

3.3 Returns and Episodes

$$G_t = \text{Expected Reward Return} \quad T = \text{final time step}$$
$$= R_{t+1} + R_{t+2} + R_{t+3} + \dots + R_T$$

S = set of all states

s^+ = Terminal state

Discounting is agent tries to select actions such that sum of discounted rewards over time is maximized.

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

γ = discount rate $0 \leq \gamma \leq 1$

$$\gamma = \begin{cases} < 1 & \text{sum has finite value as long as} \\ & R_t \text{ is bounded} \\ = 0 & \text{myopic = only cares about } R_{t+1} \end{cases}$$

As $\gamma \rightarrow 1$, takes future rewards more strongly.

$$G_t = R_{t+1} + \gamma G_{t+1}$$

$$\underline{\text{Ex 3.8}} \quad G_{t=0} = R_1 + \gamma G_{t=1} = 2$$

$$G_{t=1} = R_2 + \gamma G_{t=2} = 6$$

$$G_{t=2} = R_3 + \gamma G_{t=3} = 8$$

$$G_{t=3} = R_4 + \gamma G_{t=4} = 4$$

$$G_{t=4} = R_5 + \gamma G_{t=5} = 2$$

$$G_{t=5} = R_6 + R_7 + \dots = 0$$

3.5 Policies and Value Functions

Value Functions: How good agent to be in current state

Policy: Mapping of states to probabilities
1. can only depend on current state

$V_{\pi}(s) = \text{expected return when starting in } s \text{ and following } \pi$

$$V_{\pi}(s) = E_{\pi}[G_t | S_t=s] \quad \text{state-value function}$$

$$= E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t=s \right] \quad \forall s \in S$$

$q_{\pi}(s, a) = \text{expected return when starting in } s, \text{ taking action } a, \text{ then following policy } \pi$

action-value policy

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

$$= E_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t=s, A_t=a \right]$$

This must hold between values of s and possible successor states.

$$V_{\pi}(s) = E_{\pi}[G_t | S_t=s]$$

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

$$= \sum_a \pi(a|s) \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma E_{\pi}[G_{t+1} | S_{t+1}=s']]$$

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

Bellman Equation | look ahead

Ex:

A	B
C	D

$$V_\pi(A) = \frac{1}{2}(0 + 0.7 V_\pi(C)) + \frac{1}{2}(5 + 0.7 V_\pi(B)) \\ + \frac{1}{2}(0 + 0.7 V_\pi(A))$$

Ex 3.14

$$V_\pi(a) = \frac{1}{4}(0.9 \cdot 0.4) + \frac{1}{4}(0.9 + 0.4) + \frac{1}{4}(0.9 \cdot 0.7) + \frac{1}{4}(0.9 \cdot 2.3) \\ \approx 0.7$$

3.6 Optimal Policies and Optimum Value Functions

V_π is optimal policy - $\pi_1 \geq \pi_2$ if $V_\pi(s) \geq V_{\pi'}(s)$

$= \max_\pi V_\pi(s)$ $\forall s \in S$ - Always an optimal policy

q_π^* = optimum action-value function

$$= \max_\pi q_\pi(s,a)$$

$$= E[R_{t+1} + \gamma V_\pi^*(s_{t+1}) \mid s_t = s, A_t = a]$$

Bellman Optimum Equality

Value of a state under policy must be expected return for best action from state

$$V_\pi^*(s) = \max_{a \in A(s)} q_{\pi^*}(a,s)$$

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

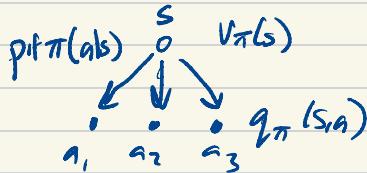
$$Q_{\pi^*}(s) = \sum_{s'} p(s'|s,a) \left[R(s,a) + \lambda \max_{a'} Q^*(s', a') \right]$$

Ex 3.11

$$E_{\pi} [R_{t+1} | s_t] = \sum_a \pi(a|s) \sum_s p(s', r | s, a) [r]$$

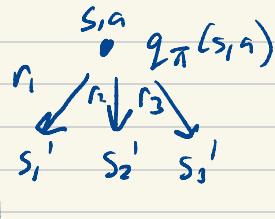
Ex 3.12

$$V_{\pi}(s) = \sum_a \pi(a|s) q_{\pi}(s, a)$$



Ex 3.13

$$q_{\pi}(s) = \sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) \cdot [r + \lambda V_{\pi}(s')]$$



$$\pi_{\pi^*}(s) = \arg \max_a \sum_{s'} \sum_r p(s', r | s, a) [r + \lambda V_{\pi^*}(s')]$$

$$\pi_{\pi^*}(s) = \arg \max_a q_{\pi^*}(s, a) \quad \text{w/ action value}$$