

Theta[s[0]] 1st sample from population

Lecture-6

- Overview:
- ① Environments
 - ② BBO for policy search
 - Simple agent that doesn't leverage MDP structure
 - ③ Value functions
 - Tool for leveraging MDP
 - not a complete agent

State-Value Function:

$$V^\pi : S \Rightarrow \mathbb{R}$$

$$V^\pi(s) \triangleq \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k} \mid S_t = s, \pi \right]$$

• Depends on policy π

$$\gamma^0 R_t + \gamma^1 R_{t+1} + \gamma^2 R_{t+2} \dots$$

= Expected discounted return from state s under policy π

$$J(\pi) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R_{t+1} \mid \pi \right]$$

$$= \sum_s d_\pi(s) \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k} \mid \pi, S_0 = s \right]$$

$$\therefore J(\pi) = \sum_s d_0(s) v^\pi(s)$$

$$J(\pi) = v^\pi(s_0)$$

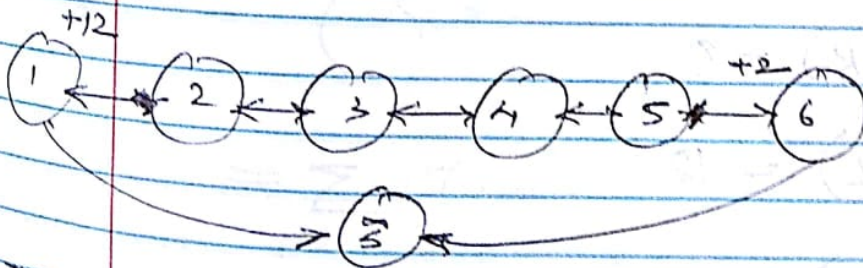
$$v^\pi(s) = \sum_{k=0}^{\infty} \gamma^k \mathbb{E} [R_{t+k} | s_t = s, \pi]$$

$$= \gamma^0 \sum_a \pi(s, a) \sum_{s'} P(a, a, s') R(s, a, s')$$

$$+ \gamma^1 \sum_a \pi(s, a) \sum_{s'} P(s, a, s') \sum_{a'} \pi(s, a') \sum_{s''} P(s', a', s'') R(s', a', s'')$$

+ ...

Example on Value function



$\gamma: 0.5$
 $\pi_1: \text{left always}$

$\pi_2: \text{Right always}$

$$v^{\pi_1}(s_1) = 0$$

$$v^{\pi_1}(s_2) = 12\gamma^0$$

$$v^{\pi_1}(s_3) = 12\gamma^1$$

$$v^{\pi_1}(s_4) = 12\gamma^2$$

$$v^{\pi_1}(s_5) = 1.5$$

$$v^{\pi_1}(s_6) = 0$$

$$v^{\pi_2}(s_1) = 0$$

$$v^{\pi_2}(s_2) = 1/4$$

$$v^{\pi_2}(s_3) = 1/2$$

$$v^{\pi_2}(s_4) = 1$$

$$v^{\pi_2}(s_5) = 2$$

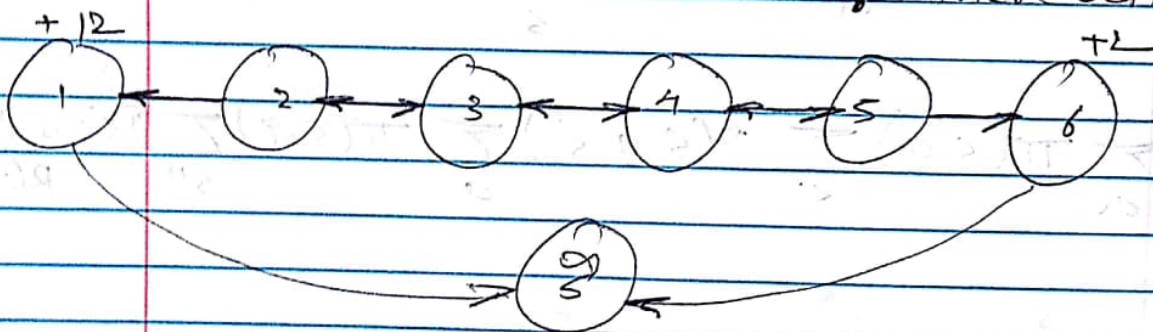
$$v^{\pi_2}(s_6) = 0$$

Action value function / state-action value $q^\pi / q-f^\pi$

$$q^\pi: S \times A \Rightarrow \mathbb{R}$$

$$q^\pi(s, a) \triangleq \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k} \mid S_t = s, A_t = a, \pi \right]$$

= Expected discounted return from state s if the agent takes action a & follows π , thereafter

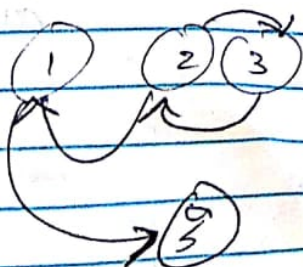


$$q^{\pi_1}(s, L) = 0, \quad q^{\pi_2}(s, R) =$$

$$q^{\pi_1}(s_1, R) = 0$$

$$q^{\pi_1}(s_2, L) = 12$$

$$q^{\pi_1}(s_2, R) = 3$$



Bellman Equation for V^π

$$V^\pi(s) \stackrel{||A}{=} \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k} \mid S_t = s, \pi \right]$$

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

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

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

$$= \sum_a \pi(s, a) \sum_{s'} P(s, a, s') \times$$

$$\left(R(s, a, s') + \gamma \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \mid S_t = s, \pi \right] \right)$$

$$= \sum_a \pi(s, a) \sum_{s'} P(s, a, s') \left[R(s, a, s') + \gamma V^\pi(s) \right]$$

$$= \mathbb{E} \left[R_t + \gamma V^\pi(S_{t+1}) \mid S_t = s, \pi \right]$$

$$= R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} \dots$$