

Ch 4 Notes

①

* Remember Bellman optimality

$$V_*(s) = \max_a \mathbb{E} [R_{t+1} + \gamma V_*(S_{t+1}) | S_t = s, A_t = a]$$

$$= \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V_*(s')]$$

or

$$Q_*(s, a) = \mathbb{E} [R_{t+1} + \gamma \max_{a'} Q_*(S_{t+1}, a') | S_t = s, A_t = a]$$
$$= \sum_{s', r} p(s', r | s, a) [r + \gamma \max_{a'} Q_*(s', a')]$$

① P algorithms are ~~learned~~ ~~into~~ obtained by
turning these into assignments. (update rules).

* Policy Iteration

(2)

Policy Iteration = Policy Evaluation

+ Policy Improvement.

→ Policy Evaluation. (also called prediction)

Recall

$$V_{\pi}(s) \doteq \mathbb{E}_{\pi} [G_t | S_t = s]$$

$$= \mathbb{E}_{\pi} [R_{t+1} + \gamma G_{t+1} | S_t = s]$$

$$= \mathbb{E}_{\pi} [R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s]$$

$$= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [\gamma + \gamma V_{\pi}(s')]$$

In principle we could solve as a giant LP!

But we prefer iterative methods.

$$\begin{aligned}
 V_{k+1}(s) &\stackrel{\circ}{=} \mathbb{E}_{\pi} [R_{t+1} + \gamma V_k(S_{t+1}) \mid S_t = s] \\
 &= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma V_k(s')] \\
 &\quad \forall s \in S.
 \end{aligned}
 \tag{3}$$

Can be shown that $\{V_k\}$ converges
as $k \rightarrow \infty$ (Iterative policy evaluation)

Policy Improvement

First let's look at Policy Improvement Theorem

Let π and π' be any pair of deterministic policies, s.t. $\forall s \in S$.

$$Q_{\pi}(s, \pi'(s)) \geq V_{\pi}(s) \rightarrow \textcircled{1.7}$$

\Rightarrow

Then π' must be $\geq \pi$.

$$\text{i.e. } V_{\pi'}(s) \geq V_{\pi}(s) \rightarrow \textcircled{1.8}$$

(9)

All they are saying is that if 4.7 is true at all states then it must be true that 4.8 is also true. This is because π' starting from any state will hit other states along the way and ~~the~~ π' is better than ~~the~~ π at all of those states so the sum must be better.

The proof of PIT proceeds by considering a policy π' that is identical to π everywhere except at a state s . So 4.7 holds everywhere except at s . So if $q_{\pi}(s, a) > v_{\pi}(s)$ then π' is indeed better.

The proof then shows by repeated application of 4.7 that $v_{\pi}(s) \leq v_{\pi'}(s)$

So if you now have a $\pi'(s)$ better everywhere — (5)

$$\pi'(s) \doteq \arg \max_a q_{\pi}(s, a)$$

$$= \arg \max_a \mathbb{E} [R_{t+1} + \gamma V_{\pi}(S_{t+1}) | S_t = s, A_t = a]$$

$$= \arg \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V_{\pi}(s')]$$

it is just greedy w.r.t. current q -values.

* Value Iteration

Drawback of PI,

→ Each iteration requires Policy Eval.
which itself is a prolonged iterative
procedure requiring multiple sweeps
through the state set.

However, we can stop early without losing convergence guarantees of policy iteration.

Important such variation =

Stop policy evaluation after 1-step.
Called as value iteration.

$$V_{k+1}(s) \triangleq \max_a \mathbb{E} [R_{t+1} + \gamma V_k(s_{t+1})]$$

$$s_t = s, A_t = a$$

$$= \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma V_k(s')]$$

$$\forall s \in \mathcal{S}.$$

$\{V_k\}$ can be shown to converge.