

強化学習入門

田中章詞

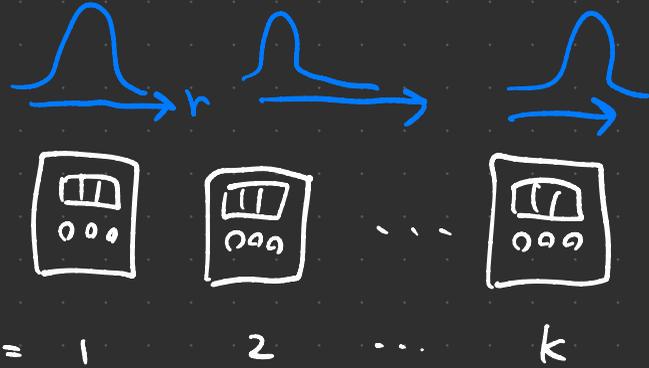
① 1. 状態がない場合

② 2. 状態がある場合

1. 状態のない場合

定義

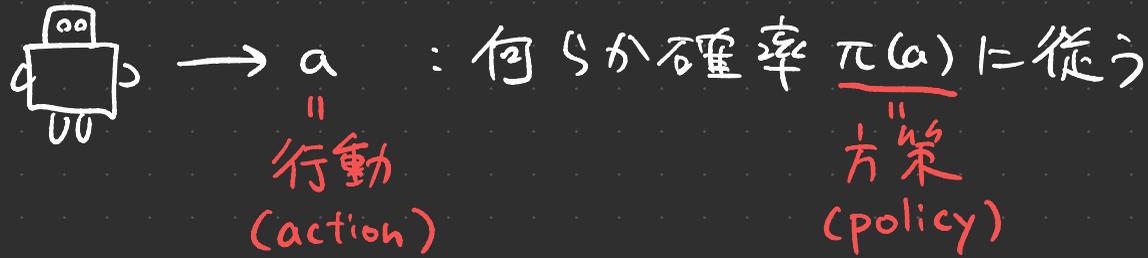
k のスロットマシン



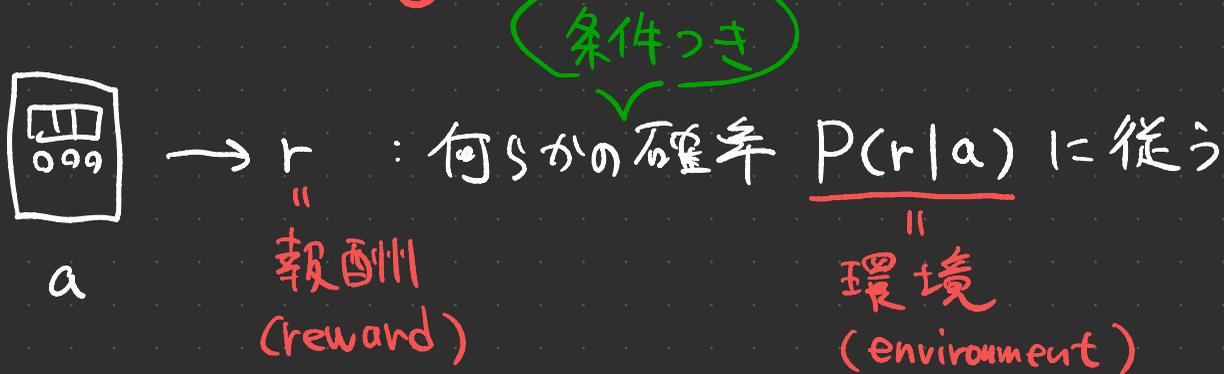
1. a を選ぶ ← ①
 2. ある $r \in \mathbb{R}$ が返される
 3. $r > 0 \Rightarrow r$ 円もらえる
 $r < 0 \Rightarrow |r|A$ L 減る
- } ②

どの a を選ぶべきか? ← ③

① 行動と方策



② 報酬と環境



③ 価値



$$\text{目的} \approx \max_{\pi} J(\pi)$$

How?

