# 나는 강화학습으로 축구한다

## 기초 이론

2025. 01. 19 - 20

# 1. 강화학습이란?

# 1. 강화학습이란? : 지도학습 vs. 비지도학습 vs. 강화학습

## 지도학습

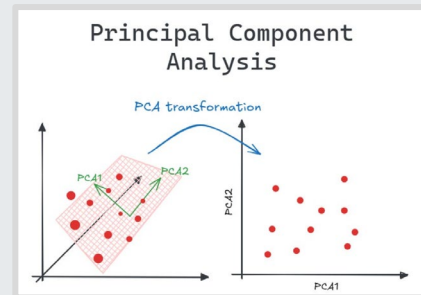$$f^* = \text{argmin}_f \, \mathbb{E} \left[ \mathcal{L}(f(\boldsymbol{x}), \boldsymbol{y}) \right]$$



**Note**

$$\mathbb{E}[X] = \sum_{i=1}^{K} P(X = x_i) \cdot x_i \approx \frac{1}{N} \sum_{i=1}^{N} x_i$$

$$\min_x (f(x)) \quad vs. \quad \text{argmin}_x(f(x))$$

https://medium.com/@dhara732002/supervised-machine-learning-a-beginners-guide-9ac0b07eccbb

## 비지도학습

$$\text{Latent structure in } \{x_1, x_2, \ldots, x_n\}$$



https://ps.mjstudio.net/clustering-methods
https://mlpills.substack.com/p/issue-91-principal-component-analysis

## 강화학습

$$\pi^*(\boldsymbol{s}_t) \in \text{argmax}_{\boldsymbol{a}_t} \, Q^*(\boldsymbol{s}_t, \boldsymbol{a}_t)$$



https://en.wikipedia.org/wiki/Reinforcement_learning

# 1. 강화학습이란? : 구성요소

# 02. 강화학습 이론

# 2. 강화학습 이론 : 개요



**State:** $\quad s_t$

**Agent**

$s_t$

$a_t$

$r_{t+1}$

$s_{t+1}$

**Action:** $\quad a_t \sim \pi(\cdot \,|\, s_t)$

**Transition:** $s_{t+1} \sim P(\cdot \,|\, s_t, a_t)$

**Reward:** $\quad r_{t+1} \sim p(\cdot \,|\, s_t, a_t, s_{t+1})$

**Environment**

**Trajectory:** $\quad \mathcal{T} = (s_0, a_0, r_1, s_1, a_1, r_2, s_2, \dots)$

**Return:** $\quad G_t = \sum_{i=0}^{\infty} \gamma^i \cdot r_{t+i+1} = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots$ **Discount rate:** $\gamma \in [0,1)$

**Q:** $\quad Q_\pi(s_t, a_t) = \mathbb{E}_\pi[G_t | s_t, a_t]$ **Optimal Q:** $Q^*(s_t, a_t) = \max_\pi Q_\pi(s_t, a_t)$

**Optimal Policy:** $\pi^*(s_t) \in \mathrm{argmax}_{a_t} Q^*(s_t, a_t)$

**Markov Decision Process:**



$$s_{t+1} \sim P(\cdot \,|\, s_t, a_t)$$

$$a_t \sim \pi(\cdot \,|\, s_t)$$

$$r_{t+1} \sim p(\cdot \,|\, s_t, a_t, s_{t+1})$$

$$r_{t+1}, s_{t+1}$$

**Trajectory:** $\mathcal{T} = (s_0, a_0, r_1, s_1, a_1, r_2, s_2, \dots)$

**Markov Property:**

$$P(s_{t+1}, r_{t+1} | s_t, a_t) = P(s_{t+1}, r_{t+1} | s_0, a_0, r_1, s_1, a_1, r_2, s_2, \dots, s_t, a_t)$$

# 02. 강화학습 이론 : 가치함수 (Value-Function)

## 상태가치함수 (State-Value Function)

- 정책 $\pi$ 를 따르는 경우, 상태 $s_t$ 에서의 기대 리턴 $G_t$
- 이 상태는 얼마나 가치가 있는가?

$$V_\pi(s_t) = \mathbb{E}_\pi[G_t|s_t]$$

## 행동가치함수 (Action-Value Function)

- 상태 $s_t$ 에서 행동 $a_t$ 를 하고, 이후 정책 $\pi$ 를 따르는 경우에 대한 기대 리턴 $G_t$
- 이 상태에서 이 행동을 얼마나 가치가 있는가?

