

DRL Lec 1

- Motivation

⇒ deep learning helps us handle unstructured environments

(it can process complex sensory inputs)

⇒ reinforcement learning provides a formalism for behavior

(it can choose complex actions)

∴ deep models are what allow reinforcement learning algorithms to solve complex problems end-to-end

- Strength of RL on real-world sequential decision making

⇒ reinforcement learning can be naturally integrated with artificial neural networks to obtain high quality generalization resulting in a significant learning speedup

⇒ reinforcement learning agents can save many learning trials by using an action model which can be learned online.

⇒ reinforcement learning agents can take advantages of instructive training instances provided by human teachers resulting in a significant learning speedup

⇒ reinforcement learning agents can significantly reduce learning time by hierarchical learning
(first solve elementary learning problems and then by combining the solutions, solve a complex one)

⇒ reinforcement learning agents can deal with non-Markovian environments by having a memory of their past

- What DRL can do

⇒ acquire high degree of proficiency in domains governed by simple, known rules

⇒ learn simple skills with raw sensory inputs, given enough experience

⇒ learn from imitating enough human provided expert behavior.

- Challenges

⇒ slow convergence, transfer learning, reward function, role of prediction

Words from Alan Turing

: Instead of trying to produce a program to simulate the adult mind,

why not rather try to produce one which simulates the child's?

If this were then subjected to an appropriate course of education,

one would obtain an adult brain

DRL Lec 2

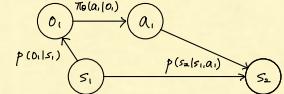
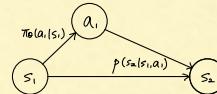
- Supervised learning of behaviors

- Notations

⇒ s_t : state, o_t : observation, a_t : action

⇒ $\pi_\theta(a_t|s_t)$: fully observed policy

⇒ $\pi_\theta(a_t|o_t)$: partially observed policy



- Imitation learning (≈ behavior cloning)

: training data $D = \{o_t, a_t\}_{t=1}^N$ is made by human experts, then the policy $\pi_\theta(a_t|o_t)$ is trained via supervised learning framework.

⇒ problems : $p_{data}(o_t) \neq p_{\pi_\theta}(o_t)$ (≈ distributional drift)

→ little mistake in the beginning results in large error at the end.

(∴ it gets to states that have never seen by the human experts)

⇒ solutions

→ DAgger (Dataset Aggregation)

⇒ train on several trajectories that involve mistakes and corrections

⇒ collect training data from $p_{\pi_\theta}(o)$ instead of $p_{data}(o)$ so that $p_{data}(o) \approx p_{\pi_\theta}(o)$

① train $\pi_\theta(a_t|o_t)$ from human data $D = \{o_t, a_t\}_{t=1}^N$

② run $\pi_\theta(a_t|o_t)$ to get dataset $D_\pi = \{o_i\}_{i=1}^M$

③ ask human to label D_π with actions $\{a_i\}_{i=1}^M$ → limitation

④ aggregate $D \leftarrow D \cup D_\pi$

→ others : hacks, better models, stable trajectory distribution

- Cause of failure

⇒ Non-markovian behavior

→ behavior depends on all the past observations

→ may be resolved by using RNN on outputs of several past observations

(but, this may also fail by causal confusion)

⇒ Multi-modal behavior

→ may be resolved by the follows.

① output mixture of gaussians → $w_1, \mu_1, \Sigma_1, \dots, w_n, \mu_n, \Sigma_n$

② latent variable models → $\xi \sim N(0, I)$

③ autoregressive discretization → discretize high dimensional action space by one dimension at a time

- Theoretical analysis

⇒ Assumption → $c(s, a) = \begin{cases} 0 & \text{if } a = \pi^*(s) \\ 1 & \text{otherwise} \end{cases}$

① $\pi_\theta(a \neq \pi^*(s) | s) \leq \varepsilon \quad \forall s \in \text{Domain}$

⇒ $E[\sum_t c(s_t, a_t)] \leq \varepsilon T + (1-\varepsilon) \cdot \varepsilon (T-1) + (1-\varepsilon)^2 \varepsilon (T-2) + \dots = O(\varepsilon T^2)$

(once it makes mistake, it ends up facing cost every time)

② $\pi_\theta(a \neq \pi^*(s) | s) \leq \varepsilon \quad \forall s \sim p_{\text{train}}(s)$ (Only from human)

⇒ $p_{\pi_\theta}(s_t) = (1-\varepsilon)^t p_{\text{train}}(s_t) + (1-(1-\varepsilon)^t) p_{\text{mistake}}(s_t)$

⇒ $\sum_t |p_{\pi_\theta}(s_t) - p_{\text{train}}(s_t)| = \sum_t (1-(1-\varepsilon)^t) |p_{\text{mistake}}(s_t) - p_{\text{train}}(s_t)|$

$\leq 2(1-(1-\varepsilon)^t) \leq 2\varepsilon t$

⇒ $E[\sum_t c(s_t, a_t)] = \sum_t \sum_{s_t} p_{\pi_\theta}(s_t) c(s_t, a_t)$

$\leq \sum_t \sum_{s_t} p_{\text{train}}(s_t) c(s_t, a_t) + |p_{\pi_\theta}(s_t) - p_{\text{train}}(s_t)| C_{\max}$

$\leq \sum_t \varepsilon + 2\varepsilon t = O(\varepsilon T^2)$

③ DAgger ($p_{\pi_\theta}(s) \rightarrow p_{\text{train}}(s)$)

⇒ $E[\sum_t c(s_t, a_t)] \leq \varepsilon T = O(\varepsilon T)$

(even it makes mistake, it immediately face correction)

DRL Lec 3

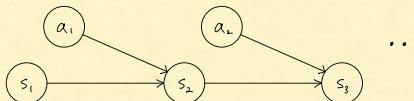
• Introduction to reinforcement learning

