

A New Natural Policy Gradient by Stationary Distribution Metric

Tetsuro Morimura^{1,2}, Eiji Uchibe¹, Junichiro Yoshimoto^{1,3}, and Kenji Doya^{1,3,4}

¹ Initial Research Project, Okinawa Institute of Science and Technology

² IBM Research, Tokyo Research Laboratory

³ Graduate School of Information Science, Nara Institute of Science and Technology

⁴ ATR Computational Neuroscience Laboratories

{morimura, uchibe, jun-y, doya}@oist.jp

Abstract. The parameter space of a statistical learning machine has a Riemannian metric structure in terms of its objective function. Amari [1] proposed the concept of “natural gradient” that takes the Riemannian metric of the parameter space into account. Kakade [2] applied it to policy gradient reinforcement learning, called a natural policy gradient (NPG). Although NPGs evidently depend on the underlying Riemannian metrics, careful attention was not paid to the alternative choice of the metric in previous studies. In this paper, we propose a Riemannian metric for the joint distribution of the state-action, which is directly linked with the average reward, and derive a new NPG named “*Natural State-action Gradient*” (NSG). Then, we prove that NSG can be computed by fitting a certain linear model into the immediate reward function. In numerical experiments, we verify that the NSG learning can handle MDPs with a large number of states, for which the performances of the existing (N)PG methods degrade.

Keywords: policy gradient reinforcement learning, natural gradient, Riemannian metric matrix, Markov decision process.

1 Introduction

Policy gradient reinforcement learning (PGRL) attempts to find a policy that maximizes the average (or time-discounted) reward, based on the gradient ascent in the policy parameter space [3,4,5]. As long as the policy is represented by a parametric statistical model that satisfies some mild conditions, PGRL can be instantly implemented in the Markov decision process (MDP). Moreover, since it is possible to treat the parameter controlling the randomness of the policy, PGRLs, rather than value-based RLs, can obtain the appropriate stochastic policy and be applied to the partially observable MDP (POMDP). Meanwhile, depending on the tasks, PGRL methods often take a huge number of learning steps. In this paper, we propose a new PGRL method that can improve the slow learning speed by focusing on the metric of the parameter space of the learning model.

It is easy to imagine that large-scale tasks suffer from a slow learning speed because the dimensionality of the policy parameters increases in conjunction with the task complexity. Besides the problem of dimensionality, the geometric structure of the parameter space also gives rise to slow learning. Ordinary PGRL methods omit the sensitivity of each element of the policy parameter and the correlation between the elements, in terms of the probability distributions of the MDP. However, most probability distributions expressed by the MDP have some manifold structures instead of Euclidean structures. Therefore, the updating direction of the policy parameter by the ordinary gradient method is different from the steepest direction on the manifold; thus, the optimization process occasionally falls into a stagnant state, commonly called a *plateau*. This is mainly due to the regions in which the geometric structure for the objective function with respect to the parameter coordinate system becomes fairly flat and its derivative becomes almost zero [6]. It was reported that a plateau was observed in a very simple MDP with only two states [2]. In order to solve such problem, Amari [1] proposed a “natural gradient” for the steepest gradient method in Riemannian space. Because the direction of the natural gradient is defined on a Riemannian metric, it is an important issue how to design the Riemannian metric. Nevertheless, the metric proposed by Kakade [2] has so far been the only metric in the application of the natural gradient for RL [7,8,9], commonly called *natural policy gradient* (NPG) reinforcement learning.

In this paper, we propose the use of the Fisher information matrix of the state-action joint distribution as the Riemannian metric for RL and derive a new robust NPG learning, “*natural state-action gradient*” (NSG) learning. It is shown that this metric considers the changes in the stationary state-action joint distribution, specifying the average reward as the objective function. In contrast, Kakade’s metric takes into account only changes in the action distribution and omits changes in the state distribution, which also depends on the policy in general. A comparison with the Hessian matrix is also given in order to confirm the adequacy of the proposed metric. We also prove that the gradient direction as computed by NSG is equal to the adjustable parameter of the linear regression model with the basis function defined on the policy when it minimizes the mean square error for the rewards. Finally, we demonstrate that the proposed NSG learning improves the performance of conventional (N)PG-based learnings by means of numerical experiments with varying scales of MDP tasks.

2 Conventional Natural Policy Gradient Method

We briefly review PGRL in section 2.1 and the natural gradient [1] and the NPG in section 2.2. In section 2.3, we introduce the controversy of NPGs.

2.1 Policy Gradient Reinforcement Learning

PGRL is modeled on a discrete-time Markov decision process (MDP) [10,11]. It is defined by the quintuplet $(\mathcal{S}, \mathcal{A}, p, r, \pi_\theta)$, where $\mathcal{S} \ni s$ and $\mathcal{A} \ni a$ are finite sets of

states and actions, respectively. Further, $p : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ is a state transition probability function of a state s_t , an action a_t , and the following state s_{t+1} at a time step t , i.e., $p(s_{t+1}|s_t, a_t) \equiv \Pr(s_{t+1}|s_t, a_t)$ ¹. $r : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathcal{R}$ is a reward function of s_t , a_t , and s_{t+1} , and is bounded, which defines an immediate reward $r_{t+1} = r(s_t, a_t, s_{t+1})$ observed by a learning agent. $\pi : \mathcal{A} \times \mathcal{S} \times \mathcal{R}^d \rightarrow [0, 1]$ is an action probability function of a_t , s_t , and a policy parameter $\theta \in \mathcal{R}^d$, and is always differentiable with respect to θ known as a policy, i.e., $\pi(a_t|s_t; \theta) \equiv \Pr(a_t|s_t, \theta)$. It defines the decision-making rule of the learning agent and is adjustable by tuning θ . We make an assumption that the Markov chain $M(\theta) = \{\mathcal{S}, \mathcal{A}, p, \pi_\theta\}$ is ergodic for all θ . Then, there exists a unique stationary state distribution $d_\theta(s) \equiv \Pr(s|M(\theta))$, which is equal to the limiting distribution and independent of the initial state, $d_\theta(s') = \lim_{t \rightarrow \infty} \Pr(S_t = s' | S_0 = s, M(\theta))$, $\forall s \in \mathcal{S}$. This distribution satisfies the balance equation:

$$d_\theta(s') = \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} p(s'|s, a) \pi(a|s; \theta) d_\theta(s). \quad (1)$$

The following equation instantly holds [10]:

$$d_\theta(s') = \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=1}^T \Pr(S_t = s' | S_0 = s, M(\theta)), \quad \forall s \in \mathcal{S}. \quad (2)$$

