f(s_{-K}, a_{-K}) | s\} &= \lim_{K \rightarrow \infty} \mathbb{E}_{M(\theta)}\{f(s_{-K}, a_{-K}) | s, d_{M(\theta)}(s_{-K})\} \\ &= \sum_{s \in S} \sum_{a \in A} \pi(s, a; \theta) d_{M(\theta)}(s) f(s, a). \end{aligned}$$

³The bold $Q_{B(\theta)}$ has no relationship with the state-action value function $Q^\pi(s, a)$.

⁴The function $\text{diag}(a)$ for a vector $a \in \mathcal{R}^d$ denotes the diagonal matrix of a , so $\text{diag}(a) \in \mathcal{R}^{d \times d}$.

Propositions 1 and 2 are significant as they indicate that the samples from the forward Markov chain $M(\theta)$ can be used directly for estimating the statistics of the backward Markov chain $B(\theta)$.

3.2 Temporal Difference Learning for LSD from the Backward to Forward Markov Chains. Using equation 3.1 the LSD, $\nabla_{\theta} \ln d_{M(\theta)}(s)$, can be decomposed into

$$\begin{aligned} \nabla_{\theta} \ln d_{M(\theta)}(s) &= \frac{1}{d_{M(\theta)}(s)} \sum_{s_{-1} \in \mathcal{S}} \sum_{a_{-1} \in \mathcal{A}} p(s \mid s_{-1}, a_{-1}) \pi(s_{-1}, a_{-1}; \theta) d_{M(\theta)}(s_{-1}) \\ &\quad \times \{ \nabla_{\theta} \ln \pi(s_{-1}, a_{-1}; \theta) + \nabla_{\theta} \ln d_{M(\theta)}(s_{-1}) \} \\ &= \sum_{s_{-1} \in \mathcal{S}} \sum_{a_{-1} \in \mathcal{A}} q_{B(\theta)}(s_{-1}, a_{-1} \mid s) \{ \nabla_{\theta} \ln \pi(s_{-1}, a_{-1}; \theta) \\ &\quad + \nabla_{\theta} \ln d_{M(\theta)}(s_{-1}) \} \\ &= \mathbb{E}_{B(\theta)} \{ \nabla_{\theta} \ln \pi(s_{-1}, a_{-1}; \theta) + \nabla_{\theta} \ln d_{M(\theta)}(s_{-1}) \mid s \}. \end{aligned} \tag{3.5}$$

There exist $\nabla_{\theta} \ln d_{M(\theta)}(s)$ and $\nabla_{\theta} \ln d_{M(\theta)}(s_{-1})$ in equation 3.5, so its recursion can be rewritten as

$$\begin{aligned} \nabla_{\theta} \ln d_{M(\theta)}(s) &= \lim_{K \rightarrow \infty} \mathbb{E}_{B(\theta)} \left\{ \sum_{k=1}^K \nabla_{\theta} \ln \pi(s_{-k}, a_{-k}; \theta) \right. \\ &\quad \left. + \nabla_{\theta} \ln d_{M(\theta)}(s_{-K}) \mid s \right\}. \end{aligned} \tag{3.6}$$

Equation 3.6 implies that the LSD of a state s is the infinite-horizon cumulation of the log policy distribution derivative (LPD), $\nabla_{\theta} \ln \pi(s, a; \theta)$, along the backward Markov chain $B(\theta)$ from state s . From equations 3.5 and 3.6, LSD could be estimated using an approach similar to temporal difference (TD) learning (Sutton, 1988) for the following backward TD-error δ on the backward Markov chain $B(\theta)$ rather than $M(\theta)$:

$$\delta(s) \equiv \nabla_{\theta} \ln \pi(s_{-1}, a_{-1}; \theta) + \nabla_{\theta} \ln d_{M(\theta)}(s_{-1}) - \nabla_{\theta} \ln d_{M(\theta)}(s),$$

where the first two terms of the right-hand side describe the one-step actual observation of the policy eligibility and the one-step ahead LSD on $B(\theta)$, respectively, while the last term is the current LSD. Interestingly, while the well-known TD error for the value-function estimation uses the reward $r(s, a, s_{+1})$ on $M(\theta)$ (Sutton & Barto, 1998), this TD error for the LSD estimation uses $\nabla_{\theta} \ln \pi(s_{-1}, a_{-1}; \theta)$ on $B(\theta)$.

While $\delta(s)$ is a random variable, $\mathbb{E}_{B(\theta)} \{ \delta(s) \mid s \} = \mathbf{0}$ holds for all states $s \in \mathcal{S}$, which comes from equation 3.5. This is a motivation for minimizing the

mean squares of the backward TD error, $\mathbb{E}_{B(\theta)}\{\|\hat{\delta}(s)\|^2\}$, for the estimation of LSDs,⁵ where $\hat{\delta}(s)$ is composed of the LSD estimate $\widehat{\nabla}_\theta \ln d_{M(\theta)}(s)$ rather than (exact) LSD $\nabla_\theta \ln d_{M(\theta)}(s)$. Here, $\|a\|$ denotes the Euclidean norm $(a^\top a)^{1/2}$.

With an eligibility decay rate $\lambda \in [0, 1]$ and a back-trace time step $K \in \mathcal{N}$, equation 3.6 is generalized, where \mathcal{N} denotes the set of natural numbers:

$$\begin{aligned} \nabla_\theta \ln d_{M(\theta)}(s) = \mathbb{E}_{B(\theta)} \left\{ \sum_{k=1}^K \lambda^{k-1} \{ \nabla_\theta \ln \pi(s_{-k}, a_{-k}; \theta) \right. \\ \left. + (1 - \lambda) \nabla_\theta \ln d_{M(\theta)}(s_{-k}) \} + \lambda^K \nabla_\theta \ln d_{M(\theta)}(s_{-K}) \mid s \right\}. \end{aligned}$$

Accordingly, the backward TD is extended into the backward TD(λ), $\delta_{\lambda,K}(s)$,

$$\begin{aligned} \delta_{\lambda,K}(s) \equiv \sum_{k=1}^K \lambda^{k-1} \{ \nabla_\theta \ln \pi(s_{-k}, a_{-k}; \theta) + (1 - \lambda) \nabla_\theta \ln d_{M(\theta)}(s_{-k}) \} \\ + \lambda^K \nabla_\theta \ln d_{M(\theta)}(s_{-K}) - \nabla_\theta \ln d_{M(\theta)}(s), \end{aligned}$$

where the unbiased property, $\mathbb{E}_{B(\theta)}\{\delta_{\lambda,K}(s) \mid s\} = \mathbf{0}$, is still retained. The minimization of $\mathbb{E}_{B(\theta)}\{\|\hat{\delta}_{\lambda,K}(s)\|^2\}$ at $\lambda = 1$ and the limit $K \rightarrow \infty$ is regarded as the Widrow-Hoff supervised learning procedure. Even if λ and K are not set in the above values, the minimization of $\mathbb{E}_{B(\theta)}\{\|\hat{\delta}_{\lambda,K}(s)\|^2\}$ in a large $\lambda \in [0, 1)$ and K would be less sensitive to a non-Markovian effect as in the case of the conventional TD(λ) learning for the value functions (Peng & Williams, 1996).

In order to minimize $\mathbb{E}_{B(\theta)}\{\|\hat{\delta}_{\lambda,K}(s)\|^2\}$ as the estimation of the LSD, we need to gather many samples drawn from the backward Markov chain $B(\theta)$. However, the actual samples are drawn from a forward Markov chain $M(\theta)$. Fortunately, by using propositions 1 and 2, we can derive the following exchangeable property:

$$\begin{aligned} \mathbb{E}_{B(\theta)}\{\|\hat{\delta}_{\lambda,K}(s)\|^2\} &= \sum_{s \in \mathcal{S}} d_{B(\theta)}(s) \mathbb{E}_{B(\theta)}\{\|\hat{\delta}_{\lambda,K}(s)\|^2 \mid s\} \\ &= \sum_{s \in \mathcal{S}} d_{M(\theta)}(s) \mathbb{E}_{M(\theta)}\{\|\hat{\delta}_{\lambda,K}(s)\|^2 \mid s, d_{M(\theta)}(s_{-K})\} \\ &= \mathbb{E}_{M(\theta)}\{\|\hat{\delta}_{\lambda,K}(s)\|^2 \mid d_{M(\theta)}(s_{-K})\}. \end{aligned} \tag{3.7}$$

⁵Actually, the classical least-squares approach to $\mathbb{E}_{B(\theta)}\{\|\hat{\delta}(s)\|^2\}$ would make the LSD estimate biased, because $\hat{\delta}(s)$ has the different time step LSDs, $\widehat{\nabla}_\theta \ln d_{M(\theta)}(s_{-1})$ and $\widehat{\nabla}_\theta \ln d_{M(\theta)}(s)$. Instead, $\mathbb{E}_{B(\theta)}\{u(s)^\top \hat{\delta}(s)\}$ is minimized for an unbiased LSD estimation, where $u(s)$ is an instrumental variable (Young, 1984; Bradtke & Barto, 1996). Such a detailed discussion is given in section 3.3. Before that, we use $\mathbb{E}_{B(\theta)}\{\|\hat{\delta}(s)\|^2\}$ to enhance readability.

In particular, the actual samples can be reused to minimize $\mathbb{E}_{B(\theta)}\{\|\hat{\delta}_{\lambda,K}(s)\|^2\}$, provided $s_{-K} \sim d_{M(\theta)}(s)$. In real problems, however, the initial state is rarely distributed according to the stationary distribution $d_{M(\theta)}(s)$. To interpolate the gap between theoretical assumption and realistic applicability, we would need to adopt either of the following two strategies: (i) K is not set at such a large integer if $\lambda \approx 1$, or (ii) λ is not set at 1 if $K \approx t$, where t is the current time step of the actual forward Markov chain $M(\theta)$.

3.3 LSD Estimation Algorithm: Least Squares on Backward TD(λ) with Constraint. In the previous sections, we introduced the theory for estimating LSD by minimizing of the mean squares of $\hat{\delta}_{\lambda,K}(s)$ on $M(\theta)$, $\mathbb{E}_{M(\theta)}\{\|\hat{\delta}_{\lambda,K}(s)\|^2 | d_{M(\theta)}(s_{-K})\}$. However, LSD also has the following constraint derived from $\sum_{s \in \mathcal{S}} d_{M(\theta)}(s) = 1$:

$$\begin{aligned} \mathbb{E}_{M(\theta)}\{\nabla_{\theta} \ln d_{M(\theta)}(s)\} &= \sum_{s \in \mathcal{S}} d_{M(\theta)}(s) \nabla_{\theta} \ln d_{M(\theta)}(s) \\ &= \nabla_{\theta} \sum_{s \in \mathcal{S}} d_{M(\theta)}(s) = \mathbf{0}. \end{aligned} \tag{3.8}$$

In this section, we propose an LSD estimation algorithm, $\mathcal{L}\text{S}\text{LSD}(\lambda)$, based on the least-squares TD approach (Young, 1984; Bradtke & Barto, 1996; Boyan, 2002; Lagoudakis & Parr, 2003), which simultaneously attempts to decrease the mean squares and satisfy the constraint. We consider the situation where the LSD estimate $\hat{\nabla}_{\theta} \ln d_{M(\theta)}(s)$ is represented by a linear vector function approximator

$$f(s; \mathbf{\Omega}) \equiv \mathbf{\Omega} \phi(s), \tag{3.9}$$

where $\phi(s) \in \mathcal{R}^e$ is a basis function and $\mathbf{\Omega} \equiv [\omega_1, \dots, \omega_d]^\top \in \mathcal{R}^{d \times e}$ is an adjustable parameter matrix, and we assume that the optimal parameter $\mathbf{\Omega}^*$ satisfies $\nabla_{\theta} \ln d_{M(\theta)}(s) = \mathbf{\Omega}^* \phi(s)$. If the estimator cannot represent the LSD exactly, $\mathcal{L}\text{S}\text{LSD}(\lambda)$ would behave as suggested by Sutton (1988) and Peng and Williams (1996), which means the estimation error for the LSD would get smaller as the value of $\lambda \in [0, 1)$ approaches 1. This will be confirmed in our numerical experiments in section 5.

