– Bellman optimality equations

(4.1)
$$V_*(s) = \max_a E\left[R_{t+1} + \gamma V_*(S_{t+1}) \mid S_t = s, A_t = a\right]$$

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

(4.2)
$$q_*(s,a) = E\left[R_{t+1} + \gamma \max_{a'} q_*(S_{t+1}, a') \mid S_t = s, A_t = a\right]$$

$$= \sum_{s',r} p(s',r \mid s,a)\left[r + \gamma \max_{a'} q_*(s', a')\right]$$

– Policy evaluation

(4.4)
$$V_\pi(s) = \sum_a \pi(a \mid s) \sum_{s',r} p(s',r \mid s,a)\left[r + \gamma V_\pi(s')\right]$$

(4.5)
$$V_{k+1}(s) = \sum_a \pi(a \mid s) \sum_{s',r} p(s',r \mid s,a)\left[r + \gamma V_k(s')\right]$$

– Iterative policy evaluation – the sequence $\{V_k\}$ converges to $V_\pi$ as $k \to \infty$

— Policy improvement

(4.6)    $q_\pi(s,a) = \sum\limits_{s',r} p(s',r \mid s,a)\left[r + \gamma V_\pi(s')\right]$

— Policy $\pi'$ is as good as or better than $\pi$ if

(4.7)    $q_\pi(s, \pi'(s)) \geq V_\pi(s)$    ✓✓

(4.8)    $V_{\pi'}(s) \geq V_\pi(s)$

— Greedy policy

(4.9)    $\pi'(s) = \underset{a}{\arg\max} \; q_\pi(s,a)$

          $= \underset{a}{\arg\max} \sum\limits_{s',r} p(s',r \mid s,a)\left[r + \gamma V_\pi(s')\right]$

— Suppose $\pi'$ is as good as but not better than $\pi$

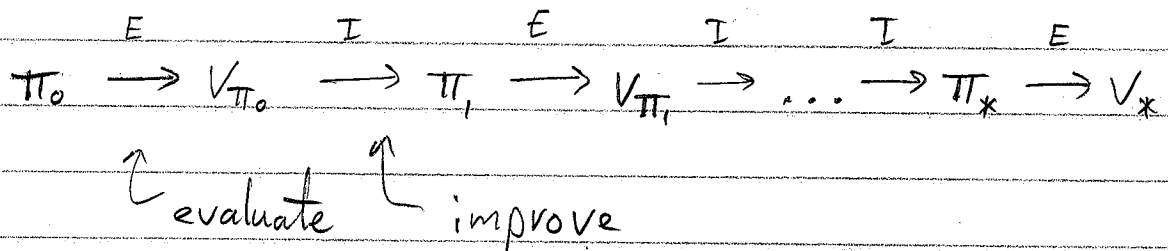$V_{\pi'}(s) = \underset{a}{\max} \sum\limits_{s',r} p(s',r \mid s,a)\left[r + \gamma V_{\pi'}(s')\right]$

∴ $V_{\pi'}$ must be $V_*$, both $\pi$ and $\pi'$ are optimal

— $q_\pi(s, \pi'(s)) = \sum\limits_a \pi'(a \mid s) \, q_\pi(s,a)$

     ↰ for stochastic policies

- Policy iteration

Once a policy, $\pi$, has been improved
using $V_\pi$ to yield a better policy, $\pi'$,
we can then compute $V_{\pi'}$ and improve
it again to yield an even better $\pi''$.

$$\pi_0 \xrightarrow{E} V_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V_{\pi_1} \xrightarrow{I} \dots \xrightarrow{I} \pi_* \xrightarrow{E} V_*$$

$\underset{\text{evaluate}}{\curvearrowleft} \underset{\text{improve}}{\curvearrowleft}$

This method can be used to find
an optimal policy.

- Value iteration

(4.10)
$$V_{k+1}(s) = \max_a \sum_{s',r} p(s',r \mid s,a)\left[r + \gamma V_k(s')\right]$$

↳ combines policy evaluation and
policy improvement

- Generalized policy evaluation (GPI) - the
general idea of letting policy evaluation and
policy improvement processes interact independently
of the details of the two processes.

— Exercise 4.1

$$q_\pi(s,a) = \sum_{s',r} p(s',r \mid s,a)\left[r + \gamma V_\pi(s')\right]$$

$$q_\pi(11, down) = -1 + \gamma V_\pi(11) = -1 - 14 = -15$$

$$q_\pi(7, down) = -1 + \gamma V_\pi(11) = -1 - 14 = -15$$

— Exercise 4.3

(4.3)
$$V_\pi(s) = E_\pi\left[R_{t+1} + \gamma V_\pi(S_{t+1}) \mid S_t = s\right]$$

$$q_\pi(s,a) = E_\pi\left[R_{t+1} + \gamma q_\pi(S_{t+1}, A_{t+1}) \mid S_t = s, A_t = a\right]$$

