

Lecture 13: Policy Gradient Actor-Critic Methods

Policy Approximation

$\pi(a|s, \theta)$  We want to learn this directly!

- Policy = a function from state to action
 - How does the agent select actions?
 - In such a way that it can be affected by learning?
 - In such a way as to assure exploration?
- Approximation: there are too many states and/or actions to represent all policies
 - To handle large/continuous action spaces

Recall: Gradient-bandit algorithm

- Store action preferences $H_t(a)$ rather than action-value estimates $Q_t(a)$
- Instead of ε -greedy, pick actions by an exponential soft-max:

$$\Pr\{A_t = a\} \doteq \frac{e^{H_t(a)}}{\sum_{b=1}^k e^{H_t(b)}} \doteq \pi_t(a)$$

- Also store the sample average of rewards as \bar{R}_t

- Then update:

$$H_{t+1}(a) = H_t(a) + \alpha (R_t - \bar{R}_t) (\mathbf{1}_{a=A_t} - \pi_t(a))$$

1 or 0, depending on whether the predicate (subscript) is true

$$\frac{\partial \pi_t(A_t)}{\partial H_t(a)} / \pi_t(A_t)$$

How can we learn $\pi(a|s, \theta)$?

- Directly from Experience?
 - REINFORCE
- From V and Q?
 - Actor Critic Algorithms
 - Deterministic Policy Gradient (DPG)

How do we parameterize π ?

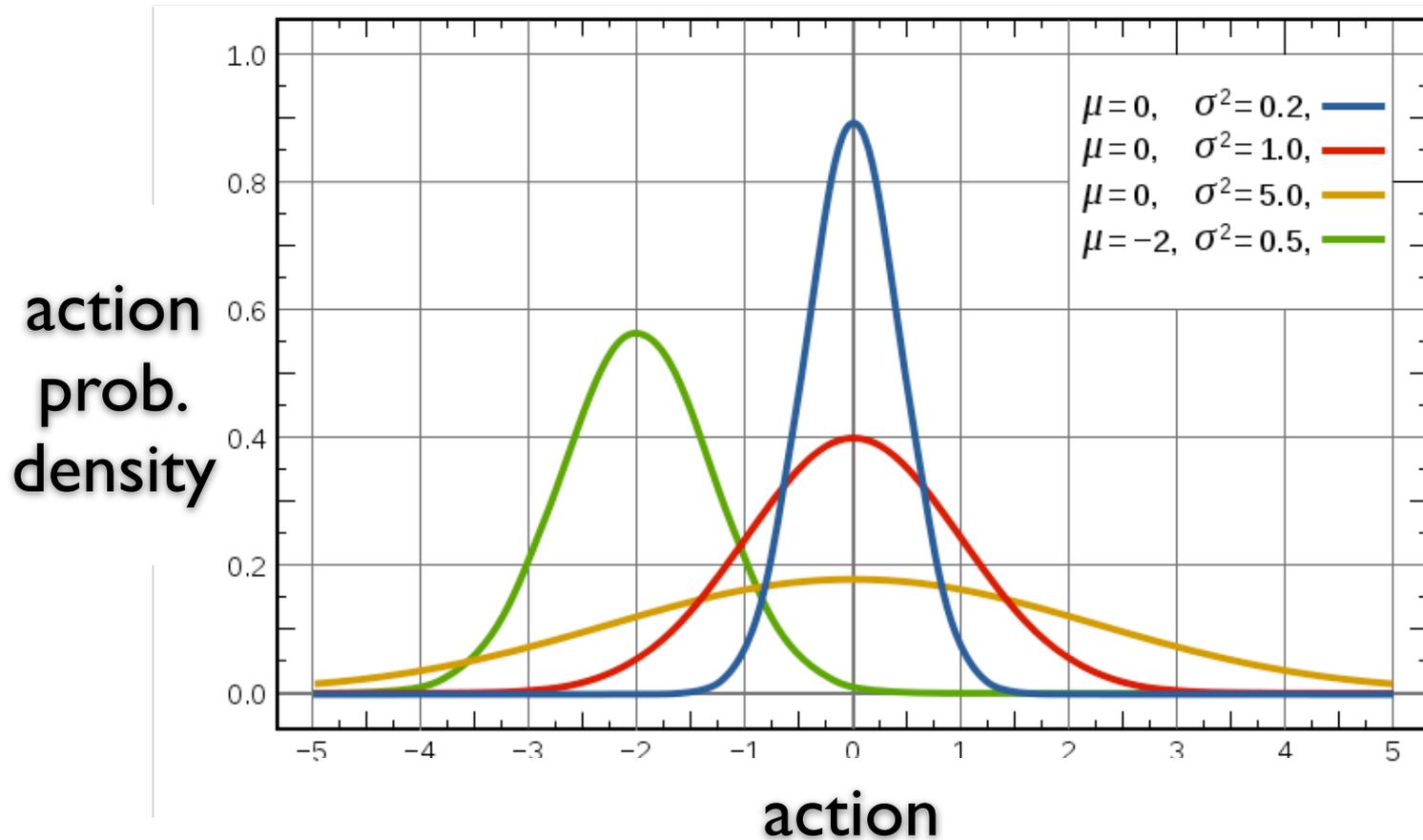
- For discrete actions?
- For continuous actions?

