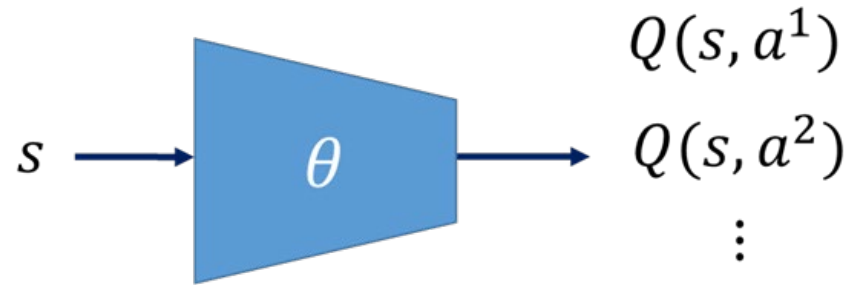


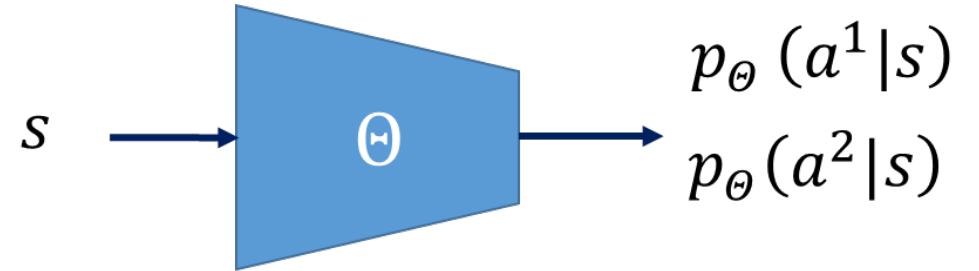
HW4 – Review and comparison of PG with DQN

- HW4 asks you to extend your class notes to review and compare policy-based vs value-based RL algorithms.
 - 1. What to learn?
 - 2. How to learn?
 - 3. How to calculate $\nabla \bar{R}_\theta$?
 - 4. How to calculate $\nabla \bar{R}_\theta$ efficiently?
 - 5. How to calculate $Q^*(s, a)$ from $Q^*(s', a')$?
 - 6. Q-learning
 - 7. How to stabilize DQN training?
 - 8. How to stabilize $\nabla \bar{R}_\theta$ calculation?
- Due: next class meeting
- Upload ppt to Teams

1. What to learn?



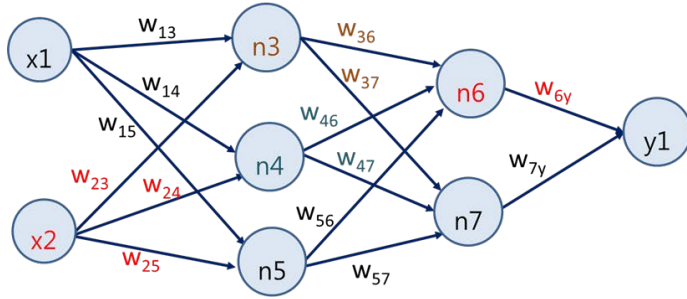
DQN



PPO

(week 7, 13 class notes)

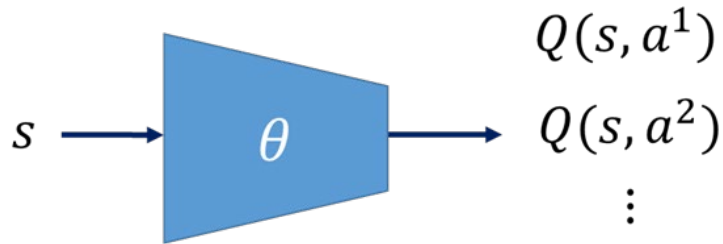
2. How to learn?



$$L = \frac{1}{N} \sum_{i=1}^N (y^i - \hat{y}^i)^2$$

$$\nabla L = \left[\frac{\partial L}{\partial w_1}, \frac{\partial L}{\partial b_1}, \dots \right]$$

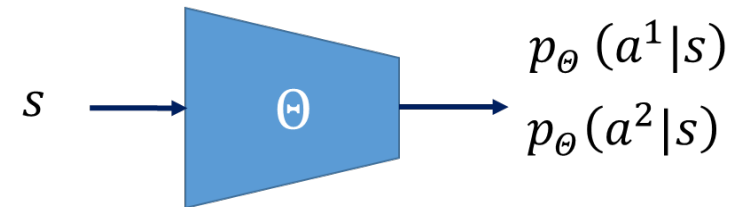
$$w_i \leftarrow w_i - \eta \frac{\partial L}{\partial w_i}$$



