

逐次的な自然方策勾配推定法の解析と 勾配推定分散の最小化による効率的な強化学習法の提案

Analysys of incremental natural policy gradient method
and minimizing the variance of gradient estimation

岩城 諒*¹ 横山 裕樹*² 浅田 稔*¹
Ryo Iwaki Hiroki Yokoyama Minoru Asada

*¹大阪大学大学院工学研究科 知能・機能創成工学専攻

Department of Adaptive Machine Systems, Graduate School of Engineering, Osaka University

*²玉川大学工学部 機械情報システム学科

Department of Intelligent Mechanical Systems, College of Engineering, Tamagawa University

We propose a natural actor critic method based on SARSA(λ). The proposed actor not only estimates unbiased natural policy gradient but also is independent of an approximated value function. Our critic explicitly minimizes the variance of actor's gradient estimation and does not need any auxiliary function, while conventional methods estimate an optimal baseline in order to reduce the variance of natural policy gradient. We apply our actor-critic algorithm to simple pendulum swing-up problem and show that our algorithm can learn more stably than the conventional method.

1. はじめに

強化学習とは、報酬という評価値を手掛かりに、システムが自律的に制御則を獲得するための手法である。学習主体であるエージェントは、将来得られる報酬を最大化する方策すなわち制御則を学習することを目的とする。エージェントは、将来得られる期待収益を価値という形式で推定し、学習に利用する。学習則や報酬の設計次第で、制御対象である環境についての事前知識がなくとも、エージェントは設計者の設定したタスクを達成するための方策を学習できる。

しかし、強化学習を実問題に適用する場合、方策や価値関数の学習に必要な時間が膨大になることが知られている。この原因として、従来の強化学習法では、入力分布に対する方策パラメータの感受性や相関を無視していることが考えられる。自然方策勾配法 [Kakade 02] は、方策空間の幾何学的構造を学習に利用することでこれらの問題を緩和し、プラトーと呼ばれる学習の停滞現象を回避することができる。しかし、オンラインでの自然方策勾配推定は、学習が不安定になりやすいことが指摘されている [木村 07]。

本研究は、自然方策勾配をオンラインで効率良く推定できる学習則を提案することを目的とする。まず、状態価値関数の推定値に依存せず自然方策勾配を不偏推定可能な学習則を提案する。次に、自然方策勾配推定の分散最小化を目的関数とする状態価値関数の学習則を提案する。提案手法を、単純化した倒立振子の振り上げ問題に適用し、頑健に学習できることを示す。

2. 自然方策勾配法

2.1 方策勾配型強化学習法

強化学習問題として、離散時間マルコフ決定過程 (MDP) を考える [Sutton 98]。環境の可能な状態の集合を S 、エージェントの選択可能な行動の集合を A 、実数の集合を \mathbb{R} とする。各離散時刻 t において、エージェントは環境の状態 $s_t \in S$ を観測し、 n 次元の方策パラメータ $\theta \in \mathbb{R}^n$ によって規定される確率の方策 $\pi(a|s; \theta) \equiv \Pr(a|s, \theta)$ に従い、行動 $a_t \in A$ を選択する。環境は、状態遷移確率 $T_{ss'}^a \equiv \Pr(s_{t+1} = s' | s_t =$

連絡先: 岩城諒, ryo.iwaki@ams.eng.osaka-u.ac.jp

$s, a_t = a)$ に従って新しい状態 s_{t+1} に遷移し、有界な報酬関数 $\mathcal{R}_{ss'}^a \equiv \mathbb{E}[r_t | s_t = s, a_t = a, s_{t+1} = s']$ によって定まる期待値に従い、エージェントに即時報酬 $r_t \in \mathbb{R}$ を与える。

エージェントは、以下で定義される期待収益 $\rho(\theta)$ を最大化する方策パラメータ θ を獲得することを学習目的とする。

$$\rho(\theta) = \mathbb{E}_{\theta} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 \right]$$

