

Bellman equation

24.01.11
정상혁

Value functions

- State-value function

$$v_{\pi}(s)$$

: state ‘ s ’ 가 **얼마나 좋은지**를 나타내는 함수.

- Action-value function

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

: state ‘ s ’ 에서 action ‘ a ’를 선택하는 것이 **얼마나 좋은지**를 나타내는 함수

‘**좋다**’ 의 기준

[Reward Hypothesis]

All goals can be described by the **maximization** of the **expected value** of the cumulative sum of rewards (=return)

Value functions

- State-value function

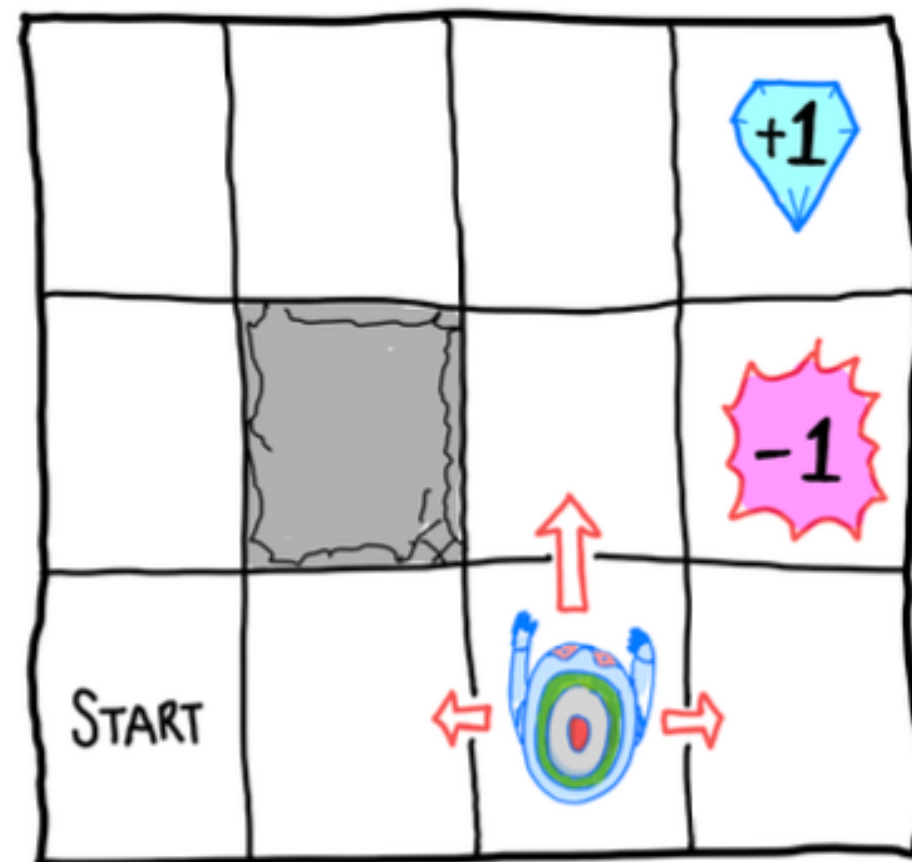
$$v_{\pi}(s) = E_{\pi} [G_t \mid S_t = s]$$

- Action-value function

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

Advantage function

$$A_{\pi}(s, a) = q_{\pi}(s, a) - v_{\pi}(s)$$



Episode $(1, 1) \xrightarrow{\text{north}} (1, 2) \xrightarrow[\text{wall}]{\text{east}} (1, 2) \xrightarrow{\text{north}} (1, 3) \xrightarrow[\text{east 10\%}]{\text{south}} (2, 3) \xrightarrow{\text{east}} (3, 3) \xrightarrow{\text{east}} (4, 3)$

EX)

[Law of total probability]

- $P(A) = \sum_n P(A | B_n) \cdot P(B_n)$
- $\mathbb{E}[X] = \sum_y \mathbb{E}[X | Y = y] \cdot P(Y = y)$
- $\mathbb{E}[X | Z = z] = \sum_y \mathbb{E}[X | Y = y, Z = z] \cdot P(Y = y | Z = z)$ \longrightarrow

$$v_\pi(s) = E_\pi [G_t | S_t = s]$$

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



$$v_\pi(s) = \sum_a \pi(a | s) q_\pi(s, a)$$

[Law of large numbers]

The average of the results obtained from a large number of trials should be close to the expected value, and will tend to become closer as more trials are performed.