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

$$Loss = \left(R_s^a + \gamma \max_{a'} Q(s', a') - Q(s, a) \right)^2$$

DQN



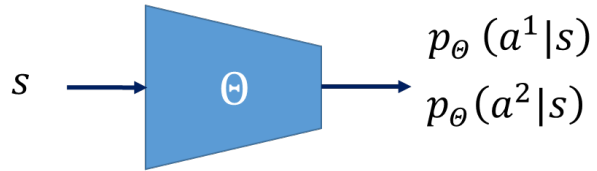
$$\bar{R}_{\theta} = E_{\tau \sim p_{\theta}(\tau)}[R(\tau)] \quad R(\tau) = \sum_{t=1}^T r_t$$

$$\theta^{\pi'} \leftarrow \theta^{\pi} + \eta \nabla \bar{R}_{\theta}$$

PPO

(week 5, 7, 13 class notes)

3. How to calculate $\nabla \bar{R}_\theta$ in PG?



$$\theta^{\pi'} \leftarrow \theta^\pi + \eta \nabla \bar{R}_\theta$$

$$\tau = (s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_T, a_T)$$

$$p_\theta(\tau) = p(s_1)p_\theta(a_1|s_1) \dots$$

$$R(\tau) = \sum_{t=1}^T r_t$$

$$\bar{R}_\theta = \sum_{\tau} R(\tau) p_\theta(\tau) = E_{\tau \sim p_\theta(\tau)}[R(\tau)]$$

$$\nabla \bar{R}_\theta = \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} R(\tau^n) \nabla \log p_\theta(a_t^n | s_t^n)$$

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} \left(\sum_{t'}^{T_n} \gamma^{t'-t} r_{t'}^n - b \right) \nabla \log p_\theta(a_t^n | s_t^n)$$

$$A^\theta(s_t, a_t) = \left(\sum_{t'}^{T_n} \gamma^{t'-t} r_{t'}^n - b \right)$$

(PG, week 7 class notes)

4. How to calculate $\nabla \bar{R}_{\theta}$ efficiently in PG?

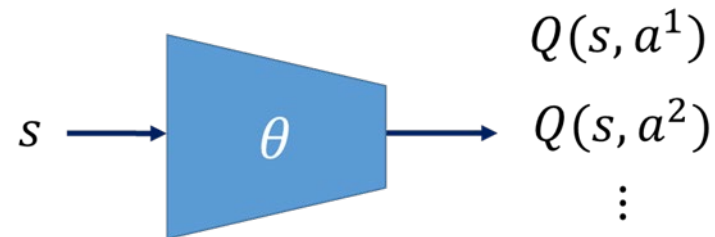
$$\nabla \bar{R}_{\theta'} \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} \left(\sum_{t'}^{T_n} \gamma^{t'-t} r_{t'}^n - b \right) \nabla \log p_{\theta'}(a_t^n | s_t^n)$$

$$\nabla \bar{R}_{\theta'} = E_{(s_t, a_t) \sim \Theta} \left[\frac{p_{\theta'}(a_t | s_t)}{p_{\Theta}(a_t | s_t)} A^{\Theta}(s_t, a_t) \nabla \log p_{\Theta}(a_t^n | s_t^n) \right]$$

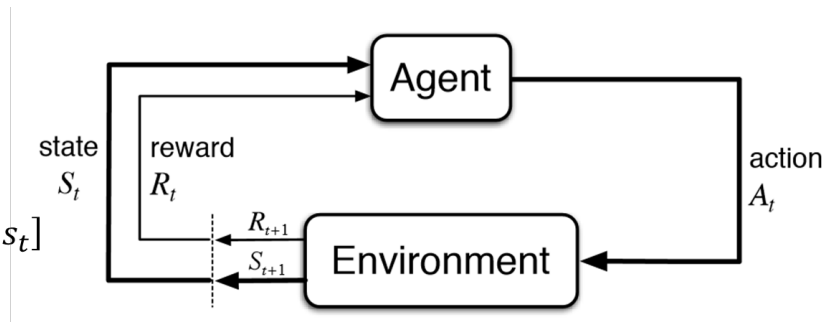
$$\text{PO2}(\theta') = \sum_{(s_t, a_t)} \min \left(\frac{p_{\theta'}(a_t | s_t)}{p_{\Theta}(a_t | s_t)} A^{\Theta}(s_t, a_t), \text{clip} \left(\frac{p_{\theta'}(a_t | s_t)}{p_{\Theta}(a_t | s_t)}, 1 - \varepsilon, 1 + \varepsilon \right) A^{\Theta}(s_t, a_t) \right)$$

(PPO, week 8 class notes)