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

$$R(\theta) \equiv \lim_{T \rightarrow \infty} \frac{1}{T} \mathbb{E} \left\{ \sum_{t=1}^T r_t \middle| s_0, M(\theta) \right\}, \quad (3)$$

where $\mathbb{E}\{\cdot\}$ denotes expectation. It is noted that, under the assumption of ergodicity (eq.2), the average reward is independent of the initial state s_0 and can be shown to equal [10]:

$$\begin{aligned} R(\theta) &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} d_\theta(s) \pi(a|s; \theta) p(s'|s, a) r(s, a, s') \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_\theta(s) \pi(a|s; \theta) \bar{r}(s, a) \\ &\equiv \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \Pr(s, a | M(\theta)) \bar{r}(s, a). \end{aligned} \quad (4)$$

where $\bar{r}(s, a) \equiv \sum_{s' \in \mathcal{S}} p(s'|s, a) r(s, a, s')$. The statistical model $\Pr(s, a | M(\theta))$ is called the stationary state-action (joint) distribution. Since $\bar{r}(s, a)$ is usually independent of the policy parameter, the derivative of the average reward with

¹ Although it should be $\Pr(S_{t+1} = s_{t+1} | S_t = s_t, A_t = a_t)$ for the random variables S_{t+1} , S_t , and A_t to be precise, we notate $\Pr(s_{t+1}|s_t, a_t)$ for simplicity. The same rule is applied to the other distributions.

respect to the policy parameter, $\nabla_{\theta} R(\theta) \equiv [\partial R(\theta)/\partial \theta_1, \dots, \partial R(\theta)/\partial \theta_d]^\top$ is given by

$$\begin{aligned} \nabla_{\theta} R(\theta) &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \nabla_{\theta} \{d_{\theta}(s)\pi(a|s;\theta)\} \bar{r}(s, a) \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_{\theta}(s)\pi(a|s;\theta) \bar{r}(s, a) \{ \nabla_{\theta} \ln \pi(a|s;\theta) + \nabla_{\theta} \ln d_{\theta}(s) \}, \end{aligned} \quad (5)$$

where \top denotes transpose and $\nabla_{\theta} a_{\theta} b_{\theta} \equiv (\nabla_{\theta} a_{\theta}) b_{\theta}$. Therefore, the average reward $R(\theta)$ increases by updating the policy parameter as follows:

$$\theta := \theta + \alpha \nabla_{\theta} R(\theta),$$

where $:=$ denotes the right-to-left substitution and α is a sufficiently small learning rate. The above framework is called the PGRL method [5].

2.2 Natural Gradient [1]

Natural gradient learning is a gradient method on a Riemannian space. The parameter space being a Riemannian space implies that the parameter $\theta \in \mathcal{R}^d$ is on the Riemannian manifold defined by the Riemannian metric matrix $\mathbf{G}(\theta) \in \mathcal{R}^{d \times d}$ (positive definite matrix) and the squared length of a small incremental vector $\Delta\theta$ connecting θ to $\theta + \Delta\theta$ is given by

$$\|\Delta\theta\|_{\mathbf{G}}^2 = \Delta\theta^\top \mathbf{G}(\theta) \Delta\theta.$$

Under the constraint $\|\Delta\theta\|_{\mathbf{G}}^2 = \varepsilon^2$ for a sufficiently small constant ε , the steepest ascent direction of a function $R(\theta)$ is given by

$$\tilde{\nabla}_{\mathbf{G},\theta} R(\theta) = \mathbf{G}(\theta)^{-1} \nabla_{\theta} R(\theta). \quad (6)$$

It is called the natural gradient of $R(\theta)$ in the Riemannian space $\mathbf{G}(\theta)$. In RL, the parameter θ is the policy parameter, the function $R(\theta)$ is the average reward, and the gradient $\tilde{\nabla}_{\mathbf{G},\theta} R(\theta)$ is called the natural policy gradient (NPG) [2]. Accordingly, in order to (locally) maximize $R(\theta)$, θ is incrementally updated by

$$\theta := \theta + \alpha \tilde{\nabla}_{\mathbf{G},\theta} R(\theta). \quad (7)$$

When we consider a statistical model of a variable x parameterized by θ , $\Pr(x|\theta)$, the Fisher information matrix (FIM) $\mathbf{F}_x(\theta)$ is often used as the Riemannian metric matrix: [12]

$$\begin{aligned} \mathbf{F}_x(\theta) &\equiv \sum_{x \in \mathcal{X}} \Pr(x|\theta) \nabla_{\theta} \ln \Pr(x|\theta) \nabla_{\theta} \ln \Pr(x|\theta)^\top \\ &= - \sum_{x \in \mathcal{X}} \Pr(x|\theta) \nabla_{\theta}^2 \ln \Pr(x|\theta), \end{aligned} \quad (8)$$

where \mathcal{X} is a set of possible values taken by x . $\nabla_{\theta}^2 a_{\theta}$ denotes $\nabla_{\theta}(\nabla_{\theta} a_{\theta})$. The reason for using $\mathbf{F}(\theta)$ as $\mathbf{G}(\theta)$ comes from the fact that $\mathbf{F}(\theta)$ is a unique metric matrix of the second-order Taylor expansion of Kullback-Leibler (KL) divergence², which is known as a (pseudo) distance between two probability distributions. That is, the KL divergence of $\Pr(x|\theta+\Delta\theta)$ from $\Pr(x|\theta)$ is represented by

$$D_{\text{KL}}\{\Pr(x|\theta)|\Pr(x|\theta+\Delta\theta)\}=\frac{1}{2}\Delta\theta^{\top}\mathbf{F}_x(\theta)\Delta\theta+O(\|\Delta\theta\|^3),$$

where $\|\mathbf{a}\|$ denotes the Euclidean norm of a vector \mathbf{a} .

2.3 Controversy of Natural Policy Gradients

PGRL is regarded as an optimizing process of the policy parameter θ on some statistical models relevant to both a stochastic policy $\pi(a|s;\theta)$ and a state transition probability $p(s'|s, a)$. If a Riemannian metric matrix $\mathbf{G}(\theta)$ can be designed on the basis of the FIM of an apposite statistical model, $\mathbf{F}^*(\theta)$, an efficient NPG $\tilde{\nabla}_{\mathbf{F}^*,\theta}R(\theta)$ is instantly derived by eq.6.

As Kakade [2] pointed out, the choice of the Riemannian metric matrix $\mathbf{G}(\theta)$ for PGRL is not unique and the question what metric is apposite to $\mathbf{G}(\theta)$ is still open. Nevertheless, all previous studies on NPG [13,14,8,7,9] did not seriously address the above problem and (naively) used the Riemannian metric matrix proposed by Kakade [2]. In the next section, we will discuss the statistical models and metric spaces for PGRL and propose a new Riemannian metric matrix.