□ 学習の方針

$J(\pi)$ の値はいろいろな理由で

アクセスするのが難しい

学習

$$\begin{array}{ccccccccc} \pi_0 & \rightarrow & \pi_1 & \rightarrow & \pi_2 & \rightarrow & \dots & \rightarrow & \pi_* \\ \downarrow & & \downarrow & & \downarrow & & \downarrow & & \downarrow \\ J(\pi_0) & \leq & J(\pi_1) & \leq & J(\pi_2) & \leq & \dots & \leq & J(\pi_*) \end{array}$$

↑
= うなごほしい

方策改善 (policy improvement)

$$\pi \rightarrow \pi'$$

$$J(\pi') - J(\pi) = E_{a \sim \pi'} [Q(a) - J(\pi)]$$

"アドバンテージ関数" $A_\pi(a)$

$$\text{方策改善} \Leftrightarrow E_{a \sim \pi'} [A_\pi(a)] \geq 0$$

□ 価値 Q - π 手法

ε -貪欲方針 (ε -greedy policy)

$$\varepsilon \in [0, 1]$$

$$\pi^\varepsilon(a) = (1-\varepsilon) \mathbb{1}_{a=\arg\max_{\tilde{a}} Q(\tilde{a})} + \varepsilon \frac{1}{K}$$

$$\dots \rightarrow \pi^\varepsilon \rightarrow \pi^{\varepsilon'} \rightarrow \dots$$

直感: $\varepsilon \geq \varepsilon' \Rightarrow$ 証明

$$A_{\pi^\varepsilon}(a) = Q(a) - J(\pi)$$

$$= Q(a) - \sum_{\tilde{a}} \pi^\varepsilon(\tilde{a}) Q(\tilde{a})$$

$$\mathbb{E}_{a \sim \pi^{\varepsilon'}} \downarrow$$

$$(1 - \varepsilon') \max_a Q(a) - \sum_{\tilde{a}} \left(\pi^\varepsilon(\tilde{a}) - \varepsilon' \frac{1}{k} \right) Q(\tilde{a}) \geq 0$$

⇒ なるには?
↓

$$\sum_{\tilde{a}} \underbrace{\left(\pi^\varepsilon(\tilde{a}) - \varepsilon' \frac{1}{k} \right)}_{\geq 0} \underbrace{\left(\max_a Q(a) - Q(\tilde{a}) \right)}_{\geq 0}$$

$$(1 - \varepsilon) \frac{1}{k} + \frac{\varepsilon}{k} - \frac{\varepsilon'}{k} \geq 0$$

□ 方策ベースの手法

パラメトリックな方策

$\theta \in \mathbb{R}^k$ でパラメータ付けられる π_θ

例

$(\theta_1, \theta_2, \dots, \theta_k) \in \mathbb{R}^k$

$$\pi_\theta(a) = \frac{e^{\theta a}}{\sum_{a'} e^{\theta a'}} \quad (\text{softmax } \frac{1}{k} \text{ 方策})$$

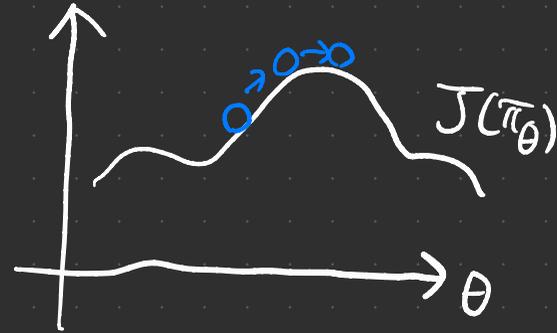
$\dots \rightarrow \pi_\theta \rightarrow \pi_{\theta'} \rightarrow \dots$ 方策改善?