$$Q_\pi(s_t, a_t) = \mathbb{E}_\pi[G_t|s_t, a_t]$$

**Note**

$$V_\pi(s_t) = \sum_{a_t} \pi(a_t|s_t)Q_\pi(s_t, a_t)$$

**Note**
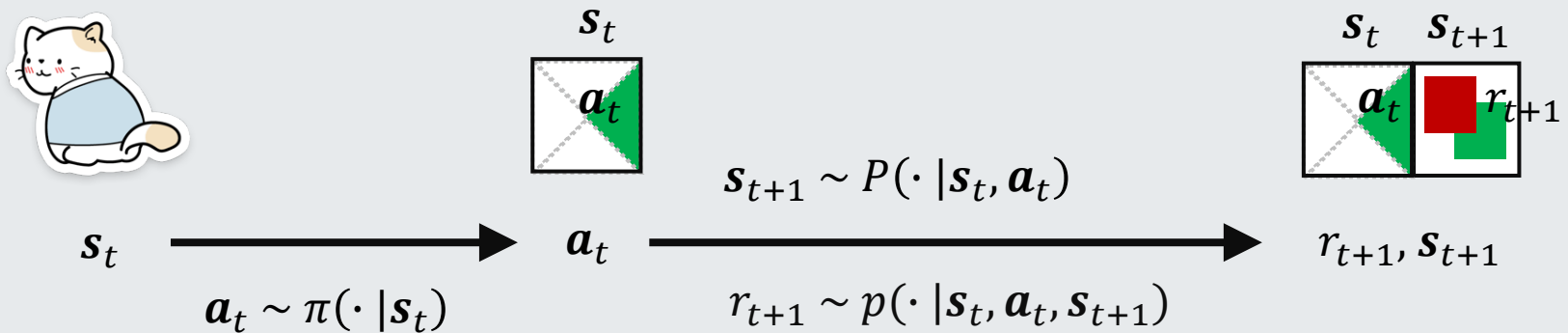
$$G_t = \sum_{i=0}^{\infty} \gamma^i \cdot r_{t+i+1} = r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots$$

# 02. 강화학습 이론 : 벨만 기대방정식 (Bellman Expectation Equation)

## 상태가치함수 (State-Value Function)

*Bellman Expectation Equation in Expectation Form*

$$V_\pi(\boldsymbol{s}_t) = \mathbb{E}_\pi[G_t|\boldsymbol{s}_t]$$

$$= \mathbb{E}_\pi[r_{t+1} + \gamma(r_{t+2} + \gamma r_{t+3} + \cdots)|\boldsymbol{s}_t]$$

$$= \mathbb{E}_\pi[r_{t+1} + \gamma G_{t+1}|\boldsymbol{s}_t]$$

$$= \mathbb{E}_\pi[r_{t+1}|\boldsymbol{s}_t] + \gamma\mathbb{E}_\pi[G_{t+1}|\boldsymbol{s}_t]$$

$$= \mathbb{E}_\pi[r_{t+1}|\boldsymbol{s}_t] + \gamma\mathbb{E}_\pi[\mathbb{E}_\pi[G_{t+1}|\boldsymbol{s}_{t+1}]|\boldsymbol{s}_t]$$

$$= \mathbb{E}_\pi[r_{t+1}|\boldsymbol{s}_t] + \gamma\mathbb{E}_\pi[V_\pi(\boldsymbol{s}_{t+1})|\boldsymbol{s}_t]$$

$$= \mathbb{E}_\pi[r_{t+1} + \gamma V_\pi(\boldsymbol{s}_{t+1})|\boldsymbol{s}_t]$$

*Bellman Expectation Equation in Summation Form*

$$V_\pi(\boldsymbol{s}_t) = \sum_{\boldsymbol{a}_t} \pi(\boldsymbol{a}_t|\boldsymbol{s}_t) \sum_{\boldsymbol{s}_{t+1}} P(\boldsymbol{s}_{t+1}|\boldsymbol{s}_t,\boldsymbol{a}_t) \sum_{r_{t+1}} p(r_{t+1}|\boldsymbol{s}_t,\boldsymbol{a}_t,\boldsymbol{s}_{t+1})\{r_{t+1} + \gamma V_\pi(\boldsymbol{s}_{t+1})\}$$

# 02. 강화학습 이론 : 벨만 기대방정식 (Bellman Expectation Equation)

## 행동가치함수 (*Action-Value Function*)

*Bellman Expectation Equation in Expectation Form*

**Note**

$$V_\pi(\boldsymbol{s}_t) = \sum_{\boldsymbol{a}_t} \pi(\boldsymbol{a}_t|\boldsymbol{s}_t) Q_\pi(\boldsymbol{s}_t, \boldsymbol{a}_t)$$

$$Q_\pi(\boldsymbol{s}_t, \boldsymbol{a}_t) = \mathbb{E}_\pi[G_t|\boldsymbol{s}_t, \boldsymbol{a}_t]$$

$$= \mathbb{E}_\pi[r_{t+1} + \gamma(r_{t+2} + \gamma r_{t+3} + \cdots)|\boldsymbol{s}_t, \boldsymbol{a}_t]$$

$$= \mathbb{E}_\pi[r_{t+1} + \gamma G_{t+1}|\boldsymbol{s}_t, \boldsymbol{a}_t]$$

$$= \mathbb{E}[r_{t+1}|\boldsymbol{s}_t, \boldsymbol{a}_t] + \gamma \mathbb{E}_\pi[G_{t+1}|\boldsymbol{s}_t, \boldsymbol{a}_t]$$

