

강화학습 맛보기

정 태 수

고려대학교 산업경영공학부
tcheong@korea.ac.kr

$$\alpha=1, \gamma=1$$

$$Q(s_t, a_t) \leftarrow (1 - \alpha) \underbrace{Q(s_t, a_t)}_{\text{기존 정보}} + \alpha \underbrace{(r_{t+1} + \gamma \max_a Q(s_{t+1}, a))}_{\hat{Q}(s_t, a_t) \text{ 새로운 정보}}$$

추정치

기존 정보

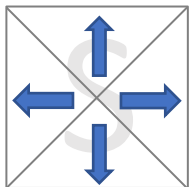
$\hat{Q}(s_t, a_t)$ 새로운 정보



강화학습

$Q(s_t, a_t) \rightarrow$ State s_t 에서 Action a_t 선택 시 얻을 수 있는 최대 누적 보상의 기대치

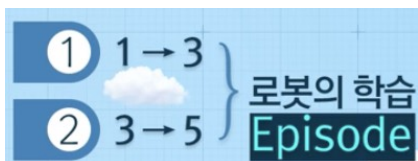
$Q(s_t, a_t)$



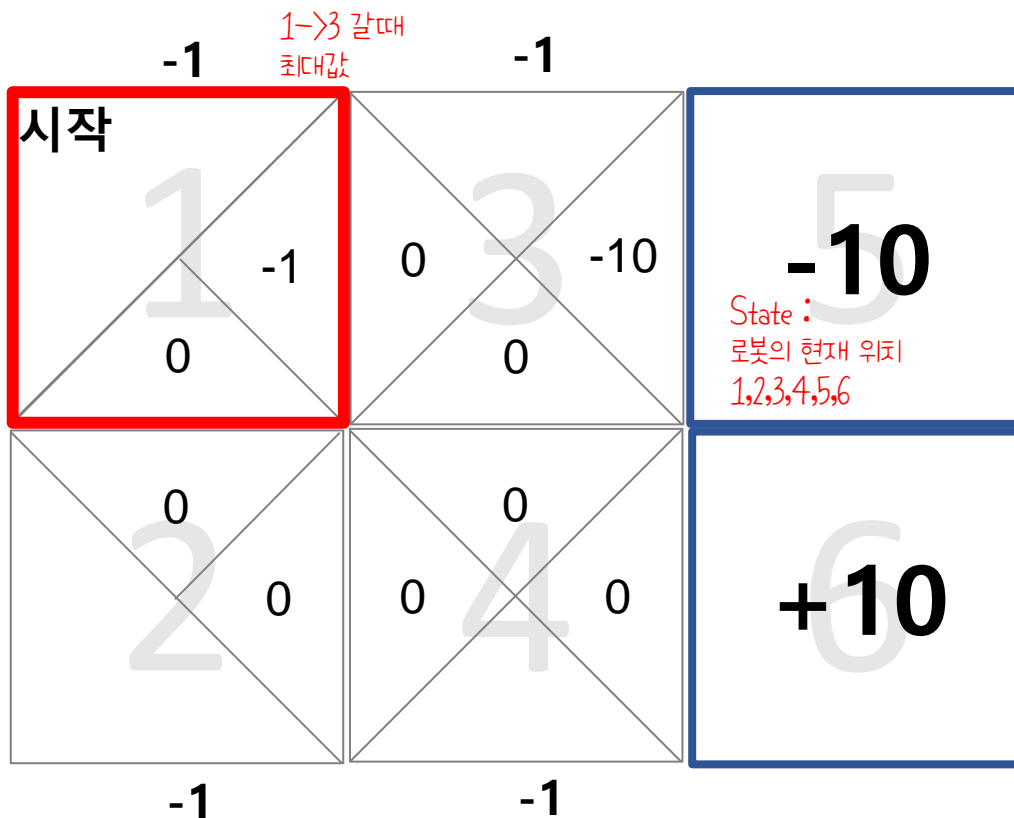
$$Q(s_t, a_t) \leftarrow (1 - \alpha)Q(s_t, a_t) + \alpha \left(r_{t+1} + \underbrace{\gamma \max_a Q(s_{t+1}, a)}_{\hat{Q}(s_t, a_t)} \right)$$

$\alpha = 1, \gamma = 1$
감가율

$(1, \rightarrow) (3, \rightarrow) 5$



러닝을 통해 어떤 정보를 업데이트 할까?