つまり $J(\pi_\theta)$ の微分を θ に関する

$$\nabla_\theta J(\pi_\theta) = \nabla_\theta \underbrace{E_{a \sim \pi_\theta} [Q(a)]}_{\sum_a Q(a) \pi_\theta(a)}$$

$$= \sum_a Q(a) \underbrace{\nabla_\theta \pi_\theta(a)}_{\nabla_\theta \log \pi_\theta(a) \cdot \pi_\theta(a)}$$

$$= E_{a \sim \pi_\theta} [Q(a) \nabla_\theta \log \pi_\theta(a)]$$



$$\theta \leftarrow \theta + \eta \underbrace{\nabla_\theta J}$$

アドバンテージから再考

$$0 \leq E[A_{\pi_{\theta+\Delta\theta}}(a)] - J(\pi_{\theta})$$

こうなると
ほしい

$\pi_{\theta+\Delta\theta}$

$\Delta\theta$ でテイラー展開 $O(\Delta\theta^2)$ まで

$$= \sum_a \left(\pi_{\theta}(a) + \Delta\theta \cdot \nabla_{\theta} \pi_{\theta}(a) \right) A_{\pi_{\theta}}(a)$$

$\leftarrow E_{\pi_{\theta}}[Q(a) - E_{\pi_{\theta}}[Q(a)]] = 0$

$$= \sum_a \Delta\theta \cdot \nabla_{\theta} \pi_{\theta}(a) A_{\pi_{\theta}}(a)$$

$\nabla_{\theta} \log \pi_{\theta}(a) \cdot \pi_{\theta}(a)$

$$= \Delta\theta E_{\pi_{\theta}}[A(a) \nabla_{\theta} \log \pi_{\theta}(a)]$$

$a \sim \pi_{\theta}$
 α

$$\Delta\theta = \eta \cdot E[\dots]$$

実は

$$\nabla_{\theta} J(\pi_{\theta}) = E_{a \sim \pi_{\theta}} [Q(a) \nabla_{\theta} \log \pi_{\theta}(a)]$$

$$= E_{a \sim \pi_{\theta}} [A_{\pi_{\theta}}(a) \nabla_{\theta} \log \pi_{\theta}(a)]$$

$Q(a) - J(\pi_{\theta})$

より一般に

$$= E_{a \sim \pi_{\theta}} [(Q(a) - B) \nabla_{\theta} \log \pi_{\theta}(a)]$$

↑
B-ライン

ふつうは $J(\pi_{\theta})$ とする

□ 実際

- ϵ -greedy π^ϵ ← $Q(a) = E[r]$ の値が必要
- parametric π_θ ← $\nabla_\theta J(\pi_\theta) = E[\dots]$ の値が必要

