Lecture 25 - April 18, 2023

- Deep Q-Network (DQN) → Finite action
  - Deep Q-Network (DQN)
  - Double DQN
  - Prioritized DQN
  - Dueling DQN
  - Noisy-Net DQN

- Deep Policy Gradients (DPG) → Large/Continuas Action
  - REINFORCE
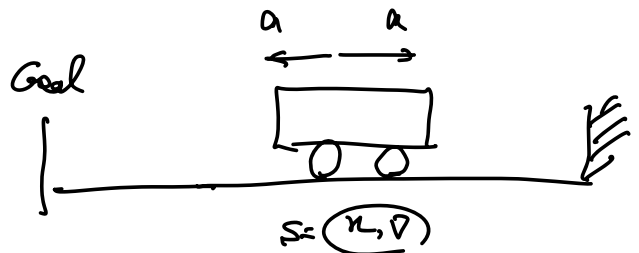  - REINFORCE with Baseline
  - Advantage Actor Critic (A2C)
  - Deep Deterministic Policy Gradient (DDPG)

HW5 ⟶ Due April 18

TA's office hour:  Wendsdays, 2Pm-3Pm (in-Person)
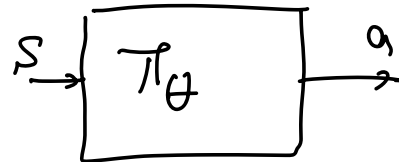                   Fridays, 2Pm-3Pm (virtual)

Goal

$S = (x, v)$

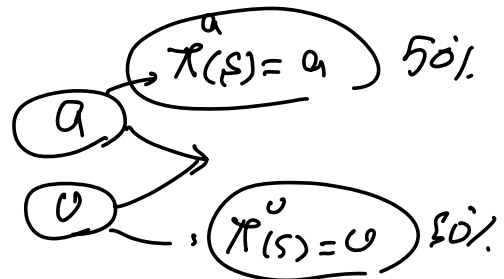Learning policy seems to be easier than $Q$-value.

$Q(s, a)$

argmax

① Smoothness

② Multi-Agent $\Rightarrow$ Stochstic

$S \rightarrow \boxed{\pi_\theta} \rightarrow a$

$\pi(a \mid s) = \begin{cases} 0.2 & a^1 \\ 0.8 & a^2 \end{cases}$

$\pi(s) = a$ 50%

$\pi(s) = v$ 50%

On-policy

$\pi_\theta(s) = a$

$\pi_\theta(a \mid s) = \begin{cases} 0.2 & a^1 \\ 0.8 & a^2 \end{cases}$

$S \rightarrow$ ... $a^1$  0.1

$a^2$  0.1

Softmax

$$s \rightarrow \boxed{\pi_\theta} \rightarrow \begin{matrix} \mu_\theta(s) \\ \sigma_\theta(s) \end{matrix} \quad \pi(\cdot \mid s)$$
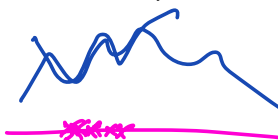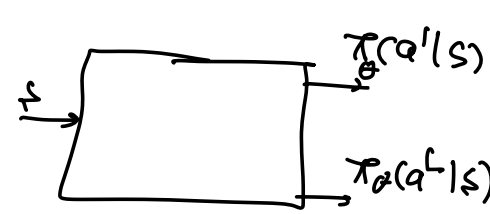
$$a \in [-1 \quad 1]$$

$$\pi_\theta(a \mid s) = \frac{1}{\sqrt{2\pi}\,\sigma_\theta(s)} \exp\left(- \frac{(a - \mu_\theta(s))}{2\,\sigma_\theta^2(s)}\right)$$



$$\pi_\theta(a \mid s) = \pi(a \mid s, \theta) = P(a_t = a \mid s_t = s, \theta)$$

$$s$$

$$s \rightarrow \boxed{\phantom{xx}} \rightarrow \begin{matrix} \pi_\theta(a' \mid s) \\ \pi_\theta(a^L \mid s) \end{matrix}$$
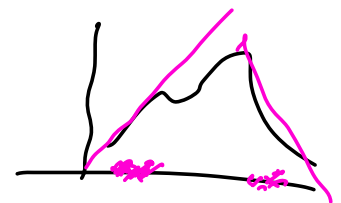
$$J(\theta) = V_{\pi_\theta}(s_0)$$

