

Policy gradient

Konpat Preechakul
Chulalongkorn University
October 2019

Recap overview

mcts



Model-free



Model-based

Environment is a black box

We know environment
(transitions, rewards)



Value-based



We learn value. Use value
to improve policy greedily

Policy-based

We directly learn policy
(from some value)



On-policy



Off-policy

Experience comes from
target policy (interactive
experience)

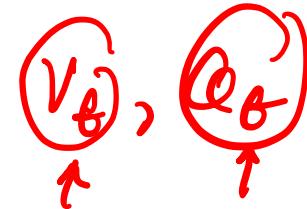
Experience comes from
behavior policy
(observative experience)

IS ↗

Recap Value-based

→ Value function approximation

- Monte Carlo
 - Value error $\hat{\theta}$
- Temporal difference \leftarrow
 - Semi gradient
 - Poor convergence
 - Tricks
 - Full gradient
 - Better convergence
 - Might not good in practice



Today topic is policy-based

- Model the policy directly (not from the value function)

$$\pi_{\theta}(a|s)$$

- Value function is used as a guide

$$\theta \leftarrow Q, V$$

Motivation of policy-based

* Continuous action space

$$\max_a Q(s, a) = ?$$

$$(-1, 1) \rightarrow 0.8$$

- Optimal policy might not be “deterministic”

$$\underset{a}{\operatorname{argmax}} \quad Q \leftarrow \text{Deterministic}$$

- Value-based RL might not give a smooth improvement of the policy

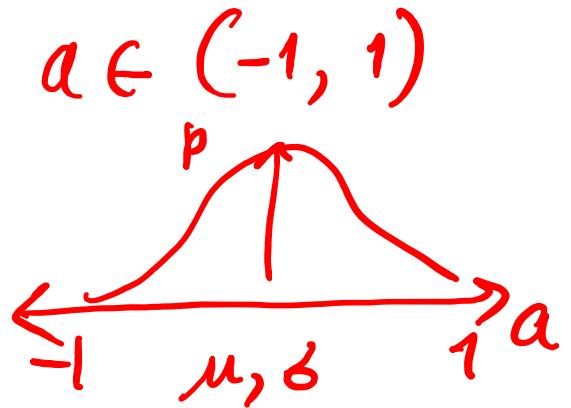
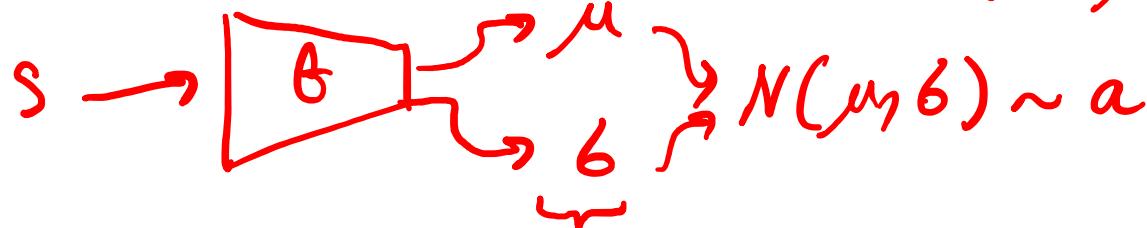
$$Q \rightarrow Q'$$

- Convergence problems in value based RL

Modeling a policy

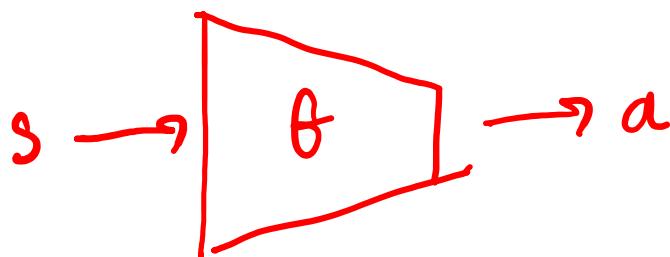
- Stochastic policy

$$\pi_{\theta}(a|s)$$



- Deterministic policy

$$\pi_{\theta}(s) \rightarrow a$$



How to improve the policy?

- We have neural nets, we want to optimize them with SGD
- What is the objective? $J(\theta)$ ←
- What is the gradient? $\text{w.r.t. } \theta$ ←

Policy gradient

Policy gradient (PG)

- Objective function

$$\underline{J}(\theta) = \sum_s P_{s_0}(s) V^\pi(s)$$

start

- Policy gradient

$$\nabla J(\theta) = \sum_s d^\pi(s) \sum_a Q^\pi(s, a) \nabla_\theta \pi_\theta(a|s)$$

d^π

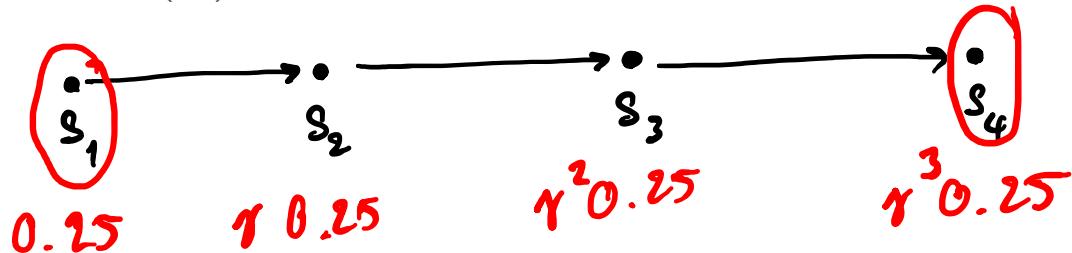
$d^\pi(s)$

Discounted state distribution

- What does it mean?

$$TV \Rightarrow r + \cancel{r\gamma} + \cancel{r^2\gamma^2} + \dots$$

$d^\pi(s) \leftarrow \text{discounted}$

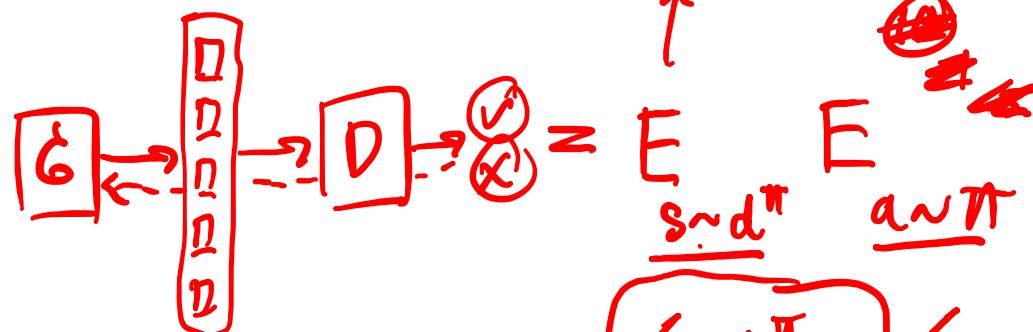


- We care less of far away states
- Why?

