

強化学習 第6章

「関数近似を用いた強化学習」

価値関数：

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a|s) \{g(s, a) + \gamma \sum_{s' \in \mathcal{S}} p_T(s'|s, a) V^\pi(s')\}, \quad \forall s \in \mathcal{S}$$

は、状態行動空間が膨大である場合、テーブルの要素数が大きくなりすぎてしまい学習が困難となる。そこで、価値関数を関数近似器を用い近似することを考える。

1 価値関数の関数近似

V^π をパラメータ $\mathbf{w} \in \mathcal{R}^d$ で規定される関数近似器 $\hat{V} : \mathcal{S} \times \mathcal{R}^d \rightarrow \mathbb{R}$ を用いて近似する。すなわち、

$$V^\pi(s) \approx \hat{V}(s; \mathbf{w}^*)$$

なるようなパラメータ \mathbf{w}^* を学習することを考える。ここで、関数近似器 $\hat{V}(s; \mathbf{w})$ の集合を

$$\mathcal{V} := \{\hat{V}_{\mathbf{w}}; \mathbf{w} \in \mathbb{R}^d\}, \quad (\hat{V}_{\mathbf{w}}(s) := \hat{V}(s; \mathbf{w}))$$

と置くと、 \mathcal{V} は状態関数空間の部分空間である。($\mathcal{V} \subset \mathbb{R}^{\mathcal{S}}$ となる。 $\mathcal{V} \equiv \mathbb{R}^{\mathcal{S}}$ でないことが重要。)

1.1 線形近似

線形関数近似器として、

$$\hat{V}_{\mathbf{w}}(s) := \mathbf{w}^\top \boldsymbol{\phi}(s)$$

がある。ただし、 $\phi : \mathcal{S} \rightarrow \mathbb{R}^d$ は状態 s の特徴ベクトル。例えば

$$\mathbf{w} := \begin{pmatrix} w_1 \\ \vdots \\ w_{|\mathcal{S}|} \end{pmatrix}, \quad \phi(s) := \begin{pmatrix} \phi_1(s) \\ \vdots \\ \phi_{|\mathcal{S}|}(s) \end{pmatrix} \in \mathbb{R}^{|\mathcal{S}| \times d}$$

とし、 $\phi_i(s) = \mathbb{I}_{\{i=s\}}$ とおけば、 $\hat{V}_{\mathbf{w}}(s) = \mathbf{w}_s$, $\forall s \in \mathcal{S}$ はテーブル形式の関数と同値。(関数近似の意味なし)

1.2 テーブル形式方法の拡張

モデルフリー型強化学習で V^π や V^* を求めたことを関数近似へと拡張する。

1.2.1 バッチ学習の場合

これまで、価値関数を求める際、なにかしらのベルマン作用素 B (履歴データから近似されるものも含む) を状態関数 v に作用させることで、

$$v'(s) := Bv(s), \quad \forall s \in \mathcal{S}$$

と更新してきた。一般に、 $\mathcal{V} \subset \mathbb{R}^{\mathcal{S}}$ より、関数近似器 $\hat{V}_{\mathbf{w}}$ にベルマン作用素を適用することで部分空間 \mathcal{V} の外に出てしまう場合がある。

$$\exists \hat{V}_{\mathbf{w}} \in \mathcal{V}, \quad B\hat{V}_{\mathbf{w}} \notin \mathcal{V}$$

このような場合、

$$\hat{V}'_{\mathbf{w}}(s) := B\hat{V}_{\mathbf{w}}(s), \quad \forall s \in \mathcal{S}$$

と従来の更新則を実行できない。そこで $B\hat{V}_{\mathbf{w}}$ を $\hat{V}_{\mathbf{w}}$ で近似するようにパラメータ \mathbf{w} を更新するアプローチを考える。すなわち、

1. 各状態 $s \in \mathcal{S}$ に対し目的関数を算出する。(fixing target)

$$V^{\text{target}}(s) := B\hat{V}_{\mathbf{w}}(s)$$

2. $\mathbf{w} \in \mathbb{R}^d$ を次のように更新する。

$$\mathbf{w} := \arg \min_{\mathbf{w} \in \mathbb{R}^d} \sum_{s \in \mathcal{S}} (V^{\text{target}}(s) - \hat{V}_{\mathbf{w}}(s))^2$$

例えば、関数近似器 \hat{V}_w が線形モデルの場合、

$$\hat{V}_w(s) = \mathbf{w}^\top \boldsymbol{\phi}(s) = (w_1 \ \dots \ w_{|\mathcal{S}|}) \begin{pmatrix} \phi_1(s) \\ \vdots \\ \phi_{|\mathcal{S}|}(s) \end{pmatrix}, \quad \forall s \in \mathcal{S}$$

$$\mathbf{v}^{\text{target}} = \begin{pmatrix} V^{\text{target}}(1) \\ \vdots \\ V^{\text{target}}(|\mathcal{S}|) \end{pmatrix}$$

とおけば、二乗誤差を微分して $\mathbf{0}$ とおけば

$$\begin{aligned} \frac{\partial}{\partial \mathbf{w}} \sum_{s \in \mathcal{S}} (V^{\text{target}}(s) - \mathbf{w}^\top \boldsymbol{\phi}(s))^2 &\propto \sum_{s \in \mathcal{S}} (V^{\text{target}}(s) - \mathbf{w}^\top \boldsymbol{\phi}(s)) \boldsymbol{\phi}(s)^\top \\ &= \sum_{s \in \mathcal{S}} V^{\text{target}}(s) \boldsymbol{\phi}(s)^\top - \mathbf{w}^\top \sum_{s \in \mathcal{S}} \boldsymbol{\phi}(s) \boldsymbol{\phi}(s)^\top = \mathbf{0} \end{aligned}$$

ここで、

$$\begin{aligned} \sum_{s \in \mathcal{S}} V^{\text{target}}(s) \boldsymbol{\phi}(s)^\top &= \sum_{s \in \mathcal{S}} V^{\text{target}}(s) (\phi_1(s) \ \dots \ \phi_{|\mathcal{S}|}(s)) \\ &= (V^{\text{target}}(1) \ \dots \ V^{\text{target}}(|\mathcal{S}|)) \begin{pmatrix} \phi_1(1) & \dots & \phi_{|\mathcal{S}|}(1) \\ \vdots & \ddots & \vdots \\ \phi_1(|\mathcal{S}|) & \dots & \phi_{|\mathcal{S}|}(|\mathcal{S}|) \end{pmatrix} \\ \sum_{s \in \mathcal{S}} \boldsymbol{\phi}(s) \boldsymbol{\phi}(s)^\top &= \sum_{s \in \mathcal{S}} \begin{pmatrix} \phi_1(s) \\ \vdots \\ \phi_{|\mathcal{S}|}(s) \end{pmatrix} (\phi_1(s) \ \dots \ \phi_{|\mathcal{S}|}(s)) \\ &= \begin{pmatrix} \phi_1(1) & \dots & \phi_1(|\mathcal{S}|) \\ \vdots & \ddots & \vdots \\ \phi_{|\mathcal{S}|}(1) & \dots & \phi_{|\mathcal{S}|}(|\mathcal{S}|) \end{pmatrix} \begin{pmatrix} \phi_1(1) & \dots & \phi_{|\mathcal{S}|}(1) \\ \vdots & \ddots & \vdots \\ \phi_1(|\mathcal{S}|) & \dots & \phi_{|\mathcal{S}|}(|\mathcal{S}|) \end{pmatrix} \end{aligned}$$

であるから

$$\Phi := \begin{pmatrix} \phi_1(1) & \dots & \phi_{|\mathcal{S}|}(1) \\ \vdots & \ddots & \vdots \\ \phi_1(|\mathcal{S}|) & \dots & \phi_{|\mathcal{S}|}(|\mathcal{S}|) \end{pmatrix}$$

とおけば (計画行列),

$$\begin{aligned} \mathbf{0} &= \mathbf{v}^{\text{target}} \Phi - \mathbf{w}^\top \Phi \Phi \\ \therefore \mathbf{w} &= (\Phi^\top \Phi)^{-1} \Phi^\top \mathbf{v}^{\text{target}} \end{aligned}$$

と解析的に w を求めることができる。(最小二乗法)
 しかし、この方法では結局全ての $s \in S$ に対し計算を行うため、関数近似の意味がない。そこで、状態の単位ではなく経験 (s, a, r, s') の単位で標本を準備して w を学習する適合価値反復法がある。履歴データを

$$\mathcal{D} := \{(s_{(1)}, a_{(1)}, r_{(1)}, s'_{(1)}), \dots, (s_{(N)}, a_{(N)}, r_{(N)}, s'_{(N)})\}$$

と置き $(1, \dots, N)$ は時間ステップとは関係なし、適合価値反復法の例を二つ示す。

- 価値関数 V^π の推定

1. 各経験の目的変数を算出

$$V_{(n)}^{\text{target}} := \hat{B}(\hat{V}; \{s_{(n)}, r_{(n)}, s'_{(n)}\})(s_{(n)}) = r_{(n)} + \gamma \hat{V}_w(s'_{(n)})$$

2. パラメータ w を更新

$$w := \arg \min_{w \in \mathbb{R}^d} \frac{1}{N} \sum_{n=1}^N (V_{(n)}^{\text{target}} - \hat{V}_w(s_{(n)}))^2$$

- 最適行動価値関数 Q^* の推定 (適合 Q 反復)