3 Riemannian Metric Matrices for PGRL

In section 3.1, we propose a new Riemannian metric matrix for RL and derive its NPG named the NSG. In sections 3.2 and 3.3, we discuss the validity of this Riemannian metric by comparing it with the Riemannian metric proposed by Kakade [2] and the Hessian matrix of the average reward.

3.1 A Proposed Riemannian Metric Matrix and NPG Based on State-Action Probability

Since the only adjustable function in PGRL is the policy function $\pi(a|s;\theta)$, previous studies on NPG focused on the policy function $\pi(a|s;\theta)$, i.e., the statistical models $\Pr(a|s, \mathbf{M}(\theta))$. However, the perturbations in the policy parameter θ also give rise to the change in the probability of the state $\Pr(s|\mathbf{M}(\theta))$. Because the average reward $R(\theta)$ as the objective function of PGRL is specified by the joint probability distribution of the state and the action $(s, a) \in \mathcal{S} \times \mathcal{A}$ (eq.4), it is natural and adequate to focus on the statistical model $\Pr(s, a|\mathbf{M}(\theta))$. For this case, the FIM of $\Pr(s, a|\mathbf{M}(\theta))$ can be used as the Riemannian metric $\mathbf{G}(\theta)$. Then, its NPG consists with the direction maximizing the average reward under

² It is same in the case of all f-divergences in general, except for scale [12].

the constraint that a measure of changes in the KL divergence of the stationary state-action distribution with respect to θ is fixed by a sufficient small constant ε : $D_{\text{KL}}\{\text{Pr}(s, a|\text{M}(\theta))|\text{Pr}(s, a|\text{M}(\theta + \Delta\theta))\} = \varepsilon^2$. The FIM of this statistical model, $\mathbf{F}_{s,a}(\theta)$, is calculated with $\text{Pr}(s, a|\text{M}(\theta)) = d_\theta(s)\pi(a|s;\theta)$ and eq.8 to be

$$\begin{aligned} \mathbf{F}_{s,a}(\theta) &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \text{Pr}(s, a|\text{M}(\theta)) \nabla_\theta \ln \text{Pr}(s, a|\text{M}(\theta)) \nabla_\theta \ln \text{Pr}(s, a|\text{M}(\theta))^\top \\ &= - \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_\theta(s) \pi(a|s;\theta) \nabla_\theta^2 \ln(d_\theta(s) \pi(a|s;\theta)) \\ &= \mathbf{F}_s(\theta) + \sum_{s \in \mathcal{S}} d_\theta(s) \mathbf{F}_a(s, \theta), \end{aligned} \quad (9)$$

where

$$\mathbf{F}_s(\theta) = \sum_{s \in \mathcal{S}} d_\theta(s) \nabla_\theta \ln d_\theta(s) \nabla_\theta \ln d_\theta(s)^\top \quad (10)$$

is the FIM defined from the statistical model comprising the state distribution, $\text{Pr}(s|\text{M}(\theta)) = d_\theta(s)$, and

$$\mathbf{F}_a(s, \theta) = \sum_{a \in \mathcal{A}} \pi(a|s;\theta) \nabla_\theta \ln \pi(a|s;\theta) \nabla_\theta \ln \pi(a|s;\theta)^\top \quad (11)$$

is the FIM of the policy comprising the action distribution given the state s , $\text{Pr}(a|s, \text{M}(\theta)) = \pi(a|s;\theta)$. Hence, the new NPG on the FIM of the stationary state-action distribution is

$$\tilde{\nabla}_{\mathbf{F}_{s,a}, \theta} R(\theta) = \mathbf{F}_{s,a}(\theta)^{-1} \nabla_\theta R(\theta).$$

We term it the ‘‘natural state-action gradient’’ (NSG).

3.2 Comparison with Kakade’s Riemannian Metric Matrix

The only Riemannian metric matrix for RL that has been proposed so far is the following matrix, which was proposed by Kakade [2] and was the weighted sum of the FIMs of the policy by the stationary state distribution $d_\theta(s)$,

$$\overline{\mathbf{F}}_a(\theta) \equiv \sum_{s \in \mathcal{S}} d_\theta(s) \mathbf{F}_a(s, \theta). \quad (12)$$

This is equal to the second term in eq.9. If it is assumed that the stationary state distribution is not changed by a variation in the policy, i.e., if $\nabla_\theta d_\theta(s) = \mathbf{0}$ holds, then $\mathbf{F}_s(\theta) = \mathbf{0}$ holds according to eq.10. While this assumption is not true in general, Kakade’s metric $\overline{\mathbf{F}}_a(\theta)$ is equivalent to $\mathbf{F}_{s,a}(\theta)$ if it holds. These facts indicate that $\overline{\mathbf{F}}_a(\theta)$ is the Riemannian metric matrix ignoring the change in the stationary state distribution $d_\theta(s)$ caused by the perturbation in the policy parameter θ in terms of the statistical model of the stationary state-action distribution $\text{Pr}(s, a|\text{M}(\theta))$.

Meanwhile, Bagnell et al. [13] and Peters et al. [14] independently, showed the relationship between the Kakade’s metric and the system trajectories $\xi_T = (s_0, a_0, s_1, \dots, a_{T-1}, s_T) \in \Xi_T$. When the FIM of the statistical model for the system trajectory ξ_T ,

$$\Pr(\xi_T | M(\theta)) = \Pr(s_0) \prod_{t=0}^{T-1} \pi(a_t | s_t; \theta) p(s_{t+1} | s_t, a_t),$$

is normalized by the time steps T with the limit $T \rightarrow \infty$, it is equivalent to the Kakade’s Riemannian metric,

$$\begin{aligned} \lim_{T \rightarrow \infty} \frac{1}{T} \mathbf{F}_{\xi_T}(\theta) &= - \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{\xi_T \in \Xi_T} \Pr(\xi_T | M(\theta)) \nabla_{\theta}^2 \left\{ \sum_{t=0}^{T-1} \ln \pi(a_t | s_t; \theta) \right\} \\ &= - \sum_{s \in \mathcal{S}} d_{\theta}(s) \sum_{a \in \mathcal{A}} \pi(a | s; \theta) \nabla_{\theta}^2 \ln \pi(a | s; \theta) \\ &= \overline{\mathbf{F}}_a(\theta) \end{aligned}$$

Since the PGRL objective, i.e., the maximization of the average reward, is reduced to the optimization of the system trajectory by eq.3 [13,14] suggested that the Kakade’s metric $\overline{\mathbf{F}}_a(\theta)$ could be a good metric. However, being equal to $\overline{\mathbf{F}}_a(\theta)$, the normalized FIM for the infinite-horizon system trajectory obviously differs with $\mathbf{F}_{s,a}(\theta)$ and is the metric that ignores the information $\mathbf{F}_s(\theta)$ about the stationary state distribution $\Pr(s | M(\theta))$. This is due to the fact that the statistical model of the system trajectory considers not only the state-action joint distribution but also the progress for the (infinite) time steps, as follows.

