

#1

Since no reward is before the terminal state of the episode, and no penalty is assigned for taking unnecessary steps in progress to the terminal state, discounting or otherwise, the robot has no incentive to minimize the steps taken to the terminal state. In fact, a random walk will eventually arrive at the same score as a direct path, so the robot never develops any path improvements.

$G_t = 1 \forall t$, regardless of action

#2

$$\begin{aligned} V_{\pi}(s) &= E_{\pi} [G_k | s_k = s] \\ &= \sum_a \pi(a|s) \sum_{s', r} p(s', r | s, a) [r + \gamma V_{\pi}(s')] \end{aligned}$$

$$\begin{aligned} V_{\pi}(\text{center}) &= (0.25)(1) [0 + 0.9(-0.4)] \\ &\quad + (0.25)(1) [0.9(0.7)] \\ &\quad + 0.25 [0.9(2.3)] \\ &\quad + 0.25 [0.9(0.4)] \\ &= -0.09 + 0.1575 + 0.5175 + 0.09 \\ &= \underline{0.675 \approx 0.7} \end{aligned}$$

#3

$$q_{\pi}(s, a) = E_{\pi}[G_t | S_t = s, A_t = a]$$

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

$$= \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma E[G_{t+1} | s_{t+1}, a_{t+1}]]$$

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

#4

$$V_{\pi} = R^{\pi} + \gamma P^{\pi} V_{\pi}$$

$$V_{\pi} - \gamma P^{\pi} V_{\pi} = R^{\pi}$$

$$(\mathbf{I} - \gamma P^{\pi}) V_{\pi} = R^{\pi}$$

$$V_{\pi} = (\mathbf{I} - \gamma P^{\pi})^{-1} R^{\pi}$$

$$V_{\pi_c} = (\mathbf{I} - \gamma P^{\pi})^{-1} (R^{\pi} + c \mathbf{J})$$

$$V_{\pi_c} - V_{\pi} = V_c$$

$$V_c = (\mathbf{I} - \gamma P^{\pi})^{-1} (R^{\pi} + c \mathbf{J}) - (\mathbf{I} - \gamma P^{\pi})^{-1} R^{\pi}$$

$$= (\mathbf{I} - \gamma P^{\pi})^{-1} [(R^{\pi} + c \mathbf{J}) - R^{\pi}]$$

$$V_c = (\mathbf{I} - \gamma P^{\pi})^{-1} (c \mathbf{J})$$

γ, P^{π}, c are constant under any
given policy π

$\Rightarrow V_c$ is constant under any given
policy π .



#5

Consider example 3.4, pole-balancing. Suppose a reward of $+1$ is given at every time interval without failure, and a reward of 0 for failure and termination of the episode.

Also consider a reward of 0 at every time interval and -1 at failure with discount factor $\gamma < 1$.

In either case, maximization of G_t would involve keeping the pole balanced as long as possible.

The result from the previous exercise holds.

#6

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

$$\max_{a'} q_*(s', a') = V_*(s')$$

\therefore

$$q_*(h, s) = p(h, r_{\text{search}} | h, s) [r_{\text{search}} + \gamma V_*(h)]$$
$$+ p(l, r_{\text{search}} | h, s) [r_{\text{search}} + \gamma V_*(l)]$$
$$= \alpha [r_{\text{search}} + \gamma V_*(h)]$$
$$+ (1 - \alpha) [r_{\text{search}} + \gamma V_*(l)]$$

$$q_*(h, w) = p(h, r_{\text{wait}} | h, w) [r_{\text{wait}} + \gamma V_*(h)]$$
$$+ p(l, r_{\text{wait}} | h, w) [r_{\text{wait}} + \gamma V_*(l)]$$
$$= r_{\text{wait}} + \gamma V_*(h)$$

$$\begin{aligned}
 q_*(l, w) &= p(h, r_{\text{wait}} | l, w) [r_{\text{wait}} + \gamma v_*(h)] \\
 &\quad + p(l, r_{\text{wait}} | l, w) [r_{\text{wait}} + \gamma v_*(l)] \\
 &= r_{\text{wait}} + \gamma v_*(l)
 \end{aligned}$$

$$\begin{aligned}
 q_*(l, s) &= p(h, r_{\text{search}} | l, s) [r_{\text{search}} + \gamma v_*(h)] \\
 &\quad + p(l, r_{\text{search}} | l, s) [r_{\text{search}} + \gamma v_*(l)] \\
 &= (1 - \beta) [-3 + \gamma v_*(h)] + \\
 &\quad \beta [r_{\text{search}} + \gamma v_*(l)]
 \end{aligned}$$

$$\begin{aligned}
 q_*(l, re) &= p(h, r_{\text{recharge}} | l, re) [r_{\text{recharge}} + \gamma v_*(h)] \\
 &\quad + p(l, r_{\text{recharge}} | l, re) [r_{\text{recharge}} + \gamma v_*(l)] \\
 &= \gamma v_*(h)
 \end{aligned}$$

where $v_*(h)$, $v_*(l)$ are defined in example 3.9, page 65.