1. Markov Decision Process ($M = \{S, A, T, r\}$)

$\Rightarrow S$: state space, A : action space, T : transition operator, r : reward

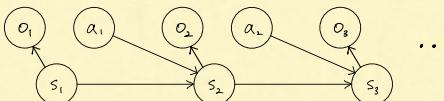
$$(M_{t,j,k} = p(s_t=j), \bar{M}_{t,k} = p(a_t=k), T_{i,j,k} = p(s_{t+1}=i | s_t=j, a_t=k), r(s_t, a_t))$$

$$(M_{t,n,i} = \sum_{j,k} T_{i,j,k} M_{t,j,k})$$



2. Partially observed Markov Decision Process ($M = \{S, A, O, T, E, R\}$)

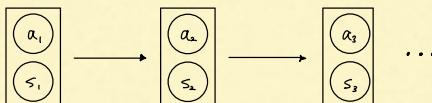
$\Rightarrow O$: observation space, E : emission probability ($p(o_t | s_t)$)



3. Reinforcement learning

$$\Rightarrow P_\theta(z) = p(s_1, a_1, \dots, s_T, a_T)$$

$$= P(s_1) \cdot \Pi_\theta(a_1 | s_1) \cdot \prod_{t=1}^{T-1} p(s_{t+1} | s_t, a_t) \Pi_\theta(a_{t+1} | s_{t+1}) \\ = P_\theta(s_1, a_1) \prod_{t=1}^{T-1} P_\theta(s_{t+1}, a_{t+1} | s_t, a_t) \rightarrow \text{markov chain on } (s_t, a_t)$$



($P_\theta(s_t, a_t)$ converges to stationary distribution $p_\theta(s, a)$ under ergodicity)

$$\Rightarrow \theta^* = \arg \max_{\theta} \frac{1}{T} \mathbb{E}_{z \sim P_\theta(z)} \left[\sum_{t=1}^T r(s_t, a_t) \right]$$

$$= \arg \max_{\theta} \frac{1}{T} \sum_{t=1}^T \mathbb{E}_{(s_t, a_t) \sim P_\theta(z)} [r(s_t, a_t)] \rightarrow \text{finite horizon case}$$

$$\approx \arg \max_{\theta} \mathbb{E}_{(s, a) \sim P_\theta(z)} [r(s, a)] \text{ as } T \rightarrow \infty \rightarrow \text{infinite horizon case}$$

4. Estimating reward

$$\left(\begin{array}{l} Q\text{-function} (Q^\pi(s_t, a_t)) \\ : Q^\pi(s_t, a_t) = \sum_{t'=t}^T \mathbb{E}_\pi [r(s_{t'}, a_{t'}) | s_t, a_t] = r(s_t, a_t) + \mathbb{E}_{s_{t+1}} [V^\pi(s_{t+1})] \\ \text{Value-function} (V^\pi(s_t)) \\ : V^\pi(s_t) = \sum_{t'=t}^T \mathbb{E}_\pi [r(s_{t'}, a_{t'}) | s_t] = \mathbb{E}_{a_t} [Q^\pi(s_t, a_t)] \\ \Rightarrow \mathbb{E}_{z \sim P_\theta(z)} \left[\sum_{t=1}^T r(s_t, a_t) \right] = \mathbb{E}_{s \sim P(z)} \left[\mathbb{E}_{a \sim \pi_\theta(a|s)} \left[r(s, a) + \mathbb{E}_{s' \sim P(s'|s,a)} [r(s', a) + \dots] \right] \right] \\ = \mathbb{E}_{s \sim P(z)} \left[\mathbb{E}_{a \sim \pi_\theta(a|s)} [Q^\pi(s, a)] \right] = \mathbb{E}_{s \sim P(z)} [V^\pi(s)] \end{array} \right)$$

\rightarrow idea 1: for all t , set $\pi'(a|s_t) = 1$ when $a = \arg \max_{a_t} Q^\pi(s_t, a_t)$

\rightarrow idea 2: modify $\pi(a|s_t)$ by computing gradients when $Q^\pi(s_t, a_t) > V^\pi(s_t)$

• Algorithms

1. Policy gradients

$$\textcircled{1} \text{ estimate } R_\pi = \mathbb{E} \left[\sum_t r(s_t, a_t) \right]$$

$$\textcircled{2} \quad \theta \leftarrow \theta + \alpha \nabla_\theta R_\pi$$

2. Value-based

$$\textcircled{1} \text{ fit } V^\pi(s) \text{ or } Q^\pi(s, a)$$

$$\textcircled{2} \text{ set } \pi'(a|s_t) = 1 \text{ when } a = \arg \max_{a_t} Q^\pi(s_t, a_t)$$

3. Actor-critic

$$\textcircled{1} \text{ fit } V^\pi(s) \text{ or } Q^\pi(s, a)$$

$$\textcircled{2} \quad \theta \leftarrow \theta + \alpha \nabla_\theta \mathbb{E} [Q^\pi(s, a)]$$

4. Model-based

$$\textcircled{1} \text{ learn } f_\phi \text{ s.t. } s_{t+1} \approx f_\phi(s_t, a_t)$$

$$\textcircled{2} \quad \left[\begin{array}{l} \text{planning (e.g. Monte Carlo Tree search)} \\ \text{learn a value function (e.g. dynamic programming)} \end{array} \right]$$

• Comparison

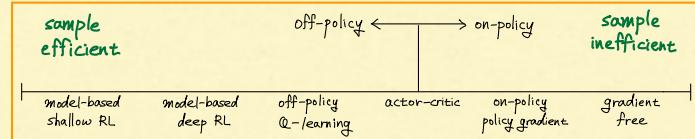
1. Sample efficiency (# wall clock time)