Here, s_{+t} and a_{+t} are the state and the action, respectively, progressed in t time steps after converging to the stationary distribution. Since the distribution of the system trajectory for T time steps from the stationary distribution, $\xi_{+T} \equiv (s, a_{+0}, s_{+1}, \dots, a_{+T-1}, s_{+T}) \in \Xi_T$, is

$$\Pr(\xi_{+T} | M(\theta)) = d_{\theta}(s) \prod_{t=0}^{T-1} \pi(a_{+t} | s_{+t}; \theta) p(s_{+t+1} | s_{+t}, a_{+t}),$$

its FIM is given by

$$\mathbf{F}_{\xi_{+T}}(\theta) = \mathbf{F}_s(\theta) + T \overline{\mathbf{F}}_a(\theta). \tag{13}$$

The derivation of which is shown in appendix A. Because of $\lim_{T \rightarrow \infty} \mathbf{F}_{\xi_{+T}}/T = \overline{\mathbf{F}}_a(\theta)$, the Kakade’s metric $\overline{\mathbf{F}}_a(\theta)$ is regarded as the limit $T \rightarrow \infty$ of the system trajectory distribution for T time steps from the stationary state distribution. Consequently, $\overline{\mathbf{F}}_a(\theta)$ omits the FIM of the state distribution, $\mathbf{F}_s(\theta)$. On the other hand, the FIM of the system trajectory distribution for one time step is obviously equivalent to the FIM of the state-action joint distribution, i.e., $\mathbf{F}_{\xi_{+1}}(\theta) = \mathbf{F}_{s,a}(\theta)$.

Now, we discuss which FIM is adequate for the average reward maximization. As discussed in section 3.1, the average reward in eq.4 is the expectation of $\bar{r}(s, a)$

over the distribution of the state-action (i.e. the +1-time-step system trajectory) and does not depend on the system trajectories after +2 time steps. It indicates that the Kakade's metric $\overline{\mathbf{F}}_a(\boldsymbol{\theta})$ supposed a redundant statistical model and the proposed metric for state-action distribution, $\mathbf{F}_{s,a}(\boldsymbol{\theta})$, would be more natural and adequate for PGRL. We give comparisons among various metrics such as $\mathbf{F}_{s,a}(\boldsymbol{\theta})$, $\overline{\mathbf{F}}_a(\boldsymbol{\theta})$, and a unit matrix \mathbf{I} through the numerical experiments in section 5.

Similarly, when the reward function is extended a function of T time steps, $r(s_t, a_t, \dots, a_{t+T-1}, s_{t+T})$, instead of one time step, $r(s_t, a_t, s_{t+1})$, the FIM of the T -time-step system trajectory distribution, $\mathbf{F}_{\xi_{+T}}(\boldsymbol{\theta})$, would be a natural metric because the average reward becomes $R(\boldsymbol{\theta}) = \sum_{\xi_{+T} \in \Xi_T} \Pr(\xi_{+T} | \mathbf{M}(\boldsymbol{\theta})) r(\xi_{+T})$.

3.3 Analogy with Hessian Matrix

We discuss the analogies between the Fisher information matrices $\mathbf{F}_{s,a}(\boldsymbol{\theta})$ and $\overline{\mathbf{F}}_a(\boldsymbol{\theta})$ and the Hessian matrix $\mathbf{H}(\boldsymbol{\theta})$, which is the second derivative of the average reward with respect to the policy parameter $\boldsymbol{\theta}$,

$$\begin{aligned}
 \mathbf{H}(\boldsymbol{\theta}) &\equiv \nabla_{\boldsymbol{\theta}}^2 R(\boldsymbol{\theta}) \\
 &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \bar{r}(s, a) d_{\boldsymbol{\theta}}(s) \pi(a|s; \boldsymbol{\theta}) \\
 &\quad \left\{ \nabla_{\boldsymbol{\theta}}^2 \ln(d_{\boldsymbol{\theta}}(s) \pi(a|s; \boldsymbol{\theta})) + \nabla_{\boldsymbol{\theta}} \ln(d_{\boldsymbol{\theta}}(s) \pi(a|s; \boldsymbol{\theta})) \nabla_{\boldsymbol{\theta}} \ln(d_{\boldsymbol{\theta}}(s) \pi(a|s; \boldsymbol{\theta}))^{\top} \right\} \\
 & \tag{14} \\
 &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \bar{r}(s, a) d_{\boldsymbol{\theta}}(s) \pi(a|s; \boldsymbol{\theta}) \\
 &\quad \left\{ \nabla_{\boldsymbol{\theta}}^2 \ln \pi(a|s; \boldsymbol{\theta}) + \nabla_{\boldsymbol{\theta}} \ln \pi(a|s; \boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}} \ln \pi(a|s; \boldsymbol{\theta})^{\top} + \nabla_{\boldsymbol{\theta}}^2 \ln d_{\boldsymbol{\theta}}(s) \right. \\
 &\quad \left. + \nabla_{\boldsymbol{\theta}} \ln d_{\boldsymbol{\theta}}(s) \nabla_{\boldsymbol{\theta}} \ln d_{\boldsymbol{\theta}}(s)^{\top} + \nabla_{\boldsymbol{\theta}} \ln d_{\boldsymbol{\theta}}(s) \nabla_{\boldsymbol{\theta}} \ln \pi(a|s; \boldsymbol{\theta})^{\top} \right. \\
 &\quad \left. + \nabla_{\boldsymbol{\theta}} \ln \pi(a|s; \boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}} \ln d_{\boldsymbol{\theta}}(s)^{\top} \right\}. \\
 & \tag{15}
 \end{aligned}$$

Comparing eq.12 of the Kakade's metric matrix $\overline{\mathbf{F}}_a(\boldsymbol{\theta})$ with eq.15 of the Hessian matrix $\mathbf{H}(\boldsymbol{\theta})$, the Kakade's metric does not have any information about the last two terms in braces $\{\cdot\}$ of eq.15, as Kakade [2] pointed out³. This is because $\overline{\mathbf{F}}_a(\boldsymbol{\theta})$ is derived under $\nabla_{\boldsymbol{\theta}} d_{\boldsymbol{\theta}}(s) = \mathbf{0}$. By eq.9 and eq.14, meanwhile, the proposed metric $\mathbf{F}_{s,a}(\boldsymbol{\theta})$ obviously has some information about all the terms of $\mathbf{H}(\boldsymbol{\theta})$. This comparison with the Hessian matrix suggests that $\mathbf{F}_{s,a}(\boldsymbol{\theta})$ should be an appropriate metric for PGRL. Additionally, $\mathbf{F}_{s,a}(\boldsymbol{\theta})$ becomes equivalent to the Hessian matrix in the cases using an atypical reward function that depends on $\boldsymbol{\theta}$ (see Appendix B).

