

Conservative Policy Iteration まとめ

北村俊徳

May 9, 2023

目次

- ① はじめに
- ② Conservative Policy Iteration
- ③ Safe Policy Iteration
- ④ Trust Region Policy Iteration
- ⑤ Constrained Policy Optimization

このスライドで説明する内容

方策更新で単調な性能の向上を保証するアルゴリズムの解説

- ① Conservative Policy Iteration [Kakade and Langford, 2002]
- ② Safe Policy Iteration [Pirootta et al., 2013]
- ③ Trust Region Policy Optimization [Schulman et al., 2015]
- ④ Proximal Policy Optimization [Schulman et al., 2017]
- ⑤ Constrained Policy Optimization [Achiam et al., 2017]

を解説するよ

- $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma, \mu \rangle$:
状態空間, 行動空間, 遷移確率, 方策関数, 割引率, 初期分布
(煩雑なので初期分布の表記は必要な箇所以外省略する)
- ベルマン方程式:

$$V^\pi(s) = \sum_{a \in \mathcal{A}} \pi(a | s) \left(\mathcal{R}(s, a) + \gamma \sum_{s' \in \mathcal{S}} \mathcal{P}(s' | s, a) V^\pi(s') \right) \quad (1)$$

- 目的関数:

$$J^\pi = \sum_{s \in \mathcal{S}} \mu(s) V^\pi(s) = \sum_{s \in \mathcal{S}} d^\pi(s) \sum_{a \in \mathcal{A}} \pi(a | s) \mathcal{R}(s, a) \quad (2)$$

ここで $d^\pi = \sum_{t=0}^{\infty} \gamma^t \Pr(s_t = s | \pi, \mu)$ は初期分布 μ から始まって π に従って動いた時の正規化されてない割引定常分布¹.

- 方策 π に従って動いたときの期待値: $\mathbb{E}_{(a_t, s_t) \sim \pi}$
- 方策 π のアドバンテージ:

$$A^\pi(s, a) = Q^\pi(s, a) - V^\pi(s) \quad (3)$$

- π のアドバンテージの π' についての期待値

$$A_{\pi}^{\pi'}(s) = \sum_{a \in \mathcal{A}} \pi'(a | s) A^\pi(s, a) \quad (4)$$

- $A_{\pi}^{\pi'}(s)$ の定常分布 d^π についての期待値

$$\mathbb{A}_{\pi, d^\pi}^{\pi'} = \sum_{s \in \mathcal{S}} d^\pi(s) A_{\pi}^{\pi'}(s) \quad (5)$$

¹正規化の有無で式展開が結構変わる。注意しよう。

目次

- ① はじめに
- ② Conservative Policy Iteration
- ③ Safe Policy Iteration
- ④ Trust Region Policy Iteration
- ⑤ Constrained Policy Optimization

基本的なアイデア

方策を $\pi' = \alpha \bar{\pi} + (1 - \alpha)\pi$ の形で更新すると、目的関数について

$$J^{\pi'} - J^{\pi} \geq \alpha \mathbb{A}_{\pi, d^{\pi}}^{\bar{\pi}} - \frac{2\alpha^2\gamma\epsilon}{(1 - \gamma(1 - \alpha))} \quad (6)$$

の不等式が成立する. ここで, $\epsilon := \frac{1}{(1-\gamma)} (\max_s A_{\pi}^{\bar{\pi}}(s))$
この左辺が常に 0 以上 & 最大になるように α を選びたい.

式 (6) の導出 I

まず, $J^{\pi'} - J^{\pi}$ をアドバンテージで表そう.

任意の方策 π' と π , π' に従って動いたときの定常分布 $d^{\pi'}$ について,

$$J^{\pi'} - J^{\pi} = \mathbb{A}_{\pi, d^{\pi'}}^{\pi'} \quad (7)$$

が成立する.