For simplicity, we focus on only the i th element θ_i of the policy parameter θ , denoting $f(s; \omega_i) \equiv \omega_i^\top \phi(s)$, $\nabla_{\theta_i} \ln \pi(s, a; \theta) \equiv \partial \ln \pi(s, a; \theta) / \partial \theta_i$, and $\hat{\delta}_{\lambda,K}(s, \omega_i)$ as the i th element of $\hat{\delta}_{\lambda,K}(s, \mathbf{\Omega})$. Accordingly, the objective function to be minimized is given by

$$\varepsilon(\omega_i) = \frac{1}{2} \mathbb{E}_{M(\theta)}\{\hat{\delta}_{\lambda,K}(s; \omega_i)^2 | d_{M(\theta)}(s_{-K})\} + \frac{1}{2} \mathbb{E}_{M(\theta)}\{f(s; \omega_i)\}^2, \tag{3.10}$$

where the second term of the right side comes from the constraint of equation 3.8.⁶ Thus, the derivative is

$$\begin{aligned} \nabla_{\omega_i} \varepsilon(\omega_i) &= \mathbb{E}_{M(\theta)} \{ \hat{\delta}_{\lambda,K}(s; \omega_i) \nabla_{\omega_i} \hat{\delta}_{\lambda,K}(s; \omega_i) \mid d_{M(\theta)}(s_{-K}) \} \\ &\quad + \nabla_{\omega_i} \mathbb{E}_{M(\theta)} \{ f(s; \omega_i) \}, \end{aligned} \tag{3.11}$$

where

$$\hat{\delta}_{\lambda,K}(s; \omega_i) = \sum_{k=1}^K \lambda^{k-1} \nabla_{\theta_i} \ln \pi(s_{-k}, a_{-k}; \theta) + \omega_i^\top \nabla_{\omega_i} \hat{\delta}_{\lambda,K}(s; \omega_i)$$

and

$$\nabla_{\omega_i} \hat{\delta}_{\lambda,K}(s; \omega_i) = (1 - \lambda) \sum_{k=1}^K \lambda^{k-1} \phi(s_{-k}) + \lambda^K \phi(s_{-K}) - \phi(s).$$

Although the conventional least-squares method aims to find the parameter satisfying $\nabla_{\omega_i} \varepsilon(\omega_i) = \mathbf{0}$ as the true parameter ω_i^* , it induces estimation bias if a correlation exists between the error $\hat{\delta}_{\lambda,K}(s; \omega_i^*)$ and its derivative $\nabla_{\omega_i} \hat{\delta}_{\lambda,K}(s; \omega_i^*)$ concerning the first term of the right-hand side in equation 3.10. That is, if

$$\mathbb{E}_{M(\theta)} \{ \hat{\delta}_{\lambda,K}(s; \omega_i^*) \nabla_{\omega_i} \hat{\delta}_{\lambda,K}(s; \omega_i^*) \mid d_{M(\theta)}(s_{-K}) \} \neq \mathbf{0},$$

$\nabla_{\omega_i} \varepsilon(\omega_i^*) \neq \mathbf{0}$ holds because $\nabla_{\omega_i} \mathbb{E}_{M(\theta)} \{ f(s; \omega_i^*) \}$ in the second term of the left-hand side in equation 3.10 is always a zero vector due to equation 3.8 and $f(s; \omega_i^*) = \nabla_{\theta_i} \ln d_{M(\theta)}(s)$. Since this correlation exists in general RL tasks, we apply the instrumental variable method to eliminate the bias (Young, 1984; Bradtke & Barto, 1996). This requires that $\nabla_{\omega_i} \hat{\delta}_{\lambda,K}(s; \omega_i)$ be replaced by the instrumental variable $\iota(s)$, which has a correlation with $\nabla_{\omega_i} \hat{\delta}_{\lambda,K}(s; \omega_i^*)$ but not $\hat{\delta}_{\lambda,K}(s; \omega_i^*)$. This condition is obviously satisfied when $\iota(s) = \phi(s)$ as well as LSTD(λ) (Bradtke & Barto, 1996; Boyan, 2002). Instead of equation 3.11 we aim to find the parameter making the equation

$$\begin{aligned} \tilde{\nabla}_{\omega_i} \varepsilon(\omega_i) &\equiv \mathbb{E}_{M(\theta)} \{ \hat{\delta}_{\lambda,K}(s; \omega_i) \phi(s) \mid d_{M(\theta)}(s_{-K}) \} \\ &\quad + \mathbb{E}_{M(\theta)} \{ \phi(s) \} \mathbb{E}_{M(\theta)} \{ \phi(s) \}^\top \omega_i \end{aligned} \tag{3.12}$$

be equal to zero in order to compute the true parameter ω_i^* , so that $\tilde{\nabla}_{\omega_i} \varepsilon(\omega_i^*) = \mathbf{0}$.

⁶While $\mathcal{L}\text{SLSD}(\lambda)$ weighs the two objectives equally, we can instantly extend it to the problem minimizing $\mathbb{E}_{M(\theta)} \{ \|\hat{\delta}_\lambda(x)\|^2 \mid d^\pi(s_{-K}) \}$ subject to the constraint of equation 3.8 with the Lagrange multiplier method.

For the remainder of this article, we denote the current state at time step t by s_t to clarify the time course on the actual Markov chain $M(\theta)$. In the proposed LSD estimation algorithm, $\mathcal{L}\text{SLSD}(\lambda)$, the back-trace time step K is set equal to the time step t , of the current state s_t , while the eligibility decay rate λ can be set in $[0, 1)$, that is,

$$\hat{\delta}_{\lambda,K}(s_t; \omega_i) = g_{\lambda,i}(s_{t-1}) + (z_{\lambda}(s_{t-1}) - \phi(s_t))^\top \omega_i,$$

where $g_{\lambda,i}(s_t) = \sum_{k=0}^t \lambda^{t-k} \nabla_{\theta_i} \ln \pi(s_k, a_k; \theta)$ and $z_{\lambda}(s_t) = (1 - \lambda) \sum_{k=1}^t \lambda^{t-k} \phi(s_k) + \lambda^t \phi(s_0)$. The expectations in equation 3.12 are estimated without bias by Bradtke and Barto (1996) and Boyan (2002):

$$\begin{aligned} & \lim_{K \rightarrow \infty} \mathbb{E}_{M(\theta)} \{ \hat{\delta}_{\lambda,K}(s; \omega_i) \phi(s) \mid d_{M(\theta)}(s_{-K}) \} \\ & \simeq \frac{1}{T} \sum_{t=1}^T \phi(s_t) \{ g_{\lambda,i}(s_{t-1}) - (\phi(s_t) - z_{\lambda}(s_{t-1}))^\top \omega_i \} \\ & = \mathbf{b}_T - \mathbf{A}_T \omega_i, \end{aligned}$$

where $\mathbf{b}_T \equiv \frac{1}{T} \sum_{t=1}^T \phi(s_t) g_{\lambda,i}(s_{t-1})$ and $\mathbf{A}_T \equiv \frac{1}{T} \sum_{t=1}^T \phi(s_t) (\phi(s_t) - z_{\lambda}(s_{t-1}))^\top$, and

$$\begin{aligned} \mathbb{E}_{M(\theta)} \{ \phi(x) \} & \simeq \frac{1}{T+1} \sum_{t=0}^T \phi(s_t) \\ & \equiv \mathbf{c}_T. \end{aligned}$$

Therefore, by substituting these estimators into equation 3.12, the estimate $\hat{\omega}_i^*$ at time step T is computed as

$$\begin{aligned} & \mathbf{b}_T - \mathbf{A}_T \hat{\omega}_i^* + \mathbf{c}_T \mathbf{c}_T^\top \hat{\omega}_i^* = \mathbf{0} \\ \Leftrightarrow & \hat{\omega}_i^* = (\mathbf{A}_T - \mathbf{c}_T \mathbf{c}_T^\top)^{-1} \mathbf{b}_T. \end{aligned}$$

The $\mathcal{L}\text{SLSD}(\lambda)$ for the matrix parameter $\hat{\Omega}^*$ rather than $\hat{\omega}_i^*$ is shown in algorithm 1, where the notation $:=$ denotes the right-to-left substitution:⁷

⁷Incidentally, although there is calculation of an inverse matrix in the algorithms, a pseudo-inverse matrix may be used instead of direct calculation of the inverse matrix so as to secure stability in numeric calculation.

Algorithm 1: $\mathcal{L}\text{SLSD}(\lambda)$: Estimation for $\nabla_{\theta} \ln d_{M(\theta)}(s)$

Given:

- a policy $\pi(s, a; \theta)$ with a fixed θ ,
- a feature vector function of state $\phi(s)$.

Initialize: $\lambda \in [0, 1)$.

Set: $c := \phi(s_0)$; $z := \phi(s_0)$; $g := \mathbf{0}$; $A := \mathbf{0}$; $B := \mathbf{0}$.

for $t = 0$ **to** $T - 1$ **do**

$c := c + \phi(s_{t+1})$;

$g := \lambda g + \nabla_{\theta} \ln \pi(s_t, a_t; \theta)$;

$A := A + \phi(s_{t+1})(\phi(s_{t+1}) - z)^{\top}$;

$B := B + \phi(s_{t+1})g^{\top}$;

$z := \lambda z + (1 - \lambda)\phi(s_{t+1})$;

end for

$\Omega := (A - cc^{\top}/T)^{-1}B$;

Return: $\widehat{\nabla}_{\theta} \ln d_{M(\theta)}(s) = \Omega^{\top} \phi(s)$.

It is intriguing that $\mathcal{L}\text{SLSD}(\lambda)$ has a relationship to a model-based method, as noted by Boyan (2002) and Lagoudakis and Parr (2003) in the references for $\text{LSTD}(\lambda)$ and $\text{LSTDQ}(\lambda)$, but $\mathcal{L}\text{SLSD}(\lambda)$ is concerned with the “backward” model $B(\theta)$ instead of the forward model $M(\theta)$. This is due to the fact that the sufficient statistics A in $\mathcal{L}\text{SLSD}(\lambda)$ can be regarded as a compressed backward model, since A is equivalent to one of the sufficient statistics to estimate the backward state transition probability $q_{B(\theta)}(s_{-1} | s)$ when $\lambda = 0$ and the feature vector ϕ corresponding to $\phi(1) = (1, 0, \dots, 0)$; $\phi(2) = (0, 1, \dots, 0)$; and so on. We give a detailed explanation of this in appendix B.

4 Policy Gradient Algorithms with the LSD Estimate

We propose a PG algorithm as a straightforward application with the LSD estimates in section 4.1. In section 4.2, we introduce baseline functions to reduce the variance of the PG as estimated by our PG algorithm.