ただし、 $\gamma \in [0, 1]$ は割引率であり、将来の報酬が現在においてどれだけの価値があるかを決定する。ここで、行動価値関数 $Q^{\theta}(s, a)$ を以下のように定義する。

$$Q^{\theta}(s, a) = \mathbb{E}_{\theta} \left[\sum_{\tau=t}^{\infty} \gamma^{\tau-t} r_{\tau} \mid s_t = s, a_t = a \right]$$

行動価値関数 $Q^{\pi}(s, a)$ は、状態 s で行動 a をとり、その後方策 π に従った場合に得られる期待収益である。一般に、行動価値関数 $Q^{\pi}(s, a)$ は、状態価値関数 $V^{\pi}(s)$ とアドバンテージ関数 $A^{\pi}(s, a)$ を用いて、

$$Q^{\pi}(s, a) = V^{\pi}(s) + A^{\pi}(s, a)$$

と表される。状態価値関数 $V^{\pi}(s)$ は、状態 s の価値を表現し、アドバンテージ関数 $A^{\pi}(s, a)$ は、ある状態 s における行動 a の相対的な価値を表現する。

ここで、近似行動価値関数 $\hat{Q}(s, a; \theta, v, w)$ を、近似状態価値関数 $\hat{V}(s; v)$ と、方策のスコア関数 (対数勾配) $\nabla_{\theta} \ln \pi(a|s; \theta)$ を用いて、以下のように表す。

$$\hat{Q}(s, a; \theta, v, w) = \hat{V}(s; v) + \left(\frac{\partial}{\partial \theta} \ln \pi(a|s; \theta) \right)^{\top} w \quad (1)$$

ただし、

$$\hat{A}(s, a; \theta, w) = \left(\frac{\partial}{\partial \theta} \ln \pi(a|s; \theta) \right)^{\top} w \quad (2)$$

は近似アドバンテージ関数, v は状態価値関数の近似パラメータ, w はアドバンテージ関数の近似パラメータである. ここで, 近似行動価値関数 $\hat{Q}(s, a; \theta, v, w)$ が局所最適解へ収束したとき, 方策勾配定理 [Sutton 00] により, 以下が成り立つ.

$$\frac{\partial \rho(\theta)}{\partial \theta} = \int_S d_\theta(s) \int_A \frac{\partial \pi(a|s; \theta)}{\partial \theta} \hat{A}(s, a) da ds \quad (3)$$

ただし, $d_\theta(s)$ は方策 $\pi(\theta)$ のもとでの状態の分布を表す.

2.2 自然方策勾配法

式 (3) に式 (2) を代入することで, 次式が成立する.

$$w = G^{-1}(\theta) \frac{\partial \rho(\theta)}{\partial \theta} \quad (4)$$

ただし, $G(\theta)$ はフィッシャー情報行列であり, 以下で定義される.

$$G(\theta) = \int_S d_\theta(s) \int_A \pi(a|s; \theta) \left(\frac{\partial}{\partial \theta} \ln \pi(a|s; \theta) \right) \left(\frac{\partial}{\partial \theta} \ln \pi(a|s; \theta) \right)^T da ds$$

フィッシャー情報行列 $G(\theta)$ は, パラメータ θ を座標系とする空間の構造を規定する. 式 (4) の右辺は自然勾配 [Amari 98] と呼ばれ, パラメータを座標系とする空間上での最急方向を表す. すなわち, 式 (1) による関数近似が局所最適を達成すれば, アドバンテージ関数の近似パラメータ w は, 方策の自然勾配に一致する [Kakade 02].

2.3 アドバンテージ関数の近似方法

木村 [木村 07] は, SARSA(λ) に基づくアドバンテージ関数の近似方法を提案した. エージェントは, 方策に従って行動を実行する Actor と呼ばれる部分と, 価値関数を用いて方策を評価する Critic と呼ばれる部分から構成される. 学習の手続きを Alg.1 に示す.