\Rightarrow on-policy \rightarrow inefficient

: each time the policy is changed, new samples need to be generated

\Rightarrow off-policy

: each time the policy is changed, new samples need not to be generated



2. Stability and ease of use

\Rightarrow no guarantee to converge when not using gradient methods

$\textcircled{1}$ policy gradient \rightarrow high variance

$\textcircled{2}$ value-based \rightarrow fitted value may not equal to expected reward

$\textcircled{3}$ model-based \rightarrow better model may not lead to better policy

3. Assumptions

$\textcircled{1}$ Full observability \rightarrow value-based

$\textcircled{2}$ Episodic learning \rightarrow policy gradient, model-based

$\textcircled{3}$ Continuity or smoothness \rightarrow value-based, model-based

DRL Lec 4

• Policy Gradients

1. Derivation

$$\begin{aligned} \Rightarrow J(\theta) &= \mathbb{E}_{z \sim p_\theta(z)} \left[\sum_{t=1}^T r(s_t, a_t) \right] = \mathbb{E}_{z \sim p_\theta(z)} [r(z)] = \int p_\theta(z) r(z) dz \\ \Rightarrow \nabla_\theta J(\theta) &= \int \nabla_\theta p_\theta(z) r(z) dz \\ &= \int p_\theta(z) \cdot \nabla_\theta \log p_\theta(z) \cdot r(z) dz \quad \rightarrow \text{score function trick} \\ &= \mathbb{E}_{z \sim p_\theta(z)} [\nabla_\theta \log p_\theta(z) \cdot r(z)] \\ &= \mathbb{E}_{z \sim p_\theta(z)} [\nabla_\theta (\log p(z) + \sum_{t=1}^T \log T_\theta(a_{i,t}|s_{i,t}) + \sum_{t=1}^{T-1} \log p(s_{i+1}|s_i, a_i)) \cdot r(z)] \\ &= \mathbb{E}_{z \sim p_\theta(z)} \left[\left(\sum_{t=1}^T \nabla_\theta \log T_\theta(a_{i,t}|s_{i,t}) \right) \left(\sum_{t=1}^T r(s_t, a_t) \right) \right] \\ &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log T_\theta(a_{i,t}|s_{i,t}) \left(\sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right) \end{aligned}$$

⇒ REINFORCE

$$\begin{aligned} \textcircled{1} \text{ sample } \{z^i\} \text{ from } \pi_\theta(a_t|s_t) \\ \textcircled{2} \quad \nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) \right) \left(\sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right) \\ \textcircled{3} \quad \theta \leftarrow \theta + \alpha \nabla_\theta J(\theta) \end{aligned}$$

episode with high reward becomes more likely

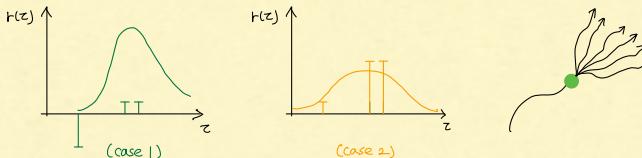
③ Maximum likelihood

$$\begin{aligned} \textcircled{1} \text{ sample } \{z^i\} \text{ from } \pi_\theta(a_t|s_t) \\ \textcircled{2} \quad \nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) \right) \quad \rightarrow \text{every episode is equally likely} \\ \textcircled{3} \quad \theta \leftarrow \theta + \alpha \nabla_\theta J(\theta) \end{aligned}$$

* Partial observability can be used without modification as Markov property is not used.

$$(\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \left(\sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) \right) \left(\sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right))$$

* Limitation : highly variable on the scale of the reward and the choice of episode



2. Variance reduction

① Causality

$$\Rightarrow \nabla_\theta J(\theta) = \mathbb{E}_{z \sim p_\theta(z)} \left[\sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) \left(\sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right) \right]$$

(future actions cannot affect past rewards)

② Subtracting baseline

$$\Rightarrow \nabla_\theta J(\theta) = \mathbb{E}_{z \sim p_\theta(z)} \left[\left(\sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) \right) \left(\sum_{t=1}^T r(s_{i,t}, a_{i,t}) - b \right) \right]$$

→ unbiasedness

$$\Rightarrow \mathbb{E}_{z \sim p_\theta(z)} [\nabla_\theta \log p_\theta(z) \cdot b] = \nabla_\theta \int p_\theta(z) dz \cdot b = 0$$

→ optimal baseline

$$\begin{aligned} \text{Var} &= \mathbb{E}_z [(\nabla_\theta \log p_\theta(z) \cdot (r(z) - b))^2] - \mathbb{E}_z [\nabla_\theta \log p_\theta(z) \cdot (r(z) - b)]^2 \\ \frac{\partial \text{Var}}{\partial b} &= -2 \mathbb{E}_z [(\nabla_\theta \log p_\theta(z))^2 \cdot (r(z) - b)] = 0 \quad \text{when } b = b^* \\ \Rightarrow b^* &= \frac{\mathbb{E}_z [(\nabla_\theta \log p_\theta(z))^2 \cdot r(z)]}{\mathbb{E}_z [(\nabla_\theta \log p_\theta(z))^2]} \end{aligned}$$

$$(b = \frac{1}{N} \sum_{i=1}^N r(z^i) \text{ works just fine, too})$$

3. Turning to off-policy

: use importance sampling to compute expected reward with updated parameter θ'