式 (6) の導出 II

式 (7) の導出

$$\begin{aligned} V^{\pi'}(s) &= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{(a_t, s_t) \sim \pi'} [\mathcal{R}(s_t, a_t)] \\ &= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{(a_t, s_t) \sim \pi'} [\mathcal{R}(s_t, a_t) + V^{\pi}(s_t) - V^{\pi}(s_t)] \\ &= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{(a_t, s_t, s_{t+1}) \sim \pi'} [\mathcal{R}(s_t, a_t) + \gamma V^{\pi}(s_{t+1}) - V^{\pi}(s_t)] + V^{\pi}(s) \\ &= V^{\pi}(s) + \mathbb{A}_{\pi, d^{\pi'}}^{\pi'} \end{aligned}$$

2 行目から 3 行目の変換: 時刻 t と $t+1$ で $\gamma V^{\pi}(s_{t+1})$ と $\gamma V^{\pi}(s_t)$ が消し合い, $-V^{\pi}(s)$ が出てくるため.

式 (6) の導出 III

次に, $A_{\pi, d^{\pi'}}^{\pi'}$ に関する不等式を考えよう.

式 (7) の π' を $\pi' = \alpha\bar{\pi} + (1 - \alpha)\pi$ のように 2 つの方策を混ぜたものと考えると, 強力な不等式が出てくる. この α は行動が $\bar{\pi}$ から出る確率とみなせる.

まず, アドバンテージを同じ方策について期待値を取ると 0 になる性質を利用して,

$$\begin{aligned} A_{\pi}^{\pi'}(s) &= \sum_a \pi'(s; a) A^{\pi}(s, a) \\ &= \sum_a ((1 - \alpha)\pi(a; s) + \alpha\bar{\pi}(a; s)) A^{\pi}(s, a) \\ &= \alpha \sum_a \bar{\pi}(a; s) A^{\pi}(s, a) \\ &= \alpha A_{\pi}^{\bar{\pi}}(s) \end{aligned}$$

式 (6) の導出 IV

ここで, c_t を時刻 t 以前に行動が $\bar{\pi}$ から選ばれた回数を示す確率変数とし, $\rho_t \equiv \Pr(c_t \geq 1) = 1 - (1 - \alpha)^t$ する. このとき,

$$\begin{aligned} & \mathbb{E}_{s \sim P(s_t; \pi')} [\alpha A_{\bar{\pi}}^{\bar{\pi}}(s)] \\ &= \alpha (1 - \rho_t) \mathbb{E}_{s \sim P(s_t | c_t=0; \pi')} [A_{\bar{\pi}}^{\bar{\pi}}(s)] \\ & \quad + \alpha \rho_t \mathbb{E}_{s \sim P(s_t | c_t \geq 1; \pi')} [A_{\bar{\pi}}^{\bar{\pi}}(s)] \\ & \geq \alpha \mathbb{E}_{s \sim P(s_t | c_t=0; \pi')} [A_{\bar{\pi}}^{\bar{\pi}}(s)] - 2\alpha \rho_t \max_s (A_{\bar{\pi}}^{\bar{\pi}}(s)) \\ &= \alpha \mathbb{E}_{s \sim P(s_t; \pi)} [A_{\bar{\pi}}^{\bar{\pi}}(s)] - 2\alpha \rho_t \max_s (A_{\bar{\pi}}^{\bar{\pi}}(s)) \end{aligned}$$

式 (6) の導出 V

式 (6) の導出

$$\begin{aligned} J^{\pi'} - J^{\pi} &= \mathbb{A}_{\pi, d^{\pi'}}^{\pi'} \\ &= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s \sim P(s_t; \pi')} \left[A_{\pi}^{\pi'}(s) \right] \\ &= \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s \sim P(s_t; \pi')} \left[\alpha A_{\pi}^{\bar{\pi}}(s) \right] \\ &\geq \alpha \sum_{t=0}^{\infty} \gamma^t \mathbb{E}_{s \sim P(s_t; \pi)} \left[A_{\pi}^{\bar{\pi}} \right] - 2\alpha \max_s (A_{\pi}^{\bar{\pi}}(s)) \sum_{t=0}^{\infty} \gamma^t (1 - (1 - \alpha)^t) \\ &= \alpha \left(\mathbb{A}_{\pi, d^{\pi}}^{\bar{\pi}} - \frac{2\alpha \gamma \max_s (A_{\pi}^{\bar{\pi}}(s))}{(1 - \gamma)(1 - \gamma(1 - \alpha))} \right) \end{aligned}$$