$X : X_1, \dots, X_n$ i.i.d random samples (independent and identically distributed)

$$\bar{X} = \frac{1}{n}(X_1 + \dots + X_n) \rightarrow \mathbb{E}[X] \text{ as } n \rightarrow \infty$$

EX)

$$\frac{1}{n} (G_1 + \dots + G_n) \longrightarrow q_\pi(s, a) = E_\pi [G_t | S_t = s, A_t = a] \quad \text{as } n \rightarrow \infty$$

Optimal value functions and policy

- Optimal state-value function $v_*(s) = \max_{\pi} v_{\pi}(s)$
 - Optimal action-value function $q_*(s, a) = \max_{\pi} q_{\pi}(s, a)$
 - Optimal policy $\pi_*(s)$ ($\pi_*(s) \geq \pi$ for all π)
- (Define a partial policies $\pi' \geq \pi$ if $v_{\pi'}(s) \geq v_{\pi}(s)$ for all s)

[Theorem] Any MDP satisfies the followings.

- There **exists** an **optimal policy** $\pi_* \geq \pi$ for all π .
- All optimal policies achieve the optimal state-value function $v_{\pi_*}(s) = v_*(s)$.
- All optimal policies achieve the optimal action-value funct. $q_{\pi_*}(s, a) = q_*(s, a)$.

Optimal value functions and policy

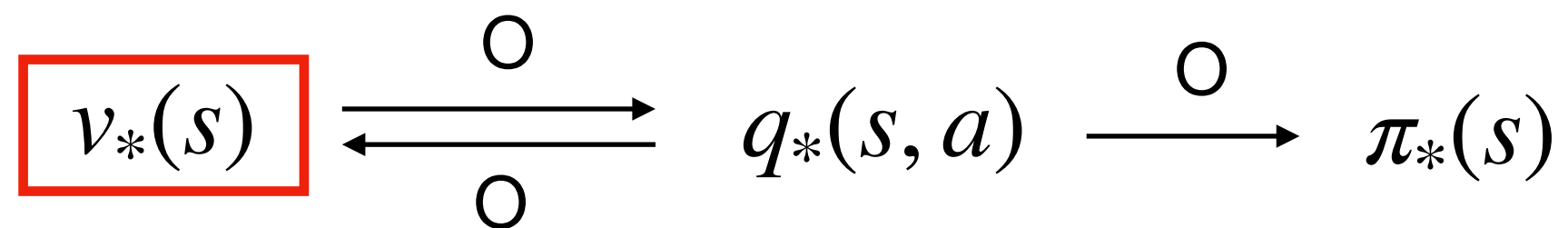
- Optimal action-value function 을 알면 Optimal policy 를 구할 수 있다.

$$\pi_*(s) = \underset{a}{\operatorname{arg\,max}} q_*(s, a) \qquad \pi_*(a | s) = \begin{cases} 1 & \text{if } a = \operatorname{arg\,max}_a q_*(s, a) \\ 0 & \text{otherwise} \end{cases}$$