$\hookrightarrow$ all $s_0$ under $\pi_\theta$

Goal: find $\theta$ that maximizes $J(\theta)$

$$J(\theta) = E\left[ R_t + \gamma R_{t+1} + \gamma^2 R_{t+2} + \cdots \mid a_{t:\infty} \sim \pi_\theta, s \right]$$

$$X = E[Q(s, \pi_\theta)]$$

$$J(\theta) = \mathop{E}_{s \sim \pi_\theta} \left[ \sum_{a \in A} Q_{\pi_\theta}(s,a) \, \pi_\theta(a|s) \mid s \right]$$

Maximizing $J(\theta) \implies \theta^* = \text{argmax} \dfrac{J(\theta)}{E_{\pi_\theta}\left[\sum_a \pi(a|s) Q_{\pi_\theta}(s,a)\right]}$

$\overset{\text{new}}{\theta} \leftarrow \overset{\text{old}}{\theta} + \alpha \, \nabla_\theta J(\theta)$

$$\nabla_\theta J(\theta) = E_{\pi_\theta}\left[ \sum_{a \in A} Q_{\pi_\theta}(s,a) \, \nabla_\theta \pi(a|s,\theta) \right]$$

$$\alpha \sum_s \mu_{(s)} \sum_{a \in A} Q_{\pi_\theta}(s,a) \, \nabla_\theta \pi(a|s,\theta)$$

Distribution of state $s$ under policy $\pi_\theta$

# REINFORCE

$$\nabla J(\theta) \propto \sum_{s} \mu_{(s)} \sum_{a} Q_{\pi_\theta}(s,a) \nabla_\theta \pi(a|s,\theta)$$

$$= E_{\pi_\theta}\left[ \sum_{a} Q_{\pi_\theta}(s,a) \nabla_\theta \pi(a|s,\theta) \right]$$

$$= E_{\pi_\theta}\left[ \sum_{a} \pi(a|s,\theta) \; Q_{\pi_\theta}(s,a) \; \frac{\nabla_\theta(a|s,\theta)}{\pi(a|s,\theta)} \right]$$

$$= E_{\pi_\theta}\left[ \sum_{a} \pi(a|s_t,\theta) Q_{\pi_\theta}(s_t,a) \nabla_\theta Ln \, \pi_\theta(a_t|s_t,\theta) \right]$$

$$\nabla J(\theta) = E_{\pi_\theta}\left[ Q_{\pi_\theta}(s_t,a_t) \nabla_\theta Ln \, \pi_\theta(a_t|s_t,\theta) \right]$$

$$a_t \sim \pi(a|s_t,\theta)$$

REINFarce Trick

$$\nabla Ln \, x = \frac{\nabla x}{x}$$

$$\overset{\pi_{\theta}}{\overline{S_0, a_0, S_1, r_1}} \longrightarrow G_1 = r_1 + \gamma r_2 + \gamma^2 r_3 \cdots \gamma^{T} r_T$$

$$S_1, a_1, S_2, r_2 \longrightarrow G_2 = r_3 + \gamma r_3 + \cdots$$

$$\vdots$$

$$S_{T-1}, a_{T-1}, S_T, r_T \longrightarrow G_t = r_T$$

$$\nabla J(\theta) = E_{\pi_\theta}\left[ Q_{\pi_\theta}(S_t, a_t) \nabla_\theta \ln \pi_\theta(a_t | S_t, \theta) \right]$$