from the samples with previous parameter θ

$$\begin{aligned} \Rightarrow J(\theta') &= \mathbb{E}_{z \sim p_\theta(z)} \left[\frac{p_\theta(z)}{p_\theta(z)} r(z) \right] = \mathbb{E}_{z \sim p_\theta(z)} \left[\frac{\prod_{i=1}^T \pi_\theta(a_{i,t}|s_{i,t})}{\prod_{i=1}^T \pi_\theta(a_{i,t}|s_{i,t})} r(z) \right] \\ \Rightarrow \nabla_\theta J(\theta') &= \mathbb{E}_{z \sim p_\theta(z)} \left[\frac{\nabla_\theta p_\theta(z)}{p_\theta(z)} r(z) \right] = \mathbb{E}_{z \sim p_\theta(z)} \left[\nabla_\theta \log p_\theta(z) \cdot \frac{p_\theta(z)}{p_\theta(z)} r(z) \right] \\ &= \mathbb{E}_{z \sim p_\theta(z)} \left[\nabla_\theta \log p_\theta(z) \cdot \frac{\prod_{i=1}^T \pi_\theta(a_{i,t}|s_{i,t})}{\prod_{i=1}^T \pi_\theta(a_{i,t}|s_{i,t})} r(z) \right] \\ &\approx \mathbb{E}_{z \sim p_\theta(z)} \left[\sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) \left(\frac{\prod_{i=1}^T \pi_\theta(a_{i,t}|s_{i,t})}{\prod_{i=1}^T \pi_\theta(a_{i,t}|s_{i,t})} \right) \left(\sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right) \right] \quad (\because \text{causality}) \end{aligned}$$

(future actions cannot affect weights on past rewards)

$$\begin{aligned} &\approx \mathbb{E}_{z \sim p_\theta(z)} \left[\sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) \frac{\pi_\theta(a_{i,t}|s_{i,t})}{\pi_\theta(a_{i,t}|s_{i,t})} \left(\sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right) \right] \quad (\because \text{first order approximation}) \\ &\approx \frac{1}{N} \sum_{i=1}^N \left[\sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) \frac{\pi_\theta(a_{i,t}|s_{i,t})}{\pi_\theta(a_{i,t}|s_{i,t})} \left(\sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right) \right] \end{aligned}$$

DRL Lec 5

• Actor-critic

1. Derivation

$$\begin{aligned} \Rightarrow \nabla_\theta J(\theta) &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) \left(\sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right) \quad \rightarrow \text{rewards to go} \\ &\quad (\text{using a single trajectory for computing gradients}) \quad \text{Variance } \uparrow, \text{ unbiased} \\ &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) \left(\mathbb{E}_{\pi_\theta} \left[\sum_{t'=t}^T r(s_{i,t'}, a_{i,t'}) \mid s_{i,t}, a_{i,t} \right] \right) \quad \rightarrow \text{True expected rewards to go} \\ &\quad (\text{using multiple trajectories for computing gradients}) \quad \text{Variance } \downarrow, \text{ unbiased} \\ \text{AC} &= \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) Q^\pi(s_{i,t}, a_{i,t}) \\ &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) \left(Q^\pi(s_{i,t}, a_{i,t}) - \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T r(s_{i,t}, a_{i,t}) \right) \quad \text{Variance } \downarrow \\ \text{A2C} &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) \left(Q^\pi(s_{i,t}, a_{i,t}) - V^\pi(s_{i,t}) \right) \quad \rightarrow \text{Advantage} \quad \text{Variance } \downarrow \end{aligned}$$

$$\begin{aligned} * V^\pi(s_t) &= \mathbb{E}_{\pi_\theta} \left[\sum_{t'=t}^T r(s_{t'}, a_{t'}) \mid s_t \right] = \mathbb{E}_{a_t} [Q^\pi(s_t, a_t)] \\ * (Q^\pi(s_t, a_t) &= r(s_t, a_t) + \mathbb{E}_{\pi_\theta} \left[\sum_{t'=t}^T r(s_{t'}, a_{t'}) \mid s_t, a_t \right] = r(s_t, a_t) + \mathbb{E}_{s_{t+1} \sim p(s_{t+1}|s_t, a_t)} [V^\pi(s_{t+1})] \\ &\approx r(s_t, a_t) + V^\pi(s_{t+1}) \quad \text{from "the" single trajectory} \\ * A^\pi(s_t, a_t) &\approx r(s_t, a_t) + V^\pi(s_{t+1}) - V^\pi(s_t) \\ \therefore \text{When } V^\pi(s_t) \text{ is fitted, } Q^\pi(s_t, a_t) \text{ and } A^\pi(s_t, a_t) \text{ can be easily computed} \end{aligned}$$

⇒ Fitting value function

① Monte Carlo policy evaluation

$$\Rightarrow V^\pi(s_t) \approx \frac{1}{T} \sum_{t=1}^T r(s_{t'}, a_{t'}) \quad \text{or} \quad \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T r(s_{i,t}, a_{i,t})$$

② Monte Carlo policy evaluation with function approximation (critic)

$$\begin{aligned} \Rightarrow L(\phi) &= \frac{1}{N} \sum_{i=1}^N \| \hat{V}_\phi^\pi(s_{i,t}) - \sum_{t'=t}^T r(s_{i,t'}, a_{i,t'}) \|^2 \quad \rightarrow \text{Monte Carlo target} \\ &= \frac{1}{N} \sum_{i=1}^N \| \hat{V}_\phi^\pi(s_{i,t}) - (r(s_{i,t}, a_{i,t}) + \hat{V}_\phi^\pi(s_{i,t+1})) \|^2 \quad \rightarrow \text{Bootstrapped target} \end{aligned}$$

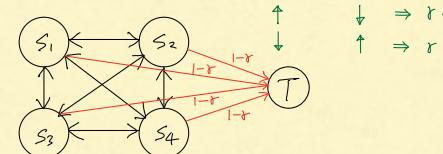
③ Policy gradient is applied with the estimated value from the critic network (actor)

$$\Rightarrow \nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) (Q^\pi(s_{i,t}, a_{i,t}) - V^\pi(s_{i,t}))$$