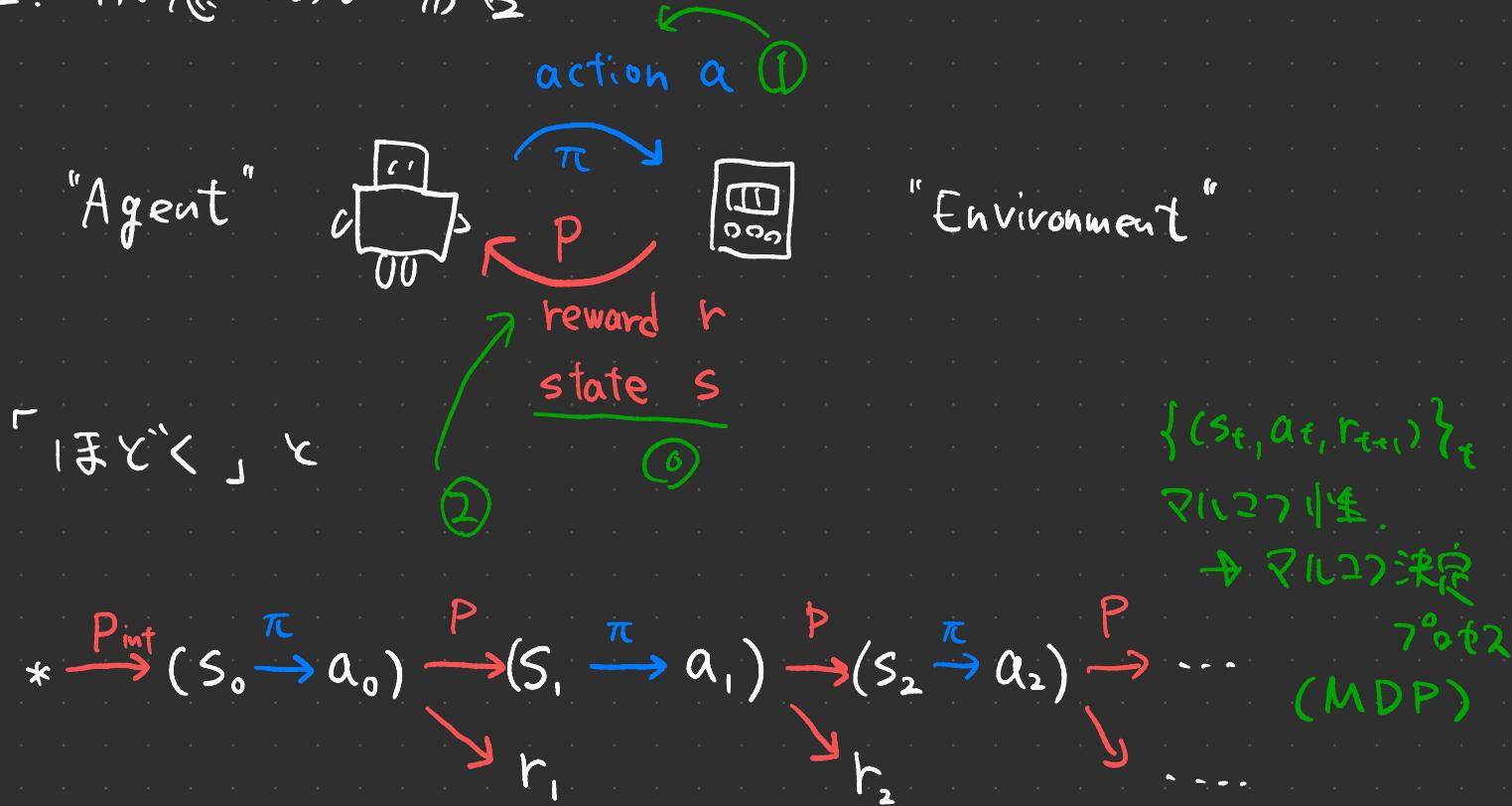


✓ $E[\dots]$ 知らず ⇒ OK.

✓ $E[\dots]$ 知らず ⇒ モンテカルロ

(次回のsection1より)

② 2. 状態のある場合



① 状態 (state)