Alg.1 はオンラインで自然方策勾配を推定できるが, [木村 07] では学習則の理論的な妥当性が示されていない. また, 学習後期に方策が決定的になるにつれて, 方策のスコア関数が極端に大きな値を取り, 学習が不安定になることが報告されている.

3. 自然方策勾配の不偏推定

本章では, [木村 07] で提案された手法をもとに, 自然方策勾配を不偏推定できるアドバンテージ関数の学習則を提案する. 提案する学習則では, TD 誤差を次のように計算する.

$$\delta_t^Q = r_t + \gamma \hat{Q}(s_{t+1}, a_{t+1}; \theta_t, v_t, w_t) - \hat{Q}(s_t, a_t; \theta_{t-1}, v_{t-1}, w_{t-1}) \quad (5)$$

既存の自然方策勾配推定法 [木村 07] と提案する更新則との違いは, TD 誤差の計算に用いる価値関数のパラメータを一時刻ずらしていることである. 提案手法を用いてアドバンテージ関数の近似パラメータを学習した場合, 固定方策のもとで漸近的に以下が成り立つ.

Algorithm 1 [木村 07]

```

1: Input:
2: • Parameterized policy
3:    $\pi(\cdot | \cdot; \theta)$ 
4: • Parameterized value function
5:    $\hat{V}(\cdot; v)$ 
6: Initialization:
7: • Policy, value function and advantage parameters
8:    $\theta = \theta_0, v = v_0, w = w_0$ 
9: • Stepsizes and discount rates
10:   $\alpha_\theta, \alpha_v, \alpha_w, \beta, \gamma$ 
11: • Draw initial state
12:   $s_0 \sim p(s_0)$ 
13: • Draw initial action
14:   $a_0 \sim \pi(a_0 | s_0; \theta_0)$ 
15: for  $t = 0, 1, 2, \dots$  do
16:   Execution:
17:   • Observe next state
18:      $s_{t+1} \sim \mathcal{T}_{ss'}^a$ 
19:   • Observe reward
20:      $r_t \sim \mathcal{R}_{ss'}^a$ 
21:   • Draw action
22:      $a_{t+1} \sim \pi(a_{t+1} | s_{t+1}; \theta_t)$ 
23:   Eligibility Traces:
24:      $\bar{e}_{t+1} = \gamma \bar{e}_t + \nabla_\theta \ln \pi(a_{t+1} | s_{t+1}; \theta_t)$ 
25:   Action Value:
26:      $\hat{Q}(s_t, a_t; \theta_{t-1}, v_t, w_t)$ 
27:      $= \hat{V}(s_t; v_t) + \left( \nabla_\theta \ln \pi(a_t | s_t; \theta_{t-1}) \right)^T w_t$ 
28:      $\hat{Q}(s_{t+1}, a_{t+1}; \theta_t, v_t, w_t)$ 
29:      $= \hat{V}(s_{t+1}; v_t) + \left( \nabla_\theta \ln \pi(a_{t+1} | s_{t+1}; \theta_t) \right)^T w_t$ 
30:   TD Error:
31:      $\delta_t^V = r_t + \gamma \hat{V}(s_{t+1}; v_t) - \hat{V}(s_t; v_t)$ 
32:      $\delta_t^Q = r_t + \gamma \hat{Q}(s_{t+1}, a_{t+1}; \theta_t, v_t, w_t)$ 
33:        $- \hat{Q}(s_t, a_t; \theta_{t-1}, v_t, w_t)$ 
34:   Advantage Update:
35:      $w_{t+1} = w_t + \alpha_w \delta_t^Q \bar{e}_t$ 
36:   Actor Update:
37:      $\theta_{t+1} = \theta_t + \alpha_\theta w_t$ 
38:   Critic Update:
39:      $\bar{c}_t = \gamma \bar{c}_{t-1} + \nabla_v \hat{V}(s_t; v_t)$ 
40:      $v_{t+1} = v_t + \alpha_v \delta_t^V \bar{c}_t$ 
41: end for