Policy gradient (PG)

- Another form (REINFORCE; likelihood ratio)

$$\Rightarrow \nabla J(\theta) = \mathbb{E}_{s,a} [Q^\pi(s, a) \nabla_\theta \log \pi_\theta(a|s)]$$



$$D(1, s_1, 15, 20)$$

\downarrow
 θ a^2

$$= E_{s \sim d^\pi} E_{a \sim \pi} \left[\sum_s d^\pi(s) \sum_a \pi(a|s) Q^\pi \right]$$

$\frac{\nabla \log \pi}{\pi}$

$$\begin{aligned} & D(-, s_1, 15, 20) \\ & \times 11 \quad \downarrow \\ & \times 11 \quad s_1 \end{aligned} = \sum_s d^\pi(s) \sum_a Q^\pi \nabla \pi$$

Making sense of policy gradient

What does it do?

$$\Rightarrow \nabla J(\theta) = \mathbb{E}_{s,a} [Q^\pi(s, a) \nabla_\theta \log \pi_\theta(a|s)]$$

- Gradient for each action is “weighted” based on the goodness

$$s; a_0 \quad Q(\cdot, a_0) = 10$$

$$a_1 \quad Q(\cdot, a_1) = 5$$

$$\nabla_{a_0} \propto 10 \cdot \nabla -$$

$$\nabla_{a_1} \propto 5 \cdot \nabla -$$

Efficiency of likelihood ratio

$$\nabla J(\theta) = \mathbb{E}_{s,a} [Q^\pi(s, a) \nabla_\theta \log \pi_\theta(a|s)]$$

- Gradient is “relative”

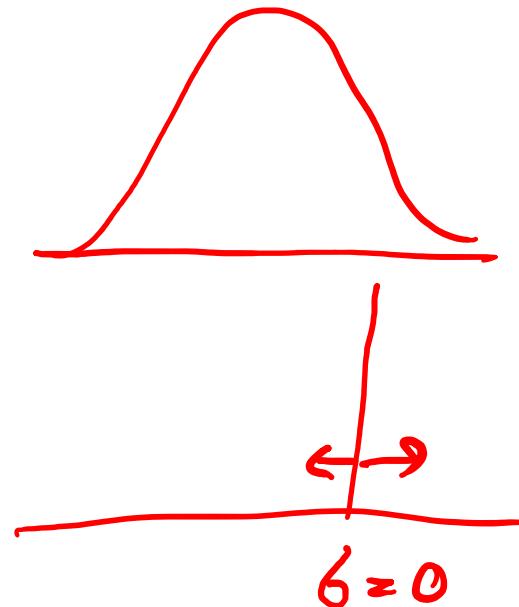
$$\begin{array}{c} \xrightarrow{\theta_0} \xrightarrow{10 \cdot \nabla} \\ \xrightarrow{\theta_1} \xrightarrow{5 \cdot \nabla} \end{array}$$

$$\nabla J(\theta) = \sum_s P_{\theta_0}(s) V(s)$$

- To get a useful gradient, you need “many samples”
- Gradient is indirect

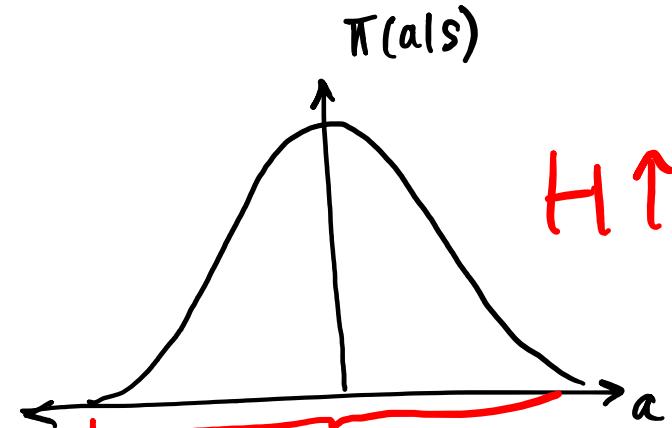
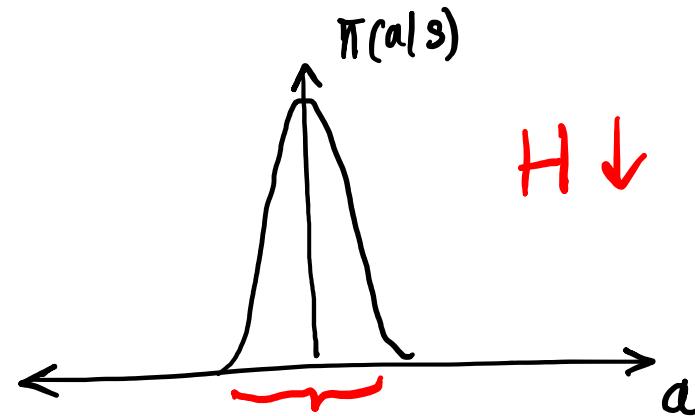
Policy collapse

- Stochastic policy



- Policy often collapses into deterministic policy
- After which, no exploration

Exploration of policy gradient



Entropy

$$H(p) = -\sum_x p(x) \log p(x)$$

$\frac{1}{N}$

- We want to discourage policy collapse

- Revised objective

$$J(\theta) = \mathbb{E}_{s,a} [Q^\pi(s, a) \log \pi_\theta(a|s) + \beta H(\pi_\theta(s))]$$

PG

Coef.

↓

entropy

Policy gradient pseudocode

for until satisfied do

→ ① collect episode (s_0, a_0, r_1, \dots) using π

$$J = \sum_i \gamma^i (g_i \log \pi_\theta(a_i | s_i) + \beta H(\pi_\theta(s_i)))$$
$$\theta \leftarrow \theta + \alpha \nabla J$$

end for

g_i

$$d^\pi \rightarrow - + r_-$$

Note: $\gamma^i \times$

$H(\cdot)$

$$Q \approx g$$

s_1

Policy gradient is on-policy

$$\nabla J(\theta) = \mathbb{E}_{s,a} [Q^\pi(s,a) \nabla \log \pi_\theta(a|s)]$$

$$\nabla J(\theta) = \sum_s d^\pi(s) \sum_a \pi(a|s) [Q^\pi(s,a) \nabla \log \pi_\theta(a|s)]$$

$E_{s,a \sim \pi}$

Proof of policy gradient

$$\begin{aligned}
 ① V^\pi(s) &= \sum_a \pi(a|s) Q^\pi(s, a) \\
 \nabla_\pi V^\pi(s) &= \sum_a Q \nabla \pi + \pi \nabla Q \\
 \nabla_\pi V^\pi(s) &= \sum_{s', r} p(s', r | s, a) [r + \gamma V^\pi(s')] \\
 \nabla_\pi V^\pi(s) &= \sum_a Q \nabla \pi + \gamma \sum_a \pi \sum_{s'} p(s' | s, a) \left[\nabla_\pi V^\pi(s') \right] \\
 \nabla_\pi V^\pi(s) &= \sum_a Q \nabla \pi + \gamma \sum_a \pi \sum_{s'} p(s' | s, a) \left[\nabla_\pi V^\pi(s') \right] \\
 &= \sum_a Q \nabla \pi + \gamma \underbrace{\sum_a \pi \sum_{s'} p(s' | s, a)}_{P(s \rightarrow s', 1)} \sum_{a'} Q \nabla \pi' + \gamma^2 \sum_a \pi \sum_{s'} p(s \rightarrow s', 2) \dots
 \end{aligned}$$

$$\cancel{\nabla_{\pi} V^{\pi}(s)} = \underbrace{\sum_a Q^{\pi} \pi_a}_{\text{Red box}} + \gamma \sum_a \pi_a \sum_{s'} \underbrace{p(s \rightarrow s', a, \pi)}_{\text{Red line}} \cdot \sum_{a'} Q^{\pi'} \pi'_{a'}$$

$$+ \gamma^2 \sum_a \pi_a \sum_{s''} \underbrace{p(s \rightarrow s'', 2, \pi)}_{\text{Red line}} \sum_{a''} Q^{\pi''} \pi''_{a''}$$

$$+ \gamma^3 \sum_a \pi_a \sum_{s'''} \underbrace{p(s \rightarrow s''', 3, \pi)}_{\text{Red line}} \sum_{a'''} Q^{\pi'''} \pi'''_{a'''}$$