1. 各経験の目的変数を算出

$$Q_{(n)}^{\text{target}} := \hat{\Upsilon}_*(\hat{Q}_w; \{s_{(n)}, a_{(n)}, r_{(n)}, s'_{(n)}\})(s_{(n)}, a_{(n)}) = r_{(n)} + \gamma \max_{a' \in \mathcal{A}} \hat{Q}_w(s_{(n+1)}, a')$$

2. パラメータ w を更新

$$w := \arg \min_{w \in \mathbb{R}^d} \frac{1}{N} \sum_{n=1}^N (Q_{(n)}^{\text{target}} - \hat{Q}_w(s_{(n)}, a_{(n)}))^2$$

以上のような価値関数の関数近似の際には、収束性に気を配る必要がある。例 6.2 において、関数近似器を

$$\hat{V}_w(s) = \begin{cases} w, & \text{if } s = 1 \\ 2w, & \text{if } s = 2 \end{cases}$$

と定義し、正しい w の値を求めることを考える。 w が $w = 0$ へと収束すれば良い。ベルマン最適 (期待) 作用素を \hat{V}_w に適用することで、目的変数は

$$V^{\text{target}} := B(\hat{V}_w)(s) = \begin{cases} 0 + \gamma \times 2w, & (\text{if } s = 1) \\ 0 + \gamma \times 2w, & (\text{if } s = 2) \end{cases}$$

なので、第 k 回目の w の更新は、

$$\begin{aligned}
w_{k+1} &= \arg \min_{w \in \mathbb{R}} \sum_{s \in \mathcal{S}} (V^{\text{target}}(s) - \hat{V}_w(s))^2 \\
&= \arg \min_{w \in \mathbb{R}} \{(V^{\text{target}}(1) - \hat{V}_w(1))^2 + (V^{\text{target}}(2) - \hat{V}_w(2))^2\} \\
&= \arg \min_{w \in \mathbb{R}} \{(2\gamma w_k - w)^2 + (2\gamma w_k - 2w)^2\} \\
&= \arg \min_{w \in \mathbb{R}} \{5(w - \frac{6}{5}\gamma w_k)^2 + \frac{4}{5}\gamma^2 w_k^2\}
\end{aligned}$$

となるので、パラメータ w の更新則は

$$w_{k+1} = \frac{6}{5}\gamma w_k$$

となる。ゆえ、 $\gamma > \frac{5}{6}$ の場合、 w は単調増加で発散してしまう。

このように、関数近似器集合が真の関数を含んでいたとしても、更新則によってはパラメータが発散してしまう場合がある。そのため、関数近似器 \mathcal{V} の表現力を評価する上で、「最悪ケースの関数近似誤差」を用いる場合がある。関数近似誤差とは、 \mathbf{w} を

$$\mathbf{w} := \arg \min_{\mathbf{w} \in \mathbb{R}^d} \sum_{s \in \mathcal{S}} (V^{\text{target}}(s) - \hat{V}_{\mathbf{w}}(s))^2$$

と更新する場合、

$$\inf_{v' \in \mathcal{V}} \frac{1}{|\mathcal{S}|} \sum_{s \in \mathcal{S}} (V^{\text{target}}(s) - v'(s))^2$$

と定義できるので、「最悪ケース」とは v' による誤差が最大になってしまうケースであり、そのときの近似誤差は

$$\epsilon^*(\mathcal{V}; B) := \sup_{v \in \mathcal{V}} \inf_{v' \in \mathcal{V}} \frac{1}{|\mathcal{S}|} \sum_{s \in \mathcal{S}} (Bv(s) - v'(s))^2$$

となる。つまり、最悪ケースの近似誤差とは、「頑張って誤差が小さくなるよう v' を設定した時生じた誤差の最大値」である。指標 ϵ^* が高いほど、 \mathcal{V} の表現力が低い (任意の V^{target} に対する誤差を小さくする能力がない) ことを意味する。例 6.2 では、近似誤差は

$$\frac{4\gamma^2}{5} w_k^2$$

より、最悪近似誤差は $w_k \rightarrow \infty$ で無限大に発散してしまうため、不良な関数近似器を用いていたことがわかる。

1.2.2 オンライン学習の場合

行動方策 π による価値 V^π を求めるため、関数近似器 \hat{V}_w を用いて近似 TD 法を考える。関数近似をしない TD 法では特定の状態 s_t に対してしか \hat{V} を更新しなかったが、近似 TD 法では w_t を更新するため、経験した状態 s_t に応じてパラメータを更新すると別の状態に対しても推定価値が変わる事になる。TD 誤差を

$$\delta_t := r_t + \gamma \hat{V}_{w_t}(s_{t+1}) - \hat{V}_{w_t}(s_t)$$

とし、 \hat{V}_w の w に関する偏微分ベクトル (定数ベクトル) を

$$\nabla_w \hat{V}_{w_t}(s) := \begin{pmatrix} \left. \frac{\partial \hat{V}_w(s)}{\partial w_1} \right|_{w=w_t} \\ \vdots \\ \left. \frac{\partial \hat{V}_w(s)}{\partial w_d} \right|_{w=w_t} \end{pmatrix}$$

と定義する。今、 $\hat{V}_{w_t+\Delta w}(s)$ を w_t 周りでテイラー展開すると、

$$\begin{aligned} \hat{V}_{w_t+\Delta w}(s) &= \hat{V}_{w_t}(s) + \top(\nabla_w \hat{V}_{w_t}(s))(\mathbf{w}_t + \Delta \mathbf{w} - \mathbf{w}_t) + o(\|\Delta(\mathbf{w})\|^2) \\ &= \hat{V}_{w_t}(s) + \top(\nabla_w \hat{V}_{w_t}(s))\Delta \mathbf{w} + o(\|\Delta(\mathbf{w})\|^2) \end{aligned}$$

となる。すなわち、 w_t を Δw だけ変化させた時の \hat{V}_w の変化量は $\top(\nabla_w \hat{V}_{w_t}(s))\Delta w$ となる。さて、近似 TD 法において、

$$\mathbf{w}_{t+1} := \mathbf{w}_t + \alpha \delta_t \nabla_w \hat{V}_{w_t}(s_t) \quad (1)$$

と w を更新することになると、この時の s_t における推定価値の変化量は、

$$\begin{aligned} \hat{V}_{w_{t+1}}(s_t) &= \hat{V}_{w_t}(s_t) + \top(\nabla_w \hat{V}_{w_t}(s_t))\alpha \delta_t \nabla_w \hat{V}_{w_t}(s_t) \\ &= \hat{V}_{w_t}(s_t) + \alpha \delta_t \|\nabla_w \hat{V}_{w_t}(s_t)\|^2 \end{aligned}$$

となり、関数近似をしない TD 法における価値関数の更新とほとんど等価になる。(1) 式の右辺第二項を $\|\nabla_w \hat{V}_{w_t}(s_t)\|^2$ で割った更新則もある。こうすれば関数近似をしない TD 法と等価。) また、 $\tilde{s} \neq s_t$ であるような \tilde{s} における価値関数は、

$$\hat{V}_{w_{t+1}}(\tilde{s}) = \hat{V}(\tilde{s}) + \alpha \delta \top(\nabla_w \hat{V}_{w_t}(s_t))\nabla_w \hat{V}_{w_t}(\tilde{s})$$

と、内積 $\top(\nabla_w \hat{V}_{w_t}(s_t))\nabla_w \hat{V}_{w_t}(\tilde{s})$ に依存して更新される。よって、状態 (s_t, \tilde{s}) 間の類似度と内積 $\top(\nabla_w \hat{V}_{w_t}(s_t))\nabla_w \hat{V}_{w_t}(\tilde{s})$ が正の相関をもつように関数近似器 \hat{V}_w を設定できれば、現時点での状態 s_t の推定価値の更新と同時に s とよ

く似た状態 \tilde{s} の推定価値も更新することができる。

TD(λ) 法においては,

$$\hat{V}(s) := \hat{V}(s) + \alpha_t \delta_t \sum_{\tau=0}^t \mathbb{I}_{s_\tau=s} (\lambda \gamma)^{t-\tau}, \quad \forall s \in \mathcal{S}$$

と価値関数を更新していた。今, 近似 TD(λ) 法において

$$\mathbf{w}_{t+1} := \mathbf{w}_t + \alpha \delta_t \sum_{\tau=0}^t (\gamma \lambda)^{t-\tau} \nabla_{\mathbf{w}} \hat{V}_{\mathbf{w}_\tau}(s_\tau)$$

と \mathbf{w} を更新することになると, 状態 s_t における推定価値の変化量は

$$\hat{V}_{\mathbf{w}_{t+1}}(s_t) = \hat{V}_{\mathbf{w}_t}(s_t) + \alpha_t \delta_t \sum_{\tau=0}^t (\gamma \lambda)^{t-\tau} \nabla_{\mathbf{w}} \hat{V}_{\mathbf{w}_t}(s_t) \nabla_{\mathbf{w}} \hat{V}_{\mathbf{w}_\tau}(s_\tau)$$

と, 関数近似をしない TD(λ) 法における価値関数の更新とほとんど等価になる。

行動価値関数の学習に関数近似器 $\hat{Q}_{\mathbf{w}}$ を用いる場合でも, TD 誤差を

$$\delta_t^{(q)} := r + \gamma \max_{a' \in \mathcal{A}} \hat{Q}_{\mathbf{w}_t}(s_{t+1}, a') - \hat{Q}_{\mathbf{w}_t}(s_t, a_t)$$

とおけば

$$\mathbf{w}_{t+1} := \mathbf{w}_t + \alpha \delta_t^{(q)} \nabla_{\mathbf{w}} \hat{Q}_{\mathbf{w}_t}(s_t, a_t)$$