```

$$\begin{aligned} & \sum_{t=0}^{\infty} \Delta w_t \\ &= \sum_{t=0}^{\infty} \delta_t^Q \sum_{\tau=0}^t \gamma^{t-\tau} e_\tau \\ &= \sum_{t=0}^{\infty} e_t \sum_{\tau=t}^{\infty} \gamma^{t-\tau} \delta_\tau^Q \\ &= \sum_{t=0}^{\infty} e_t \left(\sum_{\tau=t}^{\infty} \gamma^{\tau-t} r_\tau - \hat{Q}(s_t, a_t; \theta_{t-1}, v_{t-1}, w_{t-1}) \right) \end{aligned}$$

$$\begin{aligned}
& \mathbb{E}_\theta [\Delta w_t] \\
&= \lim_{T \rightarrow \infty} \frac{1}{T} \sum_{\tau=0}^T \Delta w_\tau \\
&= \mathbb{E}_\theta \left[\frac{\partial}{\partial \theta} \ln \pi(a|s; \theta_{t-1}) \left(Q_t^\theta(s, a) - \hat{A}(s, a; \theta_{t-1}, w_{t-1}) \right) \right] \\
&= \mathbb{E}_\theta \left[\frac{\partial}{\partial w} \hat{A}(s, a; \theta_{t-1}, w_{t-1}) \left(Q_t^\theta(s, a) - \hat{A}(s, a; \theta_{t-1}, w_{t-1}) \right) \right] \\
&= \mathbb{E}_\theta \left[-\frac{1}{2} \frac{\partial}{\partial w} \left(Q_t^\theta(s, a) - \hat{A}(s, a; \theta_{t-1}, w_{t-1}) \right)^2 \right] \quad (6)
\end{aligned}$$

式 (6) は、提案手法を用いると、アドバンテージ関数の更新方向の期待値が、未知の真の行動価値関数 Q_t^θ との二乗誤差の最急降下方向に一致すること、かつ状態価値関数の推定値には依存しないことを意味する。すなわち、価値推定の近似精度によらず、自然方策勾配を不偏推定する^{*1}。

4. 自然方策勾配推定の分散最小化

式 (6) は、アドバンテージ関数近似パラメータの更新値 $\Delta \theta$ の期待値は、状態価値関数の推定値 $\hat{V}(\cdot; v)$ に依存しないことを意味する。 $\hat{V}(\cdot; v)$ の値が不正確でも、更新を繰り返すことで Δw の期待値は最適値に近づく。しかし、各時刻の更新は $\hat{V}(\cdot; v)$ に依存しているため、 Δw は $\hat{V}(\cdot; v)$ の値に応じた分散を持っていると考えられる。従来の自然方策勾配法では、価値関数とは独立したベースライン関数 $b(s)$ を調節することで、方策勾配の分散を小さくする [Morimura 08]。一方本研究では、式 (6) より、 Δw の期待値が状態価値関数の値に依存せず最適値に近づくことが保証されているため、ベースライン関数を導入せず、 Δw の分散を最小化するように状態価値関数を学習させる。Critic は以下に示す目的関数を最小化することを学習目標とする。

$$f = \frac{1}{2} \left(\mathbb{E}_\theta \left[\left\| \sum_{t=0}^{\infty} \Delta w_t \right\|^2 \right] - \left\| \sum_{t=0}^{\infty} \mathbb{E}_\theta [\Delta w_t] \right\|^2 \right) \quad (7)$$

