

# CH3. Finite Markov Decision Process

\* MDP → 주어진 rewards, situations, state, 미래의 rewards  
 $r_t$        $x_t$        $s_t$        $r_{t+1}$

bandit 문제: 각 행동에 대해  $q_*(a)$ 를 예측

MDP: 각 상태에 대해  $q_*(s,a)$ 를 예측  
 주어진 초기의  $a$  선택에 대해 각  $s$ 의  $V_*(s)$ 를 예측

)  $s$  dependent → credit assignment

\* Agent-Environment Interface

agent: learner & decision maker

environment: agent를 향한, interaction에 모든 것

$R_t, S_t$ : discrete probability distributions

depends only on  $S, A$

$P(S', r | S, A)$ : dynamics of MDP  
 $\sum_{S' \in S} \sum_{r \in R} P(S', r | S, A) = 1$

$P(S'|S, A)$ : state-transition probabilities  
 $S \times S \times A \rightarrow [0, 1]$

$$= \sum_r P(s', r | s, a)$$

$r(s, a)$ : expected rewards for state-action pairs

$$S \times A \rightarrow R$$

$$= E[R_t | S_{t-1} = s, A_{t-1} = a] = \sum_r r P(r | s, a) = \sum_r r \sum_{S'} P(s' | r, s, a)$$

$r(s, a, s')$ : expected rewards for state-action-next-state triples

$$S \times A \times S \rightarrow R$$

$$= E[R_t | S_{t-1} = s, A_{t-1} = a, S_t = s'] = \sum_r r P(r | s, a, s') = \sum_r r \frac{P(s' | r, s, a)}{P(s' | s, a)}$$

\* Goals and Rewards

Goal = maximize total reward

= maximize the expected value of the cumulative sum of a received reward

return  $G_t = R_{t+1} + R_{t+2} + \dots + R_T$  final time step

episode: terminal state or  $\infty$

( $T \neq \infty$ : episodic task

$T = \infty$ : continuing task

discounting

discounted return  $G_t = R_{t+1} + r R_{t+2} + r^2 R_{t+3} + \dots = \sum_{k=0}^{\infty} r^k R_{t+k+1}$  ( $0 \leq r \leq 1$ : discount rate)

$$= R_{t+1} + r G_{t+1}$$

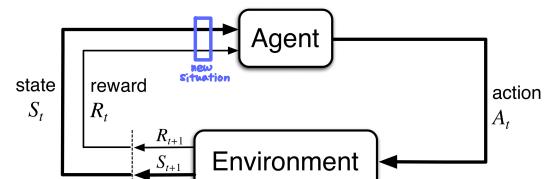
$$= \sum_{k=t+1}^{T-1} r^{k-t-1} R_k < T \neq \infty : \text{episodic}$$

$$T = \infty : \text{continuing}$$

myopic  $\rightarrow$ 长远眼光

discounting ↑

$$R_{t+1} = 0.01 \text{ 확률 } S_{t+1} = 0$$



discrete time step  $t$ ,  $(S, A, R)$ : finite number

## \* Policies and Value Functions

policy : mapping state — probabilities of selecting action  
 $= \pi(a|s)$

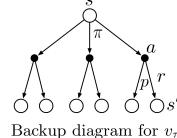
• Value function for policy  $\pi$

state-value function:  $V_{\pi}(s) = \mathbb{E}_{\pi} \{G_t | S_t = s\} = \mathbb{E}_{\pi} \left\{ \sum_{k=0}^{\infty} \delta^k R_{t+k+1} | S_t = s \right\}$

action-value function:  $q_{\pi}(s, a) = \mathbb{E}_{\pi} \{ G_t + \gamma R_{t+1} | S_t = s, A_t = a \} = \mathbb{E}_{\pi} \{ \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s, A_t = a \}$

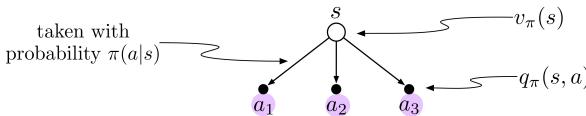
→ estimate w/ Monte-Carlo methods

$$\begin{aligned}
 V_{\pi}(s) &= \mathbb{E}[G_{t+1}|s] = \mathbb{E}[R_{t+1} + \gamma G_{t+1}|s] \\
 &= \sum_a \pi(a|s) \sum_{s'} \sum_r P(s', r | s, a) (r + \gamma \mathbb{E}_{\pi'}[G_{t+1}|S_{t+1}=s']) \\
 &= \sum_a \pi(a|s) \sum_{s'} \sum_r P(s', r | s, a) (r + \gamma V_{\pi}(s')) : \text{Bellman eq}
 \end{aligned}$$