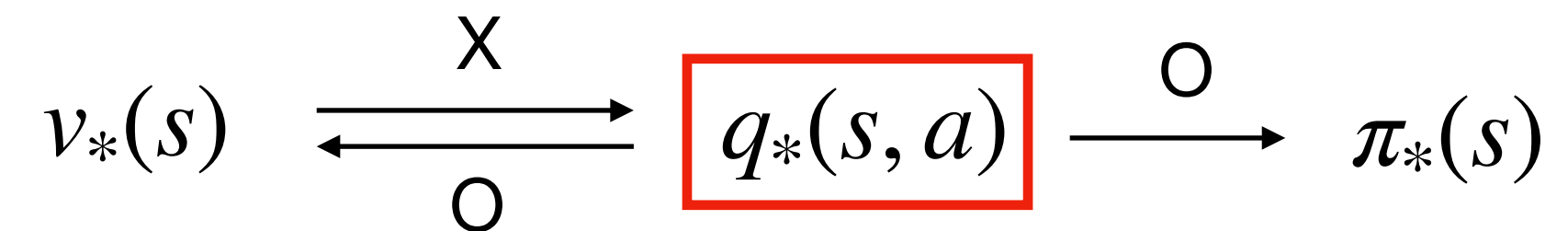
- Optimal state-value function 을 알면 Optimal policy 를 구할 수 있을까 ?

$$v_*(s) = \max_a q_*(s, a) \qquad q_*(s, a) = r(s, a) + \gamma \sum_{s'} p(s' | s, a) v_*(s')$$

Known MDP



Unknown MDP



Bellman equation

Bellman expectation equation (Markov property)

$$v_{\pi}(s) = E_{\pi} [G_t \mid S_t = s]$$

$$= E_{\pi} [R_{t+1} + \gamma v_{\pi}(S_{t+1}) \mid S_t = s]$$

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

$$= E_{\pi} [R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}) \mid S_t = s, A_t = a]$$

Bellman optimality equation

$$v_{*}(s) = \max_{\pi} v_{\pi}(s)$$

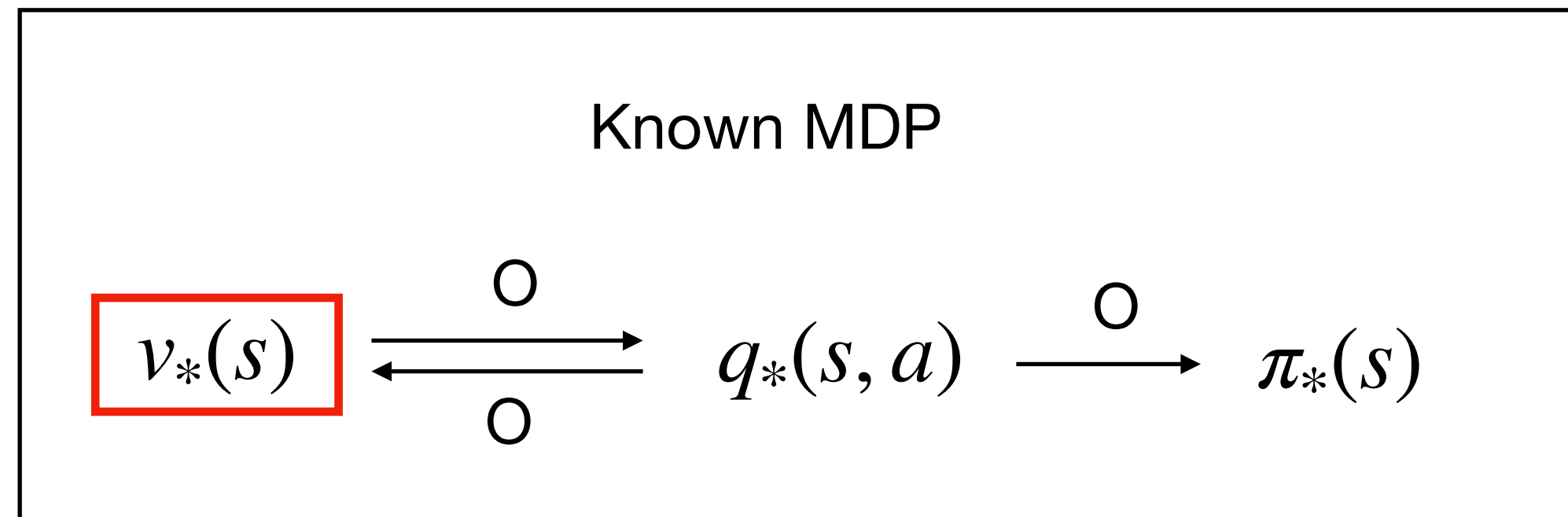
$$= \max_a \sum_{s', r} p(s', r \mid s, a) [r + \gamma v_{*}(s')]$$