+ ...

$$= \sum_{s'} \left(\sum_{k=0}^{\infty} p(s \rightarrow s', k, \pi) \right) \sum_a Q^{\pi} \pi_a$$

$$d^{\pi}(s)$$

$$= \sum_{s'} d^{\pi}(s') \sum_a Q^{\pi}(s, a) \pi(a | s)$$

$$J(B) = \sum_{s_0} P_{s_0}(s) \uparrow$$

Policy gradient extensions

Variance of policy gradient

- What is the variance?

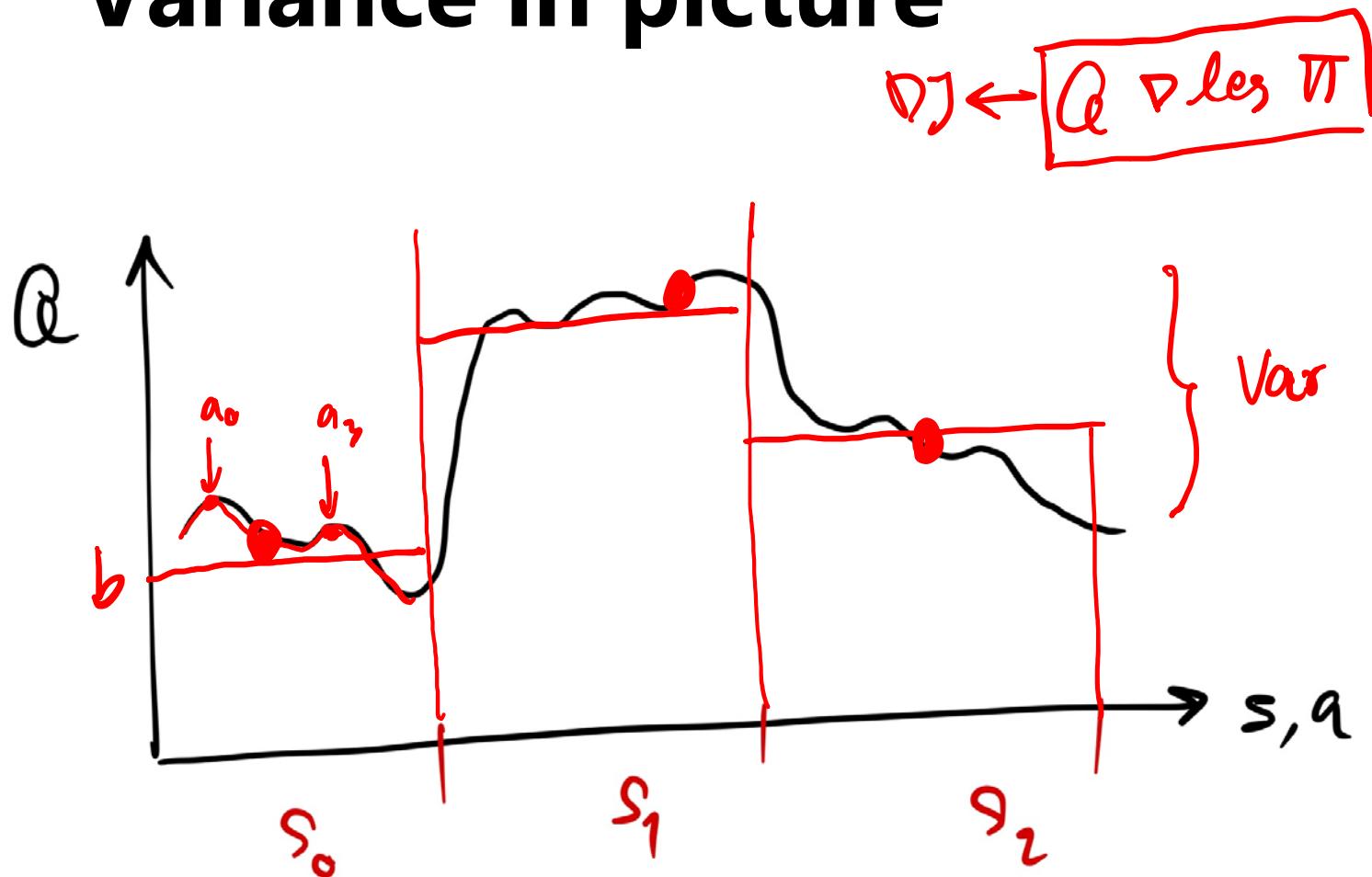
$$\color{red}{*} \text{Var}(Q) = \mathbb{E}_{s,a} [Q(s,a) - \bar{Q}]^2$$

- Policy gradient has high variance

$$\color{red}{*} \nabla J(\theta) = \mathbb{E}_{s,a} [\underbrace{Q^\pi(s,a)}_{\color{red}{Q \sim g_i}} \nabla \log \pi_\theta(a|s)]$$

- High variance slows down learning!

Variance in picture



Variance reduction

$$f(s, a) = \underbrace{Q(s, a)}_{\text{---}} - \underbrace{b(s)}_{\text{---}} - \zeta$$

- If $b(s)$ correlates well with $Q(s, a)$
- $f(s, a)$ could have “lower” variance
- ~~$b(s)$~~ is called a “baseline”



- **There is no change in gradient! Why?**

$$\rightarrow \nabla J(\theta) = \mathbb{E}_{s, a} [(\underbrace{Q^\pi(s, a)}_{\text{---}} - \underbrace{b(s)}_{\text{---}}) \nabla \log \pi_\theta(a|s)]$$

Baseline

$$\nabla J(\theta) = \mathbb{E}_{s,a}[(Q^\pi(s, a) - b(s)) \nabla \log \pi_\theta(a|s)]$$

$\xrightarrow{\quad}$

$$= \underbrace{\mathbb{E}_{s,a} Q^\pi \nabla \log}_{=0} - \underbrace{\mathbb{E}_s b \nabla \log}_{\pi}$$

$$\hookrightarrow \mathbb{E}_s \left[\sum_a \pi(a|s) \cdot b(s) \cdot \nabla \log \pi(a|s) \right] \quad \frac{\nabla \pi}{\pi}$$

$$= \mathbb{E}_s \left[\sum_a \pi \cdot b \cdot \frac{\nabla \pi}{\pi} \right]$$

$$= \mathbb{E}_s \left[\sum_a b \cdot \nabla \pi \right]$$

$$= \mathbb{E}_s \left[b \sum_a \nabla \pi \right]$$

$$= \mathbb{E}_s \left[b \nabla \sum_a \pi(a|s) \right] \Rightarrow \mathbb{E}_s \left[b \nabla \boxed{1} \right] = 0$$