のようにパラメータを更新できる。

1.3 損失関数に基づく近似価値関数学習法

テーブル形式の方法を単純に拡張した方法では, 推定値 $\hat{V}_{\mathbf{w}}$ が収束しない場合があった。そこで, 損失関数を導入し, 損失関数の最小化を考える。「ベルマン作用素を繰り返し適用して $V^\pi(s) = B_\pi V^\pi(s)$ が成り立つような V^π を近似する, というこれまでのアイデアとは異なり, 損失関数を最小化するような \mathbf{w} を勾配法や最小二乗法によって解く。

1.3.1 ベルマン残差

ベルマン方程式：

$$V^\pi(s) = B_\pi V^\pi(s)$$

より，任意の s に対し $\hat{V}_w(s)$ と $B\hat{V}_w(s)$ の差が小さくなるように w を設定したい．そこで，ベルマン残差：

$$L_{\text{BR}}(w) := \sum_{s \in \mathcal{S}} \mu(s) (\hat{V}_w(s) - B\hat{V}_w(s))^2$$

の最小化を考える．今，関数近似器が線形かつ B に真のベルマン作用素 B_π を用いる：

$$\hat{V}_w(s) = {}^\top w \phi(s), \quad B(\hat{V}_w)(s) = \mathbb{E}_{A_t | S_t \sim \pi, S_{t+1} | S_t, A_t \sim p_T} [g(S_t, A_t) + \gamma \hat{V}_w(S_{t+1}) | S_t = s]$$

とする．これを L_{BR} の式に代入することで，

$$\begin{aligned} L_{\text{BR}}(w) &= \sum_{s \in \mathcal{S}} \mu(s) (\hat{V}_w(s) - \mathbb{E}_{A_t | S_t \sim \pi, S_{t+1} | S_t, A_t \sim p_T} [g(S_t, A_t) + \gamma \hat{V}_w(S_{t+1}) | S_t = s])^2 \\ &= \sum_{s \in \mathcal{S}} \mu(s) ({}^\top w \phi(s) - \mathbb{E}_{A_t | S_t \sim \pi, S_{t+1} | S_t, A_t \sim p_T} [g(S_t, A_t) + \gamma {}^\top w \phi(S_{t+1}) | S_t = s])^2 \\ &= \sum_{s \in \mathcal{S}} \mu(s) ({}^\top w \phi(s) \\ &\quad - \mathbb{E}_{A_t | S_t \sim \pi} [g(S_t, A_t) | S_t = s] \\ &\quad - \gamma {}^\top w \mathbb{E}_{A_t | S_t \sim \pi, S_{t+1} | S_t, A_t \sim p_T} [\phi(S_{t+1}) | S_t = s])^2 \\ &= \sum_{s \in \mathcal{S}} \mu(s) \{ {}^\top w (\phi(s) - \gamma \bar{\phi}_{+1}^\pi(s)) - \bar{g}^\pi(s) \}^2 \end{aligned}$$

ただし，

$$\begin{aligned} \bar{g}^\pi(s) &= \mathbb{E}_{A_t | S_t \sim \pi} [g(S_t, A_t) | S_t = s] = \sum_{a \in \mathcal{A}} \pi(a | s) g(s, a) \\ \bar{\phi}_{+1}^\pi(s) &= \mathbb{E}_{A_t | S_t \sim \pi, S_{t+1} | S_t, A_t \sim p_T} [\phi(S_{t+1}) | S_t = s] = \sum_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \pi(a | s) p_T(s' | s, a) \phi(s') \end{aligned}$$

とおいた。従って、 L_{BR} を \mathbf{w} について微分して $\mathbf{0}$ とおくことで L_{BR} の最小値を与えるパラメータの最適値 \mathbf{w}_{BR}^* が求まる。

$$\frac{\partial L_{\text{BR}}(\mathbf{w})}{\partial \mathbf{w}} = 2 \sum_{s \in \mathcal{S}} \mu(s) \{ \mathbf{w}^\top (\phi(s) - \gamma \bar{\phi}_{+1}^\pi(s)) - \bar{g}^\pi(s) \} (\phi(s) - \gamma \bar{\phi}_{+1}^\pi(s)) \quad (1)$$

$$\begin{aligned} &= 2 \mathbf{w}^\top \sum_{s \in \mathcal{S}} \mu(s) (\phi(s) - \gamma \bar{\phi}_{+1}^\pi(s))^\top (\phi(s) - \gamma \bar{\phi}_{+1}^\pi(s)) \\ &\quad - 2 \sum_{s \in \mathcal{S}} \mu(s) \bar{g}^\pi(s) (\phi(s) - \gamma \bar{\phi}_{+1}^\pi(s)) \\ \therefore \mathbf{w}_{\text{BR}}^* &= \left\{ \sum_{s \in \mathcal{S}} \mu(s) (\phi(s) - \gamma \bar{\phi}_{+1}^\pi(s))^\top (\phi(s) - \gamma \bar{\phi}_{+1}^\pi(s)) \right\}^{-1} \\ &\quad \times \sum_{s \in \mathcal{S}} \mu(s) \bar{g}^\pi(s) (\phi(s) - \gamma \bar{\phi}_{+1}^\pi(s)) \end{aligned} \quad (2)$$

\mathbf{w}_{BR}^* は状態行動空間が膨大である場合には計算することができない、そこで、履歴データ $h_T = \{s_0, a_0, r_0, s_{t+1}, \dots, s_{T-1}, a_{T-1}, r_{T-1}, s_T\}$ から \mathbf{w}_{BR}^* の推定量 $\hat{\mathbf{w}}_{\text{BR}}^*$ を求める。 $\mu(s)$ に状態の定常分布 $p_\infty^\pi(s)$ を取ることにすれば、(2) 式は

$$\begin{aligned} \mathbf{w}_{\text{BR}}^* &= \left\{ \mathbb{E}_{S_t \sim p_\infty^\pi} [(\phi(S_t) - \gamma \bar{\phi}_{+1}^\pi(S_t))^\top (\phi(S_t) - \gamma \bar{\phi}_{+1}^\pi(S_t))] \right\}^{-1} \\ &\quad \times \mathbb{E}_{S_t \sim p_\infty^\pi} [\bar{g}^\pi(S_t) (\phi(S_t) - \gamma \bar{\phi}_{+1}^\pi(S_t))] \end{aligned}$$

と書けるので、これを標本近似して、

$$\hat{\mathbf{w}}_{\text{BR}}^* = A^{-1} \mathbf{b}$$

とすることができる。ただし、

$$\begin{aligned} A &:= \frac{1}{T} \sum_{t=0}^{T-1} (\phi(s_t) - \phi(s_{t+1}))^\top (\phi(s_t) - \phi(\tilde{s}_{t+1})) \\ \mathbf{b} &:= \frac{1}{T} \sum_{t=0}^{T-1} r_t (\phi(s_t) - \gamma \phi(s_{t+1})) \end{aligned}$$

とおいた。オンラインの場合は(1)式より勾配降下法を行えばよく、

$$\mathbf{w} := \mathbf{w} - \alpha_t \{ \mathbf{w}^\top (\phi(s_t) - \gamma \phi(s_{t+1})) - r_t \} (\phi(s_t) - \gamma \phi(\tilde{s}_{t+1}))$$

となる。しかし、以上の推定式では s_t の次状態として s_t, \tilde{s}_{t+1} を独立にサンプリングしないと不偏推定量たりえない。どういうことかという、例えば

$$\bar{\phi}_{+1}^\pi(s) = \sum_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \pi(a|s) p_T(s'|s, a) \phi(s')$$

を推定する際は、 $\phi(s')$ の実現値の平均を取ること

$$\frac{1}{\sum_{t=0}^{T-1} \mathbb{I}_{s_t=s}} \sum_{t=0}^{T-1} \phi(s_{t+1}) \mathbb{I}_{s_t=s}$$

とすればこれは標本平均なのだから不偏推定量であるが、

$$\bar{\phi}_{+1}^\pi(s)^\top \bar{\phi}_{+1}^\pi(s) = \sum_{a \in \mathcal{A}} \sum_{s \in \mathcal{S}} \pi(a|s) p_T(s'|s, a) \phi(s') \times \top \left\{ \sum_{a \in \mathcal{A}} \sum_{s \in \mathcal{S}} \pi(a|s) p_T(s'|s, a) \phi(s') \right\}$$

の推定量を

$$\frac{1}{\sum_{t=0}^{T-1} \mathbb{I}_{s_t=s}} \sum_{t=0}^{T-1} \phi(s_{t+1})^\top \phi(s_{t+1}) \mathbb{I}_{s_t=s}$$

としてしまうと、これが近似しているのは

$$\sum_{a \in \mathcal{A}} \sum_{s \in \mathcal{S}} \pi(a|s) p_T(s'|s, a) \{ \phi(s')^\top \phi(s') \}$$

であって、 $\bar{\phi}_{+1}^\pi(s)^\top \bar{\phi}_{+1}^\pi(s)$ に対する不偏推定量ではない。そのため、

$$\frac{1}{\sum_{t=0}^{T-1} \mathbb{I}_{s_t=s}} \sum_{t=0}^{T-1} \phi(s_{t+1})^\top \phi(\tilde{s}_{t+1}) \mathbb{I}_{s_t=s}$$