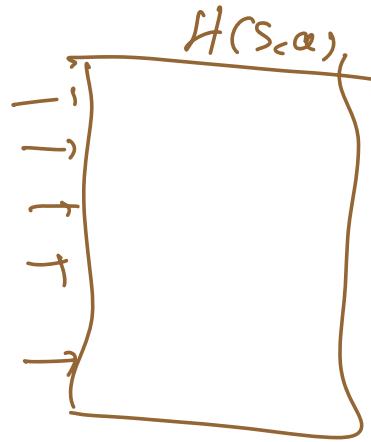$$\approx \frac{1}{N} \sum_{i=1}^{N} G_t^i \nabla_\theta \ln(a_i | S_i, \theta)$$

$$= \frac{1}{T}\left[ G_1 \nabla_\theta \ln(a_0 | S_0, \theta) + G_2 \nabla_\theta \ln(a - - |) \right]$$

$$\theta \leftarrow \theta + \alpha \nabla J(\theta)$$

# Simpler Function Approximation:

Actor critic $H(s,a) \rightarrow$ pred

$$\pi(a|s) = \frac{e^{H(s,a)}}{\sum_{a'} e^{H(s,a')}}$$



$H(s,a)$

$$H(s,a) = \theta^T \Phi(s,a)$$

$$\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$$

$$\nabla_\theta J(\theta) = E\left[ Q_{\pi_\theta}(s_t, a_t) \nabla_\theta \ln \pi_\theta(a_t|s_t, \theta) \right]$$

$$\pi_\theta(a|s) = \frac{e^{\theta^T \Phi(s,a)}}{\sum_{a'} e^{\theta^T \Phi(s,a')}}$$

$$\nabla_\theta \ln \pi_\theta(a_t|s_t, \theta) = \Phi(s_t, a_t) - \sum_{a'} \pi(a'|s_t, \theta) \Phi(s_t, a')$$

## II  Gaussian Function App

$$\pi(a \mid s, \theta) = \frac{1}{\sigma(s,\theta)\sqrt{2\pi}} \exp\left(-\frac{(a - \mu_{(s,\theta)})^2}{2\,\sigma(s,\theta)^2}\right)$$

$$\mu(s,\theta) = \theta_\mu^T \Phi_\mu(s) \qquad \qquad \theta = \begin{bmatrix} \theta_\mu \\ \theta_\sigma \end{bmatrix}$$

$$\sigma(s,\theta) = \theta_\sigma^T \phi_\sigma(s)$$

$$\nabla_\theta \ln \pi(a \mid s_t, \theta) = \qquad \checkmark$$

---

# REINFORCE with Baseline

$$\nabla_\theta J(\theta) \propto \sum_s \mu(s) \sum_a Q^\pi(s,a) \nabla_\theta \pi(a \mid s, \theta)$$

$$= \sum_s \mu(s) \sum_a \left( Q^\pi(s,a) - b(s) \right) \nabla_\theta \pi(a \mid s, \theta)$$

why? $\rightarrow \sum_a b(s) \nabla_\theta \pi(a \mid s, \theta) = b(s) \underbrace{\sum_a \nabla_\theta \pi(a \mid s, \theta)}$

$$\underbrace{\nabla_\theta \sum_a \pi(a \mid s, \theta)}_{\underbrace{1}_{0}}$$

$$\boxed{b(s) = V_{\pi_\theta}(s)}$$

$$\tau_\theta \begin{cases} s_0, a_0, r_1, s_1 & G_1 \\ s_1, a_1, r_2, s_2 & G_2 \\ \quad \vdots \\ s_{T-1}, a_{T-1}, r_T, s_T & G_T \end{cases}$$

$$V(s, w) = V_W(s)$$

$$\pi_\theta(a \mid s)$$

$$L(w) = E\left[\left(G_t - V_W(s_t)\right)^2\right]$$

Gradient
$\longrightarrow$
Decet

$$W \leftarrow W - \frac{1}{2}\alpha \nabla_W L(w)$$

$$W \leftarrow W + \alpha \left(G_t - V_W(s_t)\right) \nabla V_W(s_t)$$

$$J(\theta) \longrightarrow$$

$$\theta \leftarrow \theta + \alpha \underline{\nabla J(\theta)}$$

$$\left(G_t - V_W(s_t)\right) \nabla_\theta \ln \pi_\theta(a \mid s)$$

REINFORCE $\rightarrow$ MC

TD

$$\theta_{t+1} = \theta_t + \alpha \left( \underbrace{G_t} - \hat{V}(S_t, w) \right) \nabla_\theta \ln \pi_\theta(a_t | S_t)$$

