Learning Curve

$\mathcal{L}$

↑ accuracy
avg return $EV(s_0)$
$s_0 \sim p_0$
2,3

number of samples

## Actor - Critic Deep RL

Last time : Policy Grad



Eval

Experience

Improve π

Before : Q-learning

Eval

$\hat{Q} = \sum r_+$ ← high variance

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

Q-network params
$\phi \leftarrow \phi + \alpha \frac{dQ}{d\phi}\left(Q_\phi - r - \gamma \max_{a'} Q_{\phi'}(s', a')\right)$

$\pi(a \mid s) = \begin{cases} 1 - \varepsilon & \text{if } a = \text{argmax } Q_\phi(s, a) \\ \frac{1}{|A| - 1} & \text{o.w.} \end{cases}$

Actor - Critic
↑ policy    ↑ Value fn

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_i \sum_t \nabla_\theta \log \pi_\theta(a \mid s) \left( \underbrace{Q(\ ) - b}_{A^\pi(s,a)} \right)$$

know $Q^\pi(s,a)$

$V^\pi(s) = E\left[Q(s,a) \mid \pi(a \mid s)\right]$

$A^\pi(s,a) = Q^\pi(s,a) - V(s)$

$Q_\phi, V_\phi, $ or $A_\phi$ ?        $A_\phi$ - use it directly

$$V_\phi(s)$$

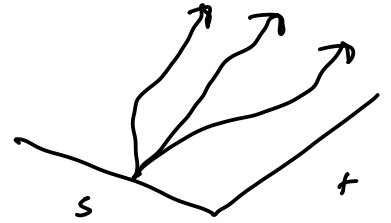$$Q(s,a) = r(s,a) +\gamma E\left[V(s')\right]$$

$$\approx r(s,a) + \gamma V(s')$$

$$A(s,a) \approx r(s,a) + \gamma V(s') - V(s)$$

targets for $V(\cdot)$

$$- \left\{\left(s_{i,t}, \sum_{t'=t}^{T} r(s_{i,t'}, a_{i,t'})\right)\right\}$$   Monte Carlo Target

$$- \left\{\left(s_{i,t}, r(s_{i,t}, a_{i,t}) + \gamma V_\phi(s_{i,t+1})\right)\right\}$$  "Bootstrapped"

target

**Batch Actor-Critic Algorithm**                              Online

loop

Collect data by running $\pi_\theta$                  $\longrightarrow$ get $(s, a, s', r)$

fit $V_\phi^{\pi_\theta}$                                                   update $\phi$ using $r + \gamma V_\phi(s')$
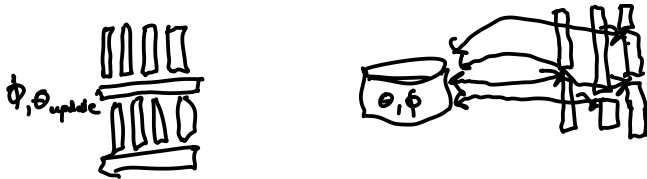
evaluate $A^\pi(s,a) = r(s,a) + \gamma V_\phi^\pi(s) - V_\phi^\pi(s)$

$\nabla_\theta J(\theta) = \sum_i \nabla_\theta \log \pi_\theta(a|s) \hat{A}(s,a)$

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

---

A3C        Asynchronous Advantage Actor-Critic

$\phi, \theta_{update}$

Design Choices:    Network

two network              shared

$s$  $\phi$  $V_\phi$            $s$  $V_\phi$

$s$  $\pi_\theta$                     $\pi_\theta$

more data,              less data, tricky to train
easier to train

# Generalized Advantage Estimation   GAE

$\hat{Q}$ ← high variance

$V_\alpha$ ← biased  ← especially in future

$$\hat{A}_n^\pi(s_t, a_t) = \sum_{t'=t}^{t+n} \gamma^{t'-t} r(s_t, a_t) - V_\phi^\pi(s_t) + \gamma^n V_\phi^\pi(s_{t+n})$$

$$\hat{A}_{GAE}^\pi = \sum w_n \hat{A}_n^\pi \qquad w_n \propto \lambda^{n-1} \quad \begin{array}{c}\text{weight} \\ \text{of} \\ \text{MC part}\end{array}$$

$$= r_t + \gamma\left((1-\lambda) V_\phi(s_{t+1}) + \lambda\left(r_{t+1} + \gamma(1-\lambda) V_\phi(s_{t+2})\right)\dots\right.$$