Baseline

- What if $b(s)$ becomes $b(s, a)$?
- Can we use the same argument?

$$b(s) = ?$$

- What is a good baseline for $Q(s, a)$?

$$\frac{1}{n(a)} \sum_a Q(a, s) \rightarrow \sum_a Q(s, a) \pi(a|s) = V(s)$$
$$E \left[\underbrace{(Q(s, a) - V(s))}_{A(s, a)} \cdot \nabla \log \pi \right] \star$$

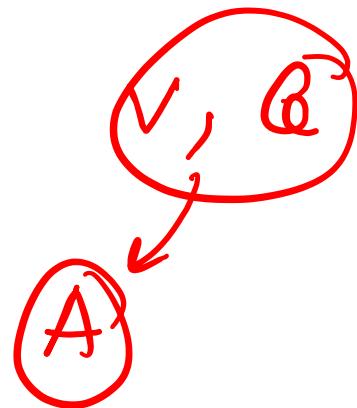
Advantage function

$$A(s, a) = Q(s, a) - V(s)$$

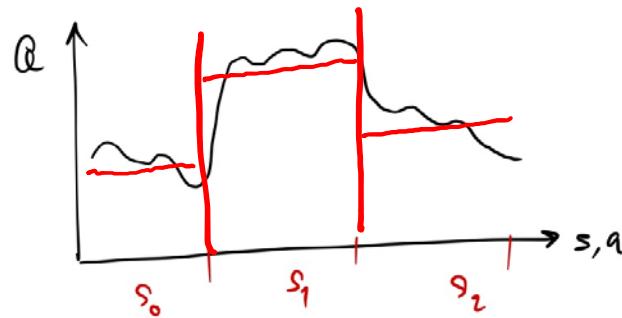
\equiv

$$\underline{A} = r + \underline{A}' X$$

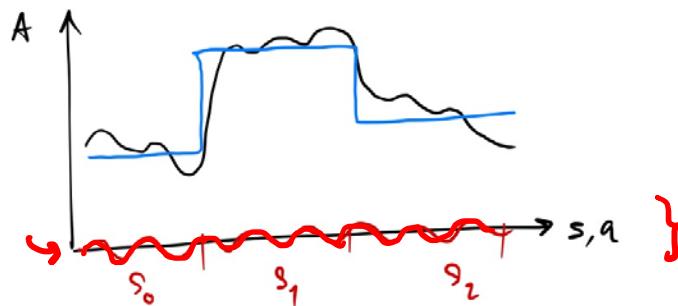
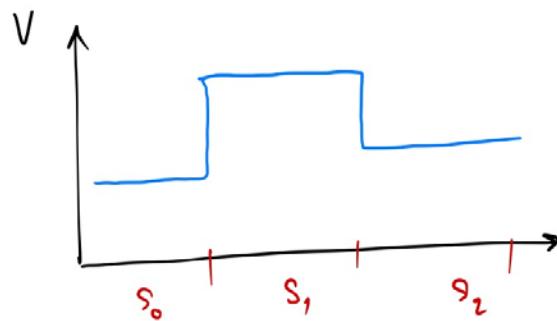
\approx



Baseline in action



$$A = Q - V$$

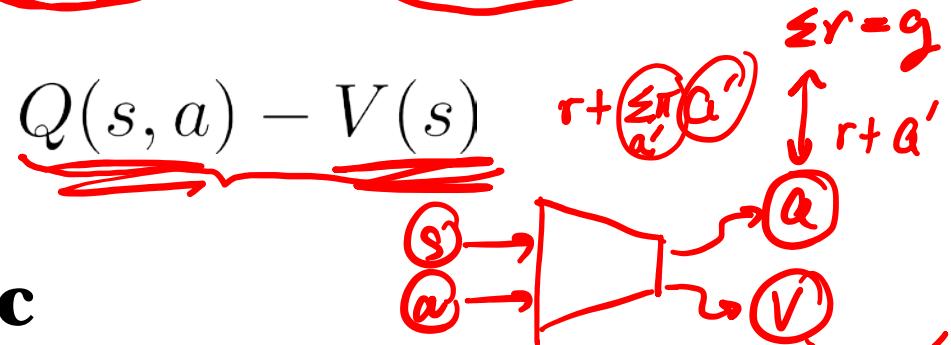


Advantage Actor-critic (A2C)

$$\nabla J(\theta) = \mathbb{E}_{s,a} [A(s,a) \nabla \log \pi_\theta(a|s)]$$

critic *Actor*

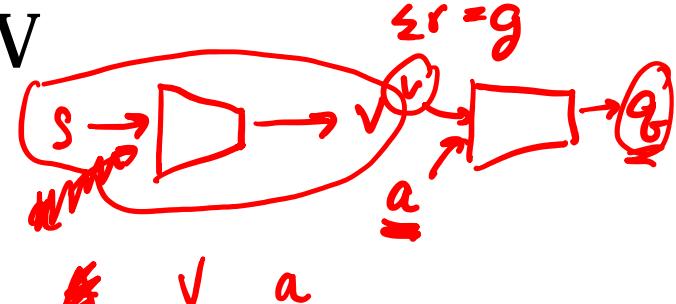
$$A(s,a) = Q(s,a) - V(s)$$

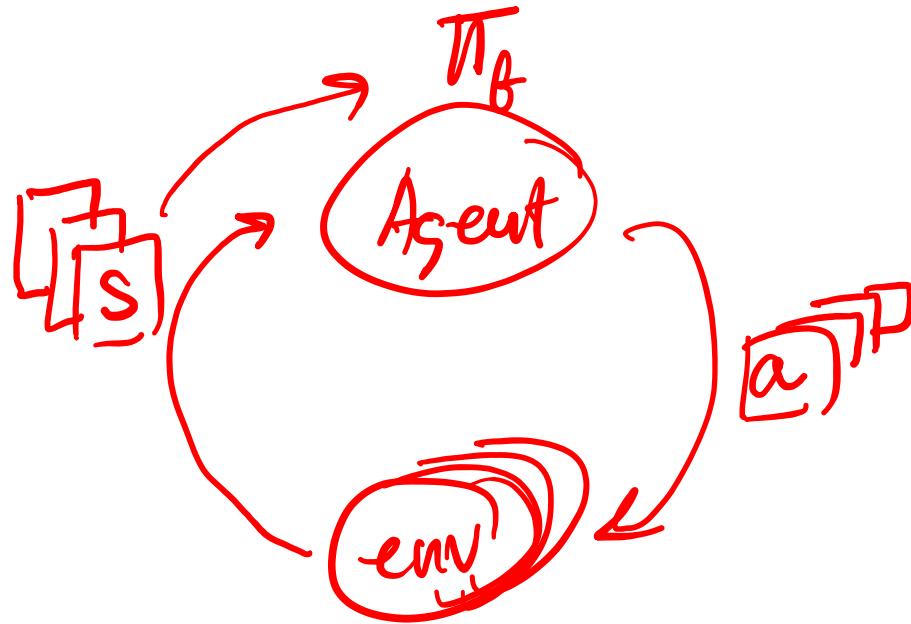


- **Advantage = Critic**
- **Policy = Actor**
- Need to model either A or Q or V
- How to do it efficiently?