4.1 Policy Update with the LSD Estimate. Now let us define the PGRL algorithm based on the LSD estimate. The realization of the estimation for $\nabla_{\theta} \ln d_{M(\theta)}(s)$ by $\mathcal{L}\text{SLSD}(\lambda)$ directly leads to the following estimate for the PG (see equation 2.3), due to its independence from the forgetting factor γ for the value functions:

$$\nabla_{\theta} \eta(\theta) \simeq \frac{1}{T} \sum_{t=0}^{T-1} (\nabla_{\theta} \ln \pi(s_t, a_t; \theta) + \nabla_{\theta} \ln d_{M(\theta)}(s_t)) r_{t+1} \quad (4.1)$$

$$\simeq \frac{1}{T} \sum_{t=0}^{T-1} (\nabla_{\theta} \ln \pi(s_t, a_t; \theta) + \widehat{\nabla}_{\theta} \ln d_{M(\theta)}(s_t)) r_{t+1}, \quad (4.2)$$

where r_{t+1} is the immediate reward defined by the reward function $r(s_t, a_t, s_{t+1})$. The policy parameter can then be updated through the stochastic gradient method with an appropriate step size α (Bertsekas and Tsitsiklis, 1996):⁸

$$\theta := \theta + \alpha(\nabla_{\theta} \ln \pi(s_t, a_t; \theta) + \widehat{\nabla}_{\theta} \ln d_{M(\theta)}(s_t))r_{t+1}.$$

\mathcal{L} SLSLSD(λ)-PG without a baseline function is shown in algorithm 2 as one of the simplest realizations of PG algorithm that uses \mathcal{L} SLSLSD(λ). In algorithm 2, the forgetting rate parameter $\beta \in [0, 1)$ is introduced to discard the past estimates given by old values of θ :

Algorithm 2: \mathcal{L} SLSLSD(λ)-PG: Optimization for the policy without baseline function

Given:

- a policy $\pi(s, a; \theta)$ with an adjustable θ ,
- a feature vector function of state $\phi(s)$.

Initialize: $\theta, \lambda \in [0, 1), \beta \in [0, 1], \alpha_t$.

Set: $c := \phi(s_0); z := \phi(s_0); g := \mathbf{0}; A := \mathbf{0}; B := \mathbf{0}$.

for $t = 0$ **to** $T - 1$ **do**

$$\theta := \theta + \alpha_t \{ \nabla_{\theta} \ln \pi(s_t, a_t; \theta) + \Omega^{\top} \phi(s_t) \} r_{t+1};$$

$$c := \beta c + \phi(s_{t+1});$$

$$g := \beta \lambda g + \nabla_{\theta} \ln \pi(s_t, a_t; \theta);$$

$$A := \beta A + \phi(s_{t+1})(\phi(s_{t+1}) - z)^{\top};$$

$$B := \beta B + \phi(s_{t+1})g^{\top};$$

$$\Omega := (A - cc^{\top} / \|c\|)^{-1} B;$$

$$z := \lambda z + (1 - \lambda)\phi(s_{t+1});$$

end for

Return: $p(a | s; \theta) = \pi(s, a; \theta)$.

Another important topic for function approximation is the choice of the basis function $\phi(s)$ of the approximator, particularly in continuous state problems. For the PG algorithm, the objective of the LSD estimate is to provide one term of the PG estimate, such as $\sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_{M(\theta)}(s) \pi(s, a; \theta) \nabla_{\theta} \ln d_{M(\theta)}(s) \bar{r}(s, a)$, but not to provide a precise estimate of the LSD $\nabla_{\theta} \ln d_{M(\theta)}(s)$, where $\bar{r}(s, a) \equiv \sum_{s_{+1} \in \mathcal{S}} p(s_{+1} | s, a) r(s, a, s_{+1})$. Therefore, the following proposition would be useful:

Proposition 3. *Let the basis function of the LSD estimator be*

$$\phi(s) = \sum_{a \in \mathcal{A}} \pi(s, a; \theta) \bar{r}(s, a),$$

⁸Alternatively, θ can also be updated through the bath gradient method: $\theta := \theta + \alpha \widehat{\nabla}_{\theta} R(\theta)$.

and then the function estimator, $f(s; \omega) = \omega \sum_{a \in \mathcal{A}} \pi(s, a; \theta) \bar{r}(s, a)$, can represent the second term of the PG, $\sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_{M(\theta)}(s) \pi(s, a; \theta) \nabla_{\theta} \ln d_{M(\theta)}(s) \bar{r}(s, a)$, where the adjustable parameter ω is a d -dimensional vector:

$$\begin{aligned} & \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_{M(\theta)}(s) \pi(s, a; \theta) \bar{r}(s, a) \nabla_{\theta} \ln d_{M(\theta)}(s) \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_{M(\theta)}(s) \pi(s, a; \theta) \bar{r}(s, a) f(s; \omega^*), \end{aligned}$$

where ω^* minimizes the mean error, $\epsilon(\omega) = \frac{1}{2} \sum_{s \in \mathcal{S}} d_{M(\theta)}(s) \{ \nabla_{\theta} \ln d_{M(\theta)}(s) - f(s; \omega) \}^2$.

Proof. The proposition follows directly from

$$\nabla_{\omega} \epsilon(\omega^*) = \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_{M(\theta)}(s) \pi(s, a; \theta) \bar{r}(s, a) \{ \nabla_{\theta} \ln d_{M(\theta)}(s) - f(s; \omega^*) \} = \mathbf{0}.$$

4.2 Baseline Function for Variance Reduction of Policy Gradient Estimates with LSD. As the variance of the PG estimates using the LSD, equation 4.2, might be large, we consider variance reduction using a baseline function for immediate reward r . The following proposition provides the kind of functions that can be used as the baseline function for PG estimation using the LSD:⁹

Proposition 4. *With the following function of the state s and the following state s_{+1} on $M(\theta)$,*

$$\rho(s, s_{+1}) = c + g(s) - g(s_{+1}), \tag{4.3}$$

where c and $g(s)$ are an arbitrary constant and an arbitrary bounded function of the state, respectively. The derivative of the average reward $\eta(\theta)$ with respect to the policy parameter θ (see equation 2.3), $\nabla_{\theta} \eta(\theta)$, is then transformed to

$$\begin{aligned} \nabla_{\theta} \eta(\theta) &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \sum_{s_{+1} \in \mathcal{S}} d_{M(\theta)}(s) \pi(s, a; \theta) p(s_{+1} | s, a) \\ &\quad \{ \nabla_{\theta} \ln \pi(s, a; \theta) + \nabla_{\theta} \ln d_{M(\theta)}(s) \} r(s, a, s_{+1}) \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \sum_{s_{+1} \in \mathcal{S}} d_{M(\theta)}(s) \pi(s, a; \theta) p(s_{+1} | s, a) \\ &\quad \{ \nabla_{\theta} \ln \pi(s, a; \theta) + \nabla_{\theta} \ln d_{M(\theta)}(s) \} \{ r(s, a, s_{+1}) - \rho(s, s_{+1}) \}. \end{aligned} \tag{4.4}$$

⁹Though a baseline might be a constant from a traditional perspective, we call the function defined in equation 4.3 a baseline function for equation 2.3 because it does not add any bias to $\nabla_{\theta} \eta(\theta)$.

Proof. See appendix C.

Proposition 4 implies that any $\rho(s, s_{+1})$ defined in equation 4.3 can be used as the baseline function of the immediate reward $r_{+1} \equiv r(s, a, s_{+1})$ for computing the PG, as in equation 4.4. Therefore, the PG can be estimated with the baseline function $\rho(s_t, s_{t+1})$ with a large time step T ,

$$\begin{aligned} \nabla_{\theta} R(\theta) &\simeq \frac{1}{T} \sum_{t=0}^{T-1} (\nabla_{\theta} \ln \pi(s_t, a_t; \theta) + \nabla_{\theta} \ln d_{M(\theta)}(s_t)) \{r(s_t, a_t, s_{t+1}) \\ &\quad - \rho(s_t, s_{t+1})\} \equiv \widehat{\nabla}_{\theta} \eta(\theta). \end{aligned} \tag{4.5}$$

In view of the form of the baseline function in equation 4.3, we use the following linear function as a representation of the baseline function,

$$\rho(s, s_{+1}; \mathbf{v}) = \begin{pmatrix} \mathbf{v}_u \\ v_d \end{pmatrix}^{\top} \begin{pmatrix} \boldsymbol{\phi}(s) - \boldsymbol{\phi}(s_{+1}) \\ 1 \end{pmatrix} \equiv \mathbf{v}^{\top} \boldsymbol{\psi}(s, s_{+1}),$$

where \mathbf{v} and $\boldsymbol{\phi}(s)$ are its coefficient parameter and feature vector function of the state.

When we consider the trace of the covariance matrix of the PG estimates $\widehat{\nabla}_{\theta} \eta(\theta)$ as the variance of $\widehat{\nabla}_{\theta} \eta(\theta)$ and use the results of Greensmith, Bartlett, and Baxter (2004), an upper bound for the variance is derived as

$$\begin{aligned} \text{Var}_{M(\theta)} [\widehat{\nabla}_{\theta} \eta(\theta)] &\leq h(\mathbb{E}_{M(\theta)} [\|\nabla_{\theta} \ln \pi(s, a; \theta) + \nabla_{\theta} \ln d_{M(\theta)}(s)\|^2 \\ &\quad \{r(s, a, s_{+1}) - \rho(s, s_{+1}; \mathbf{v})\}^2]) \\ &\equiv h(\sigma_{\widehat{\nabla}_{\theta} \eta(\theta)}^2(\mathbf{v})), \end{aligned}$$

where $h(a)$ is a monotonically increasing function of its argument a . Accordingly, since the optimal coefficient parameter \mathbf{v}^* for the optimal (linear) baseline function $b^*(s, s_{+1}) \equiv \rho(s, s_{+1}; \mathbf{v}^*)$ satisfies

$$\left. \frac{\partial \sigma_{\widehat{\nabla}_{\theta} \eta(\theta)}^2(\mathbf{v})}{\partial \mathbf{v}} \right|_{\mathbf{v}=\mathbf{v}^*} = 0,$$

the optimal coefficient parameter is computed as¹⁰

$$\mathbf{v}^* = \mathbb{E}_{M(\theta)} \left\{ \|\nabla_{\theta} \ln \pi(s, a; \theta) + \nabla_{\theta} \ln d_{M(\theta)}(s)\|^2 \boldsymbol{\psi}(s, s_{+1}) \boldsymbol{\psi}(s, s_{+1})^{\top} \right\}^{-1} \mathbb{E}_{M(\theta)} \left\{ \|\nabla_{\theta} \ln \pi(s, a; \theta) + \nabla_{\theta} \ln d_{M(\theta)}(s)\|^2 \boldsymbol{\psi}(s, s_{+1}) \mathbf{r}(s, a, s_{+1}) \right\}. \quad (4.6)$$

There is also an alternative decent baseline function $b(s, s_{+1}) \equiv \rho(s, s_{+1}; \mathbf{v}^*)$, which minimizes the residual sum of squares about the expectation of the reward function $r(s, a, s_{+1})$ given s and then works to reduce the variance of the PG estimates. If the rank of $\mathbb{E}_{M(\theta)} \{ \boldsymbol{\psi}(s) \boldsymbol{\psi}(s)^{\top} \}$ is equal to the number of states $|\mathcal{S}|$, this decent baseline function satisfies

$$\mathbb{E}_{M(\theta)} \{ b(s, s_{+1}) \mid s \} = \mathbb{E}_{M(\theta)} \{ r(s, a, s_{+1}) \mid s \}, \quad \forall s, \quad (4.7)$$

and has the following statistical interpretation: this function is a solution of Poisson's equation

$$\mathbf{v}_d^* + \mathbf{v}_u^{*\top} \boldsymbol{\phi}(s) = \mathbb{E}_{M(\theta)} \{ r(s, a, s_{+1}) - \mathbf{v}_u^* \boldsymbol{\phi}(s_{+1}) \mid s \}, \quad \forall s,$$

thus, \mathbf{v}_d^* and $\mathbf{v}_u^{*\top} \boldsymbol{\phi}(s)$ are equal to the average reward and the (undiscounted) state value function, respectively (Konda & Tsitsiklis, 2003). The parameter \mathbf{v}^* of the decent baseline function can be computed as¹¹

$$\mathbf{v}^* = \mathbb{E}_{M(\theta)} \left\{ \tilde{\boldsymbol{\phi}}(s) \boldsymbol{\psi}(s, s_{+1})^{\top} \right\}^{-1} \mathbb{E}_{M(\theta)} \left\{ \tilde{\boldsymbol{\phi}}(s) \mathbf{r}(s, a, s_{+1}) \right\}, \quad (4.8)$$

where $\tilde{\boldsymbol{\psi}}(s) \equiv (\boldsymbol{\phi}(s)^{\top}, 1)^{\top}$ (Ueno, Kawanabe, Mori, Maeda, & Ishii 2008).

