

* Model Free Control

- monte - Carlo Control
- Sarsa
- Q-learning

① monte - Carlo Control

monte - Carlo prediction을 통해 Policy Evaluation을 수행했다.

이제 Policy improvement를 해야지! → Greedy Policy 이용!

$$\dots \pi'(s) = \underset{a}{\operatorname{argmax}} Q_{\pi}(s, a)$$

$$= \underset{a}{\operatorname{argmax}} R_s^a + \gamma \sum_s P_{ss}^a V_{\pi}(s)$$

→ 어쨌든... 모르지! 다른 방법?...
이런 Evaluation은 V_{π} 가 아니라 Q_{π} 를 찾자!

• ϵ - Greedy Exploration → 다른 Action도 작은 확률로해보겠다.

$$\pi(a|s) = \begin{cases} \frac{\epsilon}{4n} + (1-\epsilon) & \text{if } a = \underset{a}{\operatorname{argmax}} Q_{\pi}(s, a) \\ \frac{\epsilon}{4n} & \text{else} \end{cases}$$

← 해당 state only

기존 Greedy Policy

$$\pi(s) = \underset{a}{\operatorname{argmax}} Q_{\pi}(s, a)$$

② Sarsa

→ MC on-line TD!

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (R_{t+1} + \gamma Q(S_{t+1}, A_{t+1}) - Q(S_t, A_t))$$

③ Q-learning

→ off-policy learning!

← $\mu(a|s)$ 을 따라 움직이며 $\pi(a|s)$ 을 평가!

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

↓ improve

μ 와 π 를 모두 update!

$$\pi(S_{t+1}) = \underset{a'}{\operatorname{argmax}} Q(S_{t+1}, a'), \quad \mu \leftarrow \epsilon\text{-greedy!}$$

Greedy!

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

$$\therefore Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha (R_{t+1} + \max_{a'} \gamma \cdot Q(S_{t+1}, a') - Q(S_t, A_t))$$