$$= \mathbb{E}[r_{t+1}|\boldsymbol{s}_t, \boldsymbol{a}_t] + \gamma \mathbb{E}[\mathbb{E}_\pi[G_{t+1}|\boldsymbol{s}_{t+1}]|\boldsymbol{s}_t, \boldsymbol{a}_t]$$

$$= \mathbb{E}[r_{t+1}|\boldsymbol{s}_t, \boldsymbol{a}_t] + \gamma \mathbb{E}[V_\pi(\boldsymbol{s}_{t+1})|\boldsymbol{s}_t, \boldsymbol{a}_t]$$

$$= \mathbb{E}[r_{t+1}|\boldsymbol{s}_t, \boldsymbol{a}_t] + \gamma \mathbb{E}\left[\sum_{\boldsymbol{a}_{t+1}} \pi(\boldsymbol{a}_{t+1}|\boldsymbol{s}_{t+1}) Q_\pi(\boldsymbol{s}_{t+1}, \boldsymbol{a}_{t+1})|\boldsymbol{s}_t, \boldsymbol{a}_t\right]$$

$$= \mathbb{E}\left[r_{t+1} + \gamma \sum_{\boldsymbol{a}_{t+1}} \pi(\boldsymbol{a}_{t+1}|\boldsymbol{s}_{t+1}) Q_\pi(\boldsymbol{s}_{t+1}, \boldsymbol{a}_{t+1})|\boldsymbol{s}_t, \boldsymbol{a}_t\right]$$

$$= \mathbb{E}[r_{t+1} + \gamma \mathbb{E}_\pi[Q_\pi(\boldsymbol{s}_{t+1}, \boldsymbol{a}_{t+1})]|\boldsymbol{s}_t, \boldsymbol{a}_t]$$

# 02. 강화학습 이론 : 벨만 기대방정식 (Bellman Expectation Equation)

## 행동가치함수 (*Action-Value Function*)

*Bellman Expectation Equation in Summation Form*

$$Q_\pi(\boldsymbol{s}_t, \boldsymbol{a}_t) = \mathbb{E}[r_{t+1} + \gamma \mathbb{E}_\pi[Q_\pi(\boldsymbol{s}_{t+1}, \boldsymbol{a}_{t+1})]|\boldsymbol{s}_t, \boldsymbol{a}_t]$$

$$= \mathbb{E}_{\boldsymbol{s}_{t+1} \sim P(\cdot|\boldsymbol{s}_t, \boldsymbol{a}_t)} \left[ \mathbb{E}_{r_{t+1} \sim p(\cdot|\boldsymbol{s}_t, \boldsymbol{a}_t, \boldsymbol{s}_{t+1})} \left[ \left[ r_{t+1} + \gamma \mathbb{E}_{\boldsymbol{a}_{t+1} \sim \pi(\cdot|\boldsymbol{s}_{t+1})}[Q_\pi(\boldsymbol{s}_{t+1}, \boldsymbol{a}_{t+1})] \right] \right] | \boldsymbol{s}_t, \boldsymbol{a}_t \right]$$

$$= \sum_{\boldsymbol{s}_{t+1}} P(\boldsymbol{s}_{t+1}|\boldsymbol{s}_t, \boldsymbol{a}_t) \sum_{r_{t+1}} p(r_{t+1}|\boldsymbol{s}_t, \boldsymbol{a}_t, \boldsymbol{s}_{t+1}) \left\{ r_{t+1} + \gamma \sum_{\boldsymbol{a}_{t+1}} \pi(\boldsymbol{a}_{t+1}|\boldsymbol{s}_{t+1}) Q_\pi(\boldsymbol{s}_{t+1}, \boldsymbol{a}_{t+1}) \right\}$$

# 02. 강화학습 이론 : 벨만 최적 방정식 (Bellman Optimality Equation)

## 상태가치함수 (State-Value Function)

$$V^*(\boldsymbol{s}_t) = \max_{\pi} V_{\pi}(\boldsymbol{s}_t)$$

$$= \max_{a_t} \mathbb{E}[r_{t+1} + \gamma V^*(\boldsymbol{s}_{t+1})|\boldsymbol{s}_t, \boldsymbol{a}_t]$$

$$= \max_{a_t} Q^*(\boldsymbol{s}_t, \boldsymbol{a}_t)$$

## 행동가치함수 (Action-Value Function)

$$Q^*(\boldsymbol{s}_t, \boldsymbol{a}_t) = \max_{\pi} Q_{\pi}(\boldsymbol{s}_t, \boldsymbol{a}_t)$$

$$= \mathbb{E}[r_{t+1} + \gamma \mathbb{E}_{\pi}[Q_{\pi}(\boldsymbol{s}_{t+1}, \boldsymbol{a}_{t+1})]|\boldsymbol{s}_t, \boldsymbol{a}_t]$$

$$= \mathbb{E}\left[r_{t+1} + \gamma \max_{\boldsymbol{a}_{t+1}}[Q^*(\boldsymbol{s}_{t+1}, \boldsymbol{a}_{t+1})]|\boldsymbol{s}_t, \boldsymbol{a}_t\right]$$