と、 s_{t+1} を観測した後で \tilde{s}_{t+1} を前状態 s_t に戻って s_{t+1} とは独立にサンプリングしなければならない。これが二重サンプリング問題である。一般の強化学習問題では、二重サンプリングは実施が難しいため、残差勾配法では \tilde{s}_{t+1} を s_{t+1} と独立にサンプリングすることを諦め、 s_{t+1} で置き換えてしまう。

$$\begin{aligned} \mathbf{w} &:= \mathbf{w} - \alpha_t \{ \top \mathbf{w} (\phi(s_t) - \gamma \phi(s_{t+1})) - r_t \} (\phi(s_t) - \gamma \phi(s_{t+1})) \\ &= \mathbf{w} + \alpha_t (r_t + \gamma \hat{V}_{\mathbf{w}}(s_{t+1}) - \hat{V}_{\mathbf{w}}(s_t)) (\phi(s_t) - \gamma \phi(s_{t+1})) \end{aligned}$$

これは、次の期待二乗 TD 誤差：

$$\begin{aligned} L_{\text{TD}}(\mathbf{w}) &:= \sum_{s \in \mathcal{S}} \mu(s) \mathbb{E}_{A_t | S_t \sim \pi, S_{t+1} | S_t, A_t \sim p_T} [\{ g(S_t, A_t) + \gamma \hat{V}_{\mathbf{w}}(S_{t+1}) - \hat{V}_{\mathbf{w}}(S_t) \}^2 | S_t] \\ &= \sum_{s \in \mathcal{S}} \mu(s) \left[\sum_{a \in \mathcal{A}} \{ \pi(a|s) g(s, a) + \gamma \sum_{s' \in \mathcal{S}} p_T(s'|s, a) \hat{V}_{\mathbf{w}}(s') \} - \hat{V}_{\mathbf{w}}(s) \right]^2 \\ &= \sum_{s \in \mathcal{S}} \mu(s) \left[\sum_{a \in \mathcal{A}} \{ \pi(a|s) g(s, a) + \gamma \sum_{s' \in \mathcal{S}} p_T(s'|s, a) \top \mathbf{w} \phi(s') \} - \top \mathbf{w} \phi(s) \right]^2 \end{aligned}$$

の重み関数を $\mu = p_\infty^\pi$ とした時の不偏な確率的勾配法に対応する (微分して $\mathbf{0}$ とおけば明らか.). 残差勾配法は対応する損失関数 L_{TD} をもつため, 安定した学習を期待できるが, 実は, L_{TD} の最小化では V^π には収束しない. 後に詳述するが, L_{TD} とは, 報酬 $r_t = g(s_t, a_t)$ を $\hat{V}(s_t) = r_t + \gamma \hat{V}(s_{t+1})$ を用いて線形モデル ${}^\top(\phi(s_t) - \gamma \phi(s_{t+1}))\mathbf{w}$ で予測する場合の

$$r_t = {}^\top(\phi(s_t) - \gamma \phi(s_{t+1}))\mathbf{w} + \sigma_t$$

の誤差項 σ_t の期待二乗値に対応する.

1.3.2 射影ベルマン残差

$B\hat{V}_\mathbf{w}$ が必ずしも関数近似器空間 \mathcal{V} に含まれているとは限らないことから, 射影ベルマン残差:

$$L_{PBR}(\mathbf{w}) := \sum_{s \in \mathcal{S}} \mu(s) (\hat{V}_\mathbf{w}(s) - \Gamma(B\hat{V}_\mathbf{w})(s))^2$$

を用いる. ただし, Γ は任意の状態関数 v を関数近似器空間 \mathcal{V} へと直交射影する作用素:

$$\Gamma(v) := \arg \min_{\hat{V}_\mathbf{w} \in \mathcal{V}} \sum_{s \in \mathcal{S}} \mu(s) (v(s) - \hat{V}_\mathbf{w}(s))^2$$

である. L_{PBR} の最小化を考える. 今, 関数近似器は線形かつ B に真のベルマン作用素 B_π を用いる:

$$\begin{aligned} \hat{V}_\mathbf{w}(s) &= {}^\top \mathbf{w} \phi(s) = w_1 \phi_1(s) + \dots + w_{|\mathcal{V}|} \phi_{|\mathcal{V}|}(s) \\ B(\hat{V}_\mathbf{w})(s) &= \mathbb{E}_{A_t | S_t \sim \pi, S_{t+1} | S_t, A_t \sim p_T} [g(S_t, A_t) + \gamma \hat{V}_\mathbf{w}(S_{t+1}) | S_t = s] \end{aligned}$$

まずは, 次のような行列表記を用意する.

- 状態遷移確率 $P^\pi \dots (i, j)$ 要素が状態 i から状態 j へと遷移する確率であるような行列.

$$(P^\pi)_{i,j} := p_{MC}^\pi(j|i) = \sum_{a \in \mathcal{A}} \pi(a|i) p_T(j|i, a)$$

- 報酬ベクトル $\mathbf{r}^\pi \dots$ 第 i 要素が状態 i における報酬の期待値であるような行列.

$$(\mathbf{r}^\pi)_i := \mathbb{E}_{A_t | S_t \sim \pi} [g(A_t, S_t) | S_t = i] = \sum_{a \in \mathcal{A}} \pi(a|i) g(i, a)$$

- 損失関数の重み行列 $U \dots i$ 番目の対角要素が L_{PBR} の状態 i の重みである行列.

$$(U)_{i,j} := \begin{cases} \mu(i), & (i = j) \\ 0, & \text{otherwise} \end{cases}$$

- 価値ベクトル \mathbf{v}^π と推定価値ベクトル $\hat{\mathbf{v}}_{\mathbf{w}} \dots$ 第 i 要素がそれぞれ状態 i の価値及び推定価値であるようなベクトル. 前者は $|\mathcal{S}|$ 次元, 後者は $|\mathcal{V}|$ 次元.

$$(\mathbf{v}^\pi) := V^\pi(i)$$

$$\hat{\mathbf{v}}_{\mathbf{w}} := \begin{pmatrix} \hat{V}_{\mathbf{w}}(1) \\ \vdots \\ \hat{V}_{\mathbf{w}}(|\mathcal{S}|) \end{pmatrix} = \Phi \mathbf{w} = \begin{pmatrix} \phi_1(1) & \dots & \phi_{|\mathcal{V}|}(1) \\ \vdots & \ddots & \vdots \\ \phi_1(|\mathcal{S}|) & \dots & \phi_{|\mathcal{V}|}(|\mathcal{S}|) \end{pmatrix} \begin{pmatrix} \mathbf{w}_1 \\ \vdots \\ \mathbf{w}_{|\mathcal{V}|} \end{pmatrix} = \begin{pmatrix} \top \phi(1) \\ \vdots \\ \top \phi(|\mathcal{S}|) \end{pmatrix} \mathbf{w}$$

今, ある状態関数 v に対し,

$$\mathbf{v} = \begin{pmatrix} v(1) \\ \vdots \\ v(|\mathcal{S}|) \end{pmatrix}$$

と置く. \mathbf{v} は $|\mathcal{S}|$ 次元のベクトルである. 対して, $|\mathcal{V}| < |\mathcal{S}|$ ならば, $\hat{\mathbf{v}}_{\mathbf{w}}$ は $|\mathcal{V}|$ 次元のベクトルである. $\Gamma(v)$ において, 二乗誤差を微分して $\mathbf{0}$ と置く.

$$\begin{aligned} \frac{\partial}{\partial \mathbf{w}} \sum_{s \in \mathcal{S}} \mu(s) (v(s) - \top \mathbf{w} \phi(s))^2 &\propto \sum_{s \in \mathcal{S}} \mu(s) (v(s) - \top \mathbf{w} \phi(s)) \top \phi(s) \\ &= \sum_{s \in \mathcal{S}} \mu(s) v(s) \top \phi(s) - \top \mathbf{w} \sum_{s \in \mathcal{S}} \mu(s) \phi(s) \top \phi(s) = \mathbf{0} \\ \therefore \mathbf{0} &= \top \mathbf{v} \top U \Phi - \top \mathbf{w} \top \Phi U \Phi \\ \therefore \hat{\mathbf{w}} &= (\top \Phi U \Phi)^{-1} \top \Phi U \mathbf{v} \end{aligned}$$

ただし, $\hat{\mathbf{w}}$ は Γ によるパラメータの推定値である. v は Γ によって \mathcal{V} 内の関数近似器 $\hat{V}_{\hat{\mathbf{w}}}$ へと射影されることになる. したがって, ある状態関数 v に対し, 任意の s について Γ を作用させることは, \mathbf{v} に対し以下のように Γ を作用させて部分空間に射影することと言える.

$$\Gamma \mathbf{v} := \begin{pmatrix} \Gamma(v)(1) \\ \vdots \\ \Gamma(v)(|\mathcal{S}|) \end{pmatrix} = \begin{pmatrix} \hat{V}_{\hat{\mathbf{w}}}(1) \\ \vdots \\ \hat{V}_{\hat{\mathbf{w}}}(|\mathcal{S}|) \end{pmatrix} = \Phi \hat{\mathbf{w}} = \Phi (\top \Phi U \Phi)^{-1} \top \Phi U \mathbf{v}$$

ここで, 行列

$$\mathbf{\Gamma} = \Phi(\mathbf{\Phi}^\top \Phi U \Phi)^{-1} \mathbf{\Phi}^\top \Phi U$$