最適な α の導出 (元論文)

報酬の最大値を R とすると, $\max_s |A_{\pi}^{\bar{\pi}}(s)| \leq R$ が成立する².
また, $0 \leq \alpha \leq 1$ なので,

$$\alpha \mathbb{A}_{\pi, d^{\pi}}^{\bar{\pi}} - \frac{2\alpha^2 \gamma \epsilon}{(1 - \gamma(1 - \alpha))} \geq \alpha \mathbb{A}_{\pi, d^{\pi}}^{\bar{\pi}} - \frac{2\alpha^2 R}{(1 - \gamma)^2}$$

右辺は α についての二次式になっているので, $\mathbb{A}_{\pi, d^{\pi}}^{\bar{\pi}} \geq 0$ である場合,

$$\alpha = \frac{\mathbb{A}_{\pi, d^{\pi}}^{\bar{\pi}} (1 - \gamma)^2}{4R} \quad (8)$$

で下界の最大値 $\frac{\mathbb{A}_{\pi, d^{\pi}}^{\bar{\pi}^2} (1 - \gamma)^2}{8R}$ を得る.

²元論文ではこの記述なのですが, どう頑張ってもこの不等式の証明ができない...
 $\max_s |A_{\pi}^{\bar{\pi}}(s)| \leq \frac{R}{1 - \gamma}$ になると思うんですよね. 実際 SPI で出てくる CPI の記述はこれになってない気がする.

最適な α の導出 (後で出てくる SPI で言及されてる版)

報酬の絶対値の最大値を R とすると, $\max_s |A_{\pi}^{\bar{\pi}}(s)| \leq \frac{2R}{1-\gamma}$ が成立する³.
また, $0 \leq \alpha \leq 1$ なので,

$$\alpha \mathbb{A}_{\pi, d^{\pi}}^{\bar{\pi}} - \frac{2\alpha^2 \gamma \epsilon}{(1 - \gamma(1 - \alpha))} \geq \alpha \mathbb{A}_{\pi, d^{\pi}}^{\bar{\pi}} - \frac{2\alpha^2 \gamma \epsilon}{1 - \gamma} \geq \alpha \mathbb{A}_{\pi, d^{\pi}}^{\bar{\pi}} - \frac{4\alpha^2 \gamma R}{(1 - \gamma)^3}$$

右辺は α についての二次式になっているので, $\mathbb{A}_{\pi, d^{\pi}}^{\bar{\pi}} \geq 0$ である場合,

$$\alpha = \frac{\mathbb{A}_{\pi, d^{\pi}}^{\bar{\pi}} (1 - \gamma)^3}{8\gamma R} \quad (9)$$

で下界の最大値 $\frac{\mathbb{A}_{\pi, d^{\pi}}^{\bar{\pi}^2} (1 - \gamma)^3}{16\gamma R}$ を得る.

³報酬の最小値が 0 なら分子の 2 は不要. SPI の論文では報酬の最小値が 0 の時で比較している.

目次

- ① はじめに
- ② Conservative Policy Iteration
- ③ Safe Policy Iteration
- ④ Trust Region Policy Iteration
- ⑤ Constrained Policy Optimization

SPI は行列を使うと証明しやすいので, 行列メインの表記:

- Π^π : $\Pi^\pi(s, (s, a)) = \pi(a \mid s)$ なる $(|\mathcal{S}| \times |\mathcal{S}||\mathcal{A}|)$ の行列.
- \mathbf{P} : $\mathbf{P}((s, a), s') = P(s' \mid s, a)$ なる $(|\mathcal{S}||\mathcal{A}| \times |\mathcal{S}|)$ の行列.
- $\mathbf{P}^\pi = \Pi^\pi \mathbf{P}$: 方策 π に従う時の状態から状態への遷移行列. サイズ $|\mathcal{S}| \times |\mathcal{S}|$.
- $\mathbf{v}^\pi, \mathbf{r}^\pi, d^\pi, A_{\pi}^{\bar{\pi}}$: サイズ $|\mathcal{S}|$ のベクトル.