⇒ Discount factor (γ)

: better to get reward sooner than later

→ allows trade-off between bias & variance



$$\begin{aligned} \Rightarrow \nabla_\theta J(\theta) &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) \left(\sum_{t'=t}^T \gamma^{t-t'} r(s_{i,t'}, a_{i,t'}) \right) \\ &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) \left(r(s_{i,t}, a_{i,t}) + \gamma \hat{V}_\phi^\pi(s_{i,t+1}) - \hat{V}_\phi^\pi(s_{i,t}) \right) \end{aligned}$$

⇒ batch actor-critic

① sample $s_{i,t}, a_{i,t}^{2:T}$ from $p_\theta(z)$

$$\textcircled{2} \text{ fit } \hat{V}_\phi^\pi(s_t) \text{ to Monte Carlo target } \left(\frac{1}{T} \sum_{t'=t}^T r(s_{t'}, a_{t'}) \right)$$

$$\textcircled{3} \text{ evaluate } \hat{A}_\phi^\pi(s_t, a_t) = r(s_t, a_t) + \gamma \hat{V}_\phi^\pi(s_{t+1}) - \hat{V}_\phi^\pi(s_t)$$

$$\textcircled{4} \quad \nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) \hat{A}_\phi^\pi(s_{i,t}, a_{i,t})$$

$$\textcircled{5} \quad \theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$$

⇒ online actor-critic

① take action $a \sim \pi_\theta(a|s)$ and get (s, a, r, s')

$$\textcircled{2} \text{ fit } \hat{V}_\phi^\pi(s) \text{ to bootstrapped target } (r(s, a) + \gamma \hat{V}_\phi^\pi(s'))$$

$$\textcircled{3} \text{ evaluate } \hat{A}_\phi^\pi(s, a) = r(s, a) + \gamma \hat{V}_\phi^\pi(s') - \hat{V}_\phi^\pi(s)$$

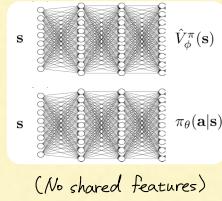
$$\textcircled{4} \quad \nabla_\theta J(\theta) \approx \nabla_\theta \log \pi_\theta(a|s) \hat{A}_\phi^\pi(s, a)$$

$$\textcircled{5} \quad \theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$$

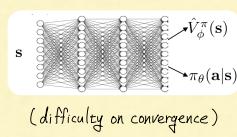
* Limitation : highly biased when the critic is not perfect, especially at the beginning

2. Architecture design

⇒ two networks



⇒ shared network



DRL Lec 6

• Value-function method

1. Knowing the transition dynamics

⇒ policy iteration

- ① evaluate $V^\pi(s) \leftarrow r(s, a) + \gamma \mathbb{E}_{s' \sim p(s'|s,a)} [V^\pi(s')]$ until convergence
- ② compute $A^\pi(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim p(s'|s,a)} [V^\pi(s')] - V^\pi(s)$
- ③ update $\pi(a|s) \leftarrow \pi'(a|s) = \begin{cases} 1 & \text{if } a = \operatorname{argmax}_a A^\pi(s, a) \\ 0 & \text{otherwise} \end{cases}$

- ④ evaluate $Q^\pi(s, a) \leftarrow r(s, a) + \gamma \mathbb{E}_{s' \sim p(s'|s,a)} [Q^\pi(s', a)]$
- ⑤ update $\pi(a|s) \leftarrow \pi'(a|s) = \begin{cases} 1 & \text{if } a = \operatorname{argmax}_a Q^\pi(s, a) \\ 0 & \text{otherwise} \end{cases}$

⇒ value iteration

: Value function is represented by a big table, one entry for each discrete s

- ① compute $Q(s, a) = r(s, a) + \gamma \mathbb{E}_{s' \sim p(s'|s,a)} [V(s')]$
- ② update $V(s) \leftarrow \max_a Q(s, a)$

⇒ fitted value iteration

: Value function is represented by a neural network with parameter ϕ

- ① set $y_i \leftarrow \max_a (r(s_i, a) + \gamma \mathbb{E}_{s' \sim p(s'|s,a)} [\hat{V}_\phi(s')])$
- ② update $\phi \leftarrow \operatorname{argmin}_\phi \frac{1}{N} \sum_i \| \hat{V}_\phi(s_i) - y_i \|^2$

2. Not knowing the transition dynamics

⇒ fitted Q iteration

: Q function is represented by a neural network with parameter ϕ

- ① collect dataset $\{(s_i, a_i, s'_i, r_i)\}$ using some policy → off-policy
- ② set $y_i \leftarrow r(s_i, a_i) + \gamma \max_a \hat{Q}_\phi(s'_i, a')$ from "the" single trajectory
- ③ update $\phi \leftarrow \operatorname{argmin}_\phi \frac{1}{N} \sum_i \| \hat{Q}_\phi(s_i, a_i) - y_i \|^2$

⇒ online Q iteration ($\approx Q$ -learning)

- ① take some action a_i and observe (s_i, a_i, s'_i, r_i) → off-policy
- ② set $y_i \leftarrow r(s_i, a_i) + \gamma \max_a \hat{Q}_\phi(s'_i, a')$ from "the" single trajectory
- ③ update $\phi \leftarrow \phi - \alpha \frac{d \hat{Q}_\phi}{d \phi} \Big|_{s=s_i, a=a_i} (\hat{Q}_\phi(s_i, a_i) - y_i)$

* Exploration in ①

$$\rightarrow \epsilon\text{-greedy} : \pi(a|s) = \begin{cases} 1-\epsilon & \text{if } a = \operatorname{argmax}_a \hat{Q}_\phi(s, a) \\ \epsilon / |A| & \text{otherwise} \end{cases}$$