式 (7) は、アドバンテージ関数近似パラメータの更新則 Δw の分散共分散行列のトレースを意味している。式 (7) を Critic のもつパラメータ v_t で偏微分することで勾配を求め、最急降下法により

$$v_{t+1} = v_t - \alpha_v \frac{\partial f}{\partial v_t}$$

として学習する。式 (6) より、式 (7) の右辺第二項は v_t に依存しないことから、

$$\begin{aligned}
\Delta v_t &= \frac{\partial f}{\partial v_t} = \mathbb{E}_\theta \left[C_t \sum_{\tau=0}^{\infty} \Delta w_\tau \right] \quad (8) \\
C_t &= \left(\frac{\partial}{\partial v_t} \hat{V}(s_{t+1}; v_t) \right) \left(\frac{\partial}{\partial \theta_t} \ln \pi(a_{t+1}|s_{t+1}; \theta_t) \right)^T
\end{aligned}$$

式 (8) の期待値計算は、時間平均で置き換える。また、式 (8) の計算には未来 ($\tau > t$) の情報が必要であるため、このまま

*1 ただし [Thomas 14] において、 $\Delta w \propto \frac{\partial}{\partial w} (Q_t^\theta - \hat{A})^2$ に従って学習した場合でも、 $\gamma \neq 1$ であればバイアスが加わり、割引収益ではなく平均報酬を最大化する自然勾配の不偏推定となることが示されている。

では計算できない。よって、以下のように変形する。

$$\sum_{t=0}^{\infty} C_t \sum_{\tau=0}^{\infty} \Delta w_\tau = \sum_{t=0}^{\infty} \left(C_t \left(\sum_{\tau=0}^{t-1} \Delta w_\tau \right) + \left(\sum_{\tau=0}^t C_\tau \right) \Delta w_t \right) \quad (9)$$

式 (9) の総和計算には、 t と τ の様々な組み合わせが含まれている。方策を逐次更新する場合、時間を大きく隔てた Δw_t と C_t は、異なる方策のもとで得られた値を保持しているため、これらの組み合わせは学習に悪影響を及ぼす可能性がある。よって、 $|t - \tau|$ の値が大きい組み合わせを割り引くための定数 $\beta \in (0, 1)$ を導入し、式 (9) を以下のように修正する。

$$\sum_{t=0}^{\infty} \left(C_t \left(\sum_{\tau=0}^{t-1} \beta^{t-\tau} \Delta w_\tau \right) + \left(\sum_{\tau=0}^t \beta^{t-\tau} C_\tau \right) \Delta w_t \right)$$

提案手法の学習手続きを Alg.2 に示す。

5. 計算機実験

提案手法の優位性を示すため、単純化した倒立振子の振り上げ問題に適用した。振子の振り上げに十分な大きさのトルクは許容されず、振動させて勢いをつけてから振り上げる必要があるため、長い行動系列を学習する必要がある。初期状態は振子が下向きに静止した状態とし、振子の先端がゴール高さへ振り上げられるまでを 1 エピソードした^{*2}。各時刻での報酬は 0 であり、リンクがゴール高さに達した時のみ報酬 $r = 1$ を与える。本研究では、Actor と Critic を Normalized Radial Basis Function Network (NRBF-Net) [Moody 89] によって構成した。Actor は、状態 $s = (q, \dot{q})^T$ を観測し、確率の方策 $\pi(a|s; \theta)$ に従って行動 $a = \tau$ を出力する。確率の方策 $\pi(a|s; \theta)$ は、正規分布を用いて以下のように定義する。

$$\pi(a|s; \theta) = \frac{1}{\theta_\sigma \sqrt{2\pi}} \exp \left(-\frac{(a - y(s, \theta_\kappa, \theta_b))^2}{2\theta_\sigma^2} \right) \quad (10)$$