3번 상태에서 획득 가능한 누적 보상합의 최대 값을 활용

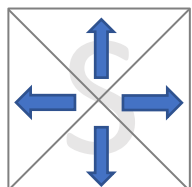
현재 1번 상태에서 우측으로 갔을 때 획득 가능한 누적 보상의 최대 값 추정

$(1, \downarrow) (2, \rightarrow) (4, \rightarrow) 6$



강화학습

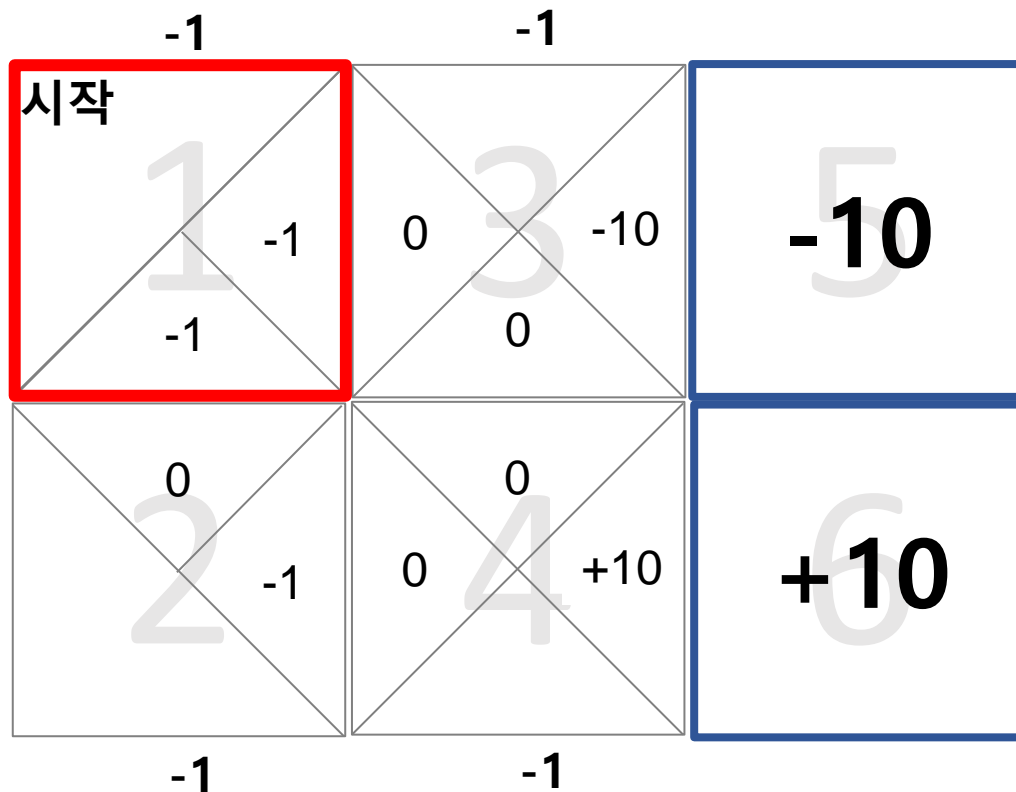
$Q(s_t, a_t)$



$$Q(s_t, a_t) \leftarrow (1 - \alpha)Q(s_t, a_t) + \alpha \left(r_{t+1} + \underbrace{\gamma \max_a Q(s_{t+1}, a)}_{\hat{Q}(s_t, a_t)} \right)$$