It is noted that the average reward would not be a quadratic form with respect to the policy parameter $\boldsymbol{\theta}$ in general. Especially when $\boldsymbol{\theta}$ is far from the optimal parameter $\boldsymbol{\theta}^*$, the Hessian matrix $\mathbf{H}(\boldsymbol{\theta})$ occasionally gets into an indefinite matrix. Meanwhile, FIM $\mathbf{F}(\boldsymbol{\theta})$ is always positive (semi-)definite, assured

³ Strictly speaking, $\mathbf{H}(\boldsymbol{\theta})$ is slightly different from the Hessian matrix used in [2]. However, the essence of argument is the same as in [2].

by its definition in eq.8. Accordingly, the natural gradient method using FIM might be a more versatile covariant gradient ascent for PGRL than the Newton-Raphson method [15], in which the gradient direction is given by $\tilde{\nabla}_{-\mathbf{H},\theta} R(\theta)$. Comparison experiments are presented in section 5.

4 Computation of Natural State-Action Gradient

In this section, we view the estimation of the NSG. It will be shown that this estimation can be reduced to the regression problem of the immediate rewards.

Consider the following linear regression model

$$f_{\theta}(s, a; \omega) \equiv \phi_{\theta}(s, a)^{\top} \omega, \quad (16)$$

where ω is the adjustable parameter and $\phi_{\theta}(s, a)$ is the basis function of the state and action, also depending on the policy parameter θ ,

$$\begin{aligned} \phi_{\theta}(s, a) &\equiv \nabla_{\theta} \ln(d_{\theta}(s)\pi(a|s;\theta)) \\ &= \nabla_{\theta} \ln d_{\theta}(s) + \nabla_{\theta} \ln \pi(a|s;\theta). \end{aligned} \quad (17)$$

Then, the following theorem holds:

Theorem 1. *Let the Markov chain $M(\theta)$ have the fixed policy parameter θ , if the objective is to minimize the mean square error $\epsilon(\omega)$ of the linear regression model $f_{\theta}(s_t, a_t; \omega)$ in eq.16 for the rewards r_{t+1} ,*

$$\epsilon(\omega) = \lim_{T \rightarrow \infty} \frac{1}{2T} \sum_{t=0}^{T-1} \{r_{t+1} - f_{\theta}(s_t, a_t; \omega)\}^2, \quad (18)$$

then the optimal adjustable parameter ω^* is equal to NSG as the natural policy gradient on $\mathbf{F}_{s,a}(\theta)$:

$$\tilde{\nabla}_{\mathbf{F}_{s,a},\theta} R(\theta) = \omega^*.$$

Proof: By the ergodic property of $M(\theta)$, eq.18 is written as

$$\epsilon(\omega) = \frac{1}{2} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_{\theta}(s)\pi(a|s;\theta) (\bar{r}(s, a) - f_{\theta}(s, a; \omega))^2.$$

Since ω^* satisfies $\nabla_{\omega} \epsilon(\omega)|_{\omega=\omega^*} = \mathbf{0}$, we have

$$\sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_{\theta}(s)\pi(a|s;\theta) \phi_{\theta}(s, a) \phi_{\theta}(s, a)^{\top} \omega^* = \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d_{\theta}(s)\pi(a|s;\theta) \phi_{\theta}(s, a) \bar{r}(s, a).$$

By the definition of the basis function (eq.17), the following equations hold,

$$\begin{aligned} \sum_{s,a} d_{\theta}(s)\pi(a|s;\theta) \phi_{\theta}(s, a) \phi_{\theta}(s, a)^{\top} &= \mathbf{F}_{s,a}(\theta), \\ \sum_{s,a} d_{\theta}(s)\pi(a|s;\theta) \phi_{\theta}(s, a) \bar{r}(s, a) &= \nabla_{\theta} R(\theta). \end{aligned}$$

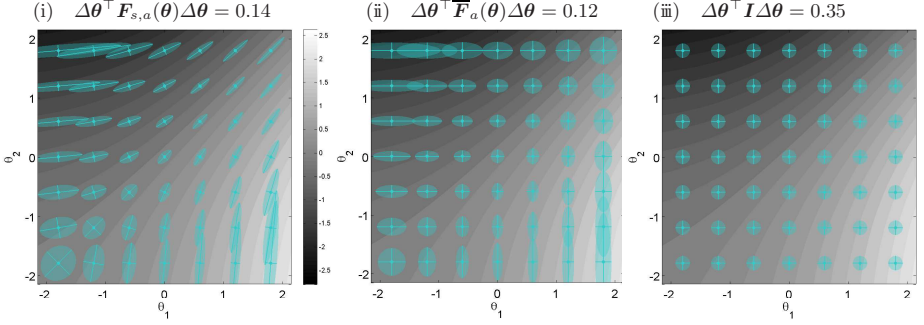


Fig. 1. Phase planes of a policy parameter in a two-state MDP: The gray level denotes $\ln d_{\theta}(1)/d_{\theta}(2)$. Each ellipsoid denotes the fixed distance spaces by each metric $G(\theta) :=$ (i) $F_{s,a}(\theta)$, (ii) $\bar{F}_a(\theta)$, or (iii) I .

Therefore, the following equation holds:

$$\omega^* = F_{s,a}(\theta)^{-1} \nabla_{\theta} R(\theta) = \widetilde{\nabla}_{F_{s,a},\theta} R(\theta). \quad \square$$

It is confirmed by theorem 1 that if the least-square regression to the immediate reward r_{t+1} by the linear function approximator $f_{\theta}(s_t, a_t; \omega)$ with the basis function $\phi_{\theta}(s, a) \equiv \nabla_{\theta} \ln(d_{\theta}(s)\pi(a|s; \theta))$ is performed, the adjustable parameter ω becomes the unbiased estimate of NSG $\widetilde{\nabla}_{F_{s,a},\theta} R(\theta)$. Therefore, since the NSG estimation problem is reduced to the regression problem of the reward function, NSG would be simply estimated by the least-square technique or by such a gradient descent technique as the method with the eligibility traces proposed by Morimura et al. [7], where the matrix inversion is not required.