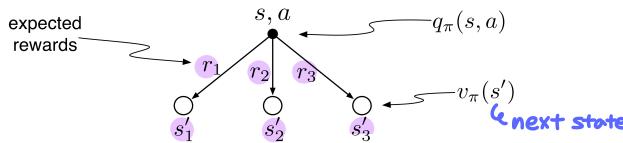


Backup diagram for  $v_7$

$\theta_{\pi}(s,a)$



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



## \* Optimal Policies and Optimal Value Functions

$$\pi \Sigma \pi' \longleftrightarrow V_\pi(s) \geq V_{\pi'}(s)$$

모든 다른 policies보다 크거나 같은 최소 하나의 policy = optimal policy

$$\pi^* \geq \text{all other } \pi \rightarrow U_{\pi^*}(s) \geq U_\pi(s)$$

$\therefore V_*(s) = \max_{\pi} V_{\pi}(s)$  : optimal state-value function

$q_{\pi^*}(s,a) = \max_{\pi} q_{\pi}(s,a)$ : optimal action-value function

$$= \{E[R_{t+1} + \delta V_t(S_{t+1}) | S_t = s, A_t = a]\}$$

greedy: 확률 고려 않고 다른 길 가보는 것

## Bellman optimality equation

$$V_{\pi}(s) = \max_a q_{\pi_{\pi}}(s, a) = \max_a \mathbb{E}_{\pi^a} [G_t | S_t = s, A_t = a]$$

$$= \max_a \mathbb{E}_{\pi^a} [R_{t+1} + \gamma G_{t+1} | S_t = s, A_t = a]$$

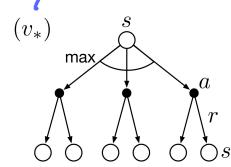
$$= \max_a [E f(R_{t+1}) + \gamma V_\pi(S_{t+1}) | S_t = s, A_t = a]$$

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

$$Q_{\pi}^*(s, a) = \mathbb{E}\{R_{t+1} + \gamma \max_{a'} Q_{\pi}(s_{t+1}, a') \mid s_t = s, A_t = a\}$$

$$= \sum_{s' \in R} p(s', r | s, a) \{ r + \gamma \max_{a'} q_{\pi}(s', a') \}$$

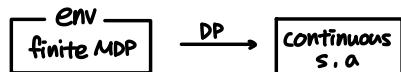
unique solution



Short-term  $\xrightarrow{a}$  [E<sup>f</sup> long-term]  
one-step-ahead search

$E_i$  {long-term return}

# CH4. Dynamic Programming



s, a, r  
dynamics  
P(s', r | s, a)

$$\begin{aligned} V_{\pi}(s) &= \max_a \{ \mathbb{E}[R_{t+1} + \gamma V_{\pi}(s_{t+1})] \mid S_t = s, A_t = a \} \\ &= \max_a \sum_{s', r} P(s', r | s, a) (r + \gamma V_{\pi}(s')) \\ q_{\pi}(s, a) &= \mathbb{E}[R_{t+1} + \gamma \max_{a'} q_{\pi}(s_{t+1}, a') \mid S_t = s, A_t = a] \\ &= \sum_{s', r} P(s', r | s, a) (r + \gamma \max_{a'} q_{\pi}(s', a')) \end{aligned}$$

\* env를 단계별로 만드는 approximation for  $V_{\pi}$

$$V_{k+1}(s) = \mathbb{E}_{\pi} [R_{t+1} + \gamma V_k(s_{t+1}) \mid S_t = s] = \sum_a \pi(a|s) \sum_{s', r} P(s', r | s, a) (r + \gamma V_k(s'))$$

↑ 가수화된 next state on each expected update

$\rightarrow \infty : V_k \rightarrow V^*$

In place · Old  $\frac{1}{2}$  new  $\frac{1}{2}$  overwrite Sweep

$$\max_s |V_{k+1}(s) - V_k(s)| > \frac{\epsilon}{2} \text{ 까지} \text{ 깜아지면 stop}$$

\* policy improvement

$$\begin{aligned} V_{\pi}(s) &\leq q_{\pi}(s, \pi'(s)) = \mathbb{E}_{\pi'} [R_{t+1} + \gamma V_{\pi}(s_{t+1}) \mid S_t = s] \\ &\leq \mathbb{E}_{\pi'} [R_{t+1} + \gamma q_{\pi}(s_{t+1}, \pi'(s_{t+1})) \mid S_t = s] = \mathbb{E}_{\pi'} [R_{t+1} + \gamma R_{t+2} + \gamma^2 V_{\pi}(s_{t+2}) \mid S_t = s] \\ &\leq \dots \leq \mathbb{E}_{\pi'} [R_{t+1} + \gamma R_{t+2} + \dots \mid S_t = s] = V_{\pi'}(s) \end{aligned}$$

$$\pi'(s) = \arg \max_a q_{\pi}(s, a) = \arg \max_a \{ \mathbb{E}[R_{t+1} + \gamma V_{\pi}(s_{t+1}) \mid S_t = s, A_t = a] \}$$

$$\begin{array}{ccccccc} \pi_0 & \xrightarrow[E]{\quad} & v_{\pi_0} & \xrightarrow[I]{\quad} & \pi_1 & \xrightarrow[E]{\quad} & v_{\pi_1} \\ & \text{Evaluation} & & & & \text{Improvement} & \\ & \max_s |V_{k+1}(s) - V_k(s)| & \text{만족: 각을 대비하기} & & & & : \pi' = \pi \text{일 때까지} \end{array}$$

\* value iteration

$$\begin{aligned} V_{k+1} &= \max_a \{ \mathbb{E}[R_{t+1} + \gamma V_k(s_{t+1}) \mid S_t = s, A_t = a] \} \\ &= \max_a \sum_{s', r} P(s', r | s, a) (r + \gamma V_k(s')) \end{aligned}$$

\* generalized policy iteration (GPI)

