

強化学習導入

西浦学

東京大学

July N, 2020

- 1 Characters in RL world
- 2 policy gradient theorem
- 3 Relationship between $Q(s, a)$ function and $V(s)$

Value function

$$\begin{aligned} V^\pi(s) &= \mathbb{E}^\pi [G_t | S_t = s] \\ &= \sum_{s' \in \mathcal{S}} \sum_{a \in \mathcal{A}(s)} P(S_{t+1} = s', A_t = a | S_t = s) r(s, a, s') \\ &= \sum_{s' \in \mathcal{S}} \sum_{a \in \mathcal{A}(s)} P(S_{t+1} = s' | S_t = s, A_t = a) \pi(a|s) r(s, a, s') \end{aligned}$$

an example of finite horizon return $T = 2$

$$G_t = R_{t+1} + R_{t+2} \quad (1)$$

therefore the value function can be calculated as following.

$$\begin{aligned} V^\pi(s) &= \mathbb{E}^\pi [G_t | S_t = s] = \mathbb{E}^\pi [R_{t+1} + R_{t+2} | S_t = s] \\ &= \sum_{s'' \in \mathcal{S}} \sum_{a' \in \mathcal{A}(s)} \sum_{s' \in \mathcal{S}} \sum_{a \in \mathcal{A}(s)} \\ &\quad P(S_{t+2} = s'', A_{t+1} = a', S_{t+1} = s', A_t = a | S_t = s) \\ &\quad \times \{r(s, a, s') + r(s', a', s'')\} \\ &= \sum_{s'' \in \mathcal{S}} \sum_{a' \in \mathcal{A}(s)} \sum_{s' \in \mathcal{S}} \sum_{a \in \mathcal{A}(s)} \\ &\quad P(S_{t+2} = s'' | S_{t+1} = s', A_{t+1} = a') \pi(a' | s') \\ &\quad \times P(S_{t+1} = s' | S_t = s, A_t = a) \pi(a | s) \{r(s, a, s') + r(s', a', s'')\} \end{aligned}$$

Theorem 2.1

定理型環境が使える。使い方は普通の \LaTeX と同じ

Theorem 2.1

定理型環境が使える。使い方は普通の \LaTeX と同じ

Proof.

証明も書ける。



Theorem 2.1

定理型環境が使える。使い方は普通の \LaTeX と同じ

Proof.

証明も書ける。



Example 2.2

example

文字の色を変えてみよう

赤

文字の色を変えてみよう

赤青

文字の色を変えてみよう

赤青緑