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

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

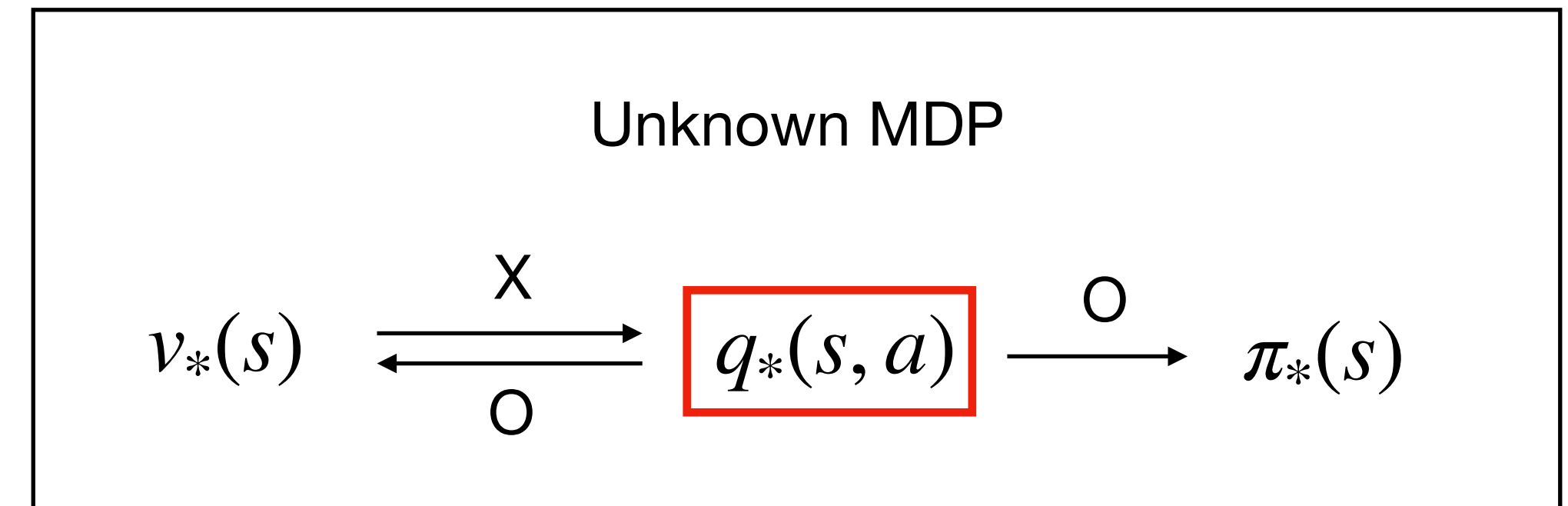
$$= \sum_{s', r} p(s', r \mid s, a) [r + \gamma \max_{a'} q_{*}(s', a')]$$

Bellman equation



Dynamic Programming 사용

- $V(s)$ (State-value function table)을 반복해서 update 하며, optimal policy를 구한다.
- 이때, $v_*(s)$ 를 구하기 위해서, Bellman equation 을 이용한다.



Reinforcement Learning 사용

- $Q(s, a)$ (Action-value function table)을 반복해서 update 하며, optimal policy를 구한다.
- 이때, $q_*(s, a)$ 를 근사하기 위해서, random samples, Bellman equation 을 이용한다.

