

[Next](#) [Up](#) [Previous](#) [Contents](#)

Next: [4.3 Policy Iteration](#) **Up:** [4. Dynamic Programming](#) **Previous:** [4.1 Policy Evaluation](#) [Contents](#)

4.2 Policy Improvement

Our reason for computing the value function for a policy is to help find better policies. Suppose we have determined the value function V^π for an arbitrary deterministic policy π . For some state s we would like to know whether or not we should change the policy to deterministically choose an action $a \neq \pi(s)$. We know how good it is to follow the current policy from s --that is $V^\pi(s)$ --but would it be better or worse to change to the new policy? One way to answer this question is to consider selecting a in s and thereafter following the existing policy, π . The value of this way of behaving is

$$\begin{aligned} Q^\pi(s, a) &= E_\pi\{r_{t+1} + \gamma V^\pi(s_{t+1}) | s_t = s, a_t = a\} \quad (4.6) \\ &= \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V^\pi(s')]. \end{aligned}$$

The key criterion is whether this is greater than or less than $V^\pi(s)$. If it is greater--that is, if it is better to select a once in s and thereafter follow π than it would be to follow π all the time--then one would expect it to be better still to select a every time s is encountered, and that the new policy would in fact be a better one overall.

That this is true is a special case of a general result called the *policy improvement theorem*. Let π and π' be any pair of deterministic policies such that, for all $s \in \mathcal{S}$,

$$Q^\pi(s, \pi'(s)) \geq V^\pi(s). \quad (4.7)$$

Then the policy π' must be as good as, or better than, π . That is, it must obtain greater or equal expected return from all states $s \in \mathcal{S}$:

$$V^{\pi'}(s) \geq V^\pi(s). \quad (4.8)$$

Moreover, if there is strict inequality of (4.7) at any state, then there must be strict inequality of (4.8) at at least one state. This result applies in particular to the two policies that we considered in the previous paragraph, an original deterministic policy, π , and a changed policy, π' , that is identical to π except that $\pi'(s) = a \neq \pi(s)$.

Obviously, (4.7) holds at all states other than s . Thus, if $Q^\pi(s, a) > V^\pi(s)$, then the changed policy is indeed better than π .

The idea behind the proof of the policy improvement theorem is easy to understand. Starting from (4.7), we keep expanding the Q^π side and reapplying (4.7) until we get $V^{\pi'}(s)$:

$$\begin{aligned}
 V^\pi(s) &\leq Q^\pi(s, \pi'(s)) \\
 &= E_{\pi'}\{r_{t+1} + \gamma V^\pi(s_{t+1}) \mid s_t = s\} \\
 &\leq E_{\pi'}\{r_{t+1} + \gamma Q^\pi(s_{t+1}, \pi'(s_{t+1})) \mid s_t = s\} \\
 &= E_{\pi'}\{r_{t+1} + \gamma E_{\pi'}\{r_{t+2} + \gamma V^\pi(s_{t+2})\} \mid s_t = s\} \\
 &= E_{\pi'}\{r_{t+1} + \gamma r_{t+2} + \gamma^2 V^\pi(s_{t+2}) \mid s_t = s\} \\
 &\leq E_{\pi'}\{r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 V^\pi(s_{t+3}) \mid s_t = s\} \\
 &\vdots \\
 &\leq E_{\pi'}\{r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \gamma^3 r_{t+4} + \cdots \mid s_t = s\} \\
 &= V^{\pi'}(s).
 \end{aligned}$$

So far we have seen how, given a policy and its value function, we can easily evaluate a change in the policy at a single state to a particular action. It is a natural extension to consider changes at *all* states and to *all* possible actions, selecting at each state the action that appears best according to $Q^\pi(s, a)$. In other words, to consider the new *greedy* policy, π' , given by

$$\begin{aligned}
 \pi'(s) &= \arg \max_a Q^\pi(s, a) \\
 &= \arg \max_a E\{r_{t+1} + \gamma V^\pi(s_{t+1}) \mid s_t = s, a_t = a\} \quad (4.9) \\
 &= \arg \max_a \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V^\pi(s')],
 \end{aligned}$$

where $\arg \max_a$ denotes the value of a at which the expression that follows is maximized (with ties broken arbitrarily). The greedy policy takes the action that looks best in the short term--after one step of lookahead--according to V^π . By construction, the greedy policy meets the conditions of the policy improvement theorem (4.7), so we know that it is as good as, or better than, the original policy. The process of making a new policy that improves on an original policy, by making it greedy with respect to the value function of the original policy, is called *policy improvement*.

Suppose the new greedy policy, π' , is as good as, but not better than, the old policy π . Then $V^\pi = V^{\pi'}$, and from (4.9) it follows that for all $s \in \mathcal{S}$:

$$\begin{aligned} V^{\pi'}(s) &= \max_a E \left\{ r_{t+1} + \gamma V^{\pi'}(s_{t+1}) \mid s_t = s, a_t = a \right\} \\ &= \max_a \sum_{s'} \mathcal{P}_{ss'}^a \left[\mathcal{R}_{ss'}^a + \gamma V^{\pi'}(s') \right]. \end{aligned}$$

But this is the same as the Bellman optimality equation (4.1), and therefore, $V^{\pi'}$ must be V^* , and both π and π' must be optimal policies. Policy improvement thus must give us a strictly better policy except when the original policy is already optimal.

So far in this section we have considered the special case of deterministic policies. In the general case, a stochastic policy π specifies probabilities, $\pi(s, a)$, for taking each action, a , in each state, s . We will not go through the details, but in fact all the ideas of this section extend easily to stochastic policies. In particular, the policy improvement theorem carries through as stated for the stochastic case, under the natural definition:

$$Q^\pi(s, \pi'(s)) = \sum_a \pi'(s, a) Q^\pi(s, a).$$

In addition, if there are ties in policy improvement steps such as (4.9)--that is, if there are several actions at which the maximum is achieved--then in the stochastic case we need not select a single action from among them. Instead, each maximizing action can be given a portion of the probability of being selected in the new greedy policy. Any apportioning scheme is allowed as long as all submaximal actions are given zero probability.

The last row of Figure 4.2 shows an example of policy improvement for stochastic policies. Here the original policy, π , is the equiprobable random policy, and the new policy, π' , is greedy with respect to V^π . The value function V^π is shown in the bottom-left diagram and the set of possible π' is shown in the bottom-right diagram. The states with multiple arrows in the π' diagram are those in which several actions achieve the maximum in (4.9); any apportionment of probability among these actions is permitted. The value function of any such policy, $V^{\pi'}(s)$, can be seen by inspection to be either -1 , -2 , or -3 at all states, $s \in \mathcal{S}$, whereas $V^\pi(s)$ is at most -14 . Thus, $V^{\pi'}(s) \geq V^\pi(s)$, for all $s \in \mathcal{S}$, illustrating policy improvement. Although in this case the new policy π' happens to be optimal, in general only an improvement is guaranteed.

[Next](#) [Up](#) [Previous](#) [Contents](#)

Next: [4.3 Policy Iteration](#) **Up:** [4. Dynamic Programming](#) **Previous:** [4.1 Policy Evaluation](#) [Contents](#)

Mark Lee 2005-01-04