を射影行列と呼ぶ. 射影行列 $\mathbf{\Gamma}$ を用いることで, 射影ベルマン残差 L_{PBR} は,

$$\begin{aligned} L_{\text{PBR}} &= \mathbf{\Phi}^\top \{ \Phi \mathbf{w} - \mathbf{\Gamma}(\mathbf{r}^\pi + \gamma P^\pi \Phi \mathbf{w}) \} U \{ \Phi \mathbf{w} - \mathbf{\Gamma}(\mathbf{r}^\pi + \gamma P^\pi \Phi \mathbf{w}) \} \\ &= \mathbf{\Phi}^\top (\Phi \mathbf{w} - \mathbf{r}^\pi - \gamma P^\pi \Phi \mathbf{w}) \mathbf{\Gamma} U \mathbf{\Gamma} (\Phi \mathbf{w} - \mathbf{r}^\pi - \gamma P^\pi \Phi \mathbf{w}) \\ &\quad (\because \mathbf{\Gamma} \Phi \mathbf{w} = \Phi \mathbf{w}) \\ &= \mathbf{\Phi}^\top \{ \mathbf{r}^\pi - (\Phi - \gamma P^\pi \Phi) \mathbf{w} \} \mathbf{\Gamma} U \mathbf{\Gamma} \{ \mathbf{r}^\pi - (\Phi - \gamma P^\pi \Phi) \mathbf{w} \} \\ &= \mathbf{\Phi}^\top (\mathbf{r}^\pi - L^\pi \mathbf{w}) U \Phi (\mathbf{\Phi}^\top \Phi U \Phi)^{-1} \mathbf{\Phi}^\top \Phi U (\mathbf{r}^\pi - L^\pi \mathbf{w}) \\ &\quad (\because \mathbf{\Gamma} = \Phi (\mathbf{\Phi}^\top \Phi U \Phi)^{-1} \mathbf{\Phi}^\top \Phi U) \end{aligned}$$

ただし, $L^\pi = \Phi - \gamma P^\pi \Phi$ と置いた. よって $L_{\text{PBR}}(\mathbf{w})$ の一階微分を $\mathbf{0}$ と置くことで,

$$\frac{\partial}{\partial \mathbf{w}} L_{\text{PBR}}(\mathbf{w}) = -2 \mathbf{\Phi}^\top L^\pi U \Phi (\mathbf{\Phi}^\top \Phi U \Phi)^{-1} \mathbf{\Phi}^\top \Phi U (\mathbf{r}^\pi - L^\pi \mathbf{w}) = \mathbf{0}$$

より, $L_{\text{PBR}}(\mathbf{w})$ を最小化する最適パラメータ $\mathbf{w}_{\text{PBR}}^*$ は,

$$\begin{aligned} \mathbf{w}_{\text{PBR}}^* &= (L^\pi)^{-1} \mathbf{r}^\pi \\ &= (L^\pi)^{-1} U^{-1} \mathbf{\Phi}^{-1} \mathbf{\Phi}^{-1} U \mathbf{r}^\pi \\ &= (\mathbf{\Phi}^\top \Phi U L^\pi)^{-1} \mathbf{\Phi}^\top \Phi U \mathbf{r}^\pi \end{aligned}$$

ここで,

$$\begin{aligned} \mathbf{\Phi}^\top \Phi U L^\pi &= \mathbf{\Phi}^\top \Phi U (\Phi - \gamma P^\pi \Phi) \\ &= \mathbf{\Phi}^\top \Phi U \Phi - \gamma \mathbf{\Phi}^\top \Phi U P^\pi \Phi \\ &= \sum_{s \in \mathcal{S}} \mu(s) \phi(s) \mathbf{\Phi}^\top \phi(s) - \gamma \sum_{s \in \mathcal{S}} \mu(s) \mathbf{\Phi}^\top \phi(s) \sum_{a \in \mathcal{A}} \sum_{s' \in \mathcal{S}} \pi(a|s) p_T(s'|s, a) \phi(s') \\ &= \sum_{s \in \mathcal{S}} \mu(s) \phi(s) \mathbf{\Phi}^\top \phi(s) - \gamma \sum_{s \in \mathcal{S}} \mu(s) \phi(s) \mathbf{\Phi}^\top \bar{\phi}_{+1}^\pi(s) \\ &= \sum_{s \in \mathcal{S}} \mu(s) \phi(s) \mathbf{\Phi}^\top (\phi(s) - \gamma \bar{\phi}_{+1}^\pi(s)) \\ \mathbf{\Phi}^\top \Phi U \mathbf{r}^\pi &= \sum_{s \in \mathcal{S}} \mu(s) \bar{g}^\pi(s) \phi(s) \end{aligned}$$

なので,

$$\mathbf{w}_{\text{PBR}}^* = \left\{ \sum_{s \in \mathcal{S}} \mu(s) \phi(s) \mathbf{\Phi}^\top (\phi(s) - \gamma \bar{\phi}_{+1}^\pi(s)) \right\}^{-1} \sum_{s \in \mathcal{S}} \mu(s) \bar{g}^\pi(s) \phi(s)$$

$\mu(s) = p_\infty^\pi$ とすれば, 上式は履歴データ $\{s_0, r_0, \dots, s_{T-1}, r_{T-1}, s_T\}$ から推定することができる,

$$\hat{w}_{\text{PBR}}^* = \left\{ \frac{1}{T} \sum_{t=0}^{T-1} \phi(s_t)^\top (\phi(s_t) - \gamma \phi(s_{t+1})) \right\}^{-1} \left\{ \frac{1}{T} \sum_{t=0}^{T-1} r_t \phi(s_t) \right\}$$

である.

2 方策の関数近似

方策パラメータ $\theta \in \mathbb{R}^d$ で確率の方策 $\pi_\theta : \mathcal{A} \times \mathcal{S} \rightarrow [0, 1]$ で直接規定して, θ を学習させることを考える. π は例えば

$$\pi_{\text{normal}}(a|s; \theta) := \frac{1}{\sqrt{2\pi(\sigma(s; \theta))^2}} \exp\left\{-\frac{(a - \mu(s; \theta))^2}{2(\sigma(s; \theta))^2}\right\}$$

のような正規分布に基づくものがある. 目的関数には, 時間ステップ $t = 0$ からのリターン C_0 の期待値

$$\begin{aligned} f_0(\theta) &:= \mathbb{E}_{S_0 \sim p_{S_0}, A_t | S_t \sim \pi_\theta, S_{t+1} | S_t, A_t \sim p_T} [C_0] \\ &= \mathbb{E}_{S_0 \sim p_{S_0}} [V^{\pi_\theta}(S_0)] \\ &= \sum_{s \in \mathcal{S}} p_{S_0}(s) V^{\pi_\theta}(s) \end{aligned}$$

や, マルコフ連鎖がエルゴート性を満たす場合のみ全時間点での期待リターンの平均

$$\begin{aligned} f_\infty(\theta) &:= \lim_{T \rightarrow \infty} \mathbb{E}_{S_0 \sim p_{S_0}, A_t | S_t \sim \pi_\theta, S_{t+1} | S_t, A_t \sim p_T} \left[\frac{1}{T} \sum_{t=0}^{T-1} C_t \right] \\ &\propto \lim_{T \rightarrow \infty} \mathbb{E}_{S_0 \sim p_{S_0}, A_t | S_t \sim \pi_\theta, S_{t+1} | S_t, A_t \sim p_T} \left[\frac{1}{T} \sum_{t=0}^{T-1} g(S_t, A_t) \right] \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} p_\infty^{\pi_\theta}(s) \pi_\theta(a|s) g(s, a) \end{aligned}$$

などを用いる. 方策勾配とは, 目的関数 f の θ の一階微分ベクトル

$$\nabla_\theta f(\theta) = \begin{pmatrix} \frac{\partial f(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial f(\theta)}{\partial \theta_d} \end{pmatrix}$$

のことであり、これを求めることができれば

$$\boldsymbol{\theta} := \boldsymbol{\theta} + \alpha_t \nabla_{\boldsymbol{\theta}} f(\boldsymbol{\theta})$$

と、目的関数が大きくなる方向にパラメータを更新するアプローチを取ることができる。しかし、 $f_0(\boldsymbol{\theta})$ および $f_{\infty}(\boldsymbol{\theta})$ の微分はこのままでは求めることができない。 $(V^{\pi\theta}$ はそのままでは微分できない)

まずは、 $\nabla_{\boldsymbol{\theta}} V^{\pi\theta}$ 、ひいては $\nabla_{\boldsymbol{\theta}} f_0$ を求めるために、いくつかの概念を導入する。

エピソード... 確率 γ で行動選択し、 $1 - \gamma$ で終了する時、確率過程が終了するまでの系列 $\{s_0, a_0, r_0, \dots, s_{t-1}, a_{t-1}, r_{t-1}, s_t\}$

期待割引累積訪問数... 1 エピソードで状態 s に訪れる回数の期待値

$$\begin{aligned} d^{\pi}(s) &= \mathbb{E}_{S_0 \sim p_{S_0}, A_t | S_t \sim \pi, S_{t+1} | A_t, S_t \sim p_T} \left[\sum_{t=0}^{\infty} \gamma^t \mathbb{I}_{S_t=s} \right] \\ &= \sum_{t=0}^{\infty} \gamma^t Pr[S_t = s | M(\pi)], \quad s \in \mathcal{S} \end{aligned}$$

エルゴート性を満たさないマルコフ連鎖に対しエピソードを定義することで、エピソードごとにパラメータ更新をすることを考えることができる。今、エピソード γ を実行するときの期待値演算子を \mathbb{E}^{γ} と置くことにすると、

$$\begin{aligned} d^{\pi}(s) &= \mathbb{E}_{S_t \sim Pr}^{\gamma} \left[\sum_{t=0}^{\infty} \mathbb{I}_{S_t=s} \right] \\ V^{\pi}(s) &= \mathbb{E}_{A_t | S_t \sim \pi, S_{t+1} | S_t, A_t \sim p_T} \left[\sum_{t=0}^{\infty} \gamma^t g(S_t, A_t) | S_0 = s \right] \\ &= \mathbb{E}_{S_t, A_t | S_0 \sim Pr}^{\gamma} \left[\sum_{t=0}^{\infty} g(S_t, A_t) | S_0 = s \right] \\ Q^{\pi}(s, a) &= \mathbb{E}_{S_{t+1} | S_t, A_t \sim p_T, A_{t+1} | S_{t+1} \sim \pi} \left[\sum_{t=0}^{\infty} \gamma^t g(S_t, A_t) | S_0 = s, A_0 = a \right] \\ &= \mathbb{E}_{S_t, A_t | S_0, A_0 \sim Pr}^{\gamma} \left[\sum_{t=0}^{\infty} g(S_t, A_t) | S_0 = s, A_0 = a \right] \end{aligned}$$

(割引率 γ と行動選択確率 γ は同値。最初はエピソードで切る場合と無限時間の場合で価値関数が同値になってしまうのが変な気がしたが、よく見ると上式はエピソードで切ることは考えていなくて、「割引率を γ とする価値関数は、エピソード γ を実行するときの期待値演算子で表現出来る」ことが言い

たい.)

