

missed
dynamic programming 6/oct/2012

$$\begin{aligned} V_{\pi} &= E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid s_t = s \right] \\ &= E_{\pi} [R_{t+1} + \gamma V_{t+1} \mid s_t = s] \\ &= J_{\pi} \end{aligned}$$

Policy Improvement

We saw earlier Iterative policy evaluation.

Find value function for a given policy

a) Given value function for a policy, can one find better policy?

policy improvement

$$q_{\pi}(s, a) \triangleq E [R_{t+1} + \gamma V_{\pi}(s, a) | S_t = s, A_t = a]$$

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

b) can we use $q_{\pi}(s, a)$ & $V_{\pi}(s)$ in order to find a better policy π' ?

ie. $V_{\pi'}(s) \geq V_{\pi}(s)$

Suppose in state s , we follow the action $\pi(s)$ s.t

$$q_{\pi}(s, \pi'(s)) \geq V_{\pi}(s)$$

\Rightarrow in state s you use π'

$$V_{\pi}(s) \leq q_{\pi}(s, \pi'(s))$$

$$= E_{\pi} [R_{t+1} + \gamma V_{\pi}(s_{t+1}) | S_t = s, A_t = a]$$

$$= E_{\pi'} [R_{t+1} + \gamma V_{\pi}(s_{t+1}) | S_t = s] + \text{action is selected according to } \pi' \text{ in the beginning it self here } A_t = a \text{ absent}$$

by definition

$$\pi'(s) = \pi(s) \quad \forall s \neq s$$

$$\pi'(s) \neq \pi(s) \quad \text{--- ①}$$

$$q_{\pi}(s, \pi'(s)) \geq V_{\pi}(s) \quad \text{--- ②}$$

Kavitha

So when in s during S_t we can use $q_\pi(s, \pi'(s))$

$$\leq E_{\pi'} [R_{t+1} + \gamma q_\pi(s_{t+1}, \pi'(s_{t+1})) | S_t = s]$$

$$= E_{\pi'} [R_{t+1} + \gamma E_{\pi'} [R_{t+2} + \gamma V_\pi(s_{t+2}) | S_{t+1}, A_{t+1} = \pi'(s_{t+1}), S_t = s]]$$

$$= E_{\pi'} [R_{t+1} + \gamma R_{t+2} + \gamma^2 V_\pi(s_{t+2}) | S_t = s]$$

$$\leq E_{\pi'} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 V_\pi(s_{t+3}) | S_t = s]$$

$$\leq E_{\pi'} [R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \gamma^3 R_{t+4} + \dots | S_t = s]$$

$$\underline{\pi = V_{\pi'}(s)}$$

$$\therefore V_\pi(s) \triangleq E \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s \right]$$

$$\pi'(s) = \operatorname{argmax}_a q_{\pi'}(s, a)$$

$$= \operatorname{argmax}_a E [R_{t+1} + \gamma V_\pi(s_{t+1}) | S_t = s, A_t = a]$$

$$= \operatorname{argmax}_a \sum_{s' \in \mathcal{S}} p(s', a | s, a) [\gamma + \gamma V_\pi(s')]$$

Since model is known
i.e. expectation are calculated. The
exploration is not used

Policy improvement:

Form a new policy which is greedy wrt value f^n of old policy. Suppose new policy is such that

$$V_\pi(s) = V_{\pi'}(s) \quad \forall s \in \mathcal{S}$$

In this case π & π' are both optimal.

Also V_π is the optimal value f^n

(2) 4

~~if~~

$$\begin{aligned} V_{\pi'}(s) &= \max_a E [R_{t+1} + \gamma V_{\pi'}(s_{t+1}) \mid s_t = s, A_t = a] \\ &= \max_a E [R_{t+1} + \gamma V_{\pi'}(s_{t+1}) \mid s_t = s, A_t = a] \\ &= \max_a \sum_{s', r} p(s', r \mid s, a) [r + \gamma V_{\pi'}(s')] \end{aligned}$$

Thus $V_{\pi'}$ is the solⁿ to bellman eqn.

$$= V_{\pi'} = V_{\pi^*} \text{ optimal reward.}$$

One can instead of deterministic policies, consider stochastic policies.

Suppose $A(s) \equiv$ set of feasible actions in state s

$$\text{Let } A(s) = \{a_1, a_2, \dots, a_n\}$$

Let's assume that a_i, a_j, a_k are the maximising actions

then a stochastic policy

$$\pi'(a \mid s) = \begin{cases} a_i & \text{w.p } p_1 \\ a_j & \text{w.p } p_2 \\ a_k & \text{w.p } p_3 \end{cases}$$

$$\text{s.t. } p_1 + p_2 + p_3 = 1, \quad p_1, p_2, p_3 \geq 0$$

Thus for non-optimal, i.e., those in ^{not} $A(s) \{a_i, a_j, a_k\}$

we assign zero prob.

* grid example

π ' greedy

Random policy
(π)

V_{π}

| | | | |
|-----|-----|-----|-----|
| 0 | -14 | -20 | -22 |
| -14 | -18 | -20 | -20 |
| -20 | -20 | -18 | -14 |
| -22 | -20 | -14 | 0 |

| | | | |
|-----|---|---|-----|
| /// | ← | ← | ← |
| ↑ | ↖ | ↖ | ↓ |
| ↑ | ↗ | ↘ | ↓ |
| ↗ | → | → | /// |

Question: What is $V_{\pi'}$?

$V_0 =$

| | | | |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 |

$V_1 =$

| | | | |
|----|----|----|----|
| 0 | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | -1 |
| -1 | -1 | -1 | 0 |

$V_2 =$

| | | | |
|----|----|----|----|
| 0 | -1 | -2 | -2 |
| -1 | -2 | -3 | -2 |
| -2 | -3 | -2 | -1 |
| -3 | -2 | -1 | 0 |

$V_3 =$

| | | | |
|----|----|----|----|
| 0 | -1 | -2 | -3 |
| -1 | -2 | -3 | -2 |
| -2 | -3 | -2 | -1 |
| -3 | -2 | -1 | 0 |

V_4

| | | | |
|----|----|----|----|
| 0 | -1 | -2 | -3 |
| -1 | -2 | -3 | -2 |
| -2 | -3 | -2 | -1 |
| -3 | -2 | -1 | 0 |

← has converged.

Find π'' - the greedy policy on $V_{\pi'}$

(notes in MBROI Book)

13/04/2018