$$R_{t+1} + \gamma \hat{V}(S_{t+1}, w)$$

A2C

**REINFORCE with Baseline (episodic), for estimating $\pi_{\boldsymbol{\theta}} \approx \pi_*$**

Input: a differentiable policy parameterization $\pi(a|s, \boldsymbol{\theta})$
Input: a differentiable state-value function parameterization $\hat{v}(s, \mathbf{w})$
Algorithm parameters: step sizes $\alpha^{\boldsymbol{\theta}} > 0$, $\alpha^{\mathbf{w}} > 0$
Initialize policy parameter $\boldsymbol{\theta} \in \mathbb{R}^{d'}$ and state-value weights $\mathbf{w} \in \mathbb{R}^d$ (e.g., to $\mathbf{0}$)

Loop forever (for each episode):
    Generate an episode $S_0, A_0, R_1, \ldots, S_{T-1}, A_{T-1}, R_T$, following $\pi(\cdot|\cdot, \boldsymbol{\theta})$
    Loop for each step of the episode $t = 0, 1, \ldots, T-1$:
        $G \leftarrow \sum_{k=t+1}^{T} \gamma^{k-t-1} R_k$                                       $(G_t)$
        $\delta \leftarrow G - \hat{v}(S_t, \mathbf{w})$
        $\mathbf{w} \leftarrow \mathbf{w} + \alpha^{\mathbf{w}} \delta \nabla \hat{v}(S_t, \mathbf{w})$
        $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \quad \alpha^{\boldsymbol{\theta}} \delta \nabla \ln \pi(A_t|S_t, \boldsymbol{\theta})$

**One-step Actor–Critic (episodic), for estimating $\pi_\theta \approx \pi_*$**

Input: a differentiable policy parameterization $\pi(a|s,\boldsymbol{\theta})$
Input: a differentiable state-value function parameterization $\hat{v}(s,\mathbf{w})$
Parameters: step sizes $\alpha^{\boldsymbol{\theta}} > 0$, $\alpha^{\mathbf{w}} > 0$
Initialize policy parameter $\boldsymbol{\theta} \in \mathbb{R}^{d'}$ and state-value weights $\mathbf{w} \in \mathbb{R}^d$ (e.g., to $\mathbf{0}$)
Loop forever (for each episode):
    Initialize $S$ (first state of episode)
    $I \leftarrow 1$
    Loop while $S$ is not terminal (for each time step):
        $A \sim \pi(\cdot|S,\boldsymbol{\theta})$
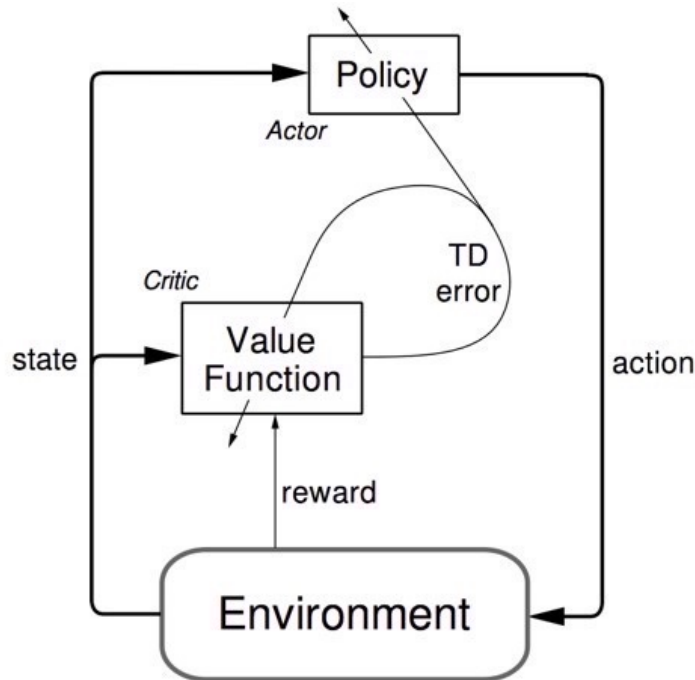        Take action $A$, observe $S', R$
        $\delta \leftarrow R + \gamma \hat{v}(S',\mathbf{w}) - \hat{v}(S,\mathbf{w})$       (if $S'$ is terminal, then $\hat{v}(S',\mathbf{w}) \doteq 0$)
        $\mathbf{w} \leftarrow \mathbf{w} + \alpha^{\mathbf{w}} I \delta \nabla \hat{v}(S,\mathbf{w})$
        $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha^{\boldsymbol{\theta}} I \delta \nabla \ln \pi(A|S,\boldsymbol{\theta})$
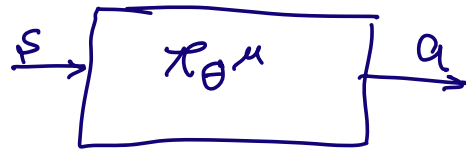        $I \leftarrow \gamma I$
        $S \leftarrow S'$