以上より, 目的関数 $f_0(\boldsymbol{\theta})$ は,

$$\begin{aligned}
f_0(\boldsymbol{\theta}) &= \sum_{s \in \mathcal{S}} p_{S_0}(s) V^{\pi_{\boldsymbol{\theta}}}(s) \\
&= \sum_{s \in \mathcal{S}} p_{S_0}(s) \mathbb{E}_{S_t, A_t | S_0 \sim Pr}^{\gamma} \left[\sum_{t=0}^{\infty} g(S_t, A_t) | S_0 = s \right] \\
&= \mathbb{E}_{S_t, A_t \sim Pr}^{\gamma} \left[\sum_{t=0}^{\infty} g(S_t, A_t) \right] \\
&= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \sum_{t=0}^{\infty} \gamma^t Pr[S_t = s] \pi_{\boldsymbol{\theta}}(a|s) g(s, a) \\
&= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d^{\pi_{\boldsymbol{\theta}}}(s) \pi_{\boldsymbol{\theta}}(a|s) g(s, a)
\end{aligned}$$

次に, 目的関数 f_{∞} を価値関数で表現するため, 割引現在価値を用いない価値関数として, 差分価値関数を導入する.

$$\begin{aligned}
Q_{\infty}^{\pi_{\boldsymbol{\theta}}}(s, a) &= \sum_{t=0}^{\infty} \mathbb{E}_{S_{t+1} | A_t, S_t \sim p_{T, A_{t+1} | S_{t+1} \sim \pi_{\boldsymbol{\theta}}}} [R_t - f_{\infty}(\boldsymbol{\theta}) | S_0 = s, A_0 = a] \\
V_{\infty}^{\pi_{\boldsymbol{\theta}}}(s) &= \sum_{a \in \mathcal{A}} \pi_{\boldsymbol{\theta}}(a|s) Q_{\infty}^{\pi_{\boldsymbol{\theta}}}(s, a)
\end{aligned}$$

割引現在価値 γ を使わないまま無限和を取ると発散してしまうため, 差分価値関数では各時点で f_{∞} を引いている, $t \rightarrow \infty$ で

$$\mathbb{E}_{S_t | A_{t-1}, S_{t-1} \sim p_{T, A_t | S_t \sim \pi_{\boldsymbol{\theta}}}} [g(S_t, A_t) | A_{t-1}, S_{t-1}] = \mathbb{E}_{S_t \sim p_{\infty}^{\pi_{\boldsymbol{\theta}}}, A_t \sim \pi_{\boldsymbol{\theta}}} [g(S_t, A_t)] = f_{\infty}(\boldsymbol{\theta})$$

となるため, 差分価値関数は有限である. 差分価値関数には, 次のような再起式が成り立つ.

$$\begin{aligned}
Q_{\infty}^{\pi_{\boldsymbol{\theta}}}(s, a) &= g(s, a) - f_{\infty}(\boldsymbol{\theta}) + \sum_{s' \in \mathcal{S}} p_T(s'|s, a) \sum_{a' \in \mathcal{A}} \pi_{\boldsymbol{\theta}}(a'|s') Q_{\infty}^{\pi_{\boldsymbol{\theta}}}(s', a') \\
V_{\infty}^{\pi_{\boldsymbol{\theta}}}(s) &= \sum_{a \in \mathcal{A}} \{g(s, a) - f_{\infty}(\boldsymbol{\theta}) + \sum_{s' \in \mathcal{S}} p_T(s'|s, a) V_{\infty}^{\pi_{\boldsymbol{\theta}}}(s')\}
\end{aligned}$$

目的関数 f の $\boldsymbol{\theta}$ による微分のことを方策勾配 $\nabla_{\boldsymbol{\theta}} f$ という. 方策勾配は, 次の性質を持つ.

方策勾配定理

a. 平均報酬 f_∞ の方策勾配は, 任意の状態関数 $b: \mathcal{S} \rightarrow \mathbb{R}$ を用いて

$$\nabla_{\theta} f_\infty(\theta) = \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} p_\infty^{\pi_\theta}(s) \pi_\theta(a|s) \nabla_{\theta} \log \pi_\theta(a|s) (Q_\infty^{\pi_\theta}(s, a) - b(s))$$

b. 価値関数の割引現在価値 γ による重みつき和 f_0 の方策勾配は,

$$\nabla_{\theta} f_0(\theta) = \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d^{\pi_\theta}(s) \pi_\theta(a|s) \nabla_{\theta} \log \pi_\theta(a|s) (Q^{\pi_\theta}(s, a) - b(s))$$

(証明 a.)

差分価値関数 $Q_\infty^{\pi_\theta}$ の再帰式:

$$Q_\infty^{\pi_\theta}(s, a) = g(s, a) - f_\infty(\theta) + \sum_{s' \in \mathcal{S}} p_T(s'|s, a) \sum_{a' \in \mathcal{A}} \pi_\theta(a'|s') Q_\infty^{\pi_\theta}(s', a')$$

の両辺を θ で微分すると,

$$\begin{aligned} \nabla_{\theta} Q_\infty^{\pi_\theta}(s, a) &= -\nabla_{\theta} f_\infty(\theta) \\ &+ \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}} p_T(s'|s, a) \{ \nabla_{\theta} \pi_\theta(a'|s') Q_\infty^{\pi_\theta}(s', a') + \pi_\theta(a'|s') \nabla_{\theta} Q_\infty^{\pi_\theta}(s', a') \} \end{aligned}$$

ここで,

$$\nabla_{\theta} \pi_\theta(a'|s') = \pi_\theta(a'|s') \nabla_{\theta} \log \pi_\theta(a'|s')$$

を用いて

$$\begin{aligned} \nabla_{\theta} f_\infty(\theta) &= -\nabla_{\theta} Q_\infty^{\pi_\theta}(s, a) \\ &+ \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}} p_T(s'|s, a) \pi_\theta(a'|s') \{ \nabla_{\theta} \log \pi_\theta(a'|s') Q_\infty^{\pi_\theta}(s', a') + \nabla_{\theta} Q_\infty^{\pi_\theta}(s', a') \} \end{aligned}$$

また, $f_\infty(\theta)$ は s, a に非依存なので,

$$f_\infty(\theta) = \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} p_\infty^{\pi_\theta}(s) \pi_\theta(a|s) \nabla_{\theta} f_\infty(\theta)$$

が成り立つのだから,

$$\begin{aligned}
\nabla_{\theta} f_{\infty}(\theta) &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} p_{\infty}^{\pi_{\theta}}(s) \pi_{\theta}(a|s) [-\nabla_{\theta} Q_{\infty}^{\pi_{\theta}}(s, a) \\
&\quad + \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}} p_T(s'|s, a) \pi_{\theta}(a'|s') \{ \nabla_{\theta} \log \pi_{\theta}(a'|s') Q_{\infty}^{\pi_{\theta}}(s', a') + \nabla_{\theta} Q_{\infty}^{\pi_{\theta}}(s', a') \}] \\
&= - \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} p_{\infty}^{\pi_{\theta}}(s) \pi_{\theta}(a|s) \nabla_{\theta} Q_{\infty}^{\pi_{\theta}}(s, a) \\
&\quad + \sum_{s' \in \mathcal{S}} \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} p_{\infty}^{\pi_{\theta}}(s) \pi_{\theta}(a|s) p_T(s'|s, a) \\
&\quad \quad \times \sum_{a' \in \mathcal{A}} \pi_{\theta}(a'|s') \{ \nabla_{\theta} \log \pi_{\theta}(a'|s') Q_{\infty}^{\pi_{\theta}}(s', a') + \nabla_{\theta} Q_{\infty}^{\pi_{\theta}}(s', a') \} \\
&= - \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} p_{\infty}^{\pi_{\theta}}(s) \pi_{\theta}(a|s) \nabla_{\theta} Q_{\infty}^{\pi_{\theta}}(s, a) \\
&\quad + \sum_{s' \in \mathcal{S}} \sum_{a' \in \mathcal{A}} p_{\infty}^{\pi_{\theta}}(s') \pi_{\theta}(a'|s') \{ \nabla_{\theta} \log \pi_{\theta}(a'|s') Q_{\infty}^{\pi_{\theta}}(s', a') + \nabla_{\theta} Q_{\infty}^{\pi_{\theta}}(s', a') \} \\
&= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} p_{\infty}^{\pi_{\theta}}(s) \pi_{\theta}(a|s) \nabla_{\theta} \log \pi_{\theta}(a|s) Q_{\infty}^{\pi_{\theta}}(s, a)
\end{aligned}$$

ここで, 任意の状態関数 $b(s)$ に対し,

$$\begin{aligned}
\sum_{a \in \mathcal{A}} \pi_{\theta}(a|s) \nabla_{\theta} \log \pi_{\theta}(a|s) b(s) &= \sum_{a \in \mathcal{A}} \nabla_{\theta} \pi_{\theta}(a|s) b(s) \\
&= b(s) \nabla_{\theta} \left\{ \sum_{a \in \mathcal{A}} \pi_{\theta}(a|s) \right\} = 0
\end{aligned}$$

なのだから,

$$\nabla_{\theta} f_{\infty}(\theta) = \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} p_{\infty}^{\pi_{\theta}}(s) \pi_{\theta}(a|s) \nabla_{\theta} \log \pi_{\theta}(a|s) (Q_{\infty}^{\pi_{\theta}}(s, a) - b(s))$$

が成り立つ.