By equations 4.6 and 4.8, both the coefficient parameters \mathbf{v}^* and \mathbf{v}^* for the optimal $b^*(s, s_{+1})$ and decent $b(s, s_{+1})$ baseline functions can be estimated by least squares and LSTD(λ), respectively, though the estimation for b^* requires LSD estimates. The $\mathcal{L}SLSD(\lambda)$ -PG algorithms with both baseline functions are shown in algorithms 3 and 4:

¹⁰The optimal baseline function is computed directly with \mathbf{v}^* as $b^*(s, s_{+1}) = \mathbf{v}^{*\top} \boldsymbol{\psi}(s, s_{+1})$. Note that there is a similarity to the optimal baseline in Peters and Schaal (2006).

¹¹ \mathbf{v}^* is given by solving the estimating function $\mathbb{E}_{M(\theta)} \{ \tilde{\boldsymbol{\psi}}(s)^{\top} (r(s, a, s_{+1}) - \rho(s, s_{+1}; \mathbf{v})) \} = 0$.

Algorithm 3: $\mathcal{L}\text{SLS}\mathcal{D}(\lambda)\text{-PG}$: Optimization for the policy with “optimal” baseline function $b^*(s, s_{+1})$

Given:

- a policy $\pi(s, a; \theta)$ with an adjustable θ ,
- a feature vector function of state $\phi(s)$.

Define: $\psi(s_t, s_{t+1}) \equiv [\phi(s_t)^\top - \phi(s_{t+1})^\top, 1]^\top$

Initialize: $\theta, \lambda \in [0, 1), \beta \in [0, 1], \alpha_t, \beta_b \in [0, 1]$.

Set: $c := \phi(s_0); z := \phi(s_0)/\beta; g := \mathbf{0}; A := \mathbf{0}; B := \mathbf{0}; X := \mathbf{0}; y := \mathbf{0}$;

for $t = 0$ **to** $T - 1$ **do**

if $t \geq 1$ **then**

$$\theta := \theta + \alpha_t \{\nabla_\theta \ln \pi(s_t, a_t; \theta) + \Omega^\top \phi(s_t)\} \{r_{t+1} - \psi(s_t, s_{t+1})^\top X^{-1} y\};$$

end if

$$c := \beta c + \phi(s_{t+1});$$

$$z := \beta \lambda z + (1 - \lambda) \phi(s_t);$$

$$g := \beta \lambda g + \nabla_\theta \ln \pi(s_t, a_t; \theta);$$

$$A := \beta A + \phi(s_{t+1})(\phi(s_{t+1}) - z)^\top;$$

$$B := \beta B + \phi(s_{t+1})g^\top;$$

$$\Omega := (A - cc^\top / \|c\|)^{-1} B;$$

$$w := \|\nabla_\theta \ln \pi(s_t, a_t; \theta) + \Omega^\top \phi(s_t)\|^2;$$

$$X := \beta_b X + w \psi(s_t, s_{t+1}) \psi(s_t, s_{t+1})^\top;$$

$$y := \beta_b y + w \psi(s_t, s_{t+1}) r_{t+1};$$

end for

Return: $p(a | s; \theta) = \pi(s, a; \theta)$.

Algorithm 4: $\mathcal{L}\text{SLS}\mathcal{D}(\lambda)\text{-PG}$: Optimization for the policy with “decent” baseline function $b^*(s, s_{+1})$

Given:

- a policy $\pi(s, a; \theta)$ with an adjustable θ ,
- a feature vector function of state $\phi(s)$.

Define: $\psi(s_t, s_{t+1}) \equiv [\phi(s_t)^\top - \phi(s_{t+1})^\top, 1]^\top, \tilde{\phi}(s_t) \equiv [\phi(s_t)^\top, 1]^\top$.

Initialize: $\theta, \lambda \in [0, 1), \beta \in [0, 1], \alpha_t, \lambda_b \in [0, 1), \beta_b \in [0, 1]$

Set: $c := \phi(s_0); z := \phi(s_0)/\beta; g := \mathbf{0}; A := \mathbf{0}; B := \mathbf{0}$;

$X := \mathbf{0}; y := \mathbf{0}; z_b := \mathbf{0}$

for $t = 0$ **to** $T - 1$ **do**

if $t \geq 1$ **then**

$$\theta := \theta + \alpha_t \{\nabla_\theta \ln \pi(s_t, a_t; \theta) + \Omega^\top \phi(s_t)\} \{r_{t+1} - \psi(s_t, s_{t+1})^\top X^{-1} y\};$$

end if

$$c := \beta c + \phi(s_{t+1});$$

$$z := \beta \lambda z + (1 - \lambda) \phi(s_t);$$

$$g := \beta \lambda g + \nabla_\theta \ln \pi(s_t, a_t; \theta);$$

$$A := \beta A + \phi(s_{t+1})(\phi(s_{t+1}) - z)^\top;$$

$$B := \beta B + \phi(s_{t+1})g^\top;$$

$$\Omega := (A - cc^\top / \|c\|)^{-1} B;$$

$$z_b := \beta_b \lambda_b z_b + \tilde{\phi}(s_t);$$

$$X := \beta_b X + z_b \psi(s_t, s_{t+1})^\top;$$

```

y :=  $\beta_b \mathbf{y} + z_b r_{t+1}$ ;
end for
Return:  $p(a \mid s; \boldsymbol{\theta}) = \pi(s, a; \boldsymbol{\theta})$ .
    
```

5 Numerical Experiments

We verify the performance of our proposed algorithms in stochastic “torus” MDPs in section 5.1. The proposed algorithms and other existing PG algorithms are also applied to a pendulum balancing problem as a continuous state-action problem in section 5.2.

5.1 Torus $|\mathcal{S}|$ -State MDP. We tested the performance of our proposed algorithms in a stochastic one-dimensional torus grid-world with a finite set of grids $\mathcal{S} = \{1, \dots, |\mathcal{S}|\}$ and a set of two possible actions $\mathcal{A} = \{L, R\}$. This is a typical $|\mathcal{S}|$ -state MDP task where the state transition probabilities p are given by

$$\begin{cases} p(s-1 \mid s, L) = q_s \\ p(s \mid s, L) = \frac{1 - q_s}{2} \\ p(s+1 \mid s, L) = \frac{1 - q_s}{2} \end{cases} \quad \begin{cases} p(s-1 \mid s, R) = \frac{1 - q_s}{2} \\ p(s \mid s, R) = \frac{1 - q_s}{2} \\ p(s+1 \mid s, R) = q_s; \end{cases}$$

otherwise, $p = 0$, where $s = 0$ and $s = |\mathcal{S}|$ ($s = 1$ and $s = |\mathcal{S}| + 1$) are the identical states and $q_s \in [0, 1]$ is a task-dependent constant. In this experiment, a stochastic policy is a so-called Boltzmann or Gibbs policy represented by a sigmoidal function,

$$\pi(s, a = L; \boldsymbol{\theta}) = 1 - \pi(s, a = R; \boldsymbol{\theta}) = \frac{1}{1 + \exp(\boldsymbol{\theta}^\top \boldsymbol{\phi}(s))}.$$

Here, all of the elements of the state-feature vectors $\boldsymbol{\phi}(1), \dots, \boldsymbol{\phi}(|\mathcal{S}|) \in \mathcal{R}^{|\mathcal{S}|}$ were independently drawn from the gaussian distribution $\mathcal{N}(\mu = 0, \sigma^2 = 1)$ for each episode (simulation run). This was to assess how the parameterization of the stochastic policy affected the performance of our algorithms. The state-feature vectors $\boldsymbol{\phi}(s)$ were also used as the basis function for the LSD estimate $\mathbf{f}(s; \boldsymbol{\Omega})$ (see equation 3.9).

5.1.1 Performance of $\mathcal{L}SLSD(\lambda)$ Algorithm. First, we verified how precisely the $\mathcal{L}SLSD(\lambda)$ algorithm estimated $\nabla_{\boldsymbol{\theta}} \ln d_{\mathcal{M}(\boldsymbol{\theta})}(s)$ without regard to the setting of q_s and the policy parameter $\boldsymbol{\theta}$. Each element of $\boldsymbol{\theta}$ and the task-dependent constant q_s were initialized according to $\mathcal{N}(\mu = 0, \sigma^2 = 0.5^2)$ and $\mathcal{U}(a = 0.7, b = 1)$, respectively, where $\mathcal{U}(a = 0.7, b = 1)$ is the uniform distribution over the interval of $[a, b]$. These values were fixed during each episode.

Figure 1A shows a typical time course of the LSD estimates $\widehat{\nabla}_{\boldsymbol{\theta}} \ln d_{\mathcal{M}(\boldsymbol{\theta})}(s)$ for $|\mathcal{S}|=3$ -state MDP, where nine different gray scales indicate all of the

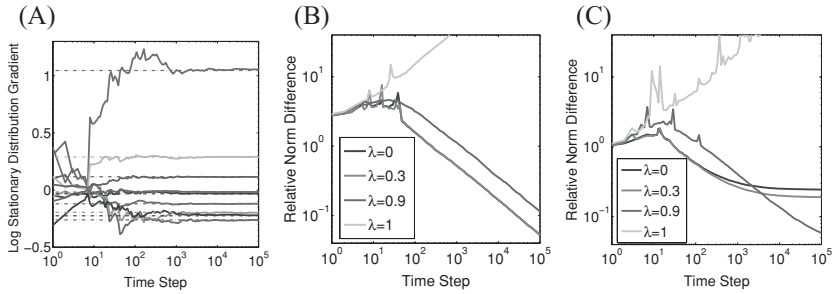


Figure 1: Performances of $\mathcal{L}SLSD(\lambda)$ for the estimation of the LSD $\nabla_{\theta} \ln d_{M(\theta)}(s)$. (A) A typical time course of LSD estimates in a three-state MDP. (B, C) The relative errors averaged over independent 200 episodes in seven-state MDPs for various λ (B) with a proper basis function $\phi(s) \in \mathcal{R}^7$ and (C) with an improper basis function $\phi(s) \in \mathcal{R}^6$.

different elements of the LSD, respectively. The solid lines denote the values estimated by $\mathcal{L}SLSD(0)$, and the dotted lines denote the analytical solution of the LSD. This result shows that the proposed algorithm, $\mathcal{L}SLSD(\lambda)$, can estimate the LSD $\nabla_{\theta} \ln d_{M(\theta)}(s)$. Besides this result, we have confirmed that the estimates by $\mathcal{L}SLSD(0)$ always converged to the analytical solution for $|\mathcal{S}| = 3$ as in Figure 1A.