It should be noted that, in order to implement this estimation, the computation of both the derivatives, $\nabla_{\theta} \ln \pi(a|s; \theta)$ and $\nabla_{\theta} \ln d_{\theta}(s)$, is required for the basis function $\phi_{\theta}(s, a)$. While $\nabla_{\theta} \ln \pi(a|s; \theta)$ can be instantly calculated, $\nabla_{\theta} \ln d_{\theta}(s)$ cannot be solved analytically because the state transition probabilities are generally unknown in RL. However, an efficient online estimation manner for $\nabla_{\theta} \ln d_{\theta}(s)$, which is similar to the method of estimating the value function, has been established by Morimura et al. [16]. However, we have not discussed the concrete implementations in this paper.

5 Numerical Experiments

5.1 Comparison of Metrics

We first looked into the differences among the Riemannian metric matrices $G(\theta)$ —the proposed metric $F_{s,a}(\theta)$, Kakade’s metric $\bar{F}_a(\theta)$, and unit matrix I —in a simple two-state MDP [2], where each state $s \in \{1, 2\}$ has self- and cross-transition actions $\mathcal{A} = \{l, m\}$ and each state transition is deterministic.

The policy with $\theta \in \mathcal{R}^2$ is represented by the sigmoidal function: $\pi(l|s; \theta) = 1/(1 + \exp(-\theta^\top \psi(s)))$, where $\psi(1) = [1, 0]^\top$ and $\psi(2) = [0, 1]^\top$. Figure 1 shows the phase planes of the policy parameter θ . The gray level denotes the log ratio of the stationary state distribution, and each ellipsoid corresponds to the set of $\Delta\theta$ satisfying a constant distance $\Delta\theta^\top \mathbf{G}(\theta) \Delta\theta = \varepsilon^2$, in which NPG looks for the steepest direction maximizing the average reward. It is confirmed that the ellipsoids by the proposed metric $\mathbf{F}_{s,a}(\theta)$ coped with the changes in the state distribution by the perturbation in θ because the alignment of the minor axis of the ellipsoid on $\mathbf{F}_{s,a}(\theta)$ complied with the direction significantly changing the $d_\theta(s)$. This indicates that the policy update with NSG does not drastically change $d_\theta(s)$. As we see theoretically, the other metrics could not grasp the changes even though $\overline{\mathbf{F}}_a(\theta)$ is the expectation of $\mathbf{F}_a(\theta)$ over $d_\theta(s)$.

5.2 Comparison of Learnings

We compared NSG with Kakade’s NPG, the ordinary PG, and the (modified) Newton PG learnings in terms of the optimizing performances for θ through randomly synthesized MDPs with a varying number of states, $|\mathcal{S}| \in \{3, 10, 20, 35, 50, 65, 80, 100\}$. Note that the only difference among these gradients is the definition of the matrix $\mathbf{G}(\theta)$ in eq.6. The Newton PG uses a modified Hessian matrix $\mathbf{H}^*(\theta)$ to assure the negative definiteness:

$$\mathbf{H}^*(\theta) = \mathbf{H}(\theta) - \max(0, \lambda_{\max} - \lambda'_{\max}) \mathbf{I},$$

where λ_{\max} and λ'_{\max} are the maximum and the largest-negative eigenvalues of $\mathbf{H}(\theta)$, respectively⁴.

It is noted that each gradient was computed analytically because we focussed on the direction of the gradients rather than the sampling issue in this paper.

Experimental Setup. We initialized the $|\mathcal{S}|$ -state MDP in each episode as follows. The set of the actions was always $|\mathcal{A}| = \{l, m\}$. The state transition probability function was set by using the Dirichlet distribution $\text{Dir}(\alpha \in \mathcal{R}^2)$ and the uniform distribution $\mathbf{U}(|\mathcal{S}|; b)$ generating an integer from 1 to $|\mathcal{S}|$ other than b : we first initialized it such that $p(s'|s, a) := 0, \forall (s', s, a)$ and then, with $\mathbf{q}(s, a) \sim \text{Dir}(\alpha = [.3, .3])$ and $x_{\setminus b} \sim \mathbf{U}(|\mathcal{S}|; b)$,

$$\begin{cases} p(s+1|s, l) := q_1(s, l) \\ p(x_{\setminus s+1}|s, l) := q_2(s, l) \end{cases} \quad \begin{cases} p(s|s, m) := q_1(s, m) \\ p(x_{\setminus s}|s, m) := q_2(s, m) \end{cases}$$

where $s' = 1$ and $s' = |\mathcal{S}| + 1$ are the identical states. The reward function $r(s, a, s')$ was temporarily set for each argument by Gaussian distribution $\mathbf{N}(\mu = 0, \sigma^2 = 1)$ and was normalized such that $\max_\theta R(\theta) = 1$ and $\min_\theta R(\theta) = 0$;

$$r(s, a, s') := \frac{r(s, a, s') - \min_\theta R(\theta)}{\max_\theta R(\theta) - \min_\theta R(\theta)}.$$

⁴ We examined various Hessian modifications [15]. The modification adopted here worked best in this task.

The policy parameterization was the same as that for previous experiment. Accordingly, in this MDP setting, there is no local optimum except for the global optimum. Each element of $\theta_0 \in \mathcal{R}^{|\mathcal{S}|}$ and $\psi(s) \in \mathcal{R}^{|\mathcal{S}|}$ for any state s was drawn from $N(0, .5)$ and $N(0, 1)$, respectively. We set the total episode time step at $T = 300$ and the initial learning rate α_0 in eq.7 for each (N)PG before each episode at the inverse of RMS,

$$\alpha_0 = \sqrt{|\mathcal{S}|} / \|\tilde{\nabla}_{G,\theta} R(\theta)|_{\theta=\theta_0}\|.$$

If the learning rate α is decent, $R(\theta)$ will always increase by the policy update of eq.7. Hence, when the policy update decreased $R(\theta)$, we tuned the learning rate “ $\alpha := \alpha/2$ ” and reattempted the update in the same time step. This tuning was kept until $\Delta R(\theta) \geq 0$. On the other hand, when $\alpha_0 > \alpha$ held true at the following time step, we also tuned “ $\alpha := 2\alpha$ ” to avoid standstills of the learning.

Results and Discussions. Figure 2 shows the learning curves for ten individual episodes in 100-state MDPs and reveals that NSG learning was able to succeed in optimizing the policy parameter uniformly and robustly though, compared with the other gradients, NSG was not infrequently slow in improving of performance at a moment. These are consistent with the results about the application of the natural gradient method to the learning of the multilayer perceptron [17].

