

Reinforcement learning - part II

Announcements

- ① Gradience 10 due on Sunday night
 - ② For finals, a practice quiz and list of topics will be posted soon
 - ③ Assignment 3 due this week on Friday
-

RL

↳ Markov Decision Processes (MDP)

$$p(s', r | s, a)$$

$$G_t = R_{t+1} + \dots$$

Agent

Policy (π)

$$\pi(a|s) \rightarrow A_t = a \mid S_t = s$$

value function at state s

$$V_{\pi}(s) \doteq \mathbb{E}_{\pi}[\underline{G_t} \mid S_t = s] \quad \textcircled{\circ}$$

$$Q_{\pi}(s, a) \doteq \mathbb{E}_{\pi}[G_T \mid S_t = s, A_t = a] \quad \textcircled{\circ}$$

complete model of the environment.

$$\underline{p(s', r \mid s, a)}$$

Dynamic Programming ^(DP) based methods

init: Start with a policy

Step 1. Find $\underline{V_{\pi}}(s) \quad \forall s \in S$

Policy prediction

Step 2 $\pi \rightarrow V_{\pi}(s) \rightarrow \pi'$

Policy Improvement

$$\underline{p(s', r \mid s, a)}$$

DP \rightarrow works if we ~~have~~ a perfect
 know the environment model
 completely $\underbrace{p(s', r | s, a)}$
 \rightarrow Data-free
 \rightarrow not very useful

Monte-carlo methods.

Policy prediction \nwarrow
 \nearrow Policy Control

first visit MC prediction

s_1 $V(s_1) = 1$ Returns(s_1) = [.]
 s_2 $V(s_2) = 1$ Returns(s_2) = [.]

$(s_1) A_0 R_1 s_1 A_1 R_2 s_2 A_2 R_3 s_3 A \dots$

$V(s_t)$
$Q(s_t, a_t)$

PlayerTotal , Ace usable , House Card

s_1	11	0	6
s_2	13	0	6
s_3	17	0	6
s_4	18	0	6
s_5	22	0	6

-1

$$\text{Returns}(22, 0, 6) = \underline{\underline{[-1]}}$$

✓

$$\text{Returns}(18, 0, 6) = \underline{\underline{[-1]}}$$

⋮

$$V(_) = -1,$$

$$\left[\begin{array}{ccc} 11 & 0 & 6 \end{array} \right]$$

{