$\alpha = 1, \gamma = 1$

(1, ↓) (2, →) (4, →) 6

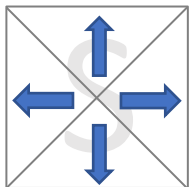


(1, ↓) (2, →) (4, ↑) (3, →) 5



강화학습

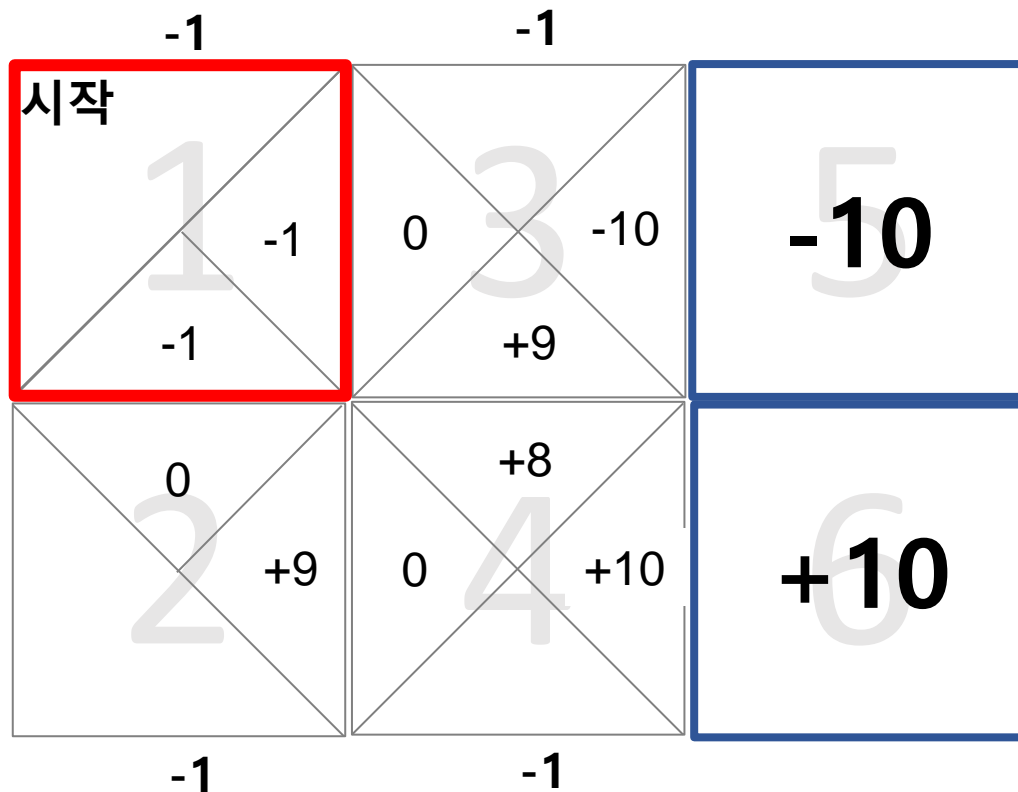
$Q(s_t, a_t)$



$$Q(s_t, a_t) \leftarrow (1 - \alpha)Q(s_t, a_t) + \alpha \left(\underbrace{r_{t+1} + \gamma \max_a Q(s_{t+1}, a)}_{\hat{Q}(s_t, a_t)} \right)$$

$\alpha = 1, \gamma = 1$

(1, →) (3, ↓) (4, ↑)
(3, →) 5

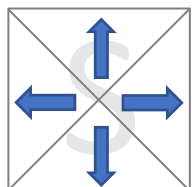


(1, ↓) (2, →) (4, ←) (2, ↑) (1, →) (3, →) 5



강화학습

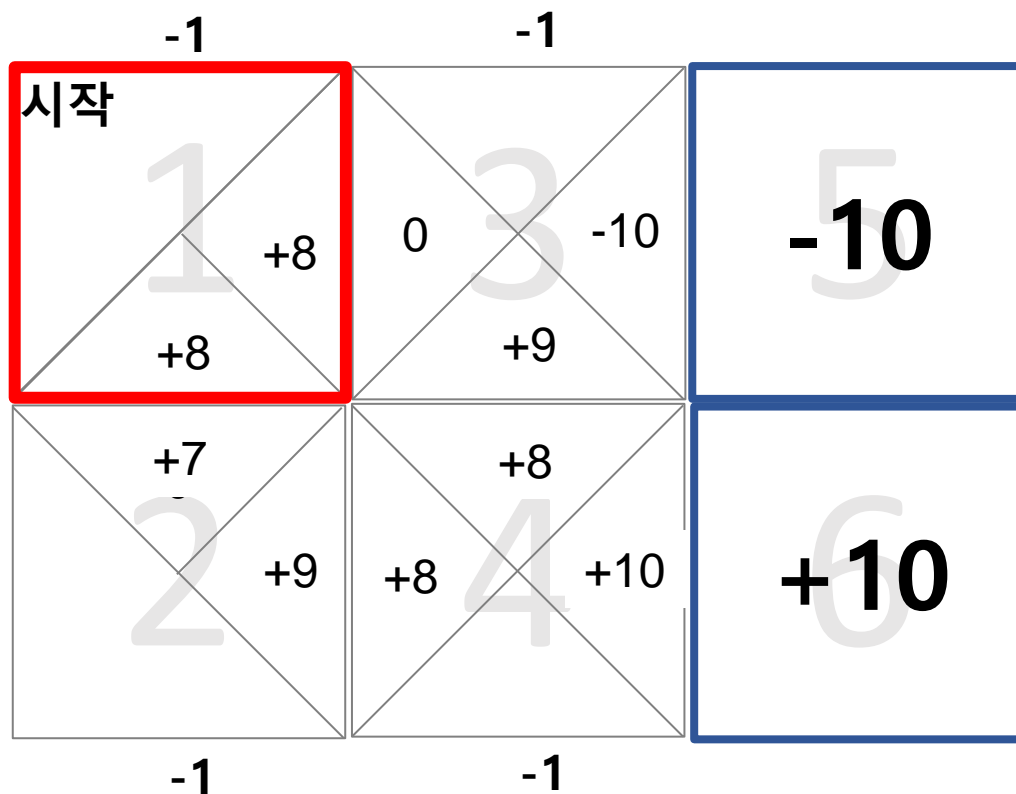
$Q(s_t, a_t)$



$$Q(s_t, a_t) \leftarrow (1 - \alpha)Q(s_t, a_t) + \alpha \left(\underbrace{r_{t+1} + \gamma \max_a Q(s_{t+1}, a)}_{\hat{Q}(s_t, a_t)} \right)$$

$\alpha = 1, \gamma = 1$

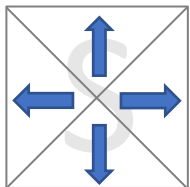
(1, ↓) (2, →) (4, ←)
(2, ↑) (1, →) (3, →) 5





강화학습

$Q(s_t, a_t)$



$$Q(s_t, a_t) \leftarrow (1 - \alpha)Q(s_t, a_t) + \alpha \left(\underbrace{r_{t+1} + \gamma \max_a Q(s_{t+1}, a)}_{\hat{Q}(s_t, a_t)} \right)$$

$\alpha = 1, \gamma = 1$