Agent が置かれている状態を表す変数



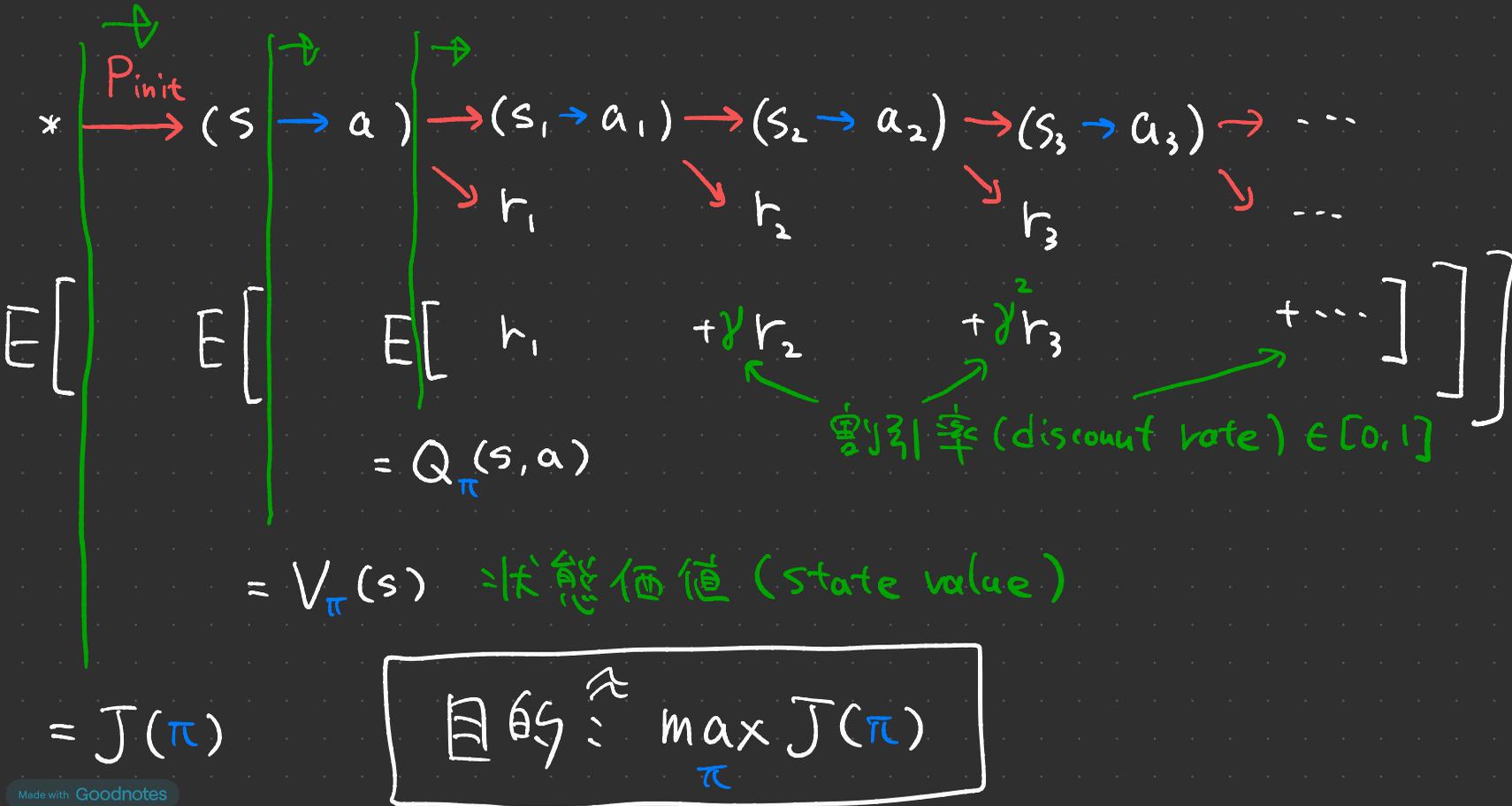
① 行動と方針

$s \xrightarrow{\pi} a$: 条件つき確率 $\pi(a|s)$

② 報酬, 状態更新

$(s, a) \rightarrow \begin{matrix} s' \\ r' \end{matrix}$: 条件つき確率 $P(s', r' | s, a)$

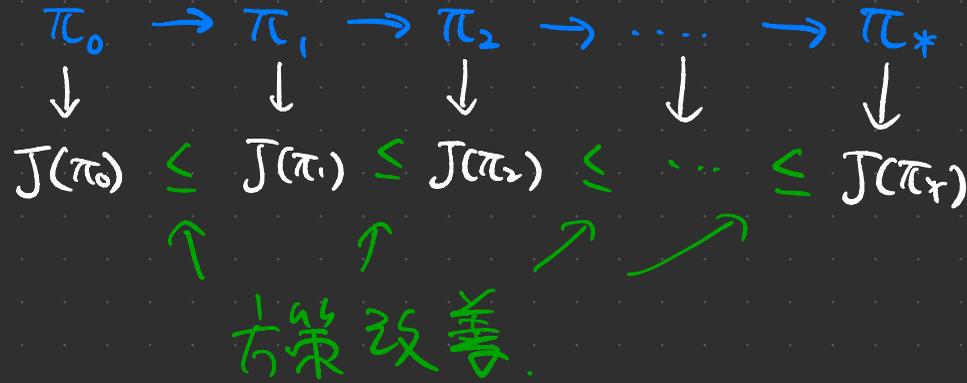
③ 價值



□ 学習の方針

$\max_{\pi} J(\pi)$ は難しいので...

学習



$$\pi \rightarrow \pi'$$

$$J(\pi') - J(\pi) = \sum_{t=0}^{\infty} \gamma^t E \left[\underbrace{Q_{\pi}(s_t, a_t) - V_{\pi}(s_t)}_{\text{アドバンテージ関数}} \right]$$

$a_0 \dots a_t \sim \pi'$
 $s_0 \dots s_t \sim P$

$A_{\pi}(s_t, a_t)$

"証明"

$$\text{方策改善} \Leftrightarrow \sum_{t=0}^{\infty} \gamma^t E \left[A_{\pi}(s_t, a_t) \right] \geq 0$$