Figure 3(A) shows the success rate of the learning by 300 episodes at each number of states. Since the maximum of the average reward was set to 1, we regarded the episodes satisfying $R(\theta_T) \geq 0.95$ as “successful” episodes. This suggests that, in the case of the MDPs with a small number of states, NSG and Kakade’s NPG methods could avoid falling into the severe plateau phenomena and robustly optimize the policy parameter θ , compared with the other methods. The reason why Kakade’s NPG could work as well as NSG would be that the Riemannian metric used in Kakade’s method has partial information about the statistical model $\Pr(s, a|M(\theta))$. Meanwhile, Kakade’s method frequently failed to improve the average reward in the cases of the MDPs with a large number of states. This could be due to the fact that Kakade’s metric omits the FIM about the state distribution, $F_s(\theta)$ unlike the proposed metric, as discussed theoretically in section 3.2. It is also confirmed that Kakade’s NPG was inferior to the modified Newton PG in the cases of many states. This could also be a result of whether the gradient has the information about the derivative of $d\theta(s)$ or not.

Finally, we analyzed how severe was the plateau in which these PG learnings were trapped. As this criterion, we utilized the smoothness of the learning curve (approximate curvature),

$$\Delta^2 R(\theta_t) = \Delta R(\theta_{t+1}) - \Delta R(\theta_t),$$

where $\Delta R(\theta_t) \equiv R(\theta_t) - R(\theta_{t-1})$. The criterion for the plateau measure of the episode was defined by

$$\text{PM} = \sum_{t=1}^{T-1} \|\Delta^2 R(\theta_t)\|.$$

Figure 3(B) represents the average of PM over all episodes for each PG and shows that NSG learning could learn very smoothly. This result indicates that the

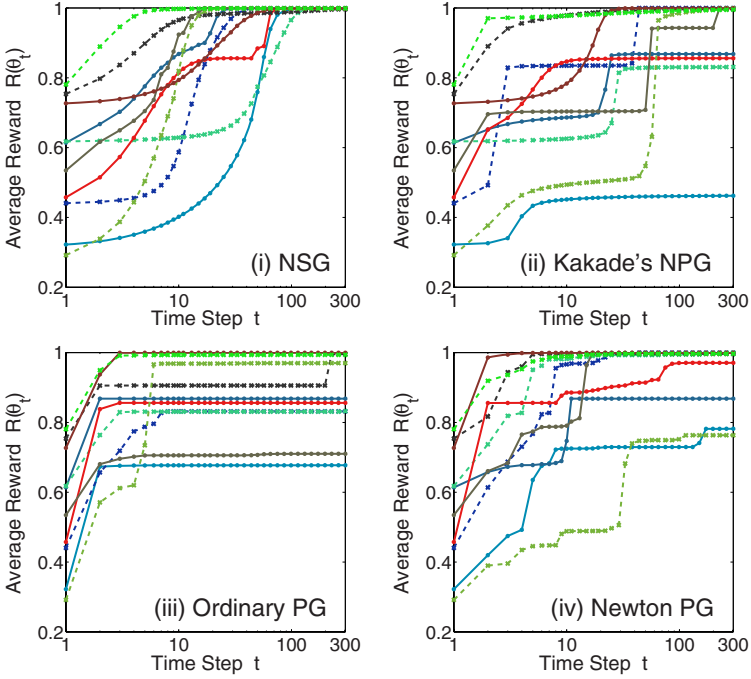


Fig. 2. Time courses of $R(\theta)$ for ten individual runs by (i) $\tilde{\nabla}_{F_{s,a},\theta} R(\theta)$, (ii) $\tilde{\nabla}_{F_a,\theta} R(\theta)$, (iii) $\tilde{\nabla}_{V,\theta} R(\theta)$, (iv) $\tilde{\nabla}_{-H^*,\theta} R(\theta)$

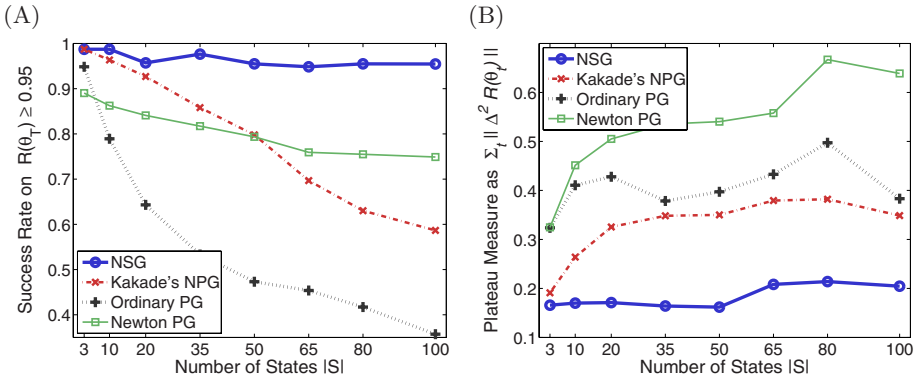


Fig. 3. (A) Learning success rates and (B) plateau measures for each number of states

learning by NSG could most successfully escape from a plateau; this is consistent with all other results.

Since NSG could avoid the plateau and robustly optimize θ without any serious effect of the setting of the MDP and the initial policy parameter, we conclude that NSG could be a more robust and natural NPG than the NPG by Kakade [2].

6 Summary and Future Work

This paper proposed a new Riemannian metric matrix for the natural gradient of the average reward, which was the Fisher information matrix of the stationary state-action distribution. We clarified that Kakade's NPG [2], which has been widely used in RL, does not consider the changes in the stationary state distribution caused by the perturbation of the policy, while our proposed NSG does. The difference was confirmed in numerical experiments where NSG learning could dramatically improve the performance and rarely fell into the plateau. Additionally, we proved that, when the immediate rewards were fitted by using the linear regression model with the basis function defined on the policy, its adjustable parameter represented the unbiased NSG estimate.

More algorithmic and experimental studies are necessary to further emphasize the effectiveness of NSG. The significant ones would be to establish an efficient Monte-Carlo estimation way of NSG along with estimating the derivative of the stationary state distribution [16], and then to clarify whether or not the proposed NSG method can still be useful even when the gradient is computed from samples. We will investigate them in future work.