$$A = Q - V = r + V' - V$$

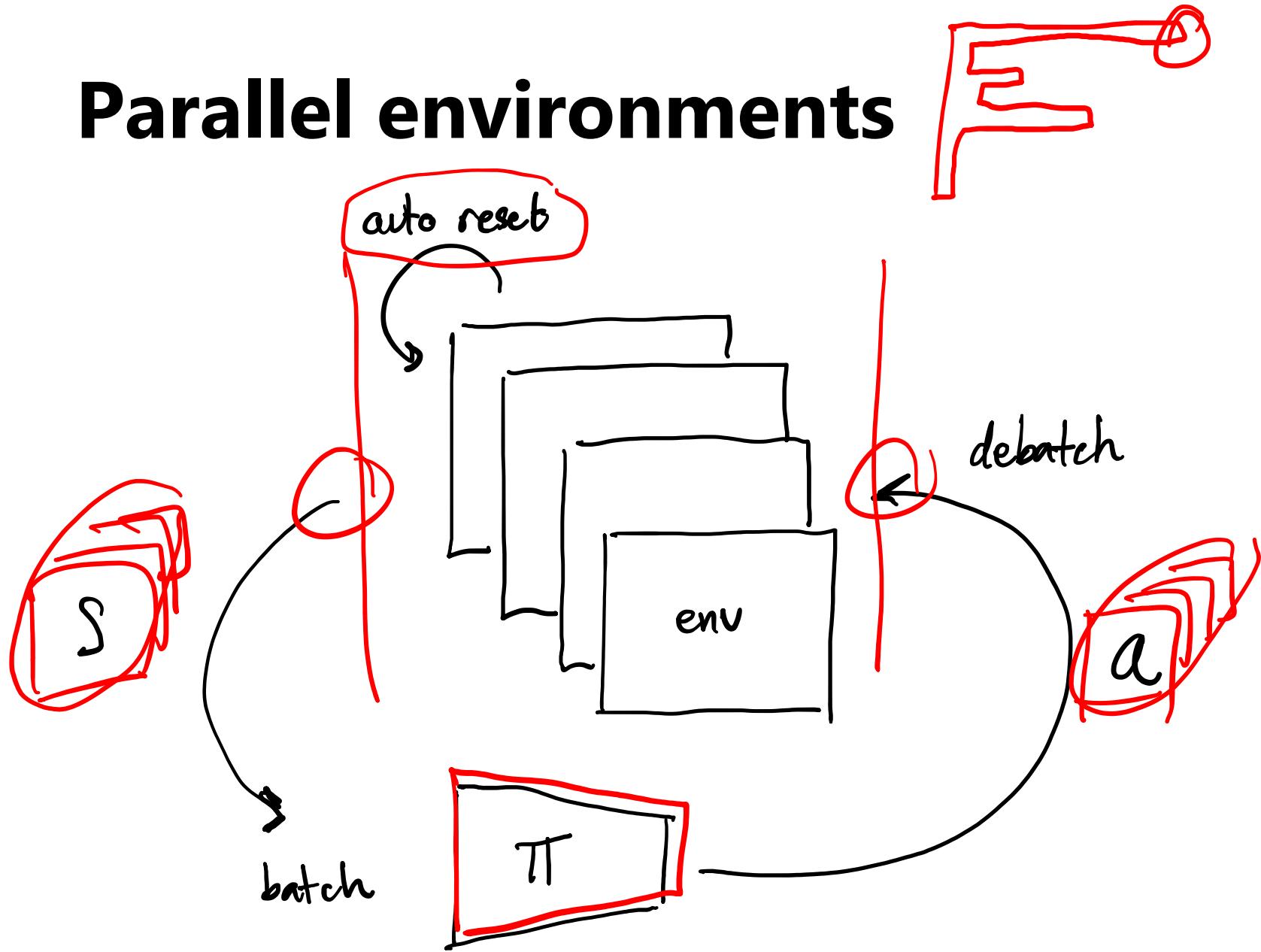
$Q \approx r + V'$



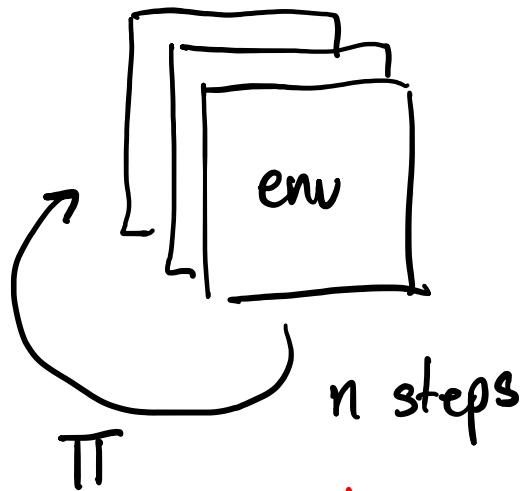


Practical considerations

Parallel environments



N-step exploration



3 step

$\left[(s, a, r, \text{done})_0, (\dots)_1, (\dots)_2 \right]$

$\left[(\dots)_0, (\dots)_1, (\dots)_2 \right]$

$\text{done} \rightarrow \text{false}, \text{true}$

$\text{done} = \text{True} \rightarrow s_0$
autoreset

[f f f f f t]
[f t f f f]

A hand-drawn diagram showing two rows of binary digits. The first row has five 'f's followed by a circled 't'. The second row has a circled 't' at the beginning, followed by four 'f's. Red arrows point from the circled 't's in both rows to the circled 't' in the first row, indicating a transition or reset condition.

A slew of returns

$$Q \rightarrow \log R$$

1. $\sum_{t=0}^{\infty} r_t$: total reward of the trajectory.
2. $\sum_{t'=t}^{\infty} r_{t'}$: reward following action a_t .
3. $\sum_{t'=t}^{\infty} r_{t'} - b(s_t)$: baselined version of previous formula.

4. $Q^{\pi}(s_t, a_t)$: state-action value function.
5. $A^{\pi}(s_t, a_t) = Q^{\pi}(s_t, a_t) - V^{\pi}(s_t)$: advantage function.
6. $r_t + V^{\pi}(s_{t+1}) - V^{\pi}(s_t)$: TD residual.
$$\underbrace{r_t}_{\sim Q} + \underbrace{V^{\pi}(s_{t+1}) - V^{\pi}(s_t)}_{\sim A} = \text{TD residual}$$

N-step TD Residual

$$A_t \approx \underbrace{r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots + \gamma^n V(s_{t+n})}_{n} - V(s_t)$$

Mnih, Volodymyr, Adrià Puigdomènech Badia, Mehdi Mirza, Alex Graves, Timothy P. Lillicrap, Tim Harley, David Silver, and Koray Kavukcuoglu. 2016. "Asynchronous Methods for Deep Reinforcement Learning." *arXiv [cs.LG]*. arXiv. <http://arxiv.org/abs/1602.01783>.