基本的なアイデア

CPI で出てきた

$$J^{\pi'} - J^{\pi} = \mathbb{A}_{\pi, d^{\pi'}}^{\pi'}$$

について, CPI とは違う下界:

$$J^{\pi'} - J^{\pi} \geq \mathbb{A}_{\pi, d^{\pi}}^{\pi'} - \frac{\gamma}{(1 - \gamma)^2} \left\| \Pi^{\pi'} - \Pi^{\pi} \right\|_{\infty} \frac{\Delta \mathbf{A}_{\pi}^{\pi'}}{2} \quad (10)$$

を最大化することを考える.

式 (10) の導出 I

次の不等式 [Haviv and Heyden, 1984] を使って下界を導く:

$$|\mathbf{c}^T \mathbf{d}| \leq \|\mathbf{c}\|_1 \frac{\Delta \mathbf{d}}{2} \quad (11)$$

ここで, \mathbf{d}, \mathbf{c} は任意の $\mathbf{c}^T \mathbf{e} = 0$ を満たすようなベクトルであり,
 $\Delta \mathbf{d} = \max_{i,j} |\mathbf{d}_i - \mathbf{d}_j|$.

これを使うと:

$$\begin{aligned} J^{\pi'} - J^{\pi} &= \mathbb{A}_{\pi, \mathbf{d}^{\pi'}}^{\pi'} \\ &= \mathbf{d}^{\pi T} \mathbf{A}_{\pi}^{\pi'} + \left(\mathbf{d}^{\pi'} - \mathbf{d}^{\pi} \right)^T \mathbf{A}_{\pi}^{\pi'} \\ &\geq \mathbf{d}^{\pi T} \mathbf{A}_{\pi}^{\pi'} - \left\| \mathbf{d}^{\pi'} - \mathbf{d}^{\pi} \right\|_1 \frac{\Delta \mathbf{A}_{\pi}^{\pi'}}{2} \end{aligned} \quad (12)$$

これは $\mathbf{d}^{\pi'} - \mathbf{d}^{\pi}$ が平均 0 のベクトルのため.

式 (10) の導出 II

さらに, $\mathbf{d}^{\pi'}$ と \mathbf{d}^{π} について, 以下の不等式が成立する:

$$\left\| \mathbf{d}^{\pi'} - \mathbf{d}^{\pi} \right\|_1 \leq \frac{\gamma}{1 - \gamma} \left\| \mathbf{P}^{\pi'} - \mathbf{P}^{\pi} \right\|_{\infty} \left\| \left(\mathbf{I} - \gamma \mathbf{P}^{\pi'} \right)^{-1} \right\|_{\infty} \quad (13)$$

上の不等式を導くため, まず $\mathbf{d}^{\pi'} - \mathbf{d}^{\pi}$ を変形する⁴

$$\begin{aligned} \mathbf{d}^{\pi'} - \mathbf{d}^{\pi} &= \gamma \mathbf{d}^{\pi'} \mathbf{P}^{\pi'} - \gamma \mathbf{d}^{\pi \top} \mathbf{P}^{\pi} \\ &= \gamma \left(\mathbf{d}^{\pi'} - \mathbf{d}^{\pi} \right) \mathbf{P}^{\pi'} + \gamma \mathbf{d}^{\pi'} \left(\mathbf{P}^{\pi'} - \mathbf{P}^{\pi} \right) \\ &= \gamma \mathbf{d}^{\pi \top} \left(\mathbf{P}^{\pi'} - \mathbf{P}^{\pi} \right) \left(\mathbf{I} - \gamma \mathbf{P}^{\pi'} \right)^{-1} \\ &\leq \gamma \left\| \mathbf{d}^{\pi \top} \right\|_{\infty} \left\| \mathbf{P}^{\pi'} - \mathbf{P}^{\pi} \right\|_{\infty} \left\| \left(\mathbf{I} - \gamma \mathbf{P}^{\pi'} \right)^{-1} \right\|_{\infty} \\ &= \frac{\gamma}{1 - \gamma} \left\| \mathbf{P}^{\pi'} - \mathbf{P}^{\pi} \right\|_{\infty} \left\| \left(\mathbf{I} - \gamma \mathbf{P}^{\pi'} \right)^{-1} \right\|_{\infty} \end{aligned}$$