References

1. Amari, S.: Natural gradient works efficiently in learning. *Neural Computation* 10(2), 251–276 (1998)
2. Kakade, S.: A natural policy gradient. In: *Advances in Neural Information Processing Systems*, vol. 14. MIT Press, Cambridge (2002)
3. Williams, R.J.: Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine Learning* 8, 229–256 (1992)
4. Kimura, H., Miyazaki, K., Kobayashi, S.: Reinforcement learning in pomdps with function approximation. In: *International Conference on Machine Learning*, pp. 152–160 (1997)
5. Baxter, J., Bartlett, P.: Infinite-horizon policy-gradient estimation. *Journal of Artificial Intelligence Research* 15, 319–350 (2001)
6. Fukumizu, K., Amari, S.: Local minima and plateaus in hierarchical structures of multilayer perceptrons. *Neural Networks* 13(3), 317–327 (2000)
7. Morimura, T., Uchibe, E., Doya, K.: Utilizing natural gradient in temporal difference reinforcement learning with eligibility traces. In: *International Symposium on Information Geometry and its Applications*, pp. 256–263 (2005)
8. Peters, J., Vijayakumar, S., Schaal, S.: Natural actor-critic. In: *European Conference on Machine Learning* (2005)

9. Richter, S., Aberdeen, D., Yu, J.: Natural actor-critic for road traffic optimisation. In: Advances in Neural Information Processing Systems. MIT Press, Cambridge (2007)
10. Bertsekas, D.P.: Dynamic Programming and Optimal Control, vol. 1, 2. Athena Scientific (1995)
11. Sutton, R.S., Barto, A.G.: Reinforcement Learning. MIT Press, Cambridge (1998)
12. Amari, S., Nagaoka, H.: Method of Information Geometry. Oxford University Press, Oxford (2000)
13. Bagnell, D., Schneider, J.: Covariant policy search. In: Proceedings of the International Joint Conference on Artificial Intelligence (July 2003)
14. Peters, J., Vijayakumar, S., Schaal, S.: Reinforcement learning for humanoid robotics. In: IEEE-RAS International Conference on Humanoid Robots (2003)
15. Nocedal, J., Wright, S.J.: Numerical Optimization. Springer, Heidelberg (2006)
16. Morimura, T., Uchibe, E., Yoshimoto, J., Doya, K.: Reinforcement learning with log stationary distribution gradient. Technical report, Nara Institute of Science and Technology (2007)
17. Amari, S., Park, H., Fukumizu, K.: Adaptive method of realizing natural gradient learning for multilayer perceptrons. Neural Computation 12(6), 1399–1409 (2000)

Appendix

A Derivation of eq.13

For simplicity, we denote $\pi_{+t} \equiv \pi(a_{+t}|s_{+t};\boldsymbol{\theta})$ and $p_{+t} \equiv p(s_{+t}|s_{+t-1}, a_{+t-1})$. Since ξ_{+T} is the system trajectory for T time steps from $d_{\boldsymbol{\theta}}(s)$, $\mathbf{F}_{\xi_{+T}}(\boldsymbol{\theta})$ is calculated to be

$$\begin{aligned} \mathbf{F}_{\xi_{+T}}(\boldsymbol{\theta}) &= - \sum_{\xi_{+T} \in \Xi_T} \Pr(\xi_{+T}) \nabla_{\boldsymbol{\theta}}^2 \left\{ \ln d_{\boldsymbol{\theta}}(s) + \sum_{t=0}^{T-1} \ln \pi(a_{+t}|s_{+t};\boldsymbol{\theta}) \right\} \\ &= - \sum_{s \in \mathcal{S}} d_{\boldsymbol{\theta}}(s) \left(\nabla_{\boldsymbol{\theta}}^2 \ln d_{\boldsymbol{\theta}}(s) + \sum_{a_{+0} \in \mathcal{A}} \pi_{+0} \left(\nabla_{\boldsymbol{\theta}}^2 \ln \pi_{+0} + \right. \right. \\ &\quad \left. \left. \sum_{s_{+1} \in \mathcal{S}} p_{+1} \sum_{a_{+1} \in \mathcal{A}} \pi_{+1} \left(\nabla_{\boldsymbol{\theta}}^2 \ln \pi_{+1} + \dots + \right. \right. \right. \\ &\quad \left. \left. \left. \sum_{s_{+T-1} \in \mathcal{S}} p_{+T-1} \sum_{a_{+T-1} \in \mathcal{A}} \pi_{+T-1} \nabla_{\boldsymbol{\theta}}^2 \ln \pi_{+T-1} \right) \dots \right) \right). \end{aligned}$$

By using the balance equation of the $d_{\boldsymbol{\theta}}(s)$ in eq.1,

$$\begin{aligned} \mathbf{F}_{\xi_{+T}}(\boldsymbol{\theta}) &= \mathbf{F}_s(\boldsymbol{\theta}) + \sum_{t=0}^{T-1} \left(\sum_{s_{+t} \in \mathcal{S}} d_{\boldsymbol{\theta}}(s_{+t}) \mathbf{F}_a(\boldsymbol{\theta}|s_{+t}) \right) \\ &= \mathbf{F}_s(\boldsymbol{\theta}) + T \overline{\mathbf{F}}_a(\boldsymbol{\theta}). \end{aligned} \quad \square$$

B Consistency of $\mathbf{F}_{s,a}(\boldsymbol{\theta})$ and $\mathbf{H}(\boldsymbol{\theta})$

If the immediate reward is dependent on $\boldsymbol{\theta}$

$$r(s, a; \boldsymbol{\theta}) = \frac{\Pr(s, a | \mathbf{M}(\boldsymbol{\theta}^*))}{\Pr(s, a | \mathbf{M}(\boldsymbol{\theta}))} \ln \Pr(s, a | \mathbf{M}(\boldsymbol{\theta})), \quad (19)$$

then the average reward becomes the negative cross entropy,

$$R(\boldsymbol{\theta}) = \sum_{s,a} \Pr(s, a | M(\boldsymbol{\theta}^*)) \ln \Pr(s, a | M(\boldsymbol{\theta})).$$

Hence, $\Pr(s, a | M(\boldsymbol{\theta}^*)) = \Pr(s, a | M(\boldsymbol{\theta}))$ holds, if the average reward is maximized. The Hessian matrix becomes $\mathbf{H}(\boldsymbol{\theta}) = \sum_{s,a} \Pr(s, a | M(\boldsymbol{\theta}^*)) \nabla_{\boldsymbol{\theta}}^2 \ln \Pr(s, a | M(\boldsymbol{\theta}))$. If the policy parameter is nearly optimal $\boldsymbol{\theta} \approx \boldsymbol{\theta}^*$, $\Pr(s, a | M(\boldsymbol{\theta})) \approx \Pr(s, a | M(\boldsymbol{\theta}^*))$ holds by the assumption of the smoothness of $\pi(a | s; \boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}$. Therefore, at this time, the Hessian matrix approximately equates the negative, proposed FIM:

$$\begin{aligned} \mathbf{H}(\boldsymbol{\theta}) &\approx \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \Pr(s, a | M(\boldsymbol{\theta})) \nabla_{\boldsymbol{\theta}}^2 \ln \Pr(s, a | M(\boldsymbol{\theta})) \\ &= -\mathbf{F}_{s,a}(\boldsymbol{\theta}). \end{aligned}$$

$\mathbf{H}(\boldsymbol{\theta}^*) = -\mathbf{F}_{s,a}(\boldsymbol{\theta}^*)$ obviously holds. Therefore, when the reward function is in eq.19 and the policy parameter is close to the optimal, NSG almost consists with the Newton direction and the NSG learning attains quadratic convergence.