Second, we investigated the effect of the eligibility decay rate λ using seven-state MDPs. In order to evaluate the average performance over various settings, we employed a relative error criterion that is defined by $\mathbb{E}_{M(\theta)}\{\|f(x; \Omega) - f(x; \Omega^*)\|^2\} / \mathbb{E}_{M(\theta)}\{\|f(x; \Omega^*)\|^2\}$, where Ω^* is the optimal parameter defined in proposition 3. Figures 1B and 1C show the time courses of the relative error averages over independent 200 episodes for $\lambda = 0, 0.3, 0.9$, and 1. The only difference between these two figures was the number of elements of the feature vectors $\phi(s)$. The feature vectors $\phi(s) \in \mathcal{R}^7$ used in Figure B were appropriate and sufficient to distinguish all of the different states, while the feature vectors $\phi(s) \in \mathcal{R}^6$ used in Figure 1C were inappropriate and insufficient. These results were consistent with the theoretical prospects. In particular, we could set λ arbitrarily in $[0, 1)$ if the basis function was appropriate (see Figure 1B); otherwise, we would need to set λ close but not equal to 1 (see Figure 1C).

5.1.2 Comparison to Other PG Methods. We compared the $\mathcal{L}SLSD(\lambda=0)$ -PG algorithm with the other PG algorithms for three-state MDPs, concerned with the estimation of PG $\nabla_{\theta} \eta(\theta)$ and the optimization of the policy parameter θ . The policy and the state transition probability were set as $\theta_i \sim \mathcal{N}(0, 0.5^2)$ and $q_i \sim \mathcal{U}(0.95, 1)$ for every $i \in \{1, 2, 3\}$, respectively.

Figure 2 shows the reward setting in the MDP. There are two types of rewards: $r = (\pm)2/Z(c)$ and $r = (\pm)c/Z(c)$, where the variable c was initialized using the uniform distribution over $[0.95, 1)$ for each episode (simulation run) and the function $Z(c)$ was the normalizing constant to ensure

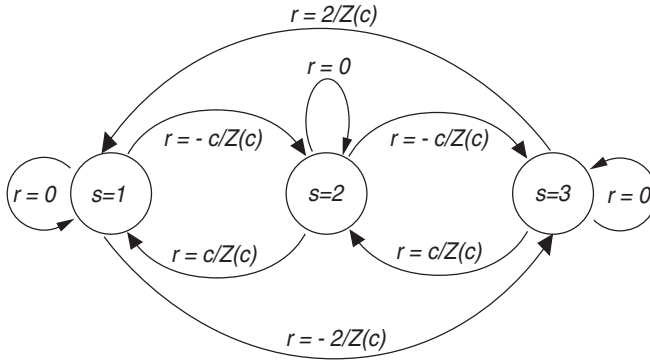


Figure 2: Reward setting of three-state MDPs used in our comparative studies. The value of c is selected from the uniform distribution $\mathcal{U}[0.95, 1)$ for each episode. $Z(c)$ is a normalizing function to ensure $\max_{\theta} \eta(\theta) = 1$.

$\max_{\theta} \eta(\theta) = 1$.¹² Note that the reward c defines the infimum value of γ to find the optimal policy: $\gamma^2 + \gamma > \frac{2c}{2-c}$. Therefore, the setting of γ is important but difficult in this task. From the performance baselines of the existing PG methods, we adopted two algorithms: GPOMDP (Baxter & Bartlett, 2001) and Konda's actor-critic (Konda & Tsitsiklis, 2003). In these algorithms, the state values were estimated by LSTD(λ) (Bradtke & Barto, 1996; Boyan, 2002; Yu & Bertsekas, 2006) and these estimates were used as the baseline functions to reduce the variance of PG estimates, while the baseline functions were not used in the original algorithms.

Figure 3 shows the results for the estimation of PG $\nabla_{\theta} \eta(\theta)$ from equation 4.1. The forgetting rate for the statistics and the eligibility decay rate were set as $\beta = 1$ and $\lambda = 0$ for all of the algorithms. Figures 3A and 3B represent the mean and the standard deviation of the angles between the estimates and the exact PG, respectively. These results show that $\mathcal{L}SLSD$ -PG, when estimating the optimal baseline function $b^*(s, s_{+1})$, termed $\mathcal{L}SLSD$ -PG: $b^*(s, s_{+1})$, worked best to estimate the PG. $\mathcal{L}SLSD$ -PG with $b(s, s_{+1})$ or $b^*(s, s_{+1})$ drastically improved the PG estimation performance for $\mathcal{L}SLSD$ -PG without a baseline function, $\mathcal{L}SLSD$ -PG:None, which was even worse than GPOMDP: $V(s)$, though $\mathcal{L}SLSD$ -PG:None worked better than GPOMDP:None. Thus, we confirmed that these baseline functions for $\mathcal{L}SLSD$ -PG are important.

Finally, we examined the optimization of the policy parameter θ , that is, the average reward, with these PG methods. In this experiment, the forgetting rate and the eligibility decay rate were set as $\beta = 0.99$ and $\lambda = 0$. In

¹²The normalizing constant $Z(c)$ was computed analytically with the unnormalized reward function as $Z(c) = \max_{\theta} \eta(\theta)$.

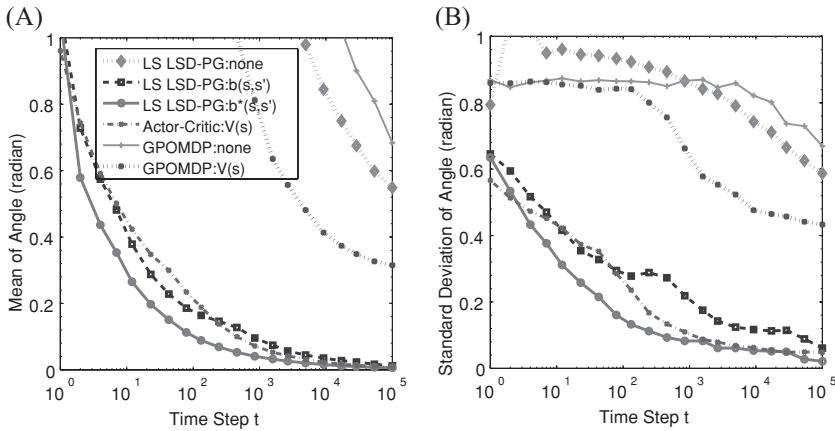


Figure 3: Comparison with various PG algorithms for the estimation of the PG over independent 2500 episodes: (A, B) The mean and the standard deviation of the angles between the estimates and the exact PG, respectively.

order to avoid the effect from poor estimations of the functions for the PG estimate, there was a prelearning period of 50 time steps, where the learning rate α was set to zero. It means that the policy remained unchanged in the first 50 time steps. Figure 4 shows the means and the standard deviations of the average rewards at an earlier stage (500 time step) and a later stage (10⁴ time step) over 1000 independent simulations of various learning rates α in order to give comparisons among the PG algorithms for the optimization of the policy parameter. It was confirmed that $\mathcal{L}SLS\mathcal{D}\text{-}PG:b^*(s, s_{+1})$ worked best, except for the high learning rate, where the learning speed of $b^*(s, s_{+1})$ could not properly follow the changes of the policy rather than that of $b(s, s_{+1})$. Figure 5 shows the time courses of the average reward, where we chose appropriate learning rates for the PG algorithms by drawing on the previous results; $\alpha = 0.16$ in $\mathcal{L}SLS\mathcal{D}\text{-}PG:b(s, s_{+1})$, $\alpha = 0.08$ in $\mathcal{L}SLS\mathcal{D}\text{-}PG:b^*(s, s_{+1})$, $\alpha = 0.08$ in Actor-Critic:V(s), and $\alpha = 0.007$ in GPOMDP:V(s). This result also indicates that our $\mathcal{L}SLS\mathcal{D}\text{-}PG$ algorithm with the optimal baseline function $b^*(s, s_{+1})$ outperformed the other PG algorithms in the sense of realizing both the highest average and the lowest standard deviation of the average rewards.

5.2 Continuous State-Action Problem. The $\mathcal{L}SLS\mathcal{D}(\lambda)\text{-}PG$ and the other existing PG algorithms were also applied to a continuous state-action problem. This task is to balance an inverted pendulum.

5.2.1 Interpretation of Continuous State Problem. Although this task obviously violates assumption 1, on which most policy gradient algorithms

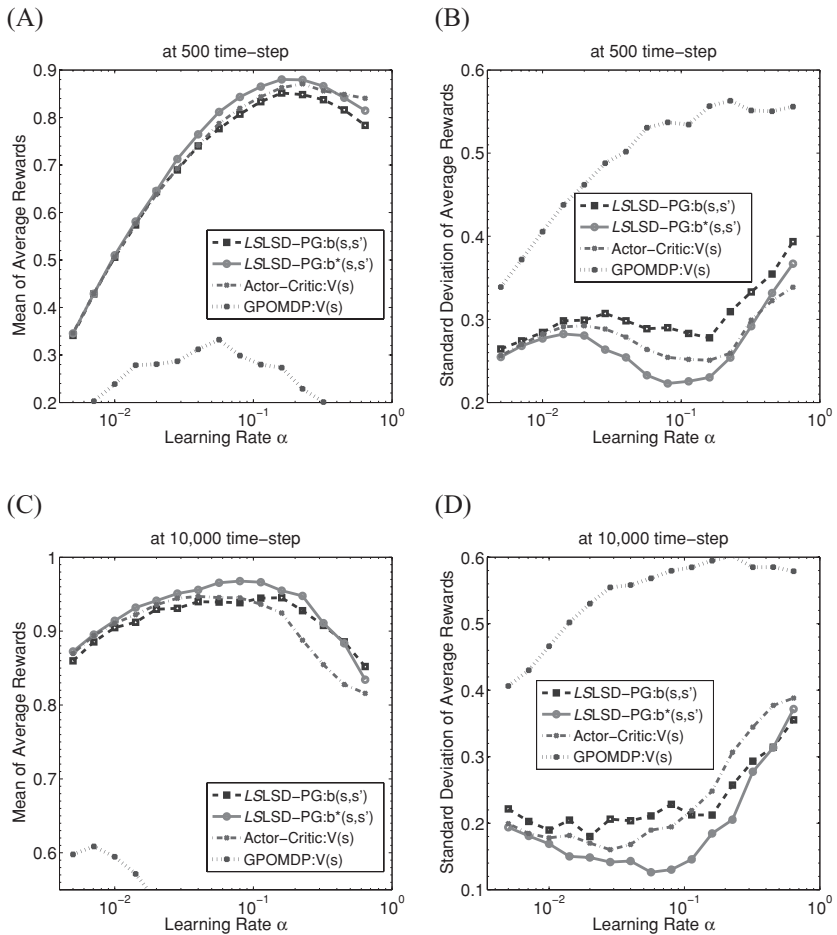


Figure 4: Comparisons with various PG algorithms for the means of the average rewards in the three-state torus MDPs with independent 1000 episodes about various learning rates. (A, B) The mean and the standard deviation of the average rewards at 500 time steps, respectively. (C, D) At 10^4 time step.

(including $\mathcal{L}\text{SLSD}(\lambda)$ -PG) have been based, it is valuable to acquire insights into the feasibility of PG algorithms. The reason is due to the following interpretations: a continuous problem to (i) a numerous-state MDP problem with some structures, or (ii) a partially observable MDP (POMDP) problem with belief states (Aberdeen, 2003). While the interpretation of (i) is straightforward, (ii) comes from the fact that when the policy (or the baseline function) in a continuous state problem is represented by a linear function with finite basis functions that output bounded activation values, then the activations

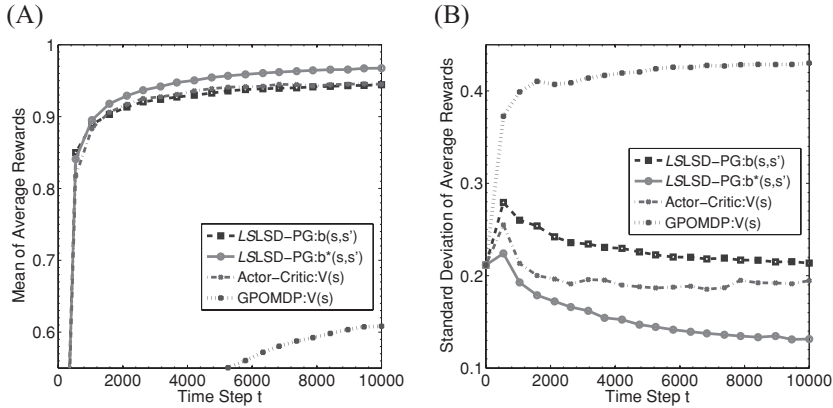


Figure 5: Comparison with various PG algorithms for the optimization of the policy parameters with the appropriate learning rate in the three-state torus MDPs over independent 1000 episodes. (A, B) Time courses of the mean and the standard deviation of the average rewards, respectively.

biased as nonnegative values and normalized can be regarded as the belief states of a finite-state POMDP (Aberdeen, 2003).