# 02. 강화학습 이론 : 최적정책 (Optimal Policy)

**결정론적 최적정책** (Deterministic Optimal Policy)

$$A^*(\boldsymbol{s}_t) := \operatorname{argmax}_{\boldsymbol{a}_t \in A} \ Q^*(\boldsymbol{s}_t, \boldsymbol{a}_t)$$

$$\pi^*(\boldsymbol{s}_t) \in A^*(\boldsymbol{s}_t)$$

**확률론적 최적정책** (Stochastic Optimal Policy)

$$\pi^*(\boldsymbol{a}_t | \boldsymbol{s}_t) = 0 \quad \forall a \notin A^*(\boldsymbol{s}_t) \ , \qquad \sum_{a \in A^*(\boldsymbol{s}_t)} \pi^*(\boldsymbol{a}_t | \boldsymbol{s}_t) = 1$$

# 03. 강화학습 알고리즘 분류

# 03. 강화학습 알고리즘 분류 : Model-based vs. Model-free

## *Model-based    vs.    Model-free*



**State:** $\boldsymbol{s}_t$

**Action:** $\boldsymbol{a}_t \sim \pi(\cdot \,|\boldsymbol{s}_t)$

**Transition:** $\boldsymbol{s}_{t+1} \sim P(\cdot \,|\boldsymbol{s}_t, \boldsymbol{a}_t)$

**Reward:** $r_{t+1} \sim p(\cdot \,|\boldsymbol{s}_t, \boldsymbol{a}_t, \boldsymbol{s}_{t+1})$

$$\pi(\boldsymbol{a}_t|\boldsymbol{s}_t) \qquad P(\boldsymbol{s}_{t+1}|\boldsymbol{s}_t, \boldsymbol{a}_t) \qquad p(r_{t+1}|\boldsymbol{s}_t, \boldsymbol{a}_t, \boldsymbol{s}_{t+1})$$

*Known*

*Unknown*

*Model-based*

*Learned by planning*

*Model-free*

*Learned by transition:*
$(\boldsymbol{s}_t, \boldsymbol{a}_t, r_t, \boldsymbol{s}_{t+1})$

# 03. 강화학습 알고리즘 분류 : Model-based vs. Model-free

**Model-based**   *vs.*   *Model-free*

*Planning*

$$V_\pi(\boldsymbol{s}_t) = \sum_{\boldsymbol{a}_t} \pi(\boldsymbol{a}_t|\boldsymbol{s}_t) \sum_{\boldsymbol{s}_{t+1}} P(\boldsymbol{s}_{t+1}|\boldsymbol{s}_t,\boldsymbol{a}_t) \sum_{r_{t+1}} p(r_{t+1}|\boldsymbol{s}_t,\boldsymbol{a}_t,\boldsymbol{s}_{t+1})\{r_{t+1} + \gamma V_\pi(\boldsymbol{s}_{t+1})\}$$

$$V^*(\boldsymbol{s}_t) = \max_{a_t} Q_\pi^*(\boldsymbol{s}_t,\boldsymbol{a}_t)$$

$$Q_\pi(\boldsymbol{s}_t,\boldsymbol{a}_t) == \sum_{\boldsymbol{s}_{t+1}} P(\boldsymbol{s}_{t+1}|\boldsymbol{s}_t,\boldsymbol{a}_t) \sum_{r_{t+1}} p(r_{t+1}|\boldsymbol{s}_t,\boldsymbol{a}_t,\boldsymbol{s}_{t+1})\left\{r_{t+1} + \gamma \sum_{\boldsymbol{a}_{t+1}} \pi(\boldsymbol{a}_{t+1}|\boldsymbol{s}_{t+1})Q_\pi(\boldsymbol{s}_{t+1},\boldsymbol{a}_{t+1})\right\}$$

$$Q^*(\boldsymbol{s}_t,\boldsymbol{a}_t) = \mathbb{E}\left[r_{t+1} + \gamma \max_{\boldsymbol{a}_{t+1}}[Q^*(\boldsymbol{s}_{t+1},\boldsymbol{a}_{t+1})]|\boldsymbol{s}_t,\boldsymbol{a}_t\right]$$

$$A^*(\boldsymbol{s}_t) := \text{argmax}_{\boldsymbol{a}_t \in A}\ Q^*(\boldsymbol{s}_t,\boldsymbol{a}_t) \qquad\qquad \pi^*(\boldsymbol{a}_t|\boldsymbol{s}_t) = 0 \quad \forall a \notin A^*(\boldsymbol{s}_t)$$

$$\pi^*(\boldsymbol{s}_t) \in A^*(\boldsymbol{s}_t) \qquad\qquad\qquad\qquad \sum_{a \in A^*(\boldsymbol{s}_t)} \pi^*(\boldsymbol{a}_t|\boldsymbol{s}_t) = 1$$