Typical example - Deep Softmax Policies for discrete actions:

$$\pi(A_i | S) = \frac{\exp(\phi(A_i, S))}{\sum_j \exp(\phi(A_j, S))}$$

where ϕ is a neural network,
or any other function approximation parametrize by some
weights.

Typical example - **Gaussian Policies** for **continuous actions**:



Typical example - **Gaussian Policies** for **continuous actions**:

$$\mu(S), \sigma(S) = \phi(S)$$

where ϕ is a neural network,
or any other function approximation
parametrize by some weights.

$$\pi(A | S) = \mathcal{N}(\mu(S), \sigma(S))$$

Typical example - **Gaussian Policies** for **continuous actions:**

These are vectors if the action has more than 1 dim,
Example: the torques for 4 different motors.

$$\mu(S), \sigma(S) = \phi(S)$$

where ϕ is a neural network,
or any other function approximation
parametrize by some weights.

$$\pi(A | S) = \mathcal{N} (\mu(S), \sigma(S))$$

Typical example - **Gaussian Policies** for **continuous actions**:

$$\mu(S), \sigma(S) = \phi(S)$$

where ϕ is a neural network, or any other function approximation parametrize by some weights.

$$\pi(A | S) = \mathcal{N}(\mu(S), \sigma(S))$$

Act by sampling from the distribution:

$$A = \mu(S) + \sigma(S)\epsilon, \quad \epsilon \sim \mathcal{N}(0,1)$$

REINFORCE ALGORITHM

Only π

V, Q, M

Gradient-bandit algorithm

- Store action preferences $H_t(a)$ rather than action-value estimates $Q_t(a)$
- Instead of ϵ -greedy, pick actions by an exponential soft-max:

$$\Pr\{A_t = a\} \doteq \frac{e^{H_t(a)}}{\sum_{b=1}^k e^{H_t(b)}} \doteq \pi_t(a)$$

- Also store the sample average of rewards as \bar{R}_t

- Then update:

$$H_{t+1}(a) = H_t(a) + \alpha (R_t - \bar{R}_t) (\mathbf{1}_{a=A_t} - \pi_t(a))$$