→ Boltzmann exploration : $\pi(a|s) \propto \exp(\hat{Q}_\phi(s, a))$

* Limitation

→ batch size is 1 → hard to optimize DNN

→ Samples are correlated → DNN needs i.i.d samples

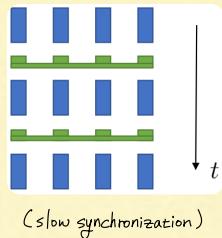
→ No gradient through the target value y_i → y_i also depends on ϕ

→ Overestimation on Q -function → action and value depend on the same network

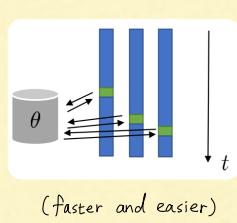
→ Biased target when Q -function is inaccurate → slow learning at the early stage

→ Taking max on continuous action space is hard → worse when Q -function is complex

⇒ synchronized parallel actor-critic



⇒ asynchronous parallel actor-critic



3. Correcting bias

① state-dependent baseline

$$\Rightarrow \nabla_\theta J(\theta) \approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) \left(\sum_{t'=t}^T r^{i,t'} - \hat{V}_\phi^{i,t}(s_{i,t}) \right) \rightarrow \text{unbiased}$$

② action-dependent baseline

$$\begin{aligned} \Rightarrow \nabla_\theta J(\theta) &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) \left(\sum_{t'=t}^T r^{i,t'} - \hat{Q}_\phi^{i,t}(s_{i,t}, a_{i,t}) \right) \rightarrow \text{biased} \\ &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) \left(\sum_{t'=t}^T r^{i,t'} - \hat{Q}_\phi^{i,t}(s_{i,t}, a_{i,t}) \right) \\ &+ \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \mathbb{E}_a [\hat{Q}_\phi^{i,t}(s_{i,t}, a)] \rightarrow \text{correction term : } Q\text{-Prop} \end{aligned}$$

③ eligibility trace

$$\begin{aligned} \Rightarrow \nabla_\theta J(\theta) &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) \left(\sum_{t'=t}^{\infty} \gamma^{t-t'} r^{i,t'} - \hat{V}_\phi^{i,t}(s_{i,t}) + \gamma^n \hat{V}_\phi^{i,t}(s_{i,\infty}) \right) \\ &\quad (\text{as } n \uparrow, \text{ bias } \uparrow, \text{ variance } \uparrow) \end{aligned}$$

④ generalized advantage estimation

$$\begin{aligned} \Rightarrow \nabla_\theta J(\theta) &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) \left(\sum_{n=1}^{\infty} w_n \sum_{t'=t}^{\infty} \gamma^{t-t'} r^{i,t'} - \hat{V}_\phi^{i,t}(s_{i,t}) + \gamma^n \hat{V}_\phi^{i,t}(s_{i,\infty}) \right) \\ &\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T \nabla_\theta \log \pi_\theta(a_{i,t}|s_{i,t}) \left(\sum_{n=1}^{\infty} (\gamma)^{t-t'} (r(s_{i,t}, a_{i,t}) - \hat{V}_\phi^{i,t}(s_{i,t}) + \gamma^n \hat{V}_\phi^{i,t}(s_{i,\infty})) \right) \end{aligned}$$

3. Learning theory

① operator B

$$BV = \max_a r_a + \gamma T_a V \quad (\text{or } BV = r + \gamma T \max_a Q)$$

\rightarrow value iteration : $V \leftarrow BV$

$\rightarrow V^*$ is a fixed point of B : $V^* = BV^*$

\rightarrow for any V and \bar{V} , we have $\|BV - B\bar{V}\|_\infty \leq \gamma \|V - \bar{V}\|_\infty \rightarrow$ contraction

\rightarrow applying B to any V would leads to V^* in terms of ∞ -norm

$$(\because \|BV - BV^*\|_\infty = \|BV - V^*\|_\infty \leq \gamma \|V - V^*\|_\infty)$$

② operator π

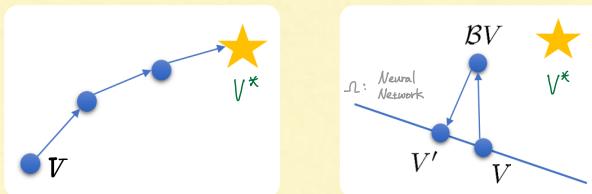
$$\pi V = \arg \min_{Q \in \mathcal{Q}} \frac{1}{2} \sum_s \|V(s) - Q(s)\|_2^2 \quad (\text{or } \pi Q = \arg \min_{V \in \mathcal{V}} \frac{1}{2} \sum_s \|Q(s, a) - V(s, a)\|_2^2)$$

\rightarrow fitted value (or Q) iteration : $V \leftarrow \pi BV$ (or $Q \leftarrow \pi BQ$)

\rightarrow for any V and \bar{V} , we have $\|\pi V - \pi \bar{V}\|_2^2 \leq \|V - \bar{V}\|_2^2 \rightarrow$ contraction

$\rightarrow \pi B$ is not contraction of any kind

(fitted bootstrapped policy evaluations are not guaranteed for convergence)



DRL Lec 7

• Deep RL with Q -functions

1. Improving Q -learning

\Rightarrow increase batch size by parallel learning

\rightarrow synchronized parallel Q -learning, asynchronous parallel Q -learning

\Rightarrow decorrelate samples

\rightarrow replay buffers (B) : dataset of transitions from some previous policies

* Q -learning with replay buffers

- ① collect dataset using some policy and add it to B
- ② sample a batch $\{(s_i, a_i, s'_i, r_i)\}_{i=1}^N$ from B
- xk ③ set $y_i \leftarrow r(s_i, a_i) + \gamma \max_a \hat{Q}_\phi(s'_i, a)$
- ④ update $\phi \leftarrow \phi - \alpha \sum_{i=1}^N \frac{d\hat{Q}_\phi}{d\phi} \Big|_{s=s_i, a=a_i} (\hat{Q}_\phi(s_i, a_i) - y_i)$