# 03. 강화학습 알고리즘 분류 : Model-based vs. Model-free

*Model-based* vs. ***Model-free***

$$V^*(\boldsymbol{s}_t) = \max_{a_t} \mathbb{E}[r_{t+1} + \gamma V^*(\boldsymbol{s}_{t+1}) | \boldsymbol{s}_t, \boldsymbol{a}_t]$$

$$(\boldsymbol{s}_t, \boldsymbol{a}_t, r_t, \boldsymbol{s}_{t+1})$$

$$V^*(\boldsymbol{s}_t) \approx r_{t+1} + \gamma V^*(\boldsymbol{s}_{t+1})$$

$$TD = r_{t+1} + \gamma \hat{V}_{\boldsymbol{\pi}}(\boldsymbol{s}_{t+1}) - \hat{V}_{\boldsymbol{\pi}}(\boldsymbol{s}_t)$$

$$\hat{V}_{\boldsymbol{\pi}}(\boldsymbol{s}_t) \leftarrow \hat{V}_{\boldsymbol{\pi}}(\boldsymbol{s}_t) + \alpha\left(r_{t+1} + \gamma\hat{\hat{V}}_{\boldsymbol{\pi}}(\boldsymbol{s}_{t+1}) - \hat{V}_{\boldsymbol{\pi}}(\boldsymbol{s}_t)\right)$$

*Q-learning, DQN, DDPG, PPO

# 03. 강화학습 알고리즘 분류 : Model-based vs. Model-free

*Model-based    vs.    **Model-free***

$$Q^*(\boldsymbol{s}_t, \boldsymbol{a}_t) = \mathbb{E}\left[r_{t+1} + \gamma \max_{\boldsymbol{a}_{t+1}}[Q^*(\boldsymbol{s}_{t+1}, \boldsymbol{a}_{t+1})]|\boldsymbol{s}_t, \boldsymbol{a}_t\right]$$

$$Q^*(\boldsymbol{s}_t, \boldsymbol{a}_t) = r_{t+1} + \gamma \max_{\boldsymbol{a}_{t+1}} Q^*(\boldsymbol{s}_{t+1}, \boldsymbol{a}_{t+1})$$

$$TD = r_{t+1} + \gamma \max_{\boldsymbol{a}_{t+1}} \hat{Q}_{\boldsymbol{\pi}}(\boldsymbol{s}_{t+1}, \boldsymbol{a}_{t+1}) - \hat{Q}_{\boldsymbol{\pi}}(\boldsymbol{s}_t, \boldsymbol{a}_t)$$

$$\hat{Q}_{\boldsymbol{\pi}}(\boldsymbol{s}_t, \boldsymbol{a}_t) \leftarrow \hat{Q}_{\boldsymbol{\pi}}(\boldsymbol{s}_t, \boldsymbol{a}_t) + \alpha\left(r_{t+1} + \gamma \max_{\boldsymbol{a}_{t+1}} \hat{Q}_{\boldsymbol{\pi}}(\boldsymbol{s}_{t+1}, \boldsymbol{a}_{t+1}) - \hat{Q}_{\boldsymbol{\pi}}(\boldsymbol{s}_t, \boldsymbol{a}_t)\right)$$

* Q-learning, DQN, DDPG, PPO

# 03. 강화학습 알고리즘 분류 : On-policy vs. Off-policy

**On-policy   vs.   Off-policy**

$$\mathcal{T} = (s_0, a_0, r_1, s_1, a_1, r_2, s_2, \dots)$$

**On-policy**

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

*\* SARSA, PPO*

**Off-policy**

$$\hat{Q}_{\pi}(s_t, a_t) \leftarrow \hat{Q}_{\pi}(s_t, a_t) + \alpha \left( r_{t+1} + \gamma \max_{a_{t+1}} \hat{Q}_{\pi}(s_{t+1}, a_{t+1}) - \hat{Q}_{\pi}(s_t, a_t) \right)$$

*\* Q-learning, DQN, DDPG*

# 03. 강화학습 알고리즘 분류 : Value-based vs. Policy-based

*Value-based    vs.    Policy-based*

## Value-based

$$\hat{\boldsymbol{\pi}}(\boldsymbol{s}_t) = \text{argmax}_{\boldsymbol{a}_t} \hat{Q}_{\boldsymbol{\pi}}(\boldsymbol{s}_t, \boldsymbol{a}_t)$$

* Q-learning, DQN

## Policy-based

Maximize $J(\boldsymbol{\theta}) = \mathbb{E}_{\tau \sim \pi_{\boldsymbol{\theta}}}[G_0]$