式 (10) の導出 III

目的の $\left\| \mathbf{d}_\mu^{\pi'} - \mathbf{d}_\mu^\pi \right\|_1$ について変形しなおして,

$$\begin{aligned} \left\| \mathbf{d}_\mu^{\pi'} - \mathbf{d}_\mu^\pi \right\|_1 &\leq \frac{\gamma}{1-\gamma} \left\| \mathbf{P}^{\pi'} - \mathbf{P}^\pi \right\|_\infty \left\| \left(\mathbf{I} - \gamma \mathbf{P}^{\pi'} \right)^{-1} \right\|_\infty \\ &\leq \frac{\gamma}{1-\gamma} \left\| \mathbf{\Pi}^{\pi'} - \mathbf{\Pi}^\pi \right\|_\infty \left\| \mathbf{P} \right\|_\infty \sum_{t=0}^{\infty} \gamma^t \left\| \mathbf{P}^{\pi'} \right\|_\infty^t \\ &= \frac{\gamma}{(1-\gamma)^2} \left\| \mathbf{\Pi}^{\pi'} - \mathbf{\Pi}^\pi \right\|_\infty \end{aligned} \tag{14}$$

ここで, $\left\| \mathbf{P} \right\|_\infty = 1$ を使ってる.
あとは式 (12) に代入して終わり.

⁴導出では $\sum_{t=0}^{\infty} (\gamma \mathbf{P}^{\pi'})^t$ がノイマン級数であるため, $\sum_{t=0}^{\infty} (\gamma \mathbf{P}^\pi)^t = \left(\mathbf{I} - \gamma \mathbf{P}^{\pi'} \right)^{-1}$ が成立することを利用.

最適な α の導出

$$J^{\pi'} - J^{\pi} \geq \mathbb{A}_{\pi, d^{\pi}}^{\pi'} - \frac{\gamma}{(1-\gamma)^2} \left\| \Pi^{\pi'} - \Pi^{\pi} \right\|_{\infty} \frac{\Delta \mathbf{A}_{\pi}^{\pi'}}{2}$$

について, π' に $\alpha \bar{\pi} + (1-\alpha)\pi$ を代入すると,

$$J^{\pi'} - J^{\pi} \geq \alpha \mathbb{A}_{\pi, d^{\pi}}^{\bar{\pi}} - \alpha^2 \frac{\gamma}{(1-\gamma)^2} \left\| \Pi^{\bar{\pi}} - \Pi^{\pi} \right\|_{\infty} \frac{\Delta \mathbf{A}_{\pi}^{\bar{\pi}}}{2} \quad (15)$$

なる不等式が出てくる. 右辺は $\alpha = \frac{(1-\gamma)^2 \mathbb{A}_{\pi, d^{\pi}}^{\bar{\pi}}}{\gamma \left\| \Pi^{\bar{\pi}} - \Pi^{\pi} \right\|_{\infty} \Delta \mathbf{A}_{\pi}^{\bar{\pi}}}$ の時最大になり, この時の下界は

$$J^{\pi'} - J^{\pi} \geq \frac{(1-\gamma)^2 \mathbb{A}_{\pi, d^{\pi}}^{\bar{\pi}^2}}{2\gamma \left\| \Pi^{\bar{\pi}} - \Pi^{\pi} \right\|_{\infty} \Delta \mathbf{A}_{\pi}^{\bar{\pi}}} \quad (16)$$

CPI との比較

SPI の下界について, $\Delta \mathbf{A}_{\pi}^{\bar{\pi}} < \frac{2R}{1-\gamma}$, $\|\Pi^{\pi} - \Pi^{\bar{\pi}}\|_{\infty} < 2$ を使って CPI と比較すると:

