# Reinforcement Learning
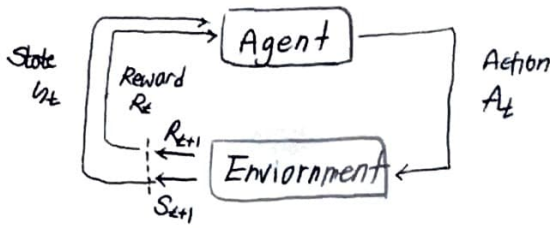


강화학습의 모든 문제는 다음과 같은 틀로 귀결된다.

Environment이 Agent에게 <u>특정상황 State</u>를 주면/Agent는 그에 대해 <u>반응 Action</u>을 하고
env는 agent에게 <u>보상 reward</u>를 준다.

- **Model** : Mathematic models of dynamics and rewards.

π · **Policy** : Function mapping agent's states to action    action을 취하는 방법론  $(s) \xrightarrow{a_t} (s')$

$g, V$ · **Value Function** : future rewards from being in a state and/or action when following a particular policy

Value
↳ Function
↳ policy

$V(s)$
ex State-value function, state-action value function

$R$ · **Reward** $R_s^a$, $R(S_t = s, a_t = a)$

$G$ · **Return** $G_t$ discounted sum of rewards from time step $t$

- **State transition Matrix**

$$P_{ss'} = P[S_{t+1} = s' \mid S_t = s]$$

$$P = \begin{pmatrix} P(S_1|S_1) & P(S_2|S_1) & \cdots & P(S_N|S_1) \\ P(S_1|S_2) & P(S_2|S_a) & \cdot & P(S_N|S_2) \\ \vdots & & & \vdots \\ P(S_1|S_N) & P(S_2|S_N) & \cdots & P(S_N|S_N) \end{pmatrix} = \begin{bmatrix} P_{11} & \cdots & & P_{1n} \\ \vdots & & & \vdots \\ P_{n1} & & \cdots & P_{nn} \end{bmatrix}$$

- **Markov Property**

State $S_t$ is Markov $\Leftrightarrow$ $P(S_{t+1} \mid S_t, a_t) = P(S_{t+1} \mid h_t, a_t)$

현재의 state가 이전의 state 에만 영향을 준다.

공식정리 **Bellman Expectation Equation**

return $G_t = R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \quad = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$

reward $R_s = E[r_t \mid S_t = s]$

$R_s^a = E[r_t \mid S_t = s, a_t = a]$

policy $\pi(a|s) = P(a_t = a, S_t = s)$

**Value Function**

state-value
$V(s) = E[G_t \mid S_t = s]$

$= E[R_t + \gamma V(S_{t+1}) \mid S_t = s]$

$= R(s) + \gamma \sum_{s' \in S} P(s'|s) V(s') \quad \checkmark$

$= R + \gamma PV$

$V_\pi(s) = E_\pi[r + \gamma E_\pi[G_{t+1} \mid S_{t+1} = s'] \mid S_t = s]$

$= \sum_a \pi(a|s) \sum_{s',r} p(s', r|s,a)[r + \gamma E_\pi[G_{t+1}|S_{t+1}=s']$

$= \sum_a \pi(a|s)[R_s^a + \gamma \sum_{s'} P_{ss'}^a V_\pi(s')] \quad \checkmark$

$= \sum_a \pi(a|s) q_\pi(s,a) \rightarrow V_\sigma = R^\pi + \gamma P^\pi V_{\pi'}$

state-action value
$q_\pi(s, a) = E_\pi[G_t \mid S_t = s, a_t = a]$

$= R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_\pi(s')$

**Bellman Optimal Equation** ( max로 차이 )

$V_*(s) = \max_\pi V_\pi(s)$

$q_*(s, a) = \max_\pi q_\pi(s, a)$

결정식 $\pi_*(a|s) = \begin{cases} 1 \\ 0 \end{cases} \quad \arg\max_a q_*(s, a)$

$V_*(s) = \max_a q_*(s, a)$

$= \max_a R_s^a + \gamma \sum_{s'} P_{ss'}^a V_*(s')$

— there is always a deterministic policy of MDP.
  ( unique 하지 않아도 가능 ).

— non-linear

## Markov Process

$\langle S, P \rangle$

No reward. No action
only state

## Markov Reward Process

$\langle S, P, R, \gamma \rangle$

$\langle S, P^{\pi}, R^{\pi}, \gamma \rangle$

No action

$P_{ss'} = P[S_{t+1} = s' \mid S_t = s]$

$V(s) = E[G_t \mid S_t = s]$

$\quad = R(s) + \gamma \sum_{s' \in S} P(s' \mid s) V(s')$

$\left( \begin{array}{l} V_{\pi} = R^{\pi} + \gamma P^{\pi} V_{\pi'} \\ q_{\pi} = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_{\pi}(s') \end{array} \right) \longrightarrow$

## Markov Decision Process

$\langle S, A, P, R, \gamma \rangle$

$P_{ss'}^a = P[S_{t+1} = s' \mid S_t = s, a_t = a]$

$P_{ss'}^{\pi} = \sum_a \pi(a \mid s) P_{ss'}^a$

$V_k^{\pi}(s) = r(s, \pi(s)) + \gamma \sum_{s \in S} P(s' \mid s, \pi(s))$
$\qquad \qquad \qquad \qquad \cdot V_{k-1}^{\pi}(s')$

$q_{\pi}(s,a) = R_s^a + \gamma \sum_{s' \in S} P_{ss'}^a V_{\pi}(s')$