$a_{0:t} \sim \pi'$
 $s_{0:t} \sim P$

$$J(\pi') - J(\pi) =$$

$$E_{a_0 \sim \pi'} \left[\underbrace{E_{\pi} [r_1 + \gamma r_2 + \dots]}_{Q_{\pi}(s_0, a_0)} \right]$$

$(s_0 \xrightarrow{\pi'} a_0) \rightarrow (s_1 \xrightarrow{\pi} a_1) \rightarrow (s_2 \xrightarrow{\pi} a_2) \rightarrow \dots$
 $\downarrow r_1 \quad \downarrow r_2 \quad \downarrow \dots$

$$E_{a_0 \sim \pi'} \left[\underbrace{E_{\pi} [r_1 + \gamma r_2 + \dots]}_{V_{\pi}(s_0)} \right]$$

$(s_0 \xrightarrow{\pi} a_0) \rightarrow (s_1 \xrightarrow{\pi} a_1) \rightarrow (s_2 \xrightarrow{\pi} a_2) \rightarrow \dots$
 $\downarrow r_1 \quad \downarrow r_2 \quad \downarrow \dots$

$$E_{a_0 \sim \pi'} [Q_{\pi}(s_0, a_0) - V_{\pi}(s_0)]$$

$$E_{a_0 a_1} \left[\underbrace{E_{\pi} [r_2 + \gamma r_3 + \dots]}_{Q_{\pi}(s_1, a_1)} \right]$$

$(s_0 \xrightarrow{\pi'} a_0) \rightarrow (s_1 \xrightarrow{\pi'} a_1) \rightarrow (s_2 \xrightarrow{\pi} a_2) \rightarrow \dots$
 $\downarrow r_1 \quad \downarrow r_2 \quad \downarrow \dots$

$$E_{a_0 a_1} \left[\underbrace{E_{\pi} [r_2 + \gamma r_3 + \dots]}_{V_{\pi}(s_1)} \right]$$

$(s_0 \xrightarrow{\pi'} a_0) \rightarrow (s_1 \xrightarrow{\pi} a_1) \rightarrow (s_2 \xrightarrow{\pi} a_2) \rightarrow \dots$
 $\downarrow r_1 \quad \downarrow r_2 \quad \downarrow \dots$

$$E_{a_0 a_1} [\gamma (Q_{\pi}(s_1, a_1) - V_{\pi}(s_1))] = 0$$

$$E_{a_0 a_1} [r_1 + \gamma r_2 + \dots]$$

$(s_0 \xrightarrow{\pi'} a_0) \rightarrow (s_1 \xrightarrow{\pi'} a_1) \rightarrow (s_2 \xrightarrow{\pi'} a_2) \rightarrow \dots$
 $\downarrow r_1 \quad \downarrow r_2 \quad \downarrow \dots$

$$E_{a_0 a_1} [r_1 + \gamma r_2 + \dots]$$

$(s_0 \xrightarrow{\pi} a_0) \rightarrow (s_1 \xrightarrow{\pi'} a_1) \rightarrow (s_2 \xrightarrow{\pi'} a_2) \rightarrow \dots$
 $\downarrow r_1 \quad \downarrow r_2 \quad \downarrow \dots$

$$0$$

□ 価値バースの手法

ε-貪欲方策

$$\frac{\pi^\varepsilon(a|s)}{\pi} = (1-\varepsilon) \mathbb{1}_{a=\arg\max_{a \in \mathcal{A}} Q_{\pi^\varepsilon}(s,a)} + \varepsilon \frac{1}{k}$$

$$\pi^{\varepsilon=1} \rightarrow \dots \rightarrow \pi^\varepsilon \rightarrow \pi^{\varepsilon'} \rightarrow \dots$$

直感: $\varepsilon: 1 \rightarrow 0$ が「...」のは?

$$(\varepsilon \geq \varepsilon')$$

$$\forall s, a, \exists \epsilon' > 0 \text{ s.t. } E[A_{\pi^{\epsilon'}}(s, a)] \geq 0$$