(4.4)
$$V_\pi(s) = \sum_a \pi(a \mid s) \sum_{s',r} p(s',r \mid s,a)\left[r + \gamma V_\pi(s')\right]$$

$$q_\pi(s,a) = \sum_{s',r} p(s',r \mid s,a)\left[r + \gamma \sum_{a'} \pi(a' \mid s) q_\pi(s', a')\right]$$

(4.5)
$$V_{k+1}(s) = \sum_a \pi(a \mid s) \sum_{s',r} p(s',r \mid s,a)\left[r + \gamma V_k(s')\right]$$

$$q_{k+1}(s,a) = \sum_{s',r} p(s',r \mid s,a)\left[r + \gamma \sum_{a'} \pi(a' \mid s') q_k(s', a')\right]$$

– Exercise 4.4

For example, if $S_{13}$ always transitions to $S_{15}$ and $S_{15}$ always transitions to $S_{13}$.

Possible Solutions:

– Identify the "trap" states and don't allow any actions that could transition into them.

– Randomly transition out of the "trap" states and into an alternate state a small percentage of the time.

– Detect if the value is not converging by checking if the difference $|V_k - V_{k+1}|$ is getting smaller on each iteration or staying the same.

– Detect if $|V_k - V_{k+1}|$ is equal to the reward on each iteration.

– Disallow transitioning into a state once it has been identified as a "trap" state.

– Exercise 4.6

$$\pi_*(s) = \underset{a}{\text{argmax}} \sum_{s',r} p(s',r|s,a)\left[r + \gamma V_*(s')\right]$$

$$\pi_*(s) = \underset{a}{\text{argmax}}\ q_*(s,a)$$

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

$$q_*(s,a) = \sum_{s',r} p(s',r|s,a)\left[r + \gamma q_*(s',\pi_*(s'))\right]$$

$$\pi'(s) = \underset{a}{\text{argmax}} \sum_{s',r} p(s',r|s,a)\left[r + \gamma V_\pi(s')\right]$$

$$\pi'(s) = \underset{a}{\text{argmax}}\ q_\pi(s,a)$$

$$V_{k+1}(s) = \sum_{s',r} p(s',r|s,\pi(s))\left[r + \gamma V_k(s')\right]$$

$$q_{k+1}(s,a) = \sum_{s',r} p(s',r|s,a)\left[r + \gamma q_k(s',\pi(s'))\right]$$

1. Initialization
   For each $s \in S$, $a \in A(s)$:
   $\quad Q(s,a) \leftarrow$ arbitrary
   $\quad \pi(s) \leftarrow$ arbitrary

2. Policy Evaluation
   Repeat
   $\quad \Delta \leftarrow 0$
   $\quad$ For each $s \in S$, $a \in A(s)$:
   $\qquad q \leftarrow Q(s,a)$

   $$Q(s,a) \leftarrow \sum_{s',r} p(s',r \mid s,a)\left[r + \gamma Q(s', \pi(s'))\right]$$

   $$\Delta \leftarrow \max\left(\Delta, |q - Q(s,a)|\right)$$
   until $\Delta < \theta$

3. Policy Improvement
   policy-stable $\leftarrow$ true
   For each $s \in S$:
   $\quad a \leftarrow \pi(s)$

   $$\pi(s) \leftarrow \overset{argmax}{a} \overbrace{\sum_{s',r} p(s',r \mid s,a)\left[r + \gamma Q(s', \pi(s').)\right]}^{Q(s,a)}$$

   $\quad$ If $a \neq \pi(s)$ then policy-stable $\leftarrow$ false
   If policy-stable then return $Q$ and $\pi$
   else go to 2.

— Exercise 4.7

A stochastic policy $\pi(a|s)$ would be used instead of a deterministic policy $\pi(s)$.

For each action $a$ where $a \neq \pi(s)$, the stochastic policy would be:

$$\pi(a|s) = \frac{\varepsilon}{|A(s)|}$$

For action $a$ where $a = \pi(s)$, the stochastic policy would be:

$$\pi(a|s) = 1 - \frac{\varepsilon(|A(s)|-1)}{|A(s)|}$$

$$= 1 - \frac{\varepsilon|A(s)|}{|A(s)|} + \frac{\varepsilon}{|A(s)|}$$

$$= 1 - \varepsilon + \frac{\varepsilon}{|A(s)|}$$

Step 3 would have to compute $\pi(a|s)$ in addition to computing $\pi(s)$

Step 2 would use $\pi(a|s)$ instead of $\pi(s)$

Step 1 would initialize $\pi(a|s)$ to satisfy $\varepsilon$-soft

Exercise 4.10

$$V_{k+1}(s) = \max_a E\left[R_{t+1} + \gamma V_k(S_{t+1}) \,\middle|\, S_t = s, A_t = a\right]$$

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

$$q_{k+1}(s,a) = E\left[R_{t+1} + \gamma \max_{a'} q_k(S_{t+1}, a') \,\middle|\, S_t = s, A_t = a\right]$$

$$= \sum_{s',r} p(s',r|s,a)\left[r + \gamma \max_{a'} q_k(s',a')\right]$$