$$= \sum_{t'=t}^{\infty} (\gamma\lambda)^{t'-t} \delta_{t'} \qquad \delta_{t'} = r(s_{t'}, a_{t'}) + \gamma V_\phi(s_{t'+1}) - V_\phi(s_{t'})$$
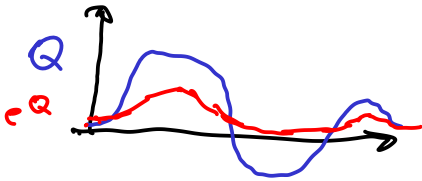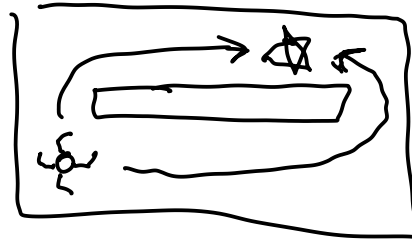
Eligibility Traces

---

# Soft - Actor - Critic

$$\pi(a|s) \propto e^{Q(s,a)}$$

$$\pi_{\text{MaxEnt}}^* = \text{argmax} \left[\sum_t r_t + \beta H(\pi(a|s))\right]$$

↑ entropy

low H         higher H

$$\mathbb{E}_{s_t}\left[D_{KL}\left(\pi_\theta(a,s) \,\Big|\Big|\, \frac{Q(s,a)}{Z(s)}\right)\right]$$

Robust Convergence

---

Didn't Cover      Model - Based Deep RL

1. Learn Dynamics Model

2. Use model Free RL / Planning / LQR

---

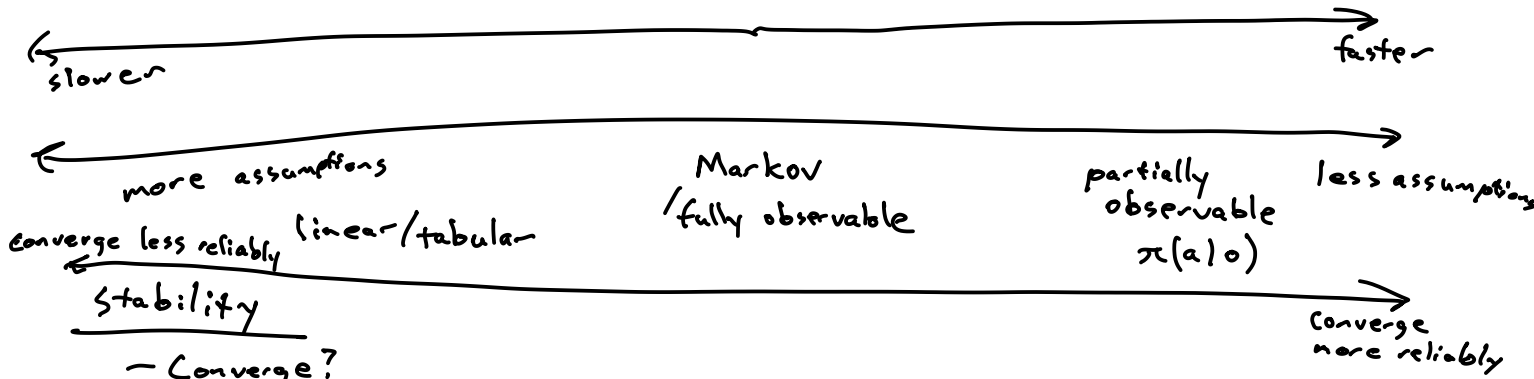How to Choose          Size S, A

<u>Sample Efficiency</u>                    <u>Stability / Ease of Use</u>

- how much experience does it take to get a good policy

- On Policy / <u>off - policy</u>

Less Samples ←————————————————————————————→ More Samples

| model-based Shallow RL | model-based Deep RL | Off-policy Q-learning | actor-critic | On-policy policy gradient | evolutionary gradientfree |

Wall Clock time

With fast simulator, roughly reversed

slower ←————————————————————————————→ faster

more assumptions ←————— Markov/fully observable ——— partially observable $\pi(a|o)$ ——→ less assumptions

converge less reliably ← linear/tabular ————————————————————→ converge more reliably

Stability

- Converge?
- to what?
- every time?

$$\|x\|^\infty = \max_i |x_i|$$

$S = \{s^1, s^2\}$



$V(s^1)$ vs $V(s^2)$

$$V \leftarrow B[V]$$

$$V(s) \approx \lambda^T \beta(s_i)$$

$$\Pi[V] = \arg\min_\lambda \sum_{i=1}^{N} \left(\lambda^T \beta(s_i) - V(s_i)\right)^2$$

$$\Pi[B[V]]$$

$$\|\ \|^{\infty}$$

$$\|f(x) - f(y)\| \leq \gamma \|x - y\|$$

Contraction mapping

Deep Q-learning        less stable

Policy Gradient        more stable