割引率改善

$$0 \leq \sum_{t=0}^{\infty} \gamma^t E[A_{\pi^{\epsilon}}(s_t, a_t)]$$

$a_0 \dots a_t$
 $s_0 \dots s_t$

E $E[\dots]$

$a_0 \dots a_{t-1}$ a_t

$s_0 \dots s_{t-1}$ s_t

↑
状態なしのときと同じ

ロ方策ベースの手法

パラメトリック方策

$\theta \in \mathbb{R}^d$ 毎に定まる $\pi_\theta(a|s)$

例

$$\{(\theta_{s1}, \theta_{s2}, \dots, \theta_{sk})\}_s \rightarrow \pi_\theta(a|s) = \frac{e^{\theta_{sa}}}{\sum_{\tilde{a}} e^{\theta_{s\tilde{a}}}}$$

$$\dots \rightarrow \pi_\theta \rightarrow \pi_{\theta+\Delta\theta} \rightarrow \dots$$

$$\sum_{t=0}^{\infty} \gamma^t E [A_{\pi_{\theta}}(s_t, a_t)] = J(\pi_{\theta+\Delta\theta}) - J(\pi_{\theta})$$

$a_0 \dots a_t \sim \pi_{\theta+\Delta\theta}$

$s_0 \dots s_t \sim P$

↓
πは-尾用 ← $E [A_{\pi_{\theta}}(s_t, a_t)] = 0$

$$= \Delta\theta \sum_{t=0}^{\infty} \gamma^t E [A_{\pi_{\theta}}(s_t, a_t) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)] + O(\Delta\theta^2)$$

$a_0 \dots a_t \sim \pi_{\theta}$

$s_0 \dots s_t \sim P$

$$\nabla_{\theta} J(\pi_{\theta})$$

||

$$\sum \dots A \rightarrow Q$$

$$Q_{\pi}(s_t, a_t) - \underbrace{V_{\pi}(s_t)}_{B(s_t)}$$

A-スライム

□ 実際.

- ϵ -greedy $\pi^\epsilon \rightarrow Q_\pi(s, a) = E[\dots]$

- param $\pi_\theta \rightarrow \nabla_\theta J(\pi_\theta) = E[\dots]$
が必須
が必須



$E[\dots]$ { モンテカルロ
 $\frac{1}{N} \sum_{i=1}^N \dots$ } Temporal Difference 法.

□ Temporal Difference に向いた.

ベルマン方程式

$$\text{ポイント: } r_1 + \underbrace{\gamma r_2 + \gamma^2 r_3 + \dots}_{\gamma(r_2 + \gamma r_3 + \dots)} = g_1$$

\downarrow g_2

$$Q_{\pi}(s_0, a_0) = E[r_1 + \gamma Q_{\pi}(s_1, a_1)]$$

$a_1 \sim \pi$
 $s_1, r_1 \sim P$

$$V_{\pi}(s_0) = E_P[r_1 + \gamma V_{\pi}(s_1)]$$

ベルマン最適方程式

$$P \pi \bar{P}: \dots \rightarrow \pi^\varepsilon \rightarrow \pi^{\varepsilon'} \rightarrow \dots \underbrace{\pi^{\varepsilon=0}} = \underset{\arg \max Q}{1}$$

$$\pi^*(a|s) = \underset{\tilde{a}}{1} a = \arg \max_{\tilde{a}} Q_{\pi^*}(s, \tilde{a})$$

ベルマン eq

$$Q_{\pi^*}(s_0, a_0) = E_{\substack{a_1 \sim \pi^* \\ s_1, r_1 \sim P}} [r_1 + \gamma Q_{\pi^*}(s_1, a_1)]$$

$$= E_{s_1, r_1 \sim P} [r_1 + \gamma \max_{\tilde{a}} Q_{\pi^*}(s_1, \tilde{a})]$$

$$F(s_0, t_0) = E[\bullet + F(s_1, a_1)]$$

時間が
12 近んだ"と=3

$$\frac{1}{N} \sum_{i=1}^N \rightarrow TD.$$

📖 参考文献

- R. Sutton, A. Barto "Reinforcement Learning
: An Introduction"

└ 訳 : ver1: 森北出版 (小ぶりな冊子)

: ver2 : (A4くらい?)

- 牧野 豊樹 さん ほか "これからの強化学習"

- 森村 哲郎 さん "強化学習"