(証明 b.)

$$\nabla_{\theta} f_0(\theta) = \sum_{s \in \mathcal{S}} p_{S_0}(s) \nabla_{\theta} V^{\pi_{\theta}}(s)$$

より, まずは $\nabla_{\theta} V^{\pi_{\theta}}(s)$ を求める.

$$V^{\pi_{\theta}}(s) = \sum_{a \in \mathcal{A}} \pi_{\theta}(a|s) Q^{\pi_{\theta}}(s, a)$$

の両辺を θ で微分することで,

$$\begin{aligned}\nabla_{\theta} V^{\pi_{\theta}}(s) &= \sum_{a \in \mathcal{A}} \{ \nabla_{\theta} \pi_{\theta}(a|s) Q^{\pi_{\theta}}(s, a) + \pi_{\theta}(a|s) \nabla_{\theta} Q^{\pi_{\theta}}(s, a) \} \\ &= \sum_{a \in \mathcal{A}} \pi_{\theta}(a|s) \{ \nabla_{\theta} \log \pi_{\theta}(a|s) Q^{\pi_{\theta}}(s, a) + \nabla_{\theta} Q^{\pi_{\theta}}(s, a) \}\end{aligned}$$

ここで,

$$Q^{\pi_{\theta}}(s, a) = g(s, a) + \gamma \sum_{s' \in \mathcal{S}} p_T(s'|s, a) V^{\pi_{\theta}}(s')$$

より

$$\nabla_{\theta} Q^{\pi_{\theta}}(s, a) = \gamma \sum_{s' \in \mathcal{S}} p_T(s'|s, a) \nabla_{\theta} V^{\pi_{\theta}}(s')$$

であるから,

$$\nabla_{\theta} V^{\pi_{\theta}}(s) = \sum_{a \in \mathcal{A}} \pi_{\theta}(a|s) \{ \nabla_{\theta} \log \pi_{\theta}(a|s) Q^{\pi_{\theta}}(s, a) + \gamma \sum_{s' \in \mathcal{S}} p_T(s'|s, a) \nabla_{\theta} V^{\pi_{\theta}}(s') \}$$

この式より, ベクトル $\nabla_{\theta} V^{\pi_{\theta}}(s)$ の第 i 要素は,

$$\begin{aligned}\frac{\partial}{\partial \theta_i} V^{\pi_{\theta}}(s) &= \sum_{a \in \mathcal{A}} \pi_{\theta}(a|s) \left\{ \frac{\partial}{\partial \theta_i} \log \pi_{\theta}(a|s) Q^{\pi_{\theta}}(s, a) + \gamma \sum_{s' \in \mathcal{S}} p_T(s'|s, a) \frac{\partial}{\partial \theta_i} V^{\pi_{\theta}}(s') \right\} \\ &= \sum_{a \in \mathcal{A}} \pi_{\theta}(a|s) \left\{ g_i(s, a) + \gamma \sum_{s' \in \mathcal{S}} p_T(s'|s, a) \frac{\partial}{\partial \theta_i} V^{\pi_{\theta}}(s') \right\}\end{aligned}$$

ただし,

$$g_i(s, a) = \frac{\partial}{\partial \theta_i} \log \pi_{\theta}(a|s) Q^{\pi_{\theta}}(s, a)$$

より, $\nabla_{\theta} V^{\pi_{\theta}}(s)$ とは, 第 i 要素が報酬関数を $g_i(s, a)$ に持つ価値関数であることがわかる. ゆえに, 第 i 要素が $g_i(s, a)$ であるようなベクトル $\mathbf{g}(s, a)$ を用いて,

$$\begin{aligned}\nabla_{\theta} V^{\pi_{\theta}}(s) &= \mathbb{E}_{A_t|S_t \sim \pi, S_{t+1}|S_t, A_t \sim p_T} \left[\sum_{t=0}^{\infty} \gamma^t \mathbf{g}(S_t, A_t) | S_0 = s \right] \\ &= \mathbb{E}_{S_t, A_t | S_0 \sim Pr}^{\gamma} \left[\sum_{t=0}^{\infty} \nabla_{\theta} \log \pi_{\theta}(A_t | S_t) Q^{\pi_{\theta}}(S_t, A_t) | S_0 = s \right]\end{aligned}$$

以上より,

$$\begin{aligned}
\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}) &= \sum_{s \in \mathcal{S}} p_{S_0}(s) \nabla_{\boldsymbol{\theta}} V^{\pi_{\boldsymbol{\theta}}}(s) \\
&= \sum_{s \in \mathcal{S}} p_{S_0}(s) \mathbb{E}_{S_t, A_t | S_0 \sim P_r}^{\gamma} \left[\sum_{t=0}^{\infty} \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(A_t | S_t) Q^{\pi_{\boldsymbol{\theta}}}(S_t, A_t) | S_0 = s \right] \\
&= \mathbb{E}_{S_t, A_t \sim P_r}^{\gamma} \left[\sum_{t=0}^{\infty} \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(A_t | S_t) Q^{\pi_{\boldsymbol{\theta}}}(S_t, A_t) \right] \\
&= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d^{\pi_{\boldsymbol{\theta}}}(s) \pi_{\boldsymbol{\theta}}(a | s) \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(a | s) Q^{\pi_{\boldsymbol{\theta}}}(s, a)
\end{aligned}$$

方策勾配定理は、空間平均を時間平均に置き換えることで次のように書き換えることができる。

$$\begin{aligned}
\nabla_{\boldsymbol{\theta}} f_{\infty}(\boldsymbol{\theta}) &= \lim_{T \rightarrow \infty} \mathbb{E}_{S_0 \sim p_{S_0}, A_t | S_t \sim \pi, S_{t+1} | A_t, S_t \sim p_T} \left[\frac{1}{T} \sum_{t=0}^{T-1} \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(A_t | S_t) (Q_{\infty}^{\pi_{\boldsymbol{\theta}}}(S_t, A_t) - b(S_t)) \right] \\
\nabla_{\boldsymbol{\theta}} f_0(\boldsymbol{\theta}) &= \mathbb{E}_{S_t, A_t \sim P_r}^{\gamma} \left[\sum_{t=0}^{\infty} \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(A_t | S_t) (Q^{\pi_{\boldsymbol{\theta}}}(S_t, A_t) - b(S_t)) \right]
\end{aligned}$$

これを用いることで、各時点で $\{s_t, a_t, r_t, s_{t+1}\}$ が得られた時,

$$\nabla_{\boldsymbol{\theta}} \log_{\boldsymbol{\theta}} \pi_{\boldsymbol{\theta}}(a_t | s_t) (\hat{Q}_t - b(s_t))$$

を計算することで勾配を推定できる。

2.1 モンテカルロ方策勾配法 (REINFORCE 法)

エピソードの履歴データ $\{s_0, a_0, r_0, \dots, s_{T-1}, a_{T-1}, r_{T-1}, s_T\}$ を用いて、 Q^{π} の推定量として

$$c_t := \sum_{k=t}^{T-1} r_k, \quad \forall t = \{0, \dots, T-1\}$$

を計算し、方策勾配を次のように更新する。

$$\boldsymbol{\theta} := \boldsymbol{\theta} + \alpha_n \frac{1}{T} \sum_{t=0}^{T-1} (c_t - b(s_t)) \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(a_t | s_t)$$

2.2 アクタークリティック方策勾配法

行動価値関数 Q_∞^π を \hat{Q}_w で近似することを考える. これにより, 価値関数と方策の両方を関数近似するアクタークリティックが得られる.

- クリティック... 近似 TD(λ) 法によって方策の価値 Q_∞^π を更新する.
 - 推定平均報酬 \hat{f} の更新

$$\hat{f}_t := \hat{f}_{t-1} + \alpha_t^{\text{critic}}$$

- TD 誤差 δ の更新

$$\delta_t := r_t - \hat{f}_t + \hat{Q}_{w_t}(s_{t+1}, a_{t+1}) - \hat{Q}_{w_t}(s_t, a_t)$$

- 関数近似器パラメータ w の更新

$$w_{t+1} := w_t + \alpha_t^{\text{critic}} \delta_t \sum_{\tau=0}^t \lambda^{t-\tau} \nabla_w \hat{Q}_{w_t}(s_t, a_t)$$

- アクター... 方策勾配法によって方策 π_θ を更新する.
 - 方策パラメータ θ の更新

$$\theta_{t+1} := \theta_t + \alpha_t^{\text{actor}} \hat{Q}_{w_t}(s_t, a_t) \nabla_\theta \log \pi_\theta(a_t | s_t)$$

アクタークリティック方策勾配法において用いる関数近似器 \hat{Q}_w について考える. \hat{Q}_w のパラメータ w を, ベクトル w_1, w_2 を用いて

$$w = \begin{pmatrix} w_1 \\ w_2 \end{pmatrix}$$

と表し, \hat{Q}_w を

$$\hat{Q}_w(s, a) := \mathbf{w}_1^\top \nabla_\theta \log \pi_\theta(a | s) + b_{w_2}(s), \quad \forall (s, a) \in \mathcal{S} \times \mathcal{A}$$

と定義する. ただし, $b_{w_2} : \mathcal{S} \rightarrow \mathbb{R}$ は w_2 で規定される任意の状態関数.