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{\tau \sim \pi_{\theta}}\left[\sum_{t=0}^{T-1} \nabla_{\boldsymbol{\theta}}\log \pi_{\theta}(\boldsymbol{a}_t | \boldsymbol{s}_t)G_t\right] = \frac{1}{N}\sum_{i=1}^{N}\sum_{t=0}^{T-1} \nabla_{\boldsymbol{\theta}}\log \pi_{\theta}(\boldsymbol{a}_t^i | \boldsymbol{s}_t^i)G_t^i$$

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha \frac{1}{N}\sum_{i=1}^{N}\sum_{t=0}^{T-1} \nabla_{\boldsymbol{\theta}}\log \pi_{\theta}(\boldsymbol{a}_t^i | \boldsymbol{s}_t^i)G_t^i$$

$$\boldsymbol{a} \sim \pi_{\theta}(\cdot | \boldsymbol{s}_t)$$

* PPO

# 03. 강화학습 알고리즘 분류 : Actor-Critic

**Actor-Critic**

**Actor**

$$a_t \sim \pi_\theta(\cdot | s_t)$$

$$\theta \ \leftarrow \ \theta + \alpha \times \delta \times \nabla_\theta \log(\pi_\theta(a_t | s_t))$$

**Critic**

$$V_w(s_t) \qquad \delta_t = r_{t+1} + \gamma V_w(s_{t+1}) - V_w(s_t)$$

$$w \leftarrow w + \beta \times \delta \times \nabla_w V_w(s_t)$$

*\* DDPG, PPO*

# 03. 강화학습 알고리즘 분류 : Deterministic vs. Stochastic

**Deterministic   vs.   Stochastic**

**Deterministic**

$$s_{t+1} = F(\boldsymbol{s}_t, \boldsymbol{a}_t)$$

$$\boldsymbol{a}_t = \ f(\boldsymbol{s}_t)$$

$$r_{t+1} = R(\boldsymbol{s}_t, \boldsymbol{a}_t, \boldsymbol{s}_{t+1})$$

*\* DDPG*

**Stochastic**

$$s_{t+1} = P(\cdot \,|\boldsymbol{s}_t, \boldsymbol{a}_t)$$

$$\boldsymbol{a}_t \ \sim \ \pi(\cdot \,|\boldsymbol{s}_t)$$

$$r_{t+1} \sim p(\cdot \,|\boldsymbol{s}_t, \boldsymbol{a}_t, \boldsymbol{s}_{t+1})$$

*\* PPO*

# 04. 강화학습 주요 알고리즘

# 04. 강화학습의 주요 알고리즘: Q-learning

*Off-policy    Value-based    Discrete Action    Model-free*

**Pseudo-code**

$\hat{Q}(s_t, a_t) = 0$

for each episode:

  $s_t$ = reset()

  while not terminal:
    $a_t$ = epsilon_greedy($Q$, $s_t$)
    $s_{t+1}$, $r_{t+1}$, done = step($a_t$)

    target = $r_{t+1} + \gamma \times \max_{a_{t+1}}[\hat{Q}(s_{t+1}, a_{t+1})] \times$(not done)

    $\hat{Q}(s_t, a_t) = \hat{Q}(s_t, a_t) + \alpha$(target $-\hat{Q}(s_t, a_t)$)

  $s_t = s_{t+1}$

$$a_t = \begin{cases} \text{argmax}_{a_t} \hat{Q}(s_t, a_t) & with\ prob\ 1 - \varepsilon \\ random\ action & with\ prob\ \varepsilon \end{cases}$$

$$s_t, a_t, r_{t+1}, s_{t+1}$$

$$Q^*(s_t, a_t) = \max_\pi Q_\pi(s_t, a_t)$$

$$= \mathbb{E}\left[ r_{t+1} + \gamma \max_{a_{t+1}}[Q^*(s_{t+1}, a_{t+1})] | s_t, a_t \right]$$

$$\approx r_{t+1} + \gamma \max_{a_{t+1}}[Q^*(s_{t+1}, a_{t+1})]$$

*(single-sample estimate of the expectation)*

$$\hat{y} = r_{t+1} + \gamma \max_{a_{t+1}}[\hat{Q}(s_{t+1}, a_{t+1})] \cdot 1[not\ done]$$

$$\hat{Q}(s_t, a_t) \leftarrow \hat{Q}(s_t, a_t) + \alpha(\hat{y} - \hat{Q}(s_t, a_t))$$

$$\hat{Q} \rightarrow Q^* \quad (as\ learning\ convergens)$$

$$\pi^*(s_t) = \text{argmax}_a Q^*(s_t, a_t)$$

# 04. 강화학습의 주요 알고리즘: SARSA

*On-policy    Value-based    Discrete Action    Model-free*

---

**Pseudo-code**

$\hat{Q}(s_t, a_t) = 0$

for each episode:

$\quad s_t$ = reset()
$\quad a_t$ = epsilon_greedy($Q$, $s_t$)

$\quad$while not terminal:
$\quad\quad s_{t+1}, r_{t+1}$, done = step($a_t$)
$\quad\quad a_{t+1}$ = epsilon_greedy($Q$, $s_t + 1$)