5. How to calculate $Q^*(s, a)$ from $Q^*(s', a')$?



$$Loss = \left(R_s^a + \gamma \max_{a'} Q(s', a') - Q(s, a) \right)^2$$



Markov property

$$p[s_{t+1}|s_t] = p[s_{t+1}|s_1, s_2, \dots, s_t]$$

Value function

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

Markov process (chain)

$$s_1, s_2, \dots, s_t$$

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots$$

$$V_\pi(s) = \mathbb{E}_\pi[R_{t+1} + \gamma V_\pi(s_{t+1}) | s_t = s]$$

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

$$V^*(s) = \max_a \sum_{s'} P(s'|s, a) (R_s^a + \gamma V^*(s'))$$

$$Q^*(s, a) = \sum_{s'} P(s'|s, a) (R_s^a + \gamma \max_{a'} Q^*(s', a'))$$

Markov decision process

Bellman equation

(week 12 class notes)

6. Q-learning

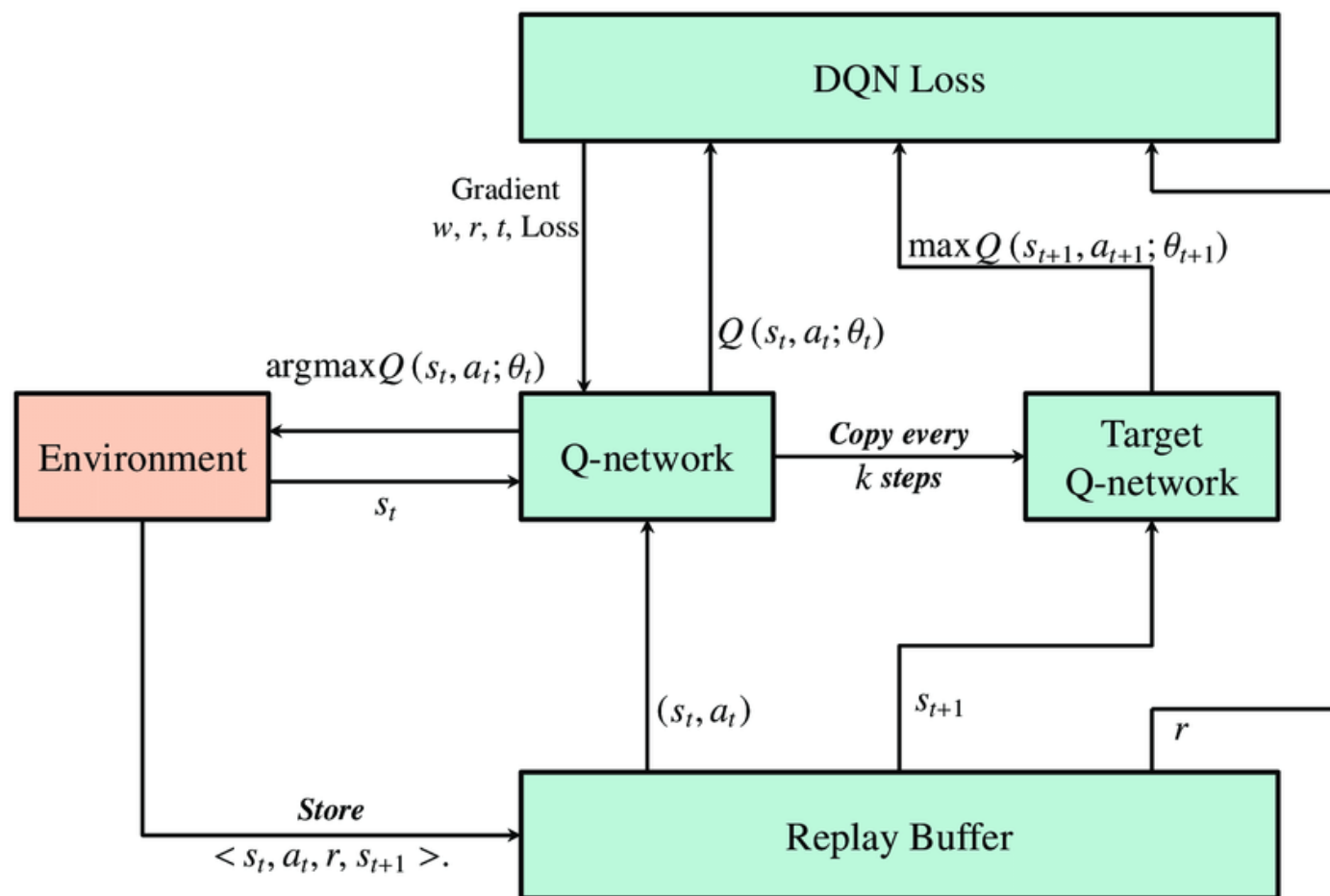
1. Initialize Table $Q(S,A)$ with random values.
2. Take a action (A) with epsilon — greedy policy and move to next state S'
3. Update the Q value of a previous state by following the update equation:

$$Q(s, a) = Q(s, a) + \alpha(R_s^a + \gamma \max_{a'} Q(s', a') - Q(s, a))$$

Q-learning

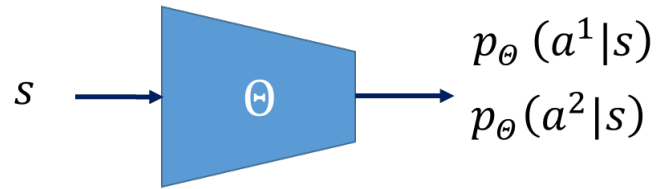
(week 13 class notes)

7. How to stabilize DQN training?



(week 13 class notes)

8. How to stabilize $\nabla \bar{R}_\theta$ calculation?



$$\theta^{\pi'} \leftarrow \theta^{\pi} + \eta \nabla \bar{R}_\theta$$

$$\nabla \bar{R}_\theta \approx \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} \left(\sum_{t'=t}^{T_n} \gamma^{t'-t} r_{t'}^n - \boxed{b} \right) \nabla \log p_\theta(a_t^n | s_t^n)$$

$\boxed{V^{\pi_\theta}(s_t^n)}$ Expected value of b

$\boxed{E[G_t^n] = Q^{\pi_\theta}(s_t^n, a_t^n)}$ Expected value of G_t^n

$G_t^n = \sum_{t'}^{T_n} \gamma^{t'-t} r_{t'}^n$ unstable when sampling amount is not large enough

$$Q^{\pi_\theta}(s_t^n, a_t^n) = \mathbb{E}[r_t^n + V^{\pi_\theta}(s_{t+1}^n)] = r_t^n + V^{\pi_\theta}(s_{t+1}^n)$$

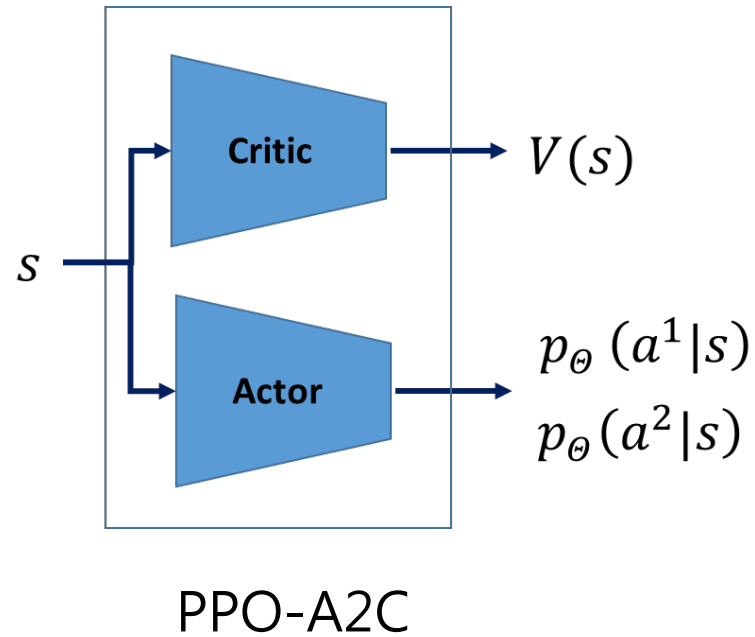
$$Q^{\pi_\theta}(s_t^n, a_t^n) - V^{\pi_\theta}(s_t^n) = r_t^n + V^{\pi_\theta}(s_{t+1}^n) - V^{\pi_\theta}(s_t^n)$$

$$A^\theta(s_t, a_t) = (r_t^n + V^{\pi_\theta}(s_{t+1}^n) - V^{\pi_\theta}(s_t^n))$$

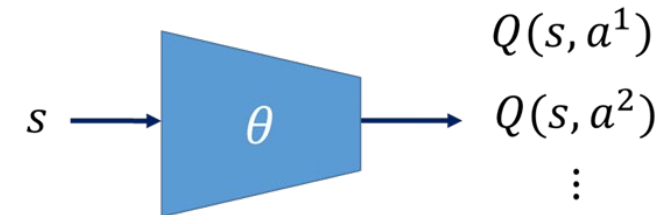
PPO-A2C

(week 9 class notes)

8. How to stabilize $\nabla \bar{R}_\theta$ calculation?



$$Loss_v = (r_t^n + \gamma V^{\pi_\theta}(s_{t+1}^n) - V^{\pi_\theta}(s_t^n))^2$$



$$Loss = (R_s^a + \gamma \max_{a'} Q(s', a') - Q(s, a))^2$$

(week 9 class notes)