1 or 0, depending on whether the predicate (subscript) is true

$$\frac{\partial \pi_t(A_t)}{\partial H_t(a)} / \pi_t(A_t)$$

Policy Gradient

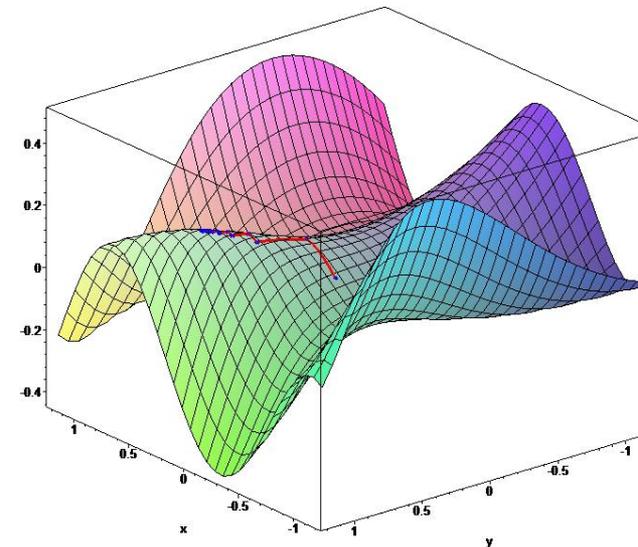
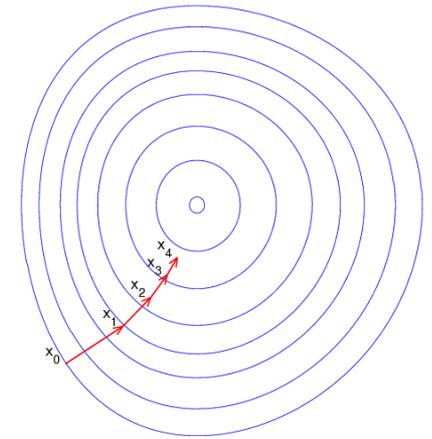
- ▶ Idea: ascent the gradient of the objective $J(\theta)$

$$\Delta\theta = \alpha \nabla_{\theta} J(\theta)$$

- ▶ Where $\nabla_{\theta} J(\theta)$ is the **policy gradient**

$$\nabla_{\theta} J(\theta) = \begin{pmatrix} \frac{\partial J(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial J(\theta)}{\partial \theta_n} \end{pmatrix}$$

- ▶ and α is a step-size parameter
- ▶ Stochastic policies help ensure $J(\theta)$ is smooth (typically/mostly)



Contextual Bandits Policy Gradient

- ▶ Consider a one-step case (a contextual bandit) such that $J(\theta) = \mathbb{E}_{\pi_\theta}[R(S, A)]$.
(Expectation is over d (states) and π (actions))
(For now, d does **not** depend on π)
- ▶ We cannot sample R_{t+1} and then take a gradient:
 R_{t+1} is just a number and does not depend on θ !
- ▶ Instead, we use the identity:

$$\nabla_{\theta} \mathbb{E}_{\pi_{\theta}}[R(S, A)] = \mathbb{E}_{\pi_{\theta}}[R(S, A) \nabla_{\theta} \log \pi(A|S)].$$

(Proof on next slide)

- ▶ The right-hand side gives an expected gradient that can be sampled
- ▶ Also known as REINFORCE (Williams, 1992)

The score function trick

Let $r_{sa} = \mathbb{E}[R(S, A) \mid S = s, A = s]$

$$\begin{aligned}\nabla_{\theta} \mathbb{E}_{\pi_{\theta}}[R(S, A)] &= \nabla_{\theta} \sum_s d(s) \sum_a \pi_{\theta}(a|s) r_{sa} \\ &= \sum_s d(s) \sum_a r_{sa} \nabla_{\theta} \pi_{\theta}(a|s) \\ &= \sum_s d(s) \sum_a r_{sa} \pi_{\theta}(a|s) \frac{\nabla_{\theta} \pi_{\theta}(a|s)}{\pi_{\theta}(a|s)} \\ &= \sum_s d(s) \sum_a \pi_{\theta}(a|s) r_{sa} \nabla_{\theta} \log \pi_{\theta}(a|s) \\ &= \mathbb{E}_{d, \pi_{\theta}}[R(S, A) \nabla_{\theta} \log \pi_{\theta}(A|S)]\end{aligned}$$