5.2.2 *Pendulum Balancing Problem.* A pendulum balancing problem is a well known benchmark in continuous RL problems (Peters et al., 2005; Morimura, Uchibe, & Doya, 2005). The state $s \equiv \{x, \dot{x}\}$ consists of the angle and the angular speed of the pendulum, which are limited to the ranges $[-\pi/6, \pi/6]$ and $[-\pi/2, \pi/2]$, respectively, as shown in Figure 6. Its dynamics is given by

$$\ddot{x}_{t+1} = \frac{-\mu\dot{x}_t + mgl \sin(x_t) + a}{ml^2},$$

where a is the torque as an action selected by a learning agent. The physical parameters are set as $m = l = 1, g = 9.8$, and $\mu = 0.01$. The reward function is

$$r(s_t, a_t, s_{t+1}) \equiv \begin{cases} -x_{t+1}^2 - 0.5\dot{x}_{t+1}^2 - 0.001a_t^2 - 1, & \text{if } |x_{t+1}| > \pi/6 \\ & \text{or } |\dot{x}_{t+1}| > \pi/2, \\ -x_{t+1}^2 - 0.5\dot{x}_{t+1}^2 - 0.001a_t^2, & \text{otherwise.} \end{cases}$$

The state s_t is initialized with the uniform distributions as $x_t \sim \mathcal{U}(-\pi/8, \pi/8)$ and $\dot{x}_t \sim \mathcal{U}(-1, 1)$, at the beginning of each episode or when the previous state s_{t-1} deviates from the permissible ranges, that is, $|x_{t-1}| > \pi/6$ or

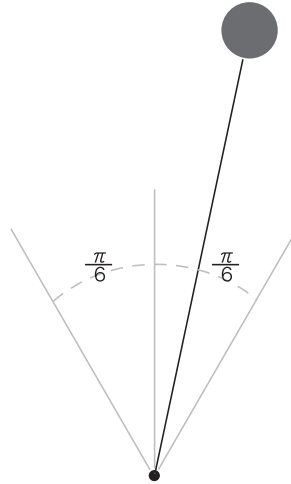


Figure 6: Pendulum balancing problem near the top ranges: $x \in [-\pi/6, \pi/6]$ and $\dot{x} \in [-\pi/2, \pi/2]$.

$|\dot{x}_{t-1}| > \pi/2$. The state is also initialized with the probability 0.01 at each time step in order to explore the state space more efficiently.

Since this problem has a potentially infinite number of states, we use a normalized radial basis function (nRBF) model (Doya, 2000) for the following policy and the linear baseline function:

$$\pi(s, a; \theta) \equiv \frac{1}{\sqrt{2\pi}} \exp \left[-\frac{\{a - \theta^\top \phi(s)\}^2}{2} \right],$$

where $\phi(s)$ is the output of the nRBF model with 3×3 RBFs as follows. The centers of the RBFs are $x \in \{-\pi/9, 0, \pi/9\}$ and $\dot{x} \in \{-\pi/3, 0, \pi/3\}$. The standard deviations of the RBFs are $\pi/9$ and $\pi/3$ for x and \dot{x} , respectively. The policy parameter $\theta \in \mathcal{R}^9$ was set to $\mathbf{0}$ at the beginning of each episode.

We used our \mathcal{L} SLSD(λ)-PGs with optimal baseline function $b^*(s, s_{+1})$ and decent baseline function $b(s, s_{+1})$, (Konda's) actor-critic: $V(s)$, and GPOMDP: $V(s)$, as in the previous experiment (see section 5.1.2). We set the meta-parameters of those algorithms appropriately: $\alpha = 0.2$ and $\beta = 0.9995$ for the \mathcal{L} SLSD-PGs and actor-critic, and $\alpha = 0.04$, $\gamma = 0.99$, and $\beta = 0.9995$ for the GPOMDP. Also, λ for the \mathcal{L} SLSD-PGs with $b^*(s, s_{+1})$ and $b(s, s_{+1})$, the actor-critic, and the GPOMDP were set to 0.7, 0.5, 0.5, and 0. There was a pre learning period of 10^3 time steps to avoid any effects from poor PG estimates, where α was set to zero.

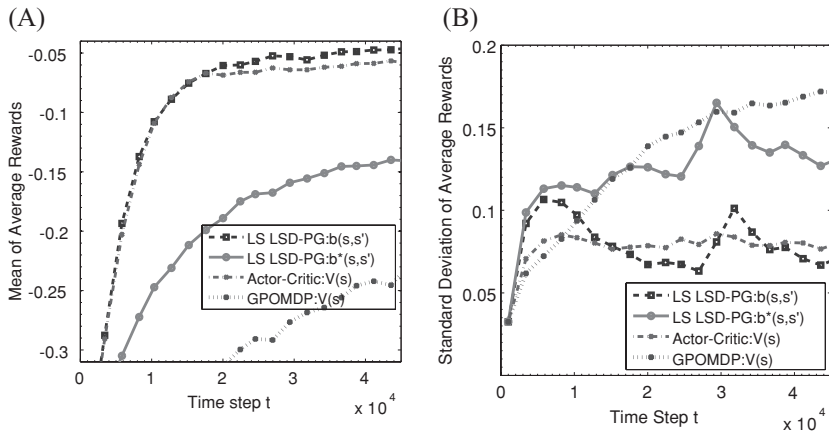


Figure 7: Comparisons with various PG algorithms for the optimization of the policy parameters in the pendulum balancing problem over independent 500 episodes. (A, B) Time courses of the mean and the standard deviation of the average awards, respectively.

Figure 7A shows the time courses of the mean and Figure 7B the standard deviations for the average rewards. We confirmed that the performance of $\mathcal{L}SLS\text{D-PG}:b(s, s_{+1})$ was better than that of $\mathcal{L}SLS\text{D-PG}:b^*(s, s_{+1})$. Because its order in the previous MDP problem (see Figure 5) was reversed, it must be caused by the difficulty in learning the parameter v^* of the “optimal” baseline function $b^*(s, s_{+1})$ in equation 4.6 in contrast to v^* of the “decent” one $b(s, s_{+1})$ in equations 4.8 with the nRBF model. This difference of the difficulties indicates that $b^*(s, s_{+1})$ could be more complex than $b(s, s_{+1})$ and $b^*(s, s_{+1})$ would need more representation capability for the baseline function approximator. From the comparison of these forms in equations 4.6 and 4.8, these observations would come from the existence of the weights $\|\nabla_{\theta} \ln \pi(s, a; \theta) + \widehat{\nabla}_{\theta} \ln d_{M(\theta)}(s)\|^2$, that is, the weights would often make the estimation of $b^*(s, s_{+1})$ with a function approximation difficult. We also confirmed from Figure 7 that the performance of the $\mathcal{L}SLS\text{D-PG}:b(s, s_{+1})$ was much better than the GPOMDP: $V(s)$ and slightly better than the actor-critic: $V(s)$ algorithm.

6 Related Work

There are two alternative methods that estimate the derivative of the (stationary) state distribution and have already been proposed in Glynn (1991) or Rubinstein (1991) and Ng et al. (2000). However, these are different from our approach and have the following problems. The method in Glynn or Rubinstein from operations research, is called the likelihood ratio gradient or the score function. This method can be problematic as to how to design

the recurrence state, since the applicability is limited to regenerative processes (see Baxter & Bartlett, 2001, in detail). The method proposed in Ng et al. is not a direct estimation of the derivative of the state distribution but is done by estimating of the state distribution with density propagation. Accordingly, both methods require knowledge of which state the agent is in, while our method needs only to observe the feature vector of the state.

Meanwhile, there is an average reward PG algorithm (Tsitsiklis & Van Roy, 1999; Konda & Tsitsiklis, 2003; Sutton et al., 2000) that eliminates the use of the forgetting rate by introducing a differential cost function as a solution of Poisson's equation (also known as the average reward Bellman equation in RL). However, since to date this is a unique PG framework proposed for maximizing the average reward,¹³ more studies of the average reward optimization would be needed and they would have to be significant. At least one possible advantage of the proposed framework over the existing average reward PG method is that a closed-form solution of an optimal baseline function for minimizing the variance bound of PG estimate can be computed by least-squares approaches, while such a solution in conventional PG frameworks has not been explored and would be intractable (Greensmith et al., 2004).

7 Conclusion

Our propositions show that the actual forward and virtual backward Markov chains are closely related and have common properties. Utilizing these properties, we proposed $\mathcal{L}SLSD(\lambda)$ as an estimation algorithm for the log stationary distribution derivative (LSD) and $\mathcal{L}SLSD(\lambda)$ -PG as a PG algorithm utilizing the LSD estimate. The experimental results also demonstrated that $\mathcal{L}SLSD(\lambda)$ worked for $\lambda \in [0, 1)$ and $\mathcal{L}SLSD(\lambda)$ -PG could learn regardless of the task's requirements for the value of γ to optimize the average reward. At the same time, it has been suggested that there is theoretically no significant difference in performances between the average-reward-based PG methods and the alternative PG methods with forgetting rate γ for the value functions close to one (Tsitsiklis & Van Roy, 2002). This might be true for the case of our proposed PG, $\mathcal{L}SLSD$ -PG, which would mean that $\mathcal{L}SLSD$ -PG might not drastically improve the performances of the PG methods with γ as does the algorithm proposed by Kimura and Kobayashi (1998) when γ was set appropriately. However, it is noted that $\mathcal{L}SLSD$ -PG is free of the setting of γ . In contrast, the learning performances of Konda's actor-critic and the $\mathcal{L}SLSD$ -PG approaches do not seem to be the significantly different, as confirmed in our numerical experiments. This would seem to be because the $\mathcal{L}SLSD$ -PG is just a dual approach compared to Konda's actor-critic approach. In dynamic programming (DP), Wang, Bowling, & Schuurmans (2007) formalize the

¹³Although R-learning also maximizes the average reward, it is based on the value function not the PG algorithm (Sutton & Barto, 1998).

dual approaches, in which probability distributions are maintained instead of the value functions (primal approaches). In contrast, in PGRL, the \mathcal{L} SLSLSD-PG approach maintains the LSD instead of the value functions to estimate the policy gradient. Considering the relationships between DP and PGRL, Konda's actor-critic and our \mathcal{L} SLSLSD-PG approaches correspond to primal and dual approaches in PGRL, respectively¹⁴. Since Wang, Lizotte, Bowling, & Schuurmans (2008) also show the advantages of the dual DP approaches in that the dual updates will not diverge even with function approximations because of the constraints caused by the probability distributions, the \mathcal{L} SLSLSD-PG approaches may share these advantages due to the constraint of equation 3.8. More theoretical and experimental work is necessary to understand the effectiveness further.

