
Supplementary Material to “A Generalized Natural Actor-Critic Algorithm”

Tetsuro Morimura[†], Eiji Uchibe[‡], Junichiro Yoshimoto[‡], Kenji Doya[‡]

[†]: IBM Research – Tokyo, Kanagawa, Japan

[‡]: Okinawa Institute of Science and Technology, Okinawa, Japan

tetsuro@jp.ibm.com, {uchibe, jun-y, doya}@oist.jp

Proof of Lemma 1

If the function $g_{\iota, \theta}(s, a; \omega)$ defined in Eq. (10) satisfies the three conditions of Eqs. (6), (7), and (8) for estimating function and its solution is equal to the gNG(ι), then Lemma 1 is proven. It is clear that it satisfies the conditions of Eqs. (7) and (8) (under Assumption 1). Thereby, if the following equation holds,

$$\begin{aligned} \mathbb{E}_{\theta} [g_{\iota, \theta}(s, a; \omega^*)] &= \mathbb{E}_{\theta} [g'_{\iota, \theta}(s, a; \omega^*) - \nabla_{\theta} \ln \{d_{\theta}(s) \pi(a|s; \theta)\} \rho(s, s_{+1})] \\ &= -\mathbb{E}_{\theta} [\nabla_{\theta} \ln \{d_{\theta}(s) \pi(a|s; \theta)\} \rho(s, s_{+1})] \\ &= \mathbf{0}, \end{aligned} \quad (\text{s-1})$$

where ω^* is the gNG(ι), $g_{\iota, \theta}(s, a; \omega)$ satisfies the condition of Eq. (8) and its unique solution is gNG(ι) and then Lemma 1 is proven. Thus, we prove Eq. (s-1) in the following.

Because of Eq. (13) and

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

we know that

$$\mathbb{E}_{\theta} [\nabla_{\theta} \ln \{d_{\theta}(s) \pi(a|s; \theta)\} \rho(s, s_{+1})] = \mathbb{E}_{\theta} [\{\nabla_{\theta} \ln \pi(a|s; \theta) + \nabla_{\theta} \ln d_{\theta}(s)\} \{b(s) - b(s_{+1})\}].$$

Since a time average is equivalent to a state-action space average with the ergodicity of $M(\theta)$ in Assumption 2, the above equation is transformed to

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

For the final transformations, we use the result

$$\mathbb{E}_{\theta}\{\nabla_{\theta}\ln\pi(a|s;\theta) + \nabla_{\theta}\ln d_{\theta}(s) \mid s_{+1}\} = \nabla_{\theta}\ln d_{\theta}(s_{+1})$$

in [1]. Therefore, Eq. (s-1) holds. \square

References

- [1] T. Morimura, E. Uchibe, J. Yoshimoto, J. Peters, and K. Doya. Derivatives of logarithmic stationary distributions for policy gradient reinforcement learning. *Neural Computation*. (in press).