Policy Gradient Theorem

- ▶ The policy gradient approach also applies to (multi-step) MDPs
- ▶ Replaces reward R with long-term return G_t or value $q_\pi(s, a)$
- ▶ There are actually two policy gradient theorems (Sutton et al., 2000):
 - average return per episode** & **average reward per step**

Policy gradient theorem (episodic)

Theorem

For any differentiable policy $\pi_{\theta}(s, a)$, let d_0 be the starting distribution over states in which we begin an episode. Then, the policy gradient of $J(\theta) = \mathbb{E}[G_0 \mid S_0 \sim d_0]$ is

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} \left[\sum_{t=0}^T \gamma^t q_{\pi_{\theta}}(S_t, A_t) \nabla_{\theta} \log \pi_{\theta}(A_t \mid S_t) \mid S_0 \sim d_0 \right]$$

where

$$\begin{aligned} q_{\pi}(s, a) &= \mathbb{E}_{\pi}[G_t \mid S_t = s, A_t = a] \\ &= \mathbb{E}_{\pi}[R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}) \mid S_t = s, A_t = a] \end{aligned}$$

Notice this is the return and not the reward,
G not r!

Policy gradient theorem (episodic)

Theorem

For any differentiable policy $\pi_{\theta}(s, a)$, let d_0 be the starting distribution over states in which we begin an episode. Then, the policy gradient of $J(\theta) = \mathbb{E}[G_0 \mid S_0 \sim d_0]$ is

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\pi_{\theta}} \left[\sum_{t=0}^T \gamma^t q_{\pi_{\theta}}(S_t, A_t) \nabla_{\theta} \log \pi_{\theta}(A_t \mid S_t) \mid S_0 \sim d_0 \right]$$

where

$$\begin{aligned} q_{\pi}(s, a) &= \mathbb{E}_{\pi}[G_t \mid S_t = s, A_t = a] \\ &= \mathbb{E}_{\pi}[R_{t+1} + \gamma q_{\pi}(S_{t+1}, A_{t+1}) \mid S_t = s, A_t = a] \end{aligned}$$

Episodic policy gradients algorithm

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^T \gamma^t q_{\pi}(S_t, A_t) \nabla_{\theta} \log \pi(A_t | S_t) \right]$$

- ▶ We can sample this, given a whole episode
- ▶ Typically, people pull out the sum, and split up this into separate gradients, e.g.,

$$\Delta \theta_t = \gamma^t G_t \nabla_{\theta} \log \pi(A_t | S_t)$$

such that $\mathbb{E}_{\pi} [\sum_t \Delta \theta_t] = \nabla_{\theta} J_{\theta}(\pi)$

- ▶ Typically, people ignore the γ^t term, use $\Delta \theta_t = G_t \nabla_{\theta} \log \pi(A_t | S_t)$
- ▶ This is actually okay-ish — we just partially pretend on each step that we could have started an episode in that state instead. Or if we use $\gamma=1$, this is also ok. (alternatively, view it as a slightly biased gradient)

REINFORCE (Monte-Carlo)

$$\nabla_{\boldsymbol{\theta}} J_{\boldsymbol{\theta}}(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^T (\gamma^t G_t) \nabla_{\boldsymbol{\theta}} \log \pi(A_t | S_t) \right]$$

REINFORCE: Monte-Carlo Policy-Gradient Control (episodic) for π_*

Input: a differentiable policy parameterization $\pi(a|s, \boldsymbol{\theta})$

Algorithm parameter: step size $\alpha > 0$

Initialize policy parameter $\boldsymbol{\theta} \in \mathbb{R}^{d'}$ (e.g., to $\mathbf{0}$)