Meanwhile, in our numerical experiments, the \mathcal{L} SLSLSD-PG approaches greatly surpassed the performances of GPOMDP. This may be because \mathcal{L} SLSLSD-PG and GPOMDP are fundamentally different approaches, since GPOMDP can be regarded as an approach based on a fully model-free estimation of the PG, while \mathcal{L} SLSLSD-PG can be regarded as a model-based estimation approach under certain conditions (see appendix B for the conditions). Although the latent difficulty of the PG estimation problem is unchanged, the \mathcal{L} SLSLSD-PG utilizes model information or prior knowledge of the model. Also, since the policy is usually a parametric model and the LSD approximator would be defined as a model similar to the policy model, the utilization of the policy model in \mathcal{L} SLSLSD-PG may be useful in many cases. Accordingly, it is important for future work to discuss and define a necessary and sufficient basis function for the (linear) LSD approximator based on the parameterization of the policy.

On the other hand, the use of LSD estimation will open up new possibilities for the natural policy gradient learning (Kakade, 2001; Peters, Vijayakumar, & Schaal, 2005; Peters & Schaal, 2008; Morimura, Uchibe, Yoshimoto, & Doya, 2008). It enables us to compute a valid Riemannian metric matrix $G(\theta)$ for the natural gradient, which is effective especially in the large-scale MDPs,

$$G(\theta) := \mathbb{E}_{M(\theta)} \left\{ \nabla_{\theta} \ln \pi(s, a; \theta) \nabla_{\theta} \ln \pi(s, a; \theta)^{\top} + \iota \nabla_{\theta} \ln d_{M(\theta)}(s) \nabla_{\theta} \ln d_{M(\theta)}(s)^{\top} \right\},$$

where $\iota \in [0, 1]$ interpolates the natural policy gradient ($\iota = 0$) and the natural state-action gradient ($\iota = 1$) (Morimura, Uchibe, Yoshimoto et al., 2008).

In addition, the use of LSD estimation might offer novel methods for addressing the trade-offs between exploration and exploitation. This is because LSD gives statistical information about how much of a change of the state stationary distribution is caused by the perturbation of each element of

¹⁴However, this is unclear, because the conversion operation on PGRL from the primal to the dual (or from the dual to the primal) is not obvious, while it should be automatic.

each policy parameter, while the stationary distribution with low entropy would make the exploration difficult.

Appendix A: Derivation of Equation 2.4

Using the relation between the forgetting (or discounted) state-value function $V_\gamma^\pi(s)$ and the average reward $\eta(\theta)$ (Singh, Jaakkola, & Jordan, 1994)

$$\eta(\theta) = (1 - \gamma) \sum_{s \in \mathcal{S}} d_{M(\theta)}(s) V_\gamma^\pi(s),$$

the derivative of average reward with respect to θ is given by

$$\nabla_\theta \eta(\theta) = (1 - \gamma) \left(\sum_{s \in \mathcal{S}} \nabla_\theta d_{M(\theta)}(s) V_\gamma^\pi(s) + \sum_{s \in \mathcal{S}} d_{M(\theta)}(s) \nabla_\theta V_\gamma^\pi(s) \right), \quad (\text{A.1})$$

where $\nabla_\alpha AB$ implies $(\nabla_\alpha A)B$. The second term is modified as follows:

$$\begin{aligned} & \sum_{s \in \mathcal{S}} d_{M(\theta)}(s) \nabla_\theta V_\gamma^\pi(s) \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_{M(\theta)}(s) \nabla_\theta \{ \pi(s, a; \theta) Q_\gamma^\pi(s, a) \} \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_{M(\theta)}(s) [\nabla_\theta \pi(s, a; \theta) Q_\gamma^\pi(s, a) + \pi(s, a; \theta) \nabla_\theta Q_\gamma^\pi(s, a)] \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_{M(\theta)}(s) \left[\nabla_\theta \pi(s, a; \theta) Q_\gamma^\pi(s, a) \right. \\ & \quad \left. + \pi(s, a; \theta) \nabla_\theta \sum_{s_{+1} \in \mathcal{S}} p(s_{+1} | s, a) \{ r(s, a, s_{+1}) + \gamma V_\gamma^\pi(s_{+1}) \} \right] \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_{M(\theta)}(s) \nabla_\theta \pi(s, a; \theta) Q_\gamma^\pi(s, a) + \gamma \sum_{s \in \mathcal{S}} d_{M(\theta)}(s) \nabla_\theta V_\gamma^\pi(s) \quad (\text{A.2}) \\ &= \frac{1}{1 - \gamma} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_{M(\theta)}(s) \nabla_\theta \pi(s, a; \theta) Q_\gamma^\pi(s, a). \quad (\text{A.3}) \end{aligned}$$

Equation A.2 is given by the property of stationary distribution, equation 2.1. Substituting equation A.3 in equation A.1, we can derive equation 2.4:

$$\begin{aligned} \nabla \eta(\theta) &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_{M(\theta)}(s) \pi(s, a; \theta) \nabla_\theta \ln \pi(s, a; \theta) Q_\gamma^\pi(s, a) \\ & \quad + (1 - \gamma) \sum_{s \in \mathcal{S}} d_{M(\theta)}(s) \nabla_\theta \ln d_{M(\theta)}(s) V_\gamma^\pi(s). \quad (\text{A.4}) \end{aligned}$$

Appendix B: $\mathcal{L}SLSD(\lambda)$ as a Model-Based Learning

Here, we show that $\mathcal{L}SLSD(\lambda)$ can be regarded as a model-based estimation of the LSD $\nabla_{\theta} \ln d_{M(\theta)}(s)$ under certain conditions.¹⁵ First, we consider a model-based LSD estimation. With the following matrix notations,

$$\begin{aligned}
 S_{\theta} &\equiv \begin{bmatrix} \nabla_{\theta} \ln d_{M(\theta)}(s=1)^{\top} \\ \vdots \\ \nabla_{\theta} \ln d_{M(\theta)}(s=|\mathcal{S}|)^{\top} \end{bmatrix}, & \Lambda_{\theta}(a) &\equiv \begin{bmatrix} \nabla_{\theta} \ln \pi(s=1, a; \theta)^{\top} \\ \vdots \\ \nabla_{\theta} \ln \pi(s=|\mathcal{S}|, a; \theta)^{\top} \end{bmatrix}, \\
 Q_{B(\theta)}^a(a) &\equiv \begin{bmatrix} q_{B(\theta)}(s=1, a | s_{+1}=1) & \dots & q_{B(\theta)}(s=|\mathcal{S}|, a | s_{+1}=1) \\ \vdots & \ddots & \vdots \\ q_{B(\theta)}(s=1, a | s_{+1}=|\mathcal{S}|) & \dots & q_{B(\theta)}(s=|\mathcal{S}|, a | s_{+1}=|\mathcal{S}|) \end{bmatrix},
 \end{aligned}$$

the recursive equation of the LSD, equation 3.5 is rewritten as

$$\begin{aligned}
 S_{\theta} &= \sum_{a \in \mathcal{A}} Q_{B(\theta)}^a(a) (\Lambda_{\theta}(a) + S_{\theta}) \\
 &= \bar{\Lambda}_{\theta} + Q_{B(\theta)} S_{\theta},
 \end{aligned}$$

where $Q_{B(\theta)} \equiv \sum_{a \in \mathcal{A}} Q_{B(\theta)}^a(a) = \{q_{B(\theta)}(s=j | s_{+1}=i)\}_{i,j}$ and $\bar{\Lambda}_{\theta} \equiv \sum_{a \in \mathcal{A}} Q_{B(\theta)}^a(a) \Lambda_{\theta}(a)$. Although this equation is transformed to

$$(I - Q_{B(\theta)}) S_{\theta} = \bar{\Lambda}_{\theta}, \tag{B.1}$$

the rank of $(I - Q_{B(\theta)})$ is equal to $|\mathcal{S}| - 1$, and thus its inverse does not exist, because the backward transition probability matrix has a unique stationary distribution from proposition 1, $d_{\theta} = Q_{B(\theta)}^{\top} d_{\theta}$, that is, $(I - Q_{B(\theta)})^{\top} d_{\theta} = \mathbf{0}$. However, $(I - Q_{B(\theta)} - d_{\theta} d_{\theta}^{\top})$ is always invertible because of $d_{\theta} d_{\theta}^{\top} d_{\theta} \neq \mathbf{0}$. By utilizing equation 3.8 as the matrix notation, $d_{\theta}^{\top} S_{\theta} = \mathbf{0}$, equation B.1 is transformed to

$$\begin{aligned}
 (I - Q_{B(\theta)} - d_{\theta} d_{\theta}^{\top}) S_{\theta} &= \bar{\Lambda}_{\theta} \\
 \Leftrightarrow S_{\theta} &= (I - Q_{B(\theta)} - d_{\theta} d_{\theta}^{\top})^{-1} \bar{\Lambda}_{\theta}.
 \end{aligned} \tag{B.2}$$

¹⁵Model based indicates that target values are estimated through the estimation of the model parameters for these target values.

Therefore, in a model-based approach for the LSD estimation, the model parameters $Q_{B(\theta)}$, d_θ , and $\bar{\Lambda}_\theta$ will be estimated as sufficient statistics $\hat{Q}_{B(\theta)}$, \hat{d}_θ , and $\hat{\Lambda}_\theta$, respectively. Then the LSD estimate \hat{S}_θ is computed as

$$\hat{S}_\theta = (I - \hat{Q}_{B(\theta)} - \hat{d}_\theta \hat{d}_\theta^\top)^{-1} \hat{\Lambda}_\theta. \tag{B.3}$$

Under the conditions that the feature vectors $\phi(s) \in \mathcal{R}^{|\mathcal{S}|}$ are $\phi(s=1) = [1, 0, \dots, 0]^\top$, $\phi(s=2) = [0, 1, \dots, 0]^\top, \dots$, and $\phi(s=|\mathcal{S}|) = [0, 0, \dots, 1]^\top$ and the eligibility decay rate λ is equal to 0, the statistics A , B , and c in the $\mathcal{L}\text{SLSLSD}(\lambda = 0)$ (see algorithm 1) can be the sufficient statistics for these model parameters as

$$\begin{aligned} T(I - \hat{Q}_{B(\theta)}) &= A, \\ T\hat{\Lambda}_\theta &= B, \\ T\hat{d}_\theta &= c, \end{aligned}$$

where T is the time-step in Algorithm 1. Then using equation B.3, the LSD estimate in this $\mathcal{L}\text{SLSLSD}$ algorithm, $\hat{S}_{\mathcal{L}\text{SLSLSD}}$, is computed as

$$\begin{aligned} \hat{S}_{\mathcal{L}\text{SLSLSD}} &= \begin{bmatrix} \phi(s=1)^\top \\ \vdots \\ \phi(s=|\mathcal{S}|)^\top \end{bmatrix} (A - cc^\top/T)^{-1} B \\ &= I \{ T(I - \hat{Q}_{B(\theta)}) - T\hat{d}_\theta \hat{d}_\theta^\top \}^{-1} T\hat{\Lambda}_\theta \\ &= \hat{S}_\theta. \end{aligned}$$

Therefore, the $\mathcal{L}\text{SLSLSD}(\lambda = 0)$ algorithm with these table-lookup features of the states is equivalent to the model-based LSD $\nabla_\theta \ln d_{M(\theta)}(s)$ estimation. It is noted that while the $\text{LSTD}(\lambda)$ by Boyan (2002) is regarded as a forward model-based approach for the value function estimation, our $\mathcal{L}\text{SLSLSD}(\lambda)$ can be regarded as a backward model-based approach for the LSD function estimation.

