# Value Function Approximation Methods $\to v, q$

맞충 생각하지 말고 근사하자. 파라미터 dimension 줄이며.

$$\begin{cases} V_\pi(s) \approx \hat{v}(s, w) & \text{where } w \in R^d \text{ with } d << |S| \\ q_\pi(s,a) \approx \hat{q}(s,a,w) & \text{where } w \in R^d \text{ with } d << |S| \end{cases}$$

- Linear Value Function Approximation

$$J(w) = MSVE(w) = E_\pi \left[ V_\pi(s) - \underbrace{x(s)^T w}_{\hat{v}(s,w)} \right)^2 \right] \quad \text{mean-squared error 사용}$$

only local optimum is possible.

$\to$ Gradient

$$w_{t+1} = w_t + \underbrace{\alpha [G_t - \hat{v}(s_t, w)] \nabla \hat{v}(s_t, w)}_{\triangle w} \quad -\frac{1}{2}\alpha \nabla [V_\pi(s_t) - \hat{v}(s_t, w)]^2$$

$$MC \quad \triangle w = \alpha (G_t - \hat{v}(s_t, w)] \nabla_w \hat{v}(s_t, w)$$

$$TD(0) = \alpha (R_{t+1} + \gamma \hat{v}(s_{t+1}, w) - \hat{v}(s_t, w)) \nabla_w \hat{v}(s_t, w)$$

$$TD(\lambda) = \alpha (\hat{G_t^\lambda} - \hat{v}(s_t, w))(\nabla_w \hat{v}(s_t, w))$$

- Linear Action - Value Function Approximation

$$J(w) = E_\pi [(q_\pi(s,A) - \hat{q}(s,A,w))^2]$$

$$\triangle w = \alpha (q_\pi(s,A) - \hat{q}(s,A,w)) x(s,A)$$

$$MC \quad \triangle w = \alpha (G_t - \hat{q}(s_t, A_t, w)) \nabla_w \hat{q}(s_t, A_t, w)$$

$$TD(0) = \alpha (R_{t+1} + \gamma \hat{q}(s_{t+1}, A_{t+1}, w) - \hat{q}(s_t, A_t, w)) \nabla_w \hat{q}(s_t, A_t, w)$$

forward-view
$$TD(\lambda) = \alpha (\hat{q_t^\lambda} - \hat{q}(s_t, A_t, w)) \nabla_w \hat{q}(s_t, A_t, w)$$

backward-view
$$TD(\lambda) = \alpha \delta_t E_t$$

$$\delta_t = R_{t+1} + \gamma \hat{q}(s_{t+1}, A_{t+1}, w) - \hat{q}(s_t, A_t, w)$$

$$E_t = \gamma \lambda E_{t-1} + \nabla_w \hat{q}(s_t, A_t, w)$$

+ DQN : Q-learning을 non-linear하게 다루기

\* Off-policy Semi-gradient Learning

$$W_{t+1} = W_t \; +$$

TD(0) $\qquad \alpha \, \rho_t \, [\, R_{t+1} + \gamma \hat{v}(b_{t+1}, w) - \hat{v}(b_t, w)\,] \, \nabla \hat{v}(b_t, w)$

SARSA(0) $\qquad \alpha \, \rho_{t+1} \, [(R_{t+1} + \gamma \hat{q}(S_{t+1}, A_{t+1}, w)) - \hat{q}(b_t, A_t, w)] \, \nabla \hat{q}(b_t, A_t, w)$

n-step
SARSA $\qquad \alpha \, \rho_{t+1:t+n} [\, G_{t:t+n} - \hat{q}(b_t, A_t, W_{t+n-1})\,] \, \nabla \hat{q}(b_t, A_t, W_{t+n-1})$

$$G_{t:t+n} = R_{t+1} + \gamma R_{t+2} + \cdots + \gamma^{m} R_{t+n} + \gamma^{n} \hat{q}(b_{t+n}, A_{t+n}, W_{t+n-1})$$

\* Batch Reinforcement Learning 끝에 업데이트

$$D = \{ <b_1, \; V_1^\pi>, <b_2, \; V_2^\pi> \; \cdots \; <b_T, \; V_T^\pi> \}$$   Batch 만틈계산

$$LS(w) = E_D \, [(V^\pi - \underbrace{\hat{v}(b, w))^2}_{W^T X(b)}]$$   의 최소가 되야 는 $W$ 구하기

① 수식적 : 계산량↑↑

② Prediction

LSMC $\longrightarrow 0 = \sum_{t=1}^{T} \alpha (G_t - \hat{v}(b, w)) X(b)$

LSTD
$\qquad \searrow 0 = \sum_{t=1}^{T} \alpha (R_{t+1} + \gamma \hat{v}(b_{t+1}, w) - \hat{v}(b_t, w)) X(S_t)$
LSTD($\lambda$)

$\qquad\qquad \searrow 0 = \sum_{t=1}^{T} \alpha \delta_t E_T$

LSTDQ $\qquad\qquad 0 = \sum_{t=1}^{T} \alpha(R_{t+1} + \gamma \hat{q}(b_{t+1}, \pi(b_{t+1}), w) - \hat{q}$

total update $\qquad (S_t, A_t, w)) X(b_t, A_t)$
= zero

# DQN (Deep Q-Networks)

Q-learning 을 Neural Network 으로 해석 (table 형태의 학습은 굉장히 느리다)

$$\langle s, v^\pi \rangle \sim D \quad \text{에피소드저장} \to \text{sample}$$

$$\Delta w = \alpha [v^\pi - \hat{v}(s, w)] \nabla_w \hat{v}(s, w)$$

$$\Rightarrow w^\pi = \arg\min_w Lh(w) \quad \text{개별 state가 중요X}$$

w가 바뀌면 전체가 바뀌기 때문에

$\to$ episode 간 correlation이 크게 발생격X

(∵ 몇개의 선별되는 state 를 이용
$\to$ sampling)

- state 이 input Data로 탄등



```
        Replay   Memory
        Update after each step        Batch < O, A, R, O' >
              |
              ↓
Reward    Agent   DQN ──── Q-value
successor
state                      action
        Environment  ←───
```

$$L_i(\theta_i) = E_{s, a \sim p(\cdot)} \left[ (y_i - Q(s, a ; \theta_i))^2 \right] \text{ 을 바탕으로 } Neural\ Network\ \text{에 넘기}$$

mini-batch data 에 대하여 bootstrap