Schulman, John, Philipp Moritz, Sergey Levine, Michael Jordan, and Pieter Abbeel. 2015. "High-Dimensional Continuous Control Using Generalized Advantage Estimation." *arXiv Preprint arXiv:1506.02438*, 1–14.

Generalized Advantage

$\text{TD}(\gamma)$

\checkmark

- Lambda-weighted infinite sum of many-step advantage functions
- Lambda for advantage functions

$$A_t^{(\gamma, \lambda)} = \sum_{i=0}^{\infty} (\gamma \lambda)^i \delta_{t+i}$$

$$\delta_t = \underline{r_{t+1}} + \gamma \underline{V(s_{t+1})} - \underline{V(s_t)}$$

- We only sum to “n” ($n \geq 5$)



A2C pseudocode + V_ϕ + n-step A

for until satisfied do

① collect n step $(s_0, a_0, r_1, \dots, s_{n-1}, a_{n-1}, r_n, s_n)$ using π

② $q_i = \sum_{j=0}^n \gamma^j r_{i+j+1} + \gamma^n V_\phi(s_n)$

③ $\nabla J(\theta) = \sum_{i=0}^n \gamma^i [q_i - V_\phi(s_i)] \nabla \log \pi_\theta(a_i | s_i)$

④ $\nabla L(\phi) = \sum_{i=0}^n [V_\phi(s_i) - q_i] \nabla V_\phi(s_i)$

⑤ $\theta \leftarrow \theta + \alpha_1 \nabla J$

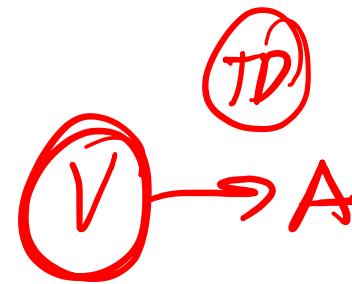
⑥ $\phi \leftarrow \phi + \alpha_2 \nabla L$

end for

\checkmark'
↓

$\alpha_2 > \alpha_1$

Stability of critics



- **Semi-gradient could diverge ✓**
- One-step return might be unstable
- You might need to use target network in such cases
- Usually $N > 5$, you don't need (but using might give you even more stable training)

Sample efficiency

- Not able to reuse data
 - Efficiency is not great
- * Off-policy is much needed
- * Off-policy actor and critic

Off-policy actor-critic

Off-policy **critic** Q one-step TD

$$\nabla J(\theta) = \mathbb{E}_{s,a} [Q_\phi(s, a) \nabla \log \pi_\theta(a|s)]$$

- If Q learns from TD (one or many steps)

$$Q \leftarrow \underline{n\text{-step}} + IS$$

- How to learn Q off-policy?

$$[r + \mathbb{E}_{\substack{a \sim \pi \\ s' \sim P}} Q'] - Q$$

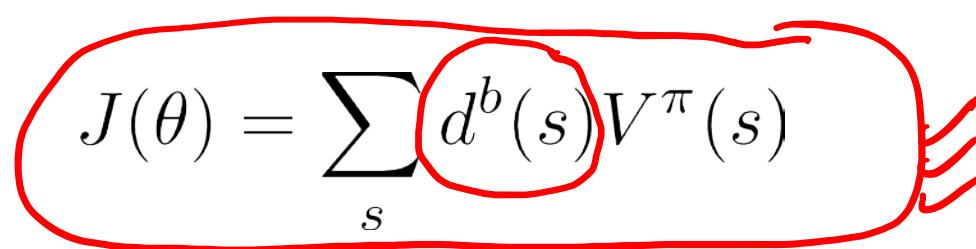
- If one step, we could use expected SARSA

actor Off-policy policy gradient

- Objective function (proposed by the paper)

$$J(\theta) = \sum_s d^b(s) V^\pi(s)$$

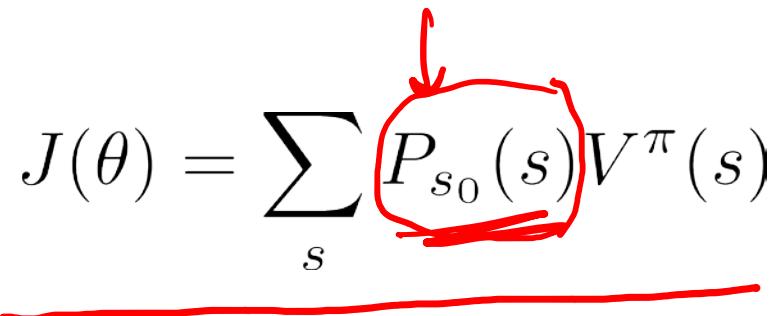
$d^b \leftarrow \text{behavior}$



- It is not the same as on-policy one!

$$J(\theta) = \sum_s P_{s_0}(s) V^\pi(s)$$

d^π



Motivation

$$\rightarrow J(\theta) = \sum_s d^b(s) V^\pi(s)$$

Replay

- We have the distribution already ↗

- Reweighting is not that bad ↑ capacity

Off-policy policy gradient

$$\nabla_{\theta} J(\theta) \rightarrow \sum_s d^b(s) \left[\sum_a Q^{\pi}(s, a) \nabla_{\theta} \pi_{\theta}(a|s) \right] + \sum_s d^{\pi}(s)$$

- Why approximation?

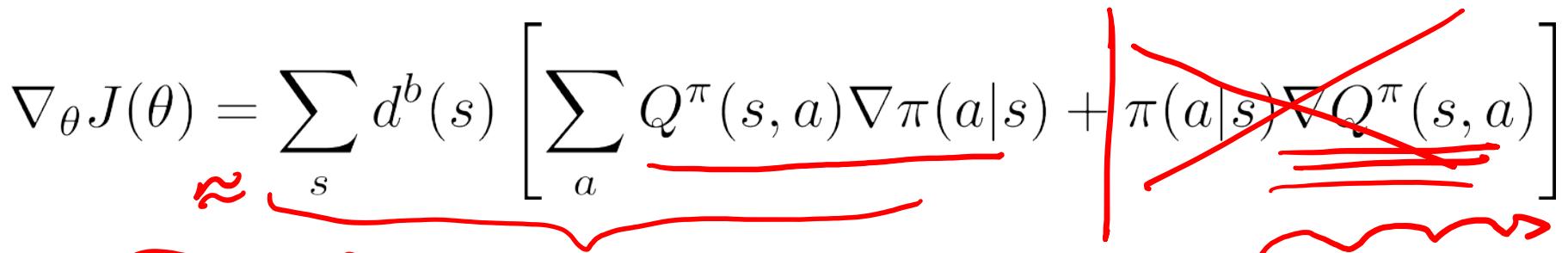
Proof of off-policy gradient

$$\nabla_{\pi} V^{\pi}(s) \Rightarrow \sum_a Q \nabla \pi + \cancel{\pi(\nabla Q)} \xrightarrow{\text{IGNORE}} \sum_a Q \nabla \pi$$
$$\approx \sum_a Q \nabla \pi$$

The diagram illustrates a sequence of states: $s \xrightarrow{\pi} s' \xrightarrow{\pi} s''$. Above the sequence, a red bracket indicates the sum of discounted rewards: $\sum \dots \sum p(r + \delta v')$. Below the sequence, another red bracket indicates the sum of discounted values: $\sum \dots$.