$\quad\quad$target = $r_{t+1} + \gamma \times \hat{Q}(s_{t+1}, a_{t+1}) \times$(not done)

$\quad\quad\quad \hat{Q}(s_t, a_t) = \hat{Q}(s_t, a_t) + \alpha$(target -$\hat{Q}(s_t, a_t)$)

$\quad\quad s_t, a_t = s_{t+1}, a_{t+1}$

$$a_t = \begin{cases} \text{argmax}_{a_t}\hat{Q}(s_t, a_t) & with\ prob\ 1 - \varepsilon \\ random\ action & with\ prob\ \varepsilon \end{cases}$$

$$s_t, a_t, r_{t+1}, s_{t+1}, a_{t+1}$$

$$Q_\pi(s_t, a_t) = \mathbb{E}[r_{t+1} + \gamma \mathbb{E}_\pi[Q_\pi(s_{t+1}, a_{t+1})]|s_t, a_t]$$

$$\approx r_{t+1} + \gamma Q_\pi(s_{t+1}, a_{t+1})$$

*(single-sample estimate of the expectation)*

$$\hat{y} = r_{t+1} + \gamma \hat{Q}(s_{t+1}, a_{t+1}) \cdot 1[not\ done]$$

$$\hat{Q}(s_t, a_t) \leftarrow \hat{Q}(s_t, a_t) + \alpha(\hat{y} - \hat{Q}(s_t, a_t))$$

$$\hat{Q} \rightarrow Q^*$$

$$(if\ GLIE: \varepsilon \rightarrow 0\ with\ infinite\ exploration)$$

$$\pi^*(s_t) = \text{argmax}_a Q^*(s_t, a_t)$$

# 04. 강화학습의 주요 알고리즘: DQN (Deep Q-Network)

*Off-policy    Value-based    Discrete Action    Model-free    Deep Neural Network    Replay Buffer    Target Network*

**Pseudo-code**

```
Replay buffer D = empty_buffer()
Online network parameters  Q_θ  = init_network()
Target network parameters  Q_θ⁻ = copy(Q_θ )

for each episode:

  s_t = reset()
  done = False

  while not done:

    a_t = epsilon_greedy(Q_θ, s_t)
    s_{t+1}, r_{t+1}, done = step(a_t)
    D.add(s_t, a_t, r_{t+1}, s_{t+1}, done)
    s_t = s_{t+1}

    if len(D) == B:
      batch = D.randomSample(B)
           # {s_t^i, a_t^i, r_{t+1}^i, s_{t+1}^i, done^i} for i =1..B
      target _i
         = r_{i+1} + γ × max_{a_t+1}[Q_θ⁻(s_{t+1}^i, a_{t+1})] ×(not done)
      L(θ) = (1/B) Σ_i(Q_θ(s_t^i, a_t^i) − target_i)²
      θ = θ − l_r × ∇_θ L(θ)

  every C steps: θ⁻ ← θ
```

$$a_t = \begin{cases} \operatorname{argmax}_{a_t} Q_\theta(s_t, a_t) & with\ prob\ 1-\varepsilon \\ random\ action & with\ prob\ \varepsilon \end{cases}$$

$$s_t, a_t, r_{t+1}, s_{t+1}$$

$$Q^*(s_t, a_t) = \max_\pi Q_\pi(s_t, a_t)$$

$$= \mathbb{E}\left[r_{t+1} + \gamma \max_{a_{t+1}}[Q^*(s_{t+1}, a_{t+1})]|s_t, a_t\right]$$

*(minibatch-sample estimate of the expectation)*

$$\hat{y}_i = r_{i+1} + \gamma \max_{a_{t+1}} Q_{\theta^-}(s_{t+1}^i, a_{t+1}) \cdot 1[not\ done]$$

$$\hat{L}(\theta) = \frac{1}{B} \sum_{i=1}^{B} \left(\hat{y}_i - \hat{Q}(s_t^i, a_t^i)\right)^2$$

$$\theta \leftarrow \theta - L_r \nabla_\theta \hat{L}(\theta) \qquad (for\ every\ steps)$$

$$\theta^- \leftarrow \theta \qquad\qquad (every\ C\ steps)$$

$$\hat{Q} \rightarrow Q^* \quad (as\ learning\ convergens)$$

$$\pi^*(s_t) = \operatorname{argmax}_a Q^*(s_t, a_t)$$

# 04. 강화학습의 주요 알고리즘: Actor–Critic

*On/Off-policy    Value/Policy-based    Discrete/Continuous Action    Model-free    Deep Neural Network*

