

1. Discounted returns

(a) Assume you observe a sequence of five rewards

$$R_1 = -1, R_2 = 2, R_3 = 6, R_4 = 3, R_5 = 2$$

until you reach a terminal state, i.e. a state that always transitions back to itself with a reward of 0. Calculate the returns G_0, \dots, G_5 for a discount factor of $\gamma = 0.5$.

$$G_0 = R_1 + \gamma R_2 + \gamma^2 R_3 + \gamma^3 R_4 + \gamma^4 R_5 + \sum_{k=0}^{\infty} \gamma^{5+k} \cdot 0$$

$$= -1 + 0.5 \times 2 + 0.25 \times 6 + 0.125 \times 3 + 0.0625 \times 2$$

$$G_4 = R_5 = 2$$

$$G_5 = \sum_{k=0}^{\infty} \gamma^k \cdot 0$$

↓
 R_{5+k+1}

$$G_3 = R_4 + \gamma R_5 = 3 + 1 = 4$$

$$G_2 = 6 + 2 = 8$$

$$G_1 = 2 + 4 = 6$$

$$G_0 = -1 + 3 = 2$$

(b) Assume an MDP produces an infinite sequence of rewards of 5, i.e.

$$R_1 = 5, R_2 = 5, R_3 = 5, \dots$$

Calculate the return G_0 for the discount factors $\gamma \in \{0, 0.2, 0.5, 0.9, 0.95, 0.99, 0.999\}$.

What would happen if the discount factor was $\gamma = 1$?

Hint: You can use the closed form of a special case of the power series.

$$G_0 = 5 \sum_{k=0}^{\infty} \gamma^k$$

~~$\gamma = 0$, $G_0 = 5$~~ only look at the next step

$$\gamma \geq 0, G_0 = \frac{5}{1-\gamma}$$

$$\gamma = 0.2 \quad G_0 = \frac{5}{0.8} = 6.25$$

$$\gamma = 0.5 \quad G_0 = \frac{5}{0.5} = 10$$

$$\gamma = 0.9 \quad G_0 = \frac{5}{0.1} = 50$$

$$\gamma = 0.95 \quad G_0 = \frac{5}{0.05} = 100$$

$$\gamma = 0.99 \quad G_0 = \frac{5}{0.01} = 500$$

$$\gamma = 0.999 \quad G_0 = \frac{5}{0.001} = 5000$$

importance of
future increases
as γ increase

↳ make sure the return is properly bounded

"one trajectory to the end,
no feedback in between"
tic-tak-toe

Real case (sparse reward environment)
e.g.

(c) Assume you observe a sequence of $T > 1$ rewards

$R_1 = 0, R_2 = 0, \dots, R_{T-1} = 0, R_T \Rightarrow$ got reward at the end of the dialogue

until you reach a terminal state. Note that all rewards except R_T are zero. How can you choose γ such that the initial return G_0 is equal to $\epsilon > 0$? Calculate these γ for the following situations:

- i. $\epsilon = 0.1, R_T = 1, T = 10$ time \uparrow factor \uparrow
- ii. $\epsilon = 0.1, R_T = 1, T = 50$ return \downarrow factor \downarrow
- iii. $\epsilon = 0.01, R_T = 1, T = 50$

$$G_0 = r^{T-1} R_T \quad r = \sqrt[T-1]{\frac{G_0}{R_T}} = \sqrt[T-1]{\frac{\epsilon}{R_T}}$$

$$G_T = \sum_{k=0}^{\infty} r^k R_{T+k+1}$$

$$G_0 = R_1 + r R_2 + \dots + r^{T-1} R_T$$

$$r = \sqrt[T-1]{\frac{\epsilon}{R_T}}$$

$$i. \quad r = \sqrt[9]{0.1} \approx 0.7743$$

$$ii. \quad r = \sqrt[49]{0.1} \approx 0.9541$$

$$iii. \quad r = \sqrt[49]{0.01} \approx 0.9103$$

2. Value functions

$$v_{\pi}(s) = \sum_a \pi(a|s) (R(s,a) +$$

- (a) For any given MDP, policy π and terminal state E , what is $v_{\pi}(E)$? All transitions from a terminal state are back to itself with a reward of 0.

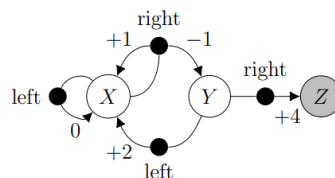
$$r \sum_{s'} P(s'|s,a) v_{\pi}(s')$$

$$p(G_t = 0 | S_t = E) = 1 \quad v_{\pi}(E) = E[G_t | S_t = E]$$

$$= 0 \cdot p(G_t = 0 | S_t = E) = 0$$

$$v_{\pi}(E) = \sum_a \pi(a|E) (R(s,a) + r p(E|E,a) v_{\pi}(E))$$

- (b) Consider the MDP and policy π_1 from the previous exercise set. Note that if action *right* is taken in state X , then the transitions to X and Y occur with probabilities 0.75 and 0.25, respectively. The deterministic policy π_1 is defined as $\pi_1(X) = \text{right}$, $\pi_1(Y) = \text{right}$.



value functions
are policy dependent

Calculate the values of states X and Y under policy π_1 , i.e. $v_{\pi_1}(X)$ and $v_{\pi_1}(Y)$, using a discount factor of $\gamma = 0.9$.

Hint: Start with the value of Y .

$$v_{\pi_1}(Y) = E[G_t | S_t = Y] = 4 \cdot p(G_t = 4 | S_t = Y) = 4$$

$$\text{for } v_{\pi_1}(X), \quad g_0 = \sum_{k=0}^{\infty} r^k \cdot (-1) \quad p = \lim_{t \rightarrow \infty} 0.75^t = 0$$

$$T = 0, 1, \dots, \quad g_T = \sum_{k=0}^{T-1} r^k (-1) + 4 r^T \quad p = 0.75^{T-1} \cdot 0.25$$

$$v_{\pi_1}(X) = \sum_{T=0}^{\infty} (- \sum_{k=0}^{T-1} r^k + 4 r^T) \cdot 0.75^{T-1} \cdot 0.25$$

$$\Downarrow$$

$$\frac{1 - r^{T-1}}{1 - r} + 4 r^T$$

$$v_{\pi_1}(X) = \sum_{k=0}^{\infty} 0.75^k \cdot 0.25 v_{\pi_1}(Y) + \sum_{k=0}^{\infty} \left(\frac{1 - r^{T-1}}{1 - r} \right) \cdot 0.75^k \cdot 0.25$$

$$S \in \{x, y, z\}$$

$$A = \{\text{left}, \text{right}\}$$

$$\pi_1(\text{right}|X) = 1, \pi_1(\text{left}|X) = 0$$

$$\pi_1(\text{right}|Y) = 1, \pi_1(\text{left}|Y) = 0$$

$$v_{\pi_1}(Z) = 0$$

$$v_{\pi_1}(Y) = \pi_1(\text{right}|Y) \left(4 + \overset{1}{\uparrow} r(P(Z|Y, \text{right}) v_{\pi_1}(Z)) \right)$$

$$= 4$$

$$v_{\pi_1}(X) = \pi_1(\text{right}|X) (0.5 + r(P(X|X, \text{right}) v_{\pi_1}(X) + P(Y|X, \text{right}) v_{\pi_1}(Y)))$$

$$R(s, a) = 0.5$$

✓ how to calculate $v_{\pi_1}(x)$
in general: we "evaluate" the policy iteratively

$$0.325 v_{\pi_1}(X) = 1.4$$

(c) We may add "-1" reward to intermediate step.