- CPI の下界⁵: $J^{\pi'} - J^{\pi} \geq \frac{(1-\gamma)^3 \mathbb{A}_{\pi, d^{\pi}}^{\bar{\pi}}{}^2}{16\gamma R}$
- SPI の下界: $J^{\pi'} - J^{\pi} \geq \frac{(1-\gamma)^2 \mathbb{A}_{\pi, d^{\pi}}^{\bar{\pi}}{}^2}{2\gamma \|\Pi^{\pi} - \Pi^{\bar{\pi}}\|_{\infty} \Delta \mathbf{A}_{\pi}^{\bar{\pi}}} \geq \frac{(1-\gamma)^3 \mathbb{A}_{\pi, d^{\pi}}^{\bar{\pi}}{}^2}{8\gamma R}$

よって SPI は CPI よりマシな下界になることがわかる。

⁵SPI の論文は報酬の最小値を 0 と仮定している. その場合分母の 16 は 8 になる。

目次

- ① はじめに
- ② Conservative Policy Iteration
- ③ Safe Policy Iteration
- ④ Trust Region Policy Iteration
- ⑤ Constrained Policy Optimization

基本的なアイデア

CPI の更新は α を使って方策同士を混ぜる: 扱いづらい...

→ TRPO はもっと汎用的な形の下界に: 以下の不等式を利用
任意の方策 π' と π について,

$$\begin{aligned} J^{\pi'} - J^{\pi} &\geq \mathbb{A}_{\pi, d^{\pi}}^{\pi'} - \frac{4\epsilon\gamma}{(1-\gamma)} D_{TV}^{\max}(\pi, \pi')^2 \\ &\geq \mathbb{A}_{\pi, d^{\pi}}^{\pi'} - \frac{4\epsilon\gamma}{(1-\gamma)} D_{KL}^{\max}(\pi, \pi') \end{aligned} \quad (17)$$

ここで, $\epsilon = \frac{1}{(1-\gamma)} \max_s \max_a |A^{\pi}(s, a)|$.

$D_{KL}^{\max}(\pi, \pi') = \max_s D_{KL}(\pi(s, \cdot) \| \pi'(s, \cdot))$,

$D_{TV}^{\max}(\pi, \pi') = \max_s D_{TV}(\pi(s, \cdot) \| \pi'(s, \cdot))$

式 (17) より, π と π' の KL ダイバージェンスを小さくしながら $\mathbb{A}_{\pi, d^{\pi}}^{\pi'}$ を最大化すると良さそう.

式 (17) の導出 I

以下, 不等式 (17) の導出を行う⁶.
導出の前に, 以下の概念を導入する.

α -coupled

2 つの方策 π と π' について, $P(a \neq a' \mid s) \leq \alpha$ が全ての状態 s について成立しているなら, 2 つの方策は α -coupled であるという.
つまり, α -coupled な方策について, α が小さいほど 2 つの方策は近くなる.

式 (17) の導出 II

α -coupled な方策についての補題 1

α -coupled な方策 π と π' で、全ての状態 s について、次が成立:

$$\left| A_{\pi'}^{\pi'}(s) \right| \leq 2\alpha \max_{s,a} |A_{\pi}(s, a)| \quad (18)$$

証明:

$$\begin{aligned} A_{\pi'}^{\pi'}(s) &= \mathbb{E}_{a' \sim \pi'} [A_{\pi}(s, a')] = \mathbb{E}_{(a,a') \sim (\pi, \pi')} [A_{\pi}(s, a') - A_{\pi}(s, a)] \\ &= P(a \neq a' \mid s) \mathbb{E}_{(a,a') \sim (\pi, \pi') \mid a \neq a'} [A_{\pi}(s, a') - A_{\pi}(s, a)] \end{aligned} \quad (19)$$

$$|A_{\pi'}^{\pi'}(s)| \leq \alpha \cdot 2 \max_{s,a} |A_{\pi}(s, a)|$$

1 行目: 方策について期待値取るとアドバンテージが 0 になる.

2 行目: $a = a'$ のときはそもそも期待値の中身が 0 になる.