\Rightarrow flow gradients to the target value

\rightarrow loosely update the target parameters (simply copying or polyak averaging)

* Q -learning with replay buffers and target network (\approx DQN if $N=k=1$)

- ① save target network parameters $\phi' \leftarrow \phi$
- ② collect dataset using some policy and add it to B
- ③ sample a batch $\{(s_i, a_i, s'_i, r_i)\}_{i=1}^N$ from B
- xk ④ set $y_i \leftarrow r(s_i, a_i) + \gamma \max_a \hat{Q}_\phi(s'_i, a)$
- ⑤ update $\phi \leftarrow \phi - \alpha \sum_{i=1}^N \frac{d\hat{Q}_\phi}{d\phi} \Big|_{s=s_i, a=a_i} (\hat{Q}_\phi(s_i, a_i) - y_i)$

\Rightarrow decorrelate the noises when computing the target value

\rightarrow use \hat{Q}_ϕ for choosing the action and $\hat{Q}_{\phi'}$ for evaluating the value

* Double Q -learning

- ① save target network parameters $\phi' \leftarrow \phi$
- ② collect dataset using some policy and add it to B
- ③ sample a batch $\{(s_i, a_i, s'_i, r_i)\}_{i=1}^N$ from B
- xk ④ set $y_i \leftarrow r(s_i, a_i) + \gamma \hat{Q}_{\phi'}(s'_i, \arg \max_a \hat{Q}_\phi(s_i, a))$
- ⑤ update $\phi \leftarrow \phi - \alpha \sum_{i=1}^N \frac{d\hat{Q}_\phi}{d\phi} \Big|_{s=s_i, a=a_i} (\hat{Q}_\phi(s_i, a_i) - y_i)$

\Rightarrow utilize multi-step returns

\rightarrow need transition history and only actually correct when learning on-policy

* Q -learning with N -step returns

- ① save target network parameters $\phi' \leftarrow \phi$
- ② collect dataset using some policy and add it to B
- ③ sample a batch $\{(s_i, a_i, s'_i, r_i)\}_{i=1}^N$ from B
- xk ④ set $y_i \leftarrow \sum_{t=t}^{t+N-1} \gamma^{t-t} r(s_{t+1}, a_{t+1}) + \gamma^N \max_a \hat{Q}_\phi(s_{t+N}, a)$
- ⑤ update $\phi \leftarrow \phi - \alpha \sum_{i=1}^N \frac{d\hat{Q}_\phi}{d\phi} \Big|_{s=s_i, a=a_i} (\hat{Q}_\phi(s_i, a_i) - y_i)$

\Rightarrow work on continuous action space

\rightarrow optimization

- (sample actions from uniform distribution)
- (cross-entropy method (CEM, CMA-ES))

\rightarrow use function class for Q -function that is easy to optimize

(Normalized Advantage functions (NAF)) : $\hat{Q}_\phi(s, a) = -\frac{1}{2} (\mu(s))^T P_\phi(s) (\mu(s) - \hat{V}_\phi(s))$

\rightarrow learn an approximate maximizer $\mu_\phi(s)$ s.t. $\mu_\phi(s) \approx \arg \max_a \hat{Q}_\phi(s, a)$

* DDPG

- ① save target network parameters $\phi' \leftarrow \phi, \theta' \leftarrow \theta$
- ② collect dataset using some policy and add it to B
- ③ sample a batch $\{(s_i, a_i, s'_i, r_i)\}_{i=1}^N$ from B
- xk ④ set $y_i \leftarrow r(s_i, a_i) + \gamma \hat{Q}_{\phi'}(s'_i, \mu_\phi(s'_i))$
- ④ update
 - $\phi \leftarrow \phi - \alpha \sum_{i=1}^N \frac{d\hat{Q}_\phi}{d\phi} \Big|_{s=s_i, a=a_i} (\hat{Q}_\phi(s_i, a_i) - y_i)$
 - $\theta \leftarrow \theta + \beta \sum_{i=1}^N \frac{d\mu_\phi}{d\theta} \Big|_{s=s_i} \cdot \frac{d\hat{Q}_\phi}{da} \Big|_{s=s_i, a=a_i}$

DRL Lec 8

- Advanced policy gradients

1. Policy gradients \approx Policy iteration

$$\begin{aligned}
\mathbb{E}[\mathcal{J}(\theta') - \mathcal{J}(\theta)] &= \mathbb{E}_{z \sim p_\theta(z)} \left[\sum_{t=0}^{\infty} \delta^t r(s_t, a_t) \right] - \mathbb{E}_{s_t \sim p_\theta(s)} \left[V^{\pi_\theta}(s_t) \right] \\
&= \mathbb{E}_{z \sim p_\theta(z)} \left[\sum_{t=0}^{\infty} \delta^t r(s_t, a_t) \right] - \mathbb{E}_{z \sim p_\theta(z)} \left[V^{\pi_\theta}(s_t) \right] \\
&= \mathbb{E}_{z \sim p_\theta(z)} \left[\sum_{t=0}^{\infty} \delta^t r(s_t, a_t) \right] - \mathbb{E}_{z \sim p_\theta(z)} \left[\sum_{t=0}^{\infty} \gamma^t V^{\pi_\theta}(s_t) - \sum_{t=1}^{\infty} \gamma^t V^{\pi_\theta}(s_t) \right] \\
&= \mathbb{E}_{z \sim p_\theta(z)} \left[\sum_{t=0}^{\infty} \delta^t r(s_t, a_t) \right] + \mathbb{E}_{z \sim p_\theta(z)} \left[\sum_{t=0}^{\infty} \gamma^t (\gamma V^{\pi_\theta}(s_{t+1}) - V^{\pi_\theta}(s_t)) \right] \\
&= \mathbb{E}_{z \sim p_\theta(z)} \left[\sum_{t=0}^{\infty} \delta^t (r(s_t, a_t) + \gamma V^{\pi_\theta}(s_{t+1}) - V^{\pi_\theta}(s_t)) \right] \\
&= \mathbb{E}_{z \sim p_\theta(z)} \left[\sum_{t=0}^{\infty} \delta^t A^\pi(s_t, a_t) \right]
\end{aligned}$$