# Deep Deterministic Policy Gradient (DDPG)

$Q$: Q-Network

$Q'$: Target Network

$\theta^\mu$: Deterministic Policy Net

$\theta^{\mu'}$: Target Policy Net



$Q$ Network: input $s$, $a$ → output $Q(s,a)$

$\pi_{\theta^\mu}$: input $s$ → output $a$

$D$

Policy Network $\theta^\mu$

$$s_0, a_0, s_1, r_1$$
$$s_1, a_1, s_1, r_2$$
$$\vdots$$

miniBatch $B$ $\Rightarrow$

Upate Q Net

$$s_0, a_0, s_1, r_1 \quad y_q = r + \gamma \max Q'(s', a') \ldots$$
$$\vdots$$

Target Net

$$w \leftarrow w + \alpha \frac{1}{B} \sum_{s,a,s'r \in B} \left( r + \gamma Q'(s', \mu'_{(s')}) - Q(s,a) \right) \nabla_w Q_{(s,a)}$$

$$J(\theta) = E\left[ Q(s,a) \mid s_t = s, a_t = \mu_{(s)} \right]$$

$$\nabla_\theta J(\theta) \approx \cdot \nabla_a Q_w(s, a) \nabla_{\theta^\mu} \mu_{\theta^\mu}(s)$$

**Algorithm 1** DDPG algorithm

---

Randomly initialize critic network $Q(s, a|\theta^Q)$ and actor $\mu(s|\theta^\mu)$ with weights $\theta^Q$ and $\theta^\mu$.
Initialize target network $Q'$ and $\mu'$ with weights $\theta^{Q'} \leftarrow \theta^Q$, $\theta^{\mu'} \leftarrow \theta^\mu$
Initialize replay buffer $R$
**for** episode = 1, M **do**
    Initialize a random process $\mathcal{N}$ for action exploration
    Receive initial observation state $s_1$
    **for** t = 1, T **do**
        Select action $a_t = \mu(s_t|\theta^\mu) + \mathcal{N}_t$ according to the current policy and exploration noise
        Execute action $a_t$ and observe reward $r_t$ and observe new state $s_{t+1}$
        Store transition $(s_t, a_t, r_t, s_{t+1})$ in $R$
        Sample a random minibatch of $N$ transitions $(s_i, a_i, r_i, s_{i+1})$ from $R$
        Set $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1}|\theta^{\mu'})|\theta^{Q'})$
        Update critic by minimizing the loss: $L = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i|\theta^Q))^2$
        Update the actor policy using the sampled policy gradient:

$$\nabla_{\theta^\mu} J \approx \frac{1}{N} \sum_i \nabla_a Q(s, a|\theta^Q)|_{s=s_i, a=\mu(s_i)} \nabla_{\theta^\mu} \mu(s|\theta^\mu)|_{s_i}$$

        Update the target networks:

$$\theta^{Q'} \leftarrow \tau\theta^Q + (1-\tau)\theta^{Q'}$$
$$\theta^{\mu'} \leftarrow \tau\theta^\mu + (1-\tau)\theta^{\mu'}$$

    **end for**
**end for**

---

*Algorithm 1 of the paper "Continuous Control with Deep Reinforcement Learning" by Timothy P. Lillicrap et al*