方策勾配の適合関数

- $M(\pi_\theta)$ がエルゴート性を満たす場合,

$$\frac{\partial}{\partial \mathbf{w}_1} \mathbb{E}_{S_0 \sim p_{S_0}, A_t | S_t \sim \pi, S_{t+1} | A_t, S_t \sim p_T} \left[\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} (Q_\infty^\pi(S_t, A_t) - \hat{Q}_w(S_t, A_t))^2 \right] \Bigg|_{\mathbf{w}=\mathbf{w}_\infty} = \mathbf{0}$$

を満たすパラメータ \mathbf{w}_∞^* を持つ関数近似器 $\hat{Q}_{\mathbf{w}_\infty^*}$ は次のように方策勾配に適合している.

$$\nabla_\theta f_\infty(\theta) = \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} p_\infty^{\pi_\theta}(s) \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s) \hat{Q}_{\mathbf{w}_\infty^*}(s, a)$$

- $\hat{Q}_{\mathbf{w}_0}$ が

$$\frac{\partial}{\partial \mathbf{w}_1} \mathbb{E}_{S_t, A_t \sim P_T}^\gamma = \left[\sum_{t=0}^{\infty} (Q^\pi(S_t, A_t) - \hat{Q}_w(S_t, A_t))^2 \right] \Bigg|_{\mathbf{w}=\mathbf{w}_0^*} = \mathbf{0}$$

を満たすパラメータ \mathbf{w}_0^* を持つとする. この時, 関数近似器 $\hat{Q}_{\mathbf{w}_0^*}$ は次のように方策勾配に適合している.

$$\nabla_\theta f_0(\theta) = \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d^{\pi_\theta}(s) \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s) \hat{Q}_{\mathbf{w}_0^*}(s, a)$$

(証明)

$M(\pi_\theta)$ はエルゴート性を満たすので, 時間平均は空間平均と一致するので

$$\begin{aligned} & \mathbb{E}_{S_0 \sim p_{S_0}, A_t | S_t \sim \pi, S_{t+1} | S_t, A_t \sim p_T} \left[\lim_{T \rightarrow \infty} \frac{1}{T} \sum_{t=0}^{T-1} (Q_\infty^\pi(S_t, A_t) - \hat{Q}_w(S_t, A_t))^2 \right] \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} p_\infty^{\pi_\theta}(s) \pi_\theta(a|s) (Q_\infty^\pi(s, a) - \hat{Q}_w(s, a))^2 \end{aligned}$$

が成り立つ。今、右辺を w_1 で偏微分して $\mathbf{0}$ となるような $w = w_\infty^*$ が存在するので、

$$\begin{aligned} & \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} p_\infty^{\pi_\theta}(s) \pi_\theta(a|s) (Q_\infty^\pi(s, a) - \hat{Q}_{w_\infty^*}(s, a)) \nabla_{w_1} \hat{Q}_w(s, a) \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} p_\infty^{\pi_\theta}(s) \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s) (Q_\infty^\pi(s, a) - \hat{Q}_{w_\infty^*}(s, a)) = \mathbf{0} \end{aligned}$$

が成り立つ。以上より、

$$\begin{aligned} \nabla_\theta f_\infty(\theta) &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} p_\infty^{\pi_\theta}(s) \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s) Q_\infty^\pi(s, a) \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} p_\infty^{\pi_\theta}(s) \pi_\theta(a|s) \nabla_\theta \log \pi_\theta(a|s) \hat{Q}_{w_\infty^*}(s, a) \end{aligned}$$

を得る。

エルゴート性を満たさない場合についても、

$$\begin{aligned} & \mathbb{E}_{S_t, A_t \sim P_r}^\gamma \left[\sum_{t=0}^{\infty} (Q^\pi(S_t, A_t) - \hat{Q}_w(S_t, A_t))^2 \right] \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \sum_{t=0}^{\infty} \gamma^t \Pr[S_t = s, A_t = a] (Q^\pi(s, a) - \hat{Q}_w(s, a))^2 \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} \sum_{t=0}^{\infty} \gamma^t \Pr[S_t = s] \pi(a|s) (Q^\pi(s, a) - \hat{Q}_w(s, a))^2 \\ &= \sum_{s \in \mathcal{S}} \sum_{a \in \mathcal{A}} d^{\pi_\theta}(s) \pi(a|s) (Q^\pi(s, a) - \hat{Q}_w(s, a))^2 \end{aligned}$$

とすれば、同様に証明できる。

(お気持ち)

アクタークリティックの際、関数 \hat{Q}_w は線形でも十分に行動価値関数 $Q_\infty^{\pi_\theta}$ を表現しうる。また、ベースライン関数 $b_{w_2}(s)$ は任意なので用いなくてもよく、この場合は $\nabla_\theta \log \pi_\theta(a|s)$ の張る空間で行動価値関数を近似すれば十分であ

ることがわかる。ただし、任意の w_1 に対し、

$$\begin{aligned}
\hat{V}_w(s) &= \sum_{a \in \mathcal{A}} \pi_\theta(a|s) \hat{Q}_w(s, a) \\
&= \sum_{a \in \mathcal{A}} \pi_\theta(a|s)^\top w_1 \nabla_\theta \log \pi_\theta(a|s) \\
&= w_1^\top \sum_{a \in \mathcal{A}} \nabla_\theta \pi_\theta(a|s) \\
&= w_1^\top \nabla_\theta \sum_{a \in \mathcal{A}} \pi_\theta(a|s) \\
&= w_1^\top \nabla_\theta 1 = 0
\end{aligned}$$

が成り立ち、 V_∞^π に対する表現能力を一切持たなくなってしまう、方策勾配法を \hat{V}_w で置き換えることができない。そこで、特徴ベクトル $\phi(s)$ を追加した

$$\varphi(s, a) := \begin{pmatrix} \nabla_\theta \log \pi_\theta(a|s) \\ \phi(s) \end{pmatrix}$$

を用いて

$$\hat{Q}_w(s, a) = w^\top \varphi(s, a)$$

とすることが考えられる。

状態価値関数 $V_\infty^{\pi_\theta}$ の TD 誤差を

$$\delta_t := g(s_t, a_t) - f_\infty(\theta) + V_\infty^{\pi_\theta}(s_{t+1}) - V_\infty^{\pi_\theta}(s_t)$$

と定義する。今、差分価値関数 $Q_\infty^{\pi_\theta}$ についての再帰式：

$$\begin{aligned}
Q_\infty^{\pi_\theta}(s, a) &= g(s, a) - f_\infty(\theta) + \sum_{s' \in \mathcal{S}} p_T(s'|s, a) \sum_{a' \in \mathcal{A}} \pi_\theta(a'|s') Q_\infty^{\pi_\theta}(s', a') \\
&= g(s, a) - f_\infty(\theta) + \sum_{s' \in \mathcal{S}} p_T(s'|s, a) V_\infty^{\pi_\theta}(s')
\end{aligned}$$

より、TD 誤差の前半 3 項の期待値を取ると、

$$\begin{aligned}
&\mathbb{E}_{S_{t+1}|S_t, A_t \sim p_T, A_{t+1}|S_{t+1} \sim \pi} [g(S_t, A_t) - f_\infty(\theta) + V_\infty^{\pi_\theta}(S_{t+1}) | S_t = s_t, A_t = a_t] \\
&= \mathbb{E}_{A_t|S_t \sim \pi, S_{t+1}|S_t, A_t \sim p_T} [g(S_t, A_t) - f_\infty(\theta) + Q_\infty^{\pi_\theta}(S_{t+1}) | S_t = s_t, A_t = a_t] \\
&= g(s_t, a_t) - f_\infty(\theta) + \sum_{s' \in \mathcal{S}} p_T(s'|s_t, a_t) \sum_{a' \in \mathcal{A}} \pi_\theta(a'|s') Q_\infty^{\pi_\theta}(s', a') \\
&= Q_\infty^{\pi_\theta}(s_t, a_t)
\end{aligned}$$

となるので,

$$\mathbb{E}[\delta_t | S_t = s, A_t = a] = Q_{\infty}^{\pi_{\theta}}(s_t, a_t) - V_{\infty}^{\pi_{\theta}}$$

である. 従って, t 時点で $\{s_t, a_t, r_t, s_{t+1}\}$ が得られた時, $V_{\infty}^{\pi_{\theta}}$ を $\hat{V}_{\mathbf{w}}$ で推定し, ベースライン関数を $\hat{V}_{\mathbf{w}}$ に設定して,

$$\nabla_{\theta} \log \pi_{\theta}(a_t | s_t) (r_t - \hat{f}_t + \hat{V}_{\mathbf{w}}(s_{t+1}) - \hat{V}_{\mathbf{w}}(s_t))$$

を計算することで勾配を推定できる.