式 (17) の導出 III

α -coupled な方策についての補題 2

α -coupled な方策 π と π' について, 次が成立:

$$\begin{aligned} \left| \mathbb{E}_{s_t \sim \pi'} \left[A_{\pi'}^{\pi'}(s_t) \right] - \mathbb{E}_{s_t \sim \pi} \left[A_{\pi'}^{\pi'}(s_t) \right] \right| &\leq 2\alpha \max_s A_{\pi'}^{\pi'}(s) \\ &\leq 4\alpha (1 - (1 - \alpha)^t) \max_s |A_{\pi}(s, a)| \end{aligned} \quad (20)$$

証明: 変数 n_t を時刻 $i < t$ で $a_i \neq a'_i$ が生じた回数とする. このとき, 方策についての期待値を n_t が 0 である場合と 0 以上の場合に分けると,

$$\begin{aligned} &\mathbb{E}_{s_t \sim \pi'} \left[A_{\pi'}^{\pi'}(s_t) \right] \\ &= P(n_t = 0) \mathbb{E}_{s_t \sim \pi' | n_t = 0} \left[A_{\pi'}^{\pi'}(s_t) \right] + P(n_t > 0) \mathbb{E}_{s_t \sim \pi' | n_t > 0} \left[A_{\pi'}^{\pi'}(s_t) \right] \\ &\mathbb{E}_{s_t \sim \pi} \left[A_{\pi'}^{\pi'}(s_t) \right] \\ &= P(n_t = 0) \mathbb{E}_{s_t \sim \pi | n_t = 0} \left[A_{\pi'}^{\pi'}(s_t) \right] + P(n_t > 0) \mathbb{E}_{s_t \sim \pi | n_t > 0} \left[A_{\pi'}^{\pi'}(s_t) \right] \end{aligned}$$

式 (17) の導出 IV

ここで, $n_t = 0$ のときの期待値は方策 π から出てると考えて良いので,
 $\mathbb{E}_{s_t \sim \pi'} [A_{\pi'}^{\pi'}(s_t)]$ と $\mathbb{E}_{s_t \sim \pi} [A_{\pi}^{\pi'}(s_t)]$ の差分は

$$\begin{aligned} & \mathbb{E}_{s_t \sim \pi'} [A_{\pi'}^{\pi'}(s_t)] - \mathbb{E}_{s_t \sim \pi} [A_{\pi}^{\pi'}(s_t)] \\ &= P(n_t > 0) \left(\mathbb{E}_{s_t \sim \pi' | n_t > 0} [A_{\pi'}^{\pi'}(s_t)] - \mathbb{E}_{s_t \sim \pi | n_t > 0} [A_{\pi}^{\pi'}(s_t)] \right) \end{aligned} \quad (21)$$

式 (21) について,

- α の定義から, $P(n_t > 0) \leq 1 - (1 - \alpha)^t$
- 式 (20) から,

$$\left| \mathbb{E}_{s_t \sim \pi' | n_t > 0} [A_{\pi'}^{\pi'}(s_t)] - \mathbb{E}_{s_t \sim \pi | n_t > 0} [A_{\pi}^{\pi'}(s_t)] \right| \leq 4\alpha \max_{s,a} |A_{\pi}(s, a)|$$

式 (17) の導出 V

これらを合わせると,

$$\begin{aligned} & \left| \mathbb{E}_{s_t \sim \pi' | n_t > 0} \left[A_{\pi}^{\pi'}(s_t) \right] - \mathbb{E}_{s_t \sim \pi | n_t > 0} \left[A_{\pi}^{\pi'}(s_t) \right] \right| \\ & \leq 4\alpha (1 - (1 - \alpha)^t) \max_{s, a} |A_{\pi}(s, a)| \end{aligned} \quad (22)$$

以上で準備が整ったので, 元々の目標である

$$\begin{aligned} J^{\pi'} - J^{\pi} & \geq \mathbb{A}_{\pi, d^{\pi}}^{\pi'} - \frac{4\epsilon\gamma}{(1 - \gamma)} D_{TV}^{\max}(\pi, \pi')^2 \\ & \geq \mathbb{A}_{\pi, d^{\pi}}^{\pi'} - \frac{4\epsilon\gamma}{(1 - \gamma)} D_{KL}^{\max}(\pi, \pi') \end{aligned} \quad (23)$$

を導出しよう.

