

Model-Free Control

→ 기본이 되는 Monte-Carlo Control \Rightarrow ϵ -greedy Exploration

$$\pi(a|s) = \begin{cases} \epsilon/m + 1-\epsilon & \text{if } a^* = \underset{a \in A}{\operatorname{argmax}} Q(s,a) \rightarrow 1-\epsilon \\ \epsilon/m & \text{otherwise} \end{cases} \quad \begin{matrix} \text{탐색 (exploration)} \\ \text{exploitation} \end{matrix}$$

(공유 이익)
 ϵ 확률로 랜덤 선택
 \rightarrow 모든 행동 탐색 가능

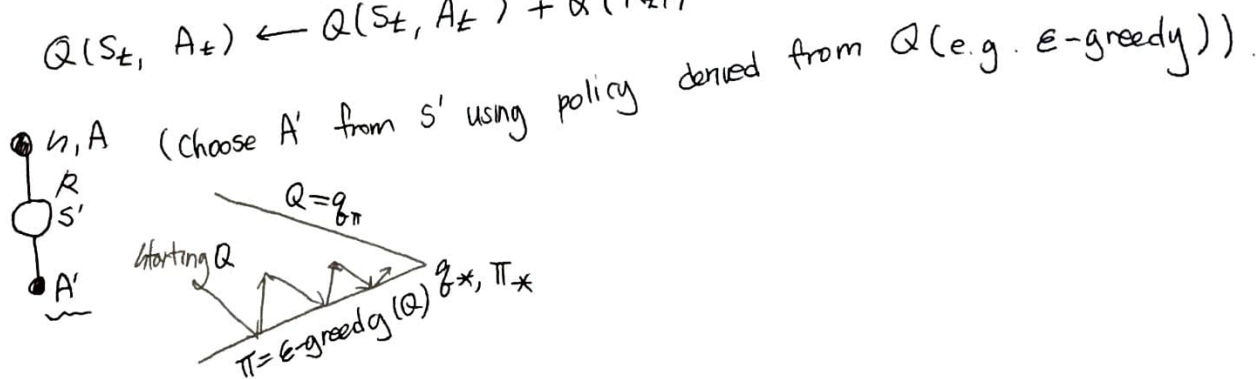
< Policy evaluation : Monte-Carlo policy evaluation $Q \approx q_\pi$
Policy Improvement : ϵ -greedy policy improvement

* Off-policy

학습 잘되는데, episode마다 ~~학습~~ learns only from the tails of episodes.
 \rightarrow learning will be slow

* SARSA : on-policy on TD Control $\forall \pi, Q$ 를 사용.

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha (R_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t))$$



n-step
SARSA

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha (q_t^{(n)} - Q(s_t, a_t))$$

Forward-view
SARSA(γ)

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha (q_t^\gamma - Q(s_t, a_t))$$

단점) 끝날 때까지 끝난게 아니다. \rightarrow 추가) n-step off-policy TD

Backward-view
SARSA(γ)

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha \delta_t E_t(s_t, a_t)$$

$$\delta_t = R_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)$$

$$E_t(s, a) = \gamma \lambda E_{t+1}(s, a) + 1 (s_t = s, a_t = a)$$

decay 추가

$$E_0(s, a) = 0$$

* Q-learning: off-policy on TD control.

$$Q(s_t, A_t) \leftarrow Q(s_t, A_t) + \alpha \left(\overbrace{R_{t+1} + \gamma \max_a Q(s_{t+1}, a)}^{\text{error}} - \underbrace{Q(s_t, A_t)}_{\text{target}} \right)$$

behavior: ϵ -greedy

target: greedy

$$\pi(s_{t+1}) = \arg\max_{a'} Q(s_{t+1}, a')$$

$$\Rightarrow R_{t+1} + \max_{a'} \gamma Q(s_{t+1}, a')$$

(S_{t+1} 까지 환경과 interaction, SARSA는 환경과 interaction)
↳ Q-value iteration



왜 importance sampling이 되어있지 않은가? 여기에 처리됨.

Expected HARGA. ^{expected value로 인해} off-policy로 작용, Q-learning 보다 효과적.

$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \dots + \gamma^{n-1} R_{t+n} + \gamma^n \bar{V}_{t+n-1}(s_{t+n})$$

$$(\bar{V}_{t+n-1}(s_{t+n}) = \sum_a \pi(a|s_{t+n}) Q_{t+n-1}(s_{t+n}, a), t+n < T)$$

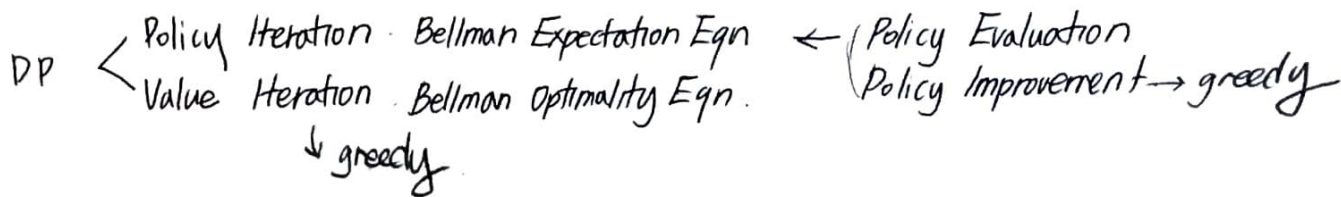
$$Q(s_t, A_t) \leftarrow Q(s_t, A_t) + \alpha (R_{t+1} + \gamma \overset{\text{Q-learning의 max와 동일하게 적용}}{E}[Q(s_{t+1}, a) | s_{t+1}] - Q(s_t, A_t))$$

$$Q(s_t, A_t) \leftarrow Q(s_t, A_t) + \alpha (R_{t+1} + \gamma \sum_a \pi(a|s_{t+1}) Q(s_{t+1}, a) - Q(s_t, A_t))$$

Prediction & Control in DP

state	1	2	3	4
action	1	2	3	4
reward	20	10	10	9
cost	8	9	10	0

- Prediction: 현재 진행하는 policy에 따라 value function 구하기 → policy evaluation 7/7 번
- Control: policy를 optimal 하게 변경. 종료. (policy를 따르게 했는지 아닌지 판단 후 update) → policy improvement 1번
↳ greedy



- Policy Improvement: state value를 바탕으로 action-value function을 이용해 policy에 좋은 action 선택하여 policy update.
 greedy max: greedy policy improvement.

- Value Iteration: 이동가능한 state s' 에 대해 max를 취해 greedy하게 improve.

$$V_{k+1}(s) = \max_{a \in A} \left(R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_k(s') \right)$$

↔ Policy evaluation: $V_{k+1}(s) = \sum \pi(a|s) \left(R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_k(s') \right)$ 합