**Pseudo-code**

```
actor: θ, π_θ(a_t|s_t)
Critic: w, V_w(s_t)

for each episode:
  a_t ~ π_θ(·|s_t)
  (r_{t+1}, s_{t+1}, done) = step(a_t)

  δ_t = r_{t+1} + γ V_w(s_{t+1}) ×(not done) - V_w(s_t)

  # critic update
  w ← w + β × δ × ∇_w V_w(s_t)

  # actor update (policy gradient)
  θ ← θ + α × δ × ∇_θ log(π_θ(a_t|s_t))

  s_t ← s_{t+1} (or reset if done)
```

$$V^*(s_t) = \max_{\pi} V_{\pi}(s_t)$$

$$= \max_{a_t} \mathbb{E}[r_{t+1} + \gamma V_{\pi}(s_{t+1})|s_t, a_t]$$

*Off-policy*          *Actor-Critic*              *Continuous Action*              *Deep Neural Network*

**Pseudo-code**

actor $\mu_{\boldsymbol{\theta}}(\boldsymbol{s_t})$          actor_target $\mu_t = copy(\mu_{\boldsymbol{\theta}})$
critic $Q_{\boldsymbol{\phi}}(\boldsymbol{s_t}, \boldsymbol{a_t})$          critic_target $\boldsymbol{Q_t} = copy(Q_{\boldsymbol{\phi}})$
Replay Buffer: $D$

for each step:
  $\boldsymbol{a_t} = \mu_{\boldsymbol{\theta}}(\boldsymbol{s_t})$ + noise
  $\boldsymbol{s_{t+1}}, r_{t+1}, done$ = step($\boldsymbol{a_t}$)
  $D$.add($\boldsymbol{s_t}, \boldsymbol{a_t}, r_{t+1}, \boldsymbol{s_{t+1}}, done$)

  batch = $D$.randomSample(B)

  $\hat{y}_i = r_{t+1,i} + \gamma \, \boldsymbol{Q_t}\left(\boldsymbol{s_{t+1,i}}, \mu_t(\boldsymbol{s_{t+1,i}})\right) \times (\text{not } done)$

  $\boldsymbol{\phi} \leftarrow \boldsymbol{\phi} - \beta\nabla_{\boldsymbol{\phi}}\frac{1}{B}\sum_i^B\left(Q_{\boldsymbol{\phi}}(\boldsymbol{s_{t,i}}, \boldsymbol{a_{t,i}}) - \hat{y}_i\right)^2$

  $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha\nabla_{\boldsymbol{\theta}}Q_{\boldsymbol{\phi}}(\boldsymbol{s_t}, \mu_{\boldsymbol{\theta}}(\boldsymbol{s_t}))$
      $\boldsymbol{\theta} + \alpha\nabla_{\boldsymbol{\theta}}Q_{\boldsymbol{\phi}}(\boldsymbol{s_t}, \boldsymbol{a})|_{a=\mu_{\boldsymbol{\theta}}(\boldsymbol{s_t})}\nabla_{\boldsymbol{\theta}}\mu_{\boldsymbol{\theta}}(\boldsymbol{s_t})$

  # soft target update
  $Q_t \leftarrow \tau \, Q_{\boldsymbol{\phi}} + (1-\tau)Q_t$
  $\mu_t \leftarrow \tau \, \mu_{\boldsymbol{\theta}} + (1-\tau)\mu_t$

  $\boldsymbol{s_t} \leftarrow \boldsymbol{s_{t+1}}$ (or reset if done)

$$Q^*(\boldsymbol{s_t}, \boldsymbol{a_t}) = \max_{\pi} Q_{\pi}(\boldsymbol{s_t}, \boldsymbol{a_t})$$

$$= \mathbb{E}\left[r_{t+1} + \gamma\max_{\boldsymbol{a_{t+1}}}[Q^*(\boldsymbol{s_{t+1}}, \boldsymbol{a_{t+1}})]|\boldsymbol{s_t}, \boldsymbol{a_t}\right]$$

# 04. 강화학습의 주요 알고리즘: PPO (Proximal Policy Optimization)

*On-policy*          *Actor-Critic*          *Discrete/Continuous Action*          *Deep Neural Network*

**Pseudo-code**

```
Initialize policy πθ and value function Vφ
for iter = 1..K:
  θold ← θ
  φold ← φ

  # Collect rollouts with old policy
  D = {(s_{t,i}, a_{t,i}, r_{t+1,i}, s_{t+1,i}, done_i)} collected by πθ,old

  # Compute return and advantages
  R_{t,i}  = r_{t+1,i} + γ V_old(s_{t+1,i}) ×(not done_i )
  Â_{t,i} = R_{t,i} - V_old(s_{t,i})

  for epoch = 1..E:
    for minibatch B ⊂ D:

      # Actor loss (clip objective)
        r_{t,i} = exp(log πθ(a_{t,i}|s_{t,i}) − log πθ.old(a_{t,i}|s_{t,i}))
      L_clip = mean_{i∈B}( min(r_{t,i} × Â_{t,i}, clip(r_{t,i},1-ε,1+ε) × Â_{t,i}))

      # Critic loss (MSE)
        L_v = mean_{i∈B}(Vφ(s_{t,i}) − R_{t,i})²

  # Update
  θ ← θ + α∇θ L_clip(θ)
  φ ← φ − β∇φ L_v(φ)
```

# 수고하셨습니다.

junsu@handong.edu