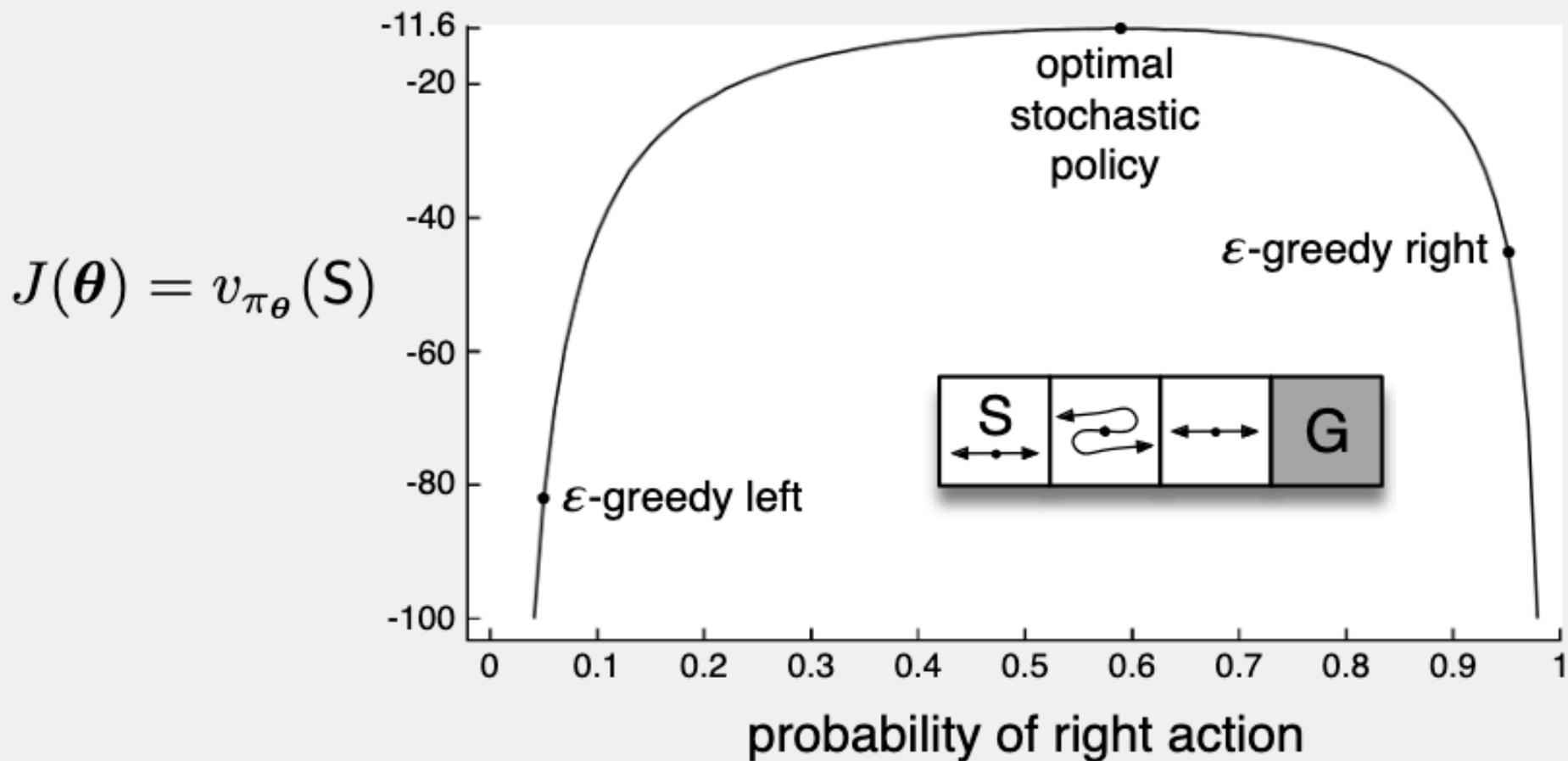
Loop forever (for each episode):

 Generate an episode $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$, following $\pi(\cdot|\cdot, \boldsymbol{\theta})$