\therefore Improvement on policy gradients \approx Improvement on policy iteration

2. Bounding the distribution change

$$\begin{aligned}
\mathbb{E}_{z \sim p_\theta(z)} \left[\sum_{t=0}^{\infty} \delta^t A^\pi(s_t, a_t) \right] &= \sum_z \mathbb{E}_{s_t \sim p_\theta(s)} \left[\mathbb{E}_{a_t \sim \pi_\theta(a|s)} \left[\frac{\pi_\theta(a_t|s_t)}{\pi_\theta(a_t|s_t)} \gamma^t A^\pi(s_t, a_t) \right] \right] \\
&\geq \sum_z \mathbb{E}_{s_t \sim p_\theta(s)} \left[\mathbb{E}_{a_t \sim \pi_\theta(a|s)} \left[\frac{\pi_\theta(a_t|s_t)}{\pi_\theta(a_t|s_t)} \gamma^t A^\pi(s_t, a_t) \right] \right] \\
&\quad - \sum_{s_t} \left| \pi_\theta(s_t) - \pi_{\theta'}(s_t) \right| \max_a \mathbb{E}_{a_t \sim \pi_\theta(a|s)} \left[\frac{\pi_\theta(a_t|s_t)}{\pi_\theta(a_t|s_t)} \gamma^t A^\pi(s_t, a_t) \right] \\
&\geq \bar{A}(\theta') - 2\varepsilon + \bar{C}
\end{aligned}$$

where $\bar{A}(\theta) = \sum_z \mathbb{E}_{s_t \sim p_\theta(s)} \left[\mathbb{E}_{a_t \sim \pi_\theta(a|s)} \left[\frac{\pi_\theta(a_t|s_t)}{\pi_\theta(a_t|s_t)} \gamma^t A^\pi(s_t, a_t) \right] \right]$

$$\bar{C} = \max_z \mathbb{E}_{a_t \sim \pi_\theta(a|s)} \left[\frac{\pi_\theta(a_t|s_t)}{\pi_\theta(a_t|s_t)} \gamma^t A^\pi(s_t, a_t) \right]$$

$$\left| \pi_\theta(a_t|s_t) - \pi_{\theta'}(a_t|s_t) \right| \leq \varepsilon$$

$$\therefore \theta' \leftarrow \operatorname{argmax}_{\theta'} \bar{A}(\theta') \text{ s.t. } |\pi_\theta(a_t|s_t) - \pi_{\theta'}(a_t|s_t)| \leq \varepsilon$$

3. Using more convenient and practical bound

$$\Rightarrow \left| \pi_\theta(a_t|s_t) - \pi_{\theta'}(a_t|s_t) \right| \leq \sqrt{\frac{1}{2} \text{KL}(\pi_{\theta'}(a_t|s_t) \| \pi_\theta(a_t|s_t))}$$

$$\therefore \theta' \leftarrow \operatorname{argmax}_{\theta'} \bar{A}(\theta') \text{ s.t. } \text{KL}(\pi_{\theta'}(a_t|s_t) \| \pi_\theta(a_t|s_t)) \leq \varepsilon$$

$$\therefore \theta' \leftarrow \operatorname{argmax}_{\theta'} \nabla_{\theta'} \bar{A}(\theta')^\top (\theta' - \theta) \text{ s.t. } \text{KL}(\pi_{\theta'}(a_t|s_t) \| \pi_\theta(a_t|s_t)) \leq \varepsilon \quad (\text{linearization})$$

$$\text{where } \nabla_{\theta'} \bar{A}(\theta) = \sum_z \mathbb{E}_{s_t \sim p_\theta(s)} \left[\mathbb{E}_{a_t \sim \pi_\theta(a|s)} \left[\frac{\pi_\theta(a_t|s_t)}{\pi_\theta(a_t|s_t)} \gamma^t \nabla_{\theta'} \log \pi_\theta(a_t|s_t) A^\pi(s_t, a_t) \right] \right]$$

$$\text{so that } \nabla_{\theta'} \bar{A}(\theta) = \sum_z \mathbb{E}_{s_t \sim p_\theta(s)} \left[\mathbb{E}_{a_t \sim \pi_\theta(a|s)} \left[\gamma^t \nabla_{\theta'} \log \pi_\theta(a_t|s_t) A^\pi(s_t, a_t) \right] \right] \quad (\text{policy gradients})$$

4. Natural gradient

$$\Rightarrow \text{KL}(\pi_\theta(a_t|s_t) \| \pi_{\theta'}(a_t|s_t)) \approx \frac{1}{2} (\theta' - \theta)^\top F (\theta' - \theta) \quad \text{where } F = \mathbb{E}_{s_t} \left[\nabla_{\theta'} \log \pi_\theta(a_t|s_t) \nabla_{\theta'} \log \pi_\theta(a_t|s_t)^\top \right]$$

$$\therefore \theta^* = \theta + \sqrt{\frac{2\varepsilon}{\nabla_{\theta'} \bar{A}(\theta)^\top F \nabla_{\theta'} \bar{A}(\theta)}} F^{-1} \nabla_{\theta'} \bar{A}(\theta) \quad (\text{Trust region policy optimization})$$