式 (17) の導出 VI

CPI の式 (7) で示したように, $J^{\pi'} = J^{\pi} + \mathbb{A}_{\pi, d^{\pi'}}^{\pi'}$ が成立.

一方, J^{π} は $L_{\pi}(\pi') = J^{\pi} + \mathbb{A}_{\pi, d^{\pi}}^{\pi'}$ と一次の項まで同じである.

以上を使って

$$\begin{aligned} \left| J^{\pi'} - L_{\pi}(\pi') \right| &= \sum_{t=0}^{\infty} \gamma^t \left| \mathbb{E}_{s_t \sim \pi'} \left[A_{\pi}^{\pi'}(s_t) \right] - \mathbb{E}_{s_t \sim \pi} \left[A_{\pi}^{\pi'}(s_t) \right] \right| \\ &\leq \sum_{t=0}^{\infty} \gamma^t \cdot 4\epsilon\alpha \left(1 - (1 - \alpha)^t \right) \\ &= \frac{4\alpha^2\gamma\epsilon}{(1 - \gamma)(1 - \gamma(1 - \alpha))} \leq \frac{4\alpha^2\gamma\epsilon}{(1 - \gamma)^2} \end{aligned} \tag{24}$$

ここで $\epsilon = \max_{s,a} |A_{\pi}(s, a)|$.

後は α を $\max_s D_{TV}(\pi(\cdot | s) || \pi'(\cdot | s)) \leq \alpha$ で置き換えたら完成!

⁶ここでの導出は元論文 [Schulman et al., 2015] に基づいてるが, 元論文よりコロンビア大学の IEOR 8100 のLecture7の方が簡単かも

目次

- ① はじめに
- ② Conservative Policy Iteration
- ③ Safe Policy Iteration
- ④ Trust Region Policy Iteration
- ⑤ Constrained Policy Optimization

CPO での performance bound

π と $\bar{\pi}$ についての performance bound

任意の関数 $f: S \rightarrow \mathbb{R}$ と任意の方策 $\bar{\pi}, \pi$ について,

$$\delta_f(s, a, s') \doteq R(s, a, s') + \gamma f(s') - f(s),$$

$$\epsilon_f^{\pi'} \doteq \max_s |\mathbb{E}_{a \sim \pi', s' \sim P} [\delta_f(s, a, s')]|,$$

$$L_{\pi, f}(\pi') \doteq \mathbb{E}_{s \sim d^\pi} \left[\left(\frac{\pi'(a | s)}{\pi(a | s)} - 1 \right) \delta_f(s, a, s') \right],$$

$$D_{\pi, f}^{\pm}(\pi') \doteq \frac{L_{\pi, f}(\pi')}{1 - \gamma} \pm \frac{2\gamma\epsilon_f^{\pi'}}{(1 - \gamma)^2} \mathbb{E}_{s \sim d^\pi} [D_{TV}(\pi' || \pi) [s]]$$

を定義する.

π と $\bar{\pi}$ についての performance bound (続き)

このとき,

$$D_{\pi, f}^+(\pi') \geq L_{\pi, f}(\pi') - \epsilon_f^{\pi'} \geq D_{\pi, f}^-(\pi')$$

(25)

References



Achiam, J., Held, D., Tamar, A., and Abbeel, P. (2017).
Constrained policy optimization.
In ICML.



Haviv, M. and Heyden, L. V. D. (1984).
Perturbation bounds for the stationary probabilities of a finite markov chain.
Advances in Applied Probability, 16(4):804–818.



Kakade, S. and Langford, J. (2002).
Approximately optimal approximate reinforcement learning.
In Proceedings of the Nineteenth International Conference on Machine Learning, ICML '02, page 267â 274, San Francisco, CA, USA.
Morgan Kaufmann Publishers Inc.



Pirotta, M., Restelli, M., Pecorino, A., and Calandriello, D. (2013).
Safe policy iteration.
In Dasgupta, S. and McAllester, D., editors, Proceedings of the 30th International Conference on Machine Learning, volume 28 of