Appendix C: Proof of Proposition 4 ---

If the following equation holds,

$$\begin{aligned} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \sum_{s_{+1} \in \mathcal{S}} d_{M(\theta)}(s) \pi(s, a; \theta) p(s_{+1} | s, a) \\ \times \{ \nabla_\theta \ln \pi(s, a; \theta) + \nabla_\theta \ln d_{M(\theta)}(s) \} \rho(s, s_{+1}) = \mathbf{0} \end{aligned} \tag{C.1}$$

then the transformation to equation 4.4 is obviously true. Because of equation 4.3 and

$$\begin{cases} \sum_{a \in \mathcal{A}} \pi(s, a; \theta) \nabla_{\theta} \ln \pi(s, a; \theta) c = \nabla_{\theta} c = \mathbf{0}, \\ \sum_{s \in \mathcal{S}} d_{M(\theta)}(s) \nabla_{\theta} \ln d_{M(\theta)}(s) c = \nabla_{\theta} c = \mathbf{0}, \end{cases}$$

we know that

$$\begin{aligned} & \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \sum_{s_{+1} \in \mathcal{S}} d_{M(\theta)}(s) \pi(s, a; \theta) p(s_{+1} | s, a) \\ & \quad \times \{ \nabla_{\theta} \ln \pi(s, a; \theta) + \nabla_{\theta} \ln d_{M(\theta)}(s) \} \rho(s, s_{+1}) \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \sum_{s_{+1} \in \mathcal{S}} d_{M(\theta)}(s) \pi(s, a; \theta) p(s_{+1} | s, a) \\ & \quad \times \{ \nabla_{\theta} \ln \pi(s, a; \theta) + \nabla_{\theta} \ln d_{M(\theta)}(s) \} \{ g(s) - g(s_{+1}) \}. \end{aligned} \tag{C.2}$$

Since a time average is equivalent to a state-action space average in an ergodic Markov chain $M(\theta)$ by assumption 1, equation C.2 is transformed to

$$\begin{aligned} & \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} \{ \nabla_{\theta} \ln \pi(s_t, a_t; \theta) + \nabla_{\theta} \ln d_{M(\theta)}(s_t) \} \{ g(s_t) - g(s_{t+1}) \} \\ &= \lim_{T \rightarrow \infty} \frac{1}{T} \left[\{ \nabla_{\theta} \ln \pi(s_0, a_0; \theta) + \nabla_{\theta} \ln d_{M(\theta)}(s_0) \} g(s_0) \right. \\ & \quad - \{ \nabla_{\theta} \ln \pi(s_T, a_T; \theta) + \nabla_{\theta} \ln d_{M(\theta)}(s_T) \} g(s_T) \\ & \quad + \sum_{t=1}^T \{ -\nabla_{\theta} \ln \pi(s_{t-1}, a_{t-1}; \theta) - \nabla_{\theta} \ln d_{M(\theta)}(s_{t-1}) + \nabla_{\theta} \ln \pi(s_t, a_t; \theta) \\ & \quad \left. + \nabla_{\theta} \ln d_{M(\theta)}(s_t) \} g(s_t) \right] \\ &= \sum_{s_{-1} \in \mathcal{S}} \sum_{a_{-1} \in \mathcal{A}} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_{M(\theta)}(s_{-1}) \pi(s_{-1}, a_{-1}; \theta) p(s | s_{-1}, a_{-1}) \pi(s, a; \theta) \\ & \quad \{ -\nabla_{\theta} \ln \pi(s_{-1}, a_{-1}; \theta) - \nabla_{\theta} \ln d_{M(\theta)}(s_{-1}) + \nabla_{\theta} \ln \pi(s, a; \theta) \\ & \quad + \nabla_{\theta} \ln d_{M(\theta)}(s) \} g(s) \\ &= \sum_{s_{-1} \in \mathcal{S}} \sum_{a_{-1} \in \mathcal{A}} \sum_{s \in \mathcal{S}} d_{M(\theta)}(s_{-1}) \pi(s_{-1}, a_{-1}; \theta) p(s | s_{-1}, a_{-1}) \\ & \quad \{ -\nabla_{\theta} \ln \pi(s_{-1}, a_{-1}; \theta) - \nabla_{\theta} \ln d_{M(\theta)}(s) + \nabla_{\theta} \ln d_{M(\theta)}(s) \} g(s) \end{aligned}$$

$$= \sum_{s \in \mathcal{S}} d_{M(\theta)}(s) g(s) \left[\mathbb{E}_{B(\theta)} \{ \nabla_{\theta} \ln \pi(s_{-1}, a_{-1}; \theta) \right. \\ \left. + \nabla_{\theta} \ln d_{M(\theta)}(s_{-1}) \mid s \} - \nabla_{\theta} \ln d_{M(\theta)}(s) \right] \quad (\text{C.3})$$

$$= \mathbf{0}, \quad (\text{C.4})$$

where equation 3.4 (proposition 2) and equation 3.5 are used for the transformations to equations C.3 and C.4. Therefore, equation C.1 holds.

References

- Aberdeen, D. (2003). *Policy-gradient algorithms for partially observable Markov decision processes*. Unpublished doctoral dissertation, Australian National University.
- Baird, L., & Moore, A. (1999). Gradient descent for general reinforcement learning. In M. S. Kearns, S. Solla, & D. Cohn (Eds.), *Advances in neural information processing systems, 11*. Cambridge, MA: MIT Press.
- Baxter, J., & Bartlett, P. (2001). Infinite-horizon policy-gradient estimation. *Journal of Artificial Intelligence Research, 15*, 319–350.
- Baxter, J., Bartlett, P., & Weaver, L. (2001). Experiments with infinite-horizon policy-gradient estimation. *Journal of Artificial Intelligence Research, 15*, 351–381.
- Bertsekas, D. P. (1995). *Dynamic programming and optimal control*. Belmont, MA: Athena Scientific.
- Bertsekas, D. P., & Tsitsiklis, J. N. (1996). *Neuro-dynamic programming*. Belmont, MA: Athena Scientific.
- Boyan, J. A. (2002). Technical update: Least-squares temporal difference Learning. *Machine Learning, 49*, 233–246.
- Bradtke, S. J., & Barto, A. G. (1996). Linear least-squares algorithms for temporal difference learning. *Machine Learning, 22*, 33–57.
- Doya, K. (2000). Reinforcement learning in continuous time and space. *Neural Computation, 12*, 219–245.
- Glynn, P. W. (1991). Likelihood ratio gradient estimation for stochastic systems. *Communications of the ACM, 33*(10), 75–84.
- Greensmith, E., Bartlett, P., & Baxter, J. (2004). Variance reduction techniques for gradient estimates in reinforcement learning. *Journal of Machine Learning Research, 5*, 1471–1530.
- Gullapalli, V. (1990). A stochastic reinforcement learning algorithm for learning real-valued functions. *Neural Networks, 3*, 671–692.
- Kakade, S. (2001). Optimizing average reward using discounted rewards. In D. Helmbold & B. Williamson (Eds.), *Annual Conference on Computational Learning Theory, 14*. Cambridge, MA: MIT Press.
- Kimura, H., & Kobayashi, S. (1998). An analysis of actor/critic algorithms using eligibility traces: Reinforcement learning with imperfect value function. In *International Conference on Machine Learning* (pp. 278–286). San Francisco: Morgan Kaufmann.
- Konda, V. S., & Tsitsiklis, J. N. (2003). On actor-critic algorithms. *SIAM Journal on Control and Optimization, 42*, 1143–1166.

- Lagoudakis, M. G., & Parr, R. (2003). Least-squares policy iteration. *Journal of Machine Learning Research*, 4, 1107–1149.
- MacKay, D. (2003). *Information theory, inference, and learning algorithms*. Cambridge: Cambridge University Press.
- Morimura, T., Uchibe, E., & Doya, K. (2005). Utilizing natural gradient in temporal difference reinforcement learning with eligibility traces. In *International Symposium on Information Geometry and Its Applications* (pp. 256–263).
- Morimura, T., Uchibe, E., & Doya, K. (2008). Natural actor-critic with baseline adjustment for variance reduction. *Artificial Life and Robotics*, 13, 275–279.
- Morimura, T., Uchibe, E., Yoshimoto, J., & Doya, K. (2007). Reinforcement learning with log stationary distribution gradient (Tech. Rep.). Nara: Nara Institute of Science and Technology.
- Morimura, T., Uchibe, E., Yoshimoto, J., & Doya, K. (2008). A new natural policy gradient by stationary distribution metric. In *European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases*. Berlin: Springer-Verlag.
- Morimura, T., Uchibe, E., Yoshimoto, J., & Doya, K. (in press). A generalized natural actor-critic algorithm. In *Advances in neural information processing systems*. Cambridge, MA: MIT Press.
- Ng, A. Y., Parr, R., & Koller, D. (2000). Policy search via density estimation. In S. Solla, T. K. Leen, & K.-R. Müller (Eds.), *Advances in neural information processing systems*, Cambridge, MA: MIT Press.
- Peng, J., & Williams, R. J. (1996). Incremental multi-step Q-learning. *Machine Learning*, 22, 283–290.
- Peters, J., & Schaal, S. (2006). Policy gradient methods for robotics. In *Proceedings of the IEEE International Conference on Intelligent Robots and Systems*. Piscataway, NJ: IEEE Press.
- Peters, J., & Schaal, S. (2008). Natural actor-critic. *Neurocomputing*, 71, 1180–1190.
- Peters, J., Vijayakumar, S., & Schaal, S. (2005). Natural actor-critic. In *European Conference on Machine Learning*. Berlin: Springer-Verlag.
- Rubinstein, R. Y. (1991). How to optimize discrete-event system from a single sample path by the score function method. *Annals of Operations Research*, 27, 175–212.
- Schinazi, R. B. (1999). *Classical and spatial stochastic processes*. Boston: Birkhauser.
- Singh, S. P., Jaakkola, T., & Jordan, M. I. (1994). Learning without state-estimation in partially observable Markovian decision processes. In *International Conference on Machine Learning* (pp. 284–292). San Francisco: Morgan Kaufmann.
- Sutton, R. S. (1988). Learning to predict by the methods of temporal differences. *Machine Learning*, 3, 9–44.
- Sutton, R. S., & Barto, A. G. (1998). *Reinforcement learning*. Cambridge, MA: MIT Press.
- Sutton, R. S., McAllester, D., Singh, S., & Mansour, Y. (2000). Policy gradient methods for reinforcement learning with function approximation. In S. A. Solla, T. K. Leen, & K.-R. Müller (Eds.), *Advances in neural information processing systems*, 12. Cambridge, MA: MIT Press.
- Tsitsiklis, J. N., & Van Roy, B. (1999). Average cost temporal-difference learning. *Automatica*, 35, 1799–1808.

- Tsitsiklis, J. N., & Van Roy, B. (2002). On average versus discounted reward temporal-difference learning. *Machine Learning*, 49, 179–191.
- Ueno, T., Kawanabe, M., Mori, T., Maeda, S., & Ishii, S. (2008). A semiparametric statistical approach to model-free policy evaluation. In *International Conference on Machine Learning* (pp. 857–864). New York: ACM.
- Wang, T., Bowling, M., & Schuurmans, D. (2007). Dual representations for dynamic programming and reinforcement learning. In *Proceedings of the IEEE International Symposium on Approximate Dynamic Programming and Reinforcement Learning* (pp. 44–51). Piscataway, NJ: IEEE.
- Wang, T., Lizotte, D., Bowling, M., & Schuurmans, D. (2008). Stable dual dynamic programming. In J. C. Platt, D. Koller, Y. Singer, & S. Roweis (Eds.), *Advances in neural information processing systems*, 20. Cambridge, MA: MIT Press.
- Williams, R. J. (1992). Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning*, 8, 229–256.
- Young, P. (1984). *Recursive estimation and time-series analysis*. Berlin: Springer-Verlag.
- Yu, H., & Bertsekas, D. P. (2006). *Convergence results for some temporal difference methods based on least squares* (Tech. Rep. LIDS 2697). Cambridge, MA: MIT Press.

Received December 12, 2008; accepted May 20, 2009.