The diagram shows two circles labeled $d^{\pi}(s)$ and $d^b(s)$, with arrows pointing from them to the state s .

Approximation vs Oracle

$$\nabla_{\theta} J(\theta) = \sum_s d^b(s) \left[\sum_a Q^{\pi}(s, a) \nabla \pi(a|s) + \cancel{\pi(a|s) \nabla Q^{\pi}(s, a)} \right]$$


- **Oracle:**

- Start from S ✓
- Update for S ✓
- Unroll, continue on-policy ✓

- **Approximation:**

- Start from S ✓
- Update for S ✓

Off-policy gradient

- It gives a “good enough” gradient
- ✓ It guarantees to improve the policy
- But the fixed point for policy might not be as good (a local minima is not as good)

Deterministic policy gradient (DPG)

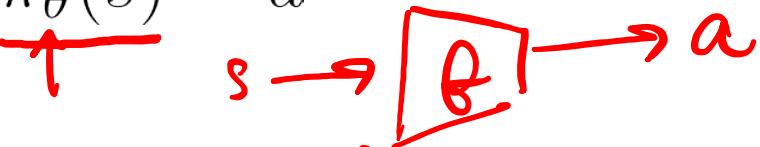
Motivation

- ✓ **Deterministic policy** simplifies much of the policy gradient
- ✓ **Gradient has very low variance!**
- ✓ Off-policy is trivial
- ✓ Potential efficiency gain

Deterministic Policy Gradient

- Deterministic policy

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



- Objective function

$$V = Q(s, \pi(s))$$

$$\nabla J(\theta) = \sum_s P_{s_0}(s) \underline{\underline{V^\pi(s)}} \Rightarrow \sum_s P_{s_0}(s) Q^\pi(s, \underline{\underline{a}}) \Big|_{a=\pi(s)}$$

- Deterministic policy gradient

$$\nabla J(\theta) = \sum_s d^\pi(s) \nabla_a Q^\pi(s, \underline{\underline{a}}) \Big|_{a=\pi(s)} \nabla_\theta \pi_\theta(s)$$

Policy improvement = backprop on critic

$$\nabla_{\theta} J(\theta) \approx \sum_s d^{\pi}(s) \nabla_a Q_{\phi}(s, a)|_{a=\pi(s)} \nabla_{\theta} \pi_{\theta}(s)$$

- We use a neural net as Q (critic)
 - Critic “tells” what is a better action

✓ Very low variance
 - Single sample can make progress

✗ Easy overfit, critic needs to be “ahead”

Needs exploration

- Policy during exploration

$$a = \pi_{\theta}(s) + \underline{\epsilon}$$

$$\epsilon \sim \mathcal{N}(0, \sigma)$$

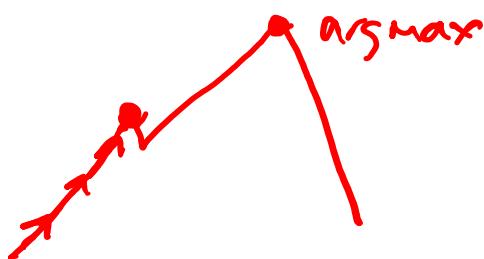
Connection with DQN

- DQN = Explicit max over discrete actions

$$\pi(s) = \underset{a}{\operatorname{argmax}} Q_{\phi}(s, a)$$


- DPG = Climb to the local max action

$$J(\theta) = \sum_s d^\pi(s) \nabla_a Q_{\phi}(s, a)|_{a=\pi(s)} \nabla_\theta \pi_\theta(s)$$



Off-policy DPG

- Objective function

$$J(\theta) = \sum_s d^b(s) Q^\pi(s, a)|_{a=\pi(s)}$$

- Gradient

$$J(\theta) \approx \sum_s d^b(s) \nabla_a Q^\pi(s, a)|_{a=\pi(s)} \nabla_\theta \pi_\theta(s)$$

- Same argument with off-policy gradient

Off-policy critic

$$r + Q'(s', \underline{a})$$

- Deterministic policy could evaluate from off-policy using SARSA-like
- (s, a, r, s')

$$\delta = \left(\underline{r} + \gamma \underline{Q}_\phi(\underline{s}', \underline{\pi}(s')) \right) - Q_\phi(s, a)$$
$$L(\phi) = \frac{1}{2} \delta^2$$

Off-policy

$$\nabla_\phi L(\phi) = -\delta \nabla_\phi Q_\phi(s, a)$$

Proof of DPG

$$J(\theta) = \sum_s d^\pi(s) \nabla_a Q^\pi(s, a)|_{a=\pi(s)} \nabla_\theta \pi_\theta(s)$$

Deterministic policy gradient extensions

Deep Deterministic Policy Gradient (DDPG) ← DQN

- DPG uses one-step bootstrapping which is unstable
- DDPG introduces:
 - ✓ Target network both actor and critic
 - ✓ Replay
- DDPG = DQN for continuous control

Algorithm 1 DDPG algorithm

Randomly initialize critic network $Q(s, a|\theta^Q)$ and actor $\mu(s|\theta^\mu)$ with weights θ^Q and θ^μ .

Initialize target network Q' and μ' 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**

 ① 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

 ② Set $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1}|\theta^{\mu'})|\theta^{Q'})$

 ③ 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:

 ④
$$\nabla_{\theta^\mu} J \approx \frac{1}{N} \sum_i \nabla_a Q(s, a|\theta^Q)|_{s=s_i, a=\mu(s_i)} \nabla_{\theta^\mu} \mu(s|\theta^\mu)|_{s_i}$$

 Update the target networks:

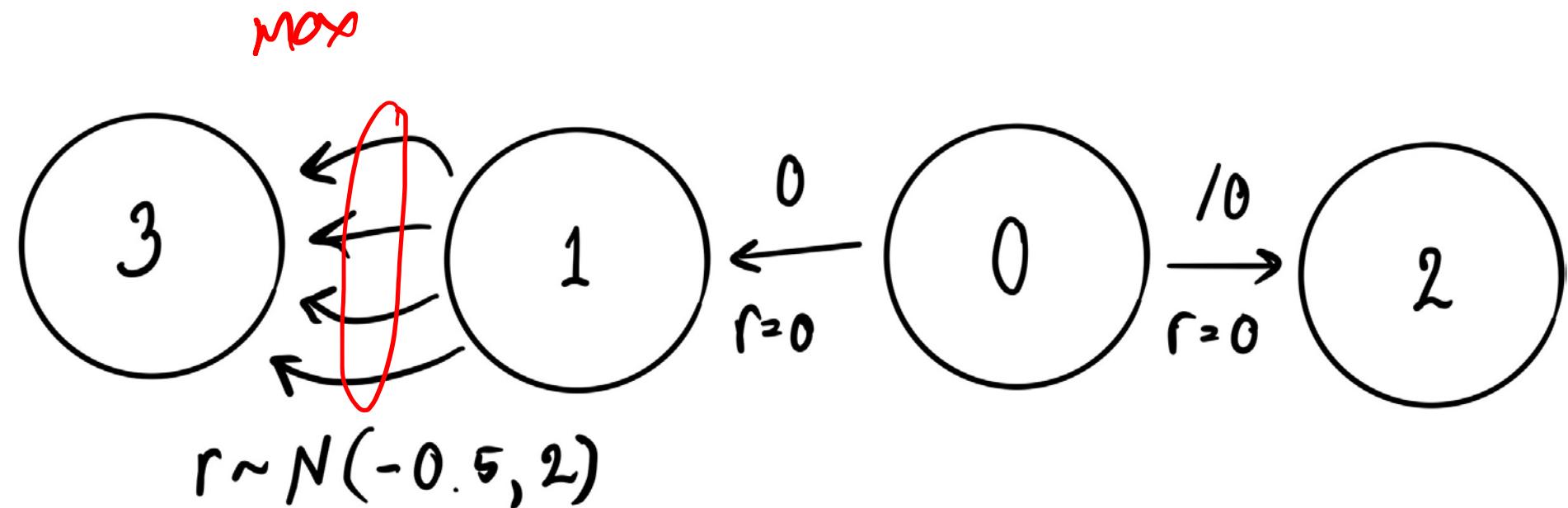
 ⑤
$$\theta^{Q'} \leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'}$$

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

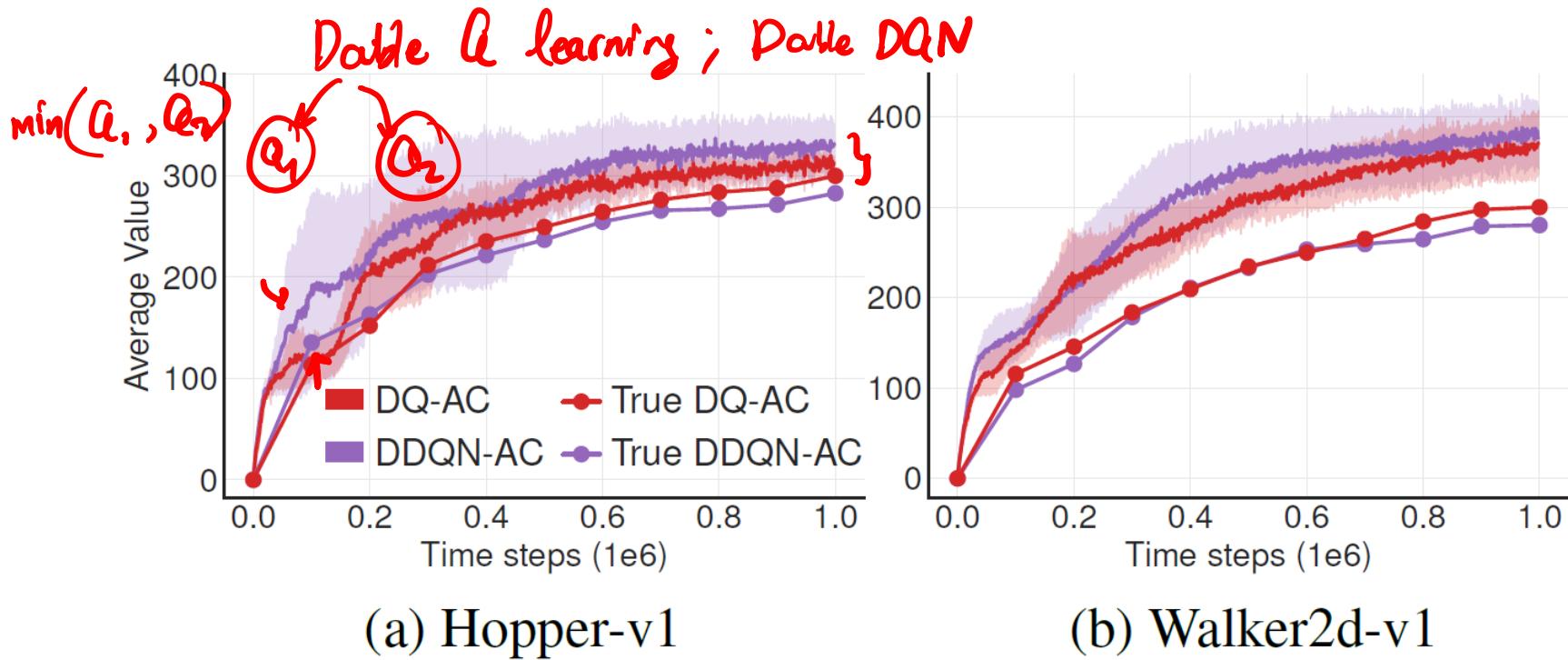
end for

end for

Maximization bias



Maximization bias in DDPG



- Double DQN doesn't work as well
- Because the policy slowly changes, the value functions are not “independent” enough

MISC.

Action dependent baseline

- PG has high variance because it uses “indirect gradient”

$$\nabla J(\theta) = \mathbb{E}_{s,a} [Q^\pi(s, a) \nabla_\theta \log \pi_\theta(a|s)]$$

- DPG has lower variance because it can “backprop”

$$\nabla J(\theta) = \sum_s d^\pi(s) \nabla_a Q_\phi(s, a)|_{a=\pi(s)} \nabla_\theta \pi_\theta(s)$$

Action dependent baseline

- Can we combine the two?
- Q-Prop

$$\begin{aligned}\nabla J(\theta) = & \mathbb{E}_{s,a} \left[(Q^\pi(s, a) - \bar{Q}(s, a)) \nabla_\theta \log \pi_\theta(a|s) \right] \\ & + \mathbb{E}_s \left[\nabla_a Q(s, a)|_{a=\pi(s)} \nabla_\theta \pi_\theta(s) \right]\end{aligned}$$

- Taylor expansion (first order)

$$\bar{Q}(s, a) = Q(s, \pi(s)) + \nabla_a Q(s, \pi(s))(a - \pi(s))$$

Policy gradient from minimizing KL

- If we look at Q as “unnormalized” policy
 - A little bit sharper of Q is $\exp(Q)$
 - This is our target policy
- We could use a KL:

$$\pi = \operatorname{argmin}_{\pi \in \Pi} D_{KL} \left(\pi(\cdot|s) \middle\| \frac{\exp(Q(\cdot, s))}{Z} \right)$$

- Minimizing KL is an optimization task

Policy gradient from minimizing KL

- KL policy gradient:

$$\nabla_{\theta} J(\theta) = \mathbb{E}_s \left[\nabla_{\theta} D_{KL} \left(\pi(\cdot|s) \middle\| \frac{\exp(Q(\cdot, s))}{Z} \right) \right]$$

- Z is a constant, ignored
- Policy improve to Q
- Policy eval: Q gets even sharper
- Repeat

Assignments

- Pull chula_rl (beware of conflicts; back up)
- A2C with continuous Cartpole
- DDPG with continuous Cartpole