ただし、式 (10) における $y(s, \theta_\kappa, \theta_b) = \theta_\kappa^T \phi(s) + \theta_b$ は NRBF-Net の出力であり、 $\phi(s)$ は RBF を並べたベクトル、 θ_κ は結合荷重、 θ_b はバイアスである。各 RBF は、状態空間に格子状に配置される。すなわち、エージェントが学習する方策パラメータは $\theta = (\theta_\kappa, \theta_b, \theta_\sigma)$ である。同様に、Critic も NRBF-Net を用いて $\hat{V}(s, v_\kappa, v_b) = v_\kappa^T \phi(s) + v_b$ と表現される。すなわち、エージェントが学習する Critic のパラメータは $v = (v_\kappa, v_b)$ である。

エージェントの学習則を以下の 3 条件として性能を比較した。(a) アドバンテージ関数の学習に文献 [木村 07] で提案された学習則を用いる条件 (Alg.1)。(b) アドバンテージ関数の学習に提案手法を用いる条件。(c) アドバンテージ関数と Critic の学習の両方に提案手法を用いる条件 (Alg.2)。乱数の初期値を変えながら 100 試行の学習を行った。各学習則を用いた場合について、乱数の初期値は同じ値を用いた。学習の 1 試行は 3000 エピソードとした。5000 ステップ経過してもリンク先端がゴール高さまで振り上げられない場合、そのエピソードを終了し、次のエピソードを開始した。(a-c) の学習則全てにおいて、 v と w の学習にトレースを利用しており、更新のスケールに大きな違いはないと考えられる。よって、学習のハイパーパラメータは、全ての条件で同じ値を用いた。

実験結果として、初期から方策の改善が見られないものを除いた学習試行の平均の学習曲線を図 1 に示す。図 1 の横軸

*2 振子先端がゴール高さまで振り上げられるとそのエピソードは終了するため、Mountain Car に近い実験設定である。

Algorithm 2 提案手法

```

1: Input:
2: • Parameterized policy
3:    $\pi(\cdot | \cdot; \theta)$ 
4: • Parameterized value function
5:    $\hat{V}(\cdot; v)$ 
6: Initialization:
7: • Policy, value function and advantage parameters
8:    $\theta = \theta_0, v = v_0, w = w_0$ 
9: • Stepsizes and discount rates
10:   $\alpha_\theta, \alpha_v, \alpha_w, \beta, \gamma$ 
11: • Draw initial state
12:   $s_0 \sim p(s_0)$ 
13: • Draw initial action
14:   $a_0 \sim \pi(a_0 | s_0; \theta_0)$ 
15: for  $t = 0, 1, 2, \dots$  do
16:   Execution:
17:   • Observe next state
18:      $s_{t+1} \sim \mathcal{T}_{ss'}^a$ 
19:   • Observe reward
20:      $r_t \sim \mathcal{R}_{ss'}^a$ 
21:   • Draw action
22:      $a_{t+1} \sim \pi(a_{t+1} | s_{t+1}; \theta_t)$ 
23:   Eligibility Traces:
24:      $\bar{e}_{t+1} = \gamma \bar{e}_t + \nabla_\theta \ln \pi(a_{t+1} | s_{t+1}; \theta_t)$ 
25:   Action Value:
26:      $\hat{Q}(s_{t+1}, a_{t+1}; \theta_t, v_t, w_t)$ 
27:      $= \hat{V}(s_{t+1}; v_t) + \left( \nabla_\theta \ln \pi(a_{t+1} | s_{t+1}; \theta_t) \right)^\top w_t$ 
28:   TD Error:
29:      $\delta_t^Q = r_t + \gamma \hat{Q}(s_{t+1}, a_{t+1}; \theta_t, v_t, w_t)$ 
30:      $- \hat{Q}(s_t, a_t; \theta_{t-1}, v_{t-1}, w_{t-1})$ 
31:   Advantage Update:
32:      $w_{t+1} = w_t + \alpha_w \delta_t^Q \bar{e}_t$ 
33:   Actor Update:
34:      $\theta_{t+1} = \theta_t + \alpha_\theta w_t$ 
35:   Critic Update:
36:      $C_t = \nabla_v \hat{V}(s_{t+1}; v_t) e_{t+1}^\top$ 
37:      $\bar{C}_t = \beta \bar{C}_{t-1} + C_t$ 
38:      $\bar{\Delta} w_t = \beta \bar{\Delta} w_{t-1} + \Delta w_t$ 
39:      $v_{t+1} = v_t + \alpha_v (\beta \bar{C}_t \bar{\Delta} w_{t-1} + \bar{C}_t \Delta w_t)$ 
40: end for

