

## # Lecture 5.

Model-free  $\rightarrow$  MDP X

Policy evaluation (O)

다음 state를 알지 X  $\rightarrow$  greedy policy X

Action-Value: Q-function

Value-function  $\rightarrow$  evaluation X

Q-function  $\rightarrow$  evaluation O

$\hookrightarrow$  state에서 할 수 있는 action 중 Q값이 가장 높은 것을 선택.

greedy하게만 선택하면 문제가 발생

Improvement

- $\epsilon$ -Greedy Exploration

$\epsilon$ 의 확률로 random하게 다른 선택을 해봄.

- $\epsilon$ -Greedy Improvement

$\pi'$ 이  $\pi$ 에 대해서 improvement된 policy  $\rightarrow V_{\pi'}(s) \geq V_{\pi}(s)$

$$q_{\pi}(s, \pi'(s)) = \sum_{a \in A} \pi'(a|s) q_{\pi}(s, a)$$

$$= \epsilon/m \sum_{a \in A} q_{\pi}(s, a) + (1-\epsilon) \max_{a \in A} q_{\pi}(s, a)$$

- Sarsa



$$Q(S, A) \leftarrow Q(S, A) + \alpha(R + rQ(S', A') - Q(S, A))$$

→ policy evaluation에 사용

↳ 한 step으로 Q를 업데이트 좀 극단적!!

- GLIE : 모든 state-action pair를 방문, 결국 Greedy policy로 수렴

- Robbins-Monro :

$$\circ \sum_{t=1}^{\infty} \alpha_t = \infty$$

$$\circ \sum_{t=1}^{\infty} \alpha_t^2 < \infty$$



- $n$ -step Sarsa

$$n=1 \text{ (Sarsa)} \quad q_t^{(1)} = R_{t+1} + \gamma Q(S_{t+1})$$

$$n=2 \quad q_t^{(2)} = R_{t+1} + R_{t+2} + \gamma Q(S_{t+2})$$

$\vdots$

$$n=\infty \text{ (MC)} \quad q_t^{(\infty)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{T-1} R_T$$

- Define the  $n$ -step Q-return

$$q_t^{(n)} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n Q(S_{t+n})$$

- $n$ -step Sarsa updates  $Q(s, a)$  towards the  $n$ -step Q-return

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (q_t^{(n)} - Q(S_t, A_t))$$

- 1-step Sarsa

- Sarsa( $\lambda$ ): 지난 모든 경로를 update 함

- Importance Sampling

$$\begin{aligned} E_{x \sim p} [f(x)] &= \sum p(x) f(x) \\ &= \sum q(x) \frac{p(x)}{q(x)} f(x) \\ &= E_{x \sim q} \left[ \frac{p(x)}{q(x)} f(x) \right] \end{aligned}$$

- Important Sampling for Off-policy Monte-Carlo

- Use returns generated from  $\mu$  to evaluate  $\pi$
- Weight return  $G_t$  according to similarity between policies
- Multiply importance sampling corrections along whole episode

$$G_t^{\pi/\mu} = \frac{\pi(A_t|S_t)}{\mu(A_t|S_t)} \frac{\pi(A_{t+1}|S_{t+1})}{\mu(A_{t+1}|S_{t+1})} \cdots \frac{\pi(A_T|S_T)}{\mu(A_T|S_T)} G_t$$

- Update value towards *corrected* return

$$V(S_t) \leftarrow V(S_t) + \alpha \left( \underbrace{G_t^{\pi/\mu}}_{\text{corrected return}} - V(S_t) \right)$$

- Cannot use if  $\mu$  is zero when  $\pi$  is non-zero
- Importance sampling can dramatically increase variance

- Important Sampling for Off-policy TD

- Use TD targets generated from  $\mu$  to evaluate  $\pi$
- Weight TD target  $R + \gamma V(S')$  by importance sampling
- Only need a single importance sampling correction

$$V(S_t) \leftarrow V(S_t) + \alpha \left( \frac{\pi(A_t|S_t)}{\mu(A_t|S_t)} (R_{t+1} + \gamma V(S_{t+1})) - V(S_t) \right)$$

- Much lower variance than Monte-Carlo importance sampling
- Policies only need to be similar over a single step

- Q-Learning

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (\underbrace{R_{t+1} + \gamma Q(S_{t+1}, A')}_{\text{추측치}} - Q(S_t, A_t))$$

- Off-policy Control with Q-Learning

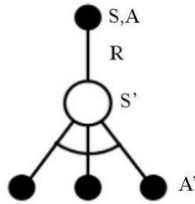
- We now allow both behaviour and target policies to **improve**
- The target policy  $\pi$  is **greedy** w.r.t.  $Q(s, a)$

$$\pi(S_{t+1}) = \operatorname{argmax}_{a'} Q(S_{t+1}, a')$$

- The behaviour policy  $\mu$  is e.g.  **$\epsilon$ -greedy** w.r.t.  $Q(s, a)$
- The Q-learning target then simplifies:

$$\begin{aligned} & R_{t+1} + \gamma Q(S_{t+1}, A') \\ &= R_{t+1} + \gamma Q(S_{t+1}, \operatorname{argmax}_{a'} Q(S_{t+1}, a')) \\ &= R_{t+1} + \max_{a'} \gamma Q(S_{t+1}, a') \end{aligned}$$

- Q-Learning Control Algorithm



$$Q(S, A) \leftarrow Q(S, A) + \alpha \left( R + \gamma \max_{a'} Q(S', a') - Q(S, A) \right)$$

### Theorem

*Q-learning control converges to the optimal action-value function,  $Q(s, a) \rightarrow q_*(s, a)$*

## Relationship Between DP and TD

	Full Backup (DP)	Sample Backup (TD)
Bellman Expectation Equation for $v_\pi(s)$	<p>Iterative Policy Evaluation</p>	<p>TD Learning</p>
Bellman Expectation Equation for $q_\pi(s, a)$	<p>Q-Policy Iteration</p>	<p>Sarsa</p>
Bellman Optimality Equation for $q_*(s, a)$	<p>Q-Value Iteration</p>	<p>Q-Learning</p>