```

はエピソード数，縦軸は1エピソード中に倒立振子の振り上げに要したステップ数を表す．縦軸の値が低いほどより良い方策を獲得していると考えられる．赤い線が条件 (a)，水色の線が条件 (b)，藍色の線が条件 (c) の場合を表す．条件 (a) で学習した場合，多くの学習試行が，初期の方策から改善した後にそれぞれ異なるエピソード数で発散した．それらの独立した学習試行間の平均をとっているため，一見方策が徐々に悪化しているように見受けられる．学習率をより小さな値に設定した条件でも実験したが，局所解に陥ったためか，多くの学習試行で初期から方策の改善が見られなかった．一方で，提案手法を用いることで，学習後期での発散を防ぎ，効率よく頑健に学習できていることがわかる．

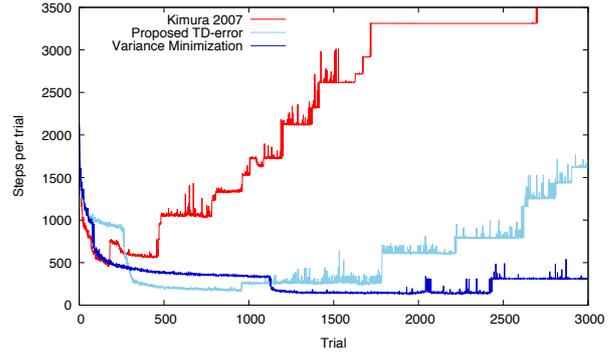


図 1: 平均の学習曲線

6. おわりに

本研究では，自然方策勾配を不偏推定する学習則と，自然方策勾配推定の分散を最小化する学習則を提案した．計算機実験により，提案手法が従来手法よりも学習後期で特に頑健であることを示した．しかし，今回提案した自然方策勾配の不偏推定・分散最小化とともに，実際は固定方策のもとでのみ成り立つ．よって今後の課題として，方策の違いによる影響を緩和可能な方策オフ学習に提案手法を拡張することが考えられる．

参考文献

- [Amari 98] S. Amari, Natural Gradient Works Efficiently in Learning, *Neural Computation*, Vol. 10, No. 2, pp. 251-276 (1998).
- [Kakade 02] S. Kakade, A Natural Policy Gradient, *Advances in Neural Information Processing Systems*, Vol. 14, pp. 227-242 (2002).
- [Moody 89] J. Moody and C. J. Darken, Fast learning in networks of locally-tuned processing units, *Neural Computation*, Vol. 1, No. 2, pp. 281-294 (1989).
- [Morimura 08] T. Morimura, E. Uchibe, and K. Doya. Natural actor-critic with baseline adjustment for variance reduction. *Artificial Life and Robotics*, Vol. 13, pp. 275-279 (2008).
- [Sutton 98] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An introduction*, MIT Press (1998).
- [Sutton 00] R. S. Sutton, D. McAllester, S. Singh and Y. Mansour, Policy Gradient Methods for Reinforcement Learning with Function Approximation, *Advances in Neural Information Processing Systems*, Vol. 12, pp. 1057-1063 (2000).
- [Thomas 14] P. Thomas, Bias in Natural Actor-Critic Algorithms, *Proceedings of The 31st International Conference on Machine Learning*, pp. 441-448 (2014).
- [木村 07] 木村 元, 適正度の履歴を用いた自然勾配 Actor-Critic 法, 第 19 回自律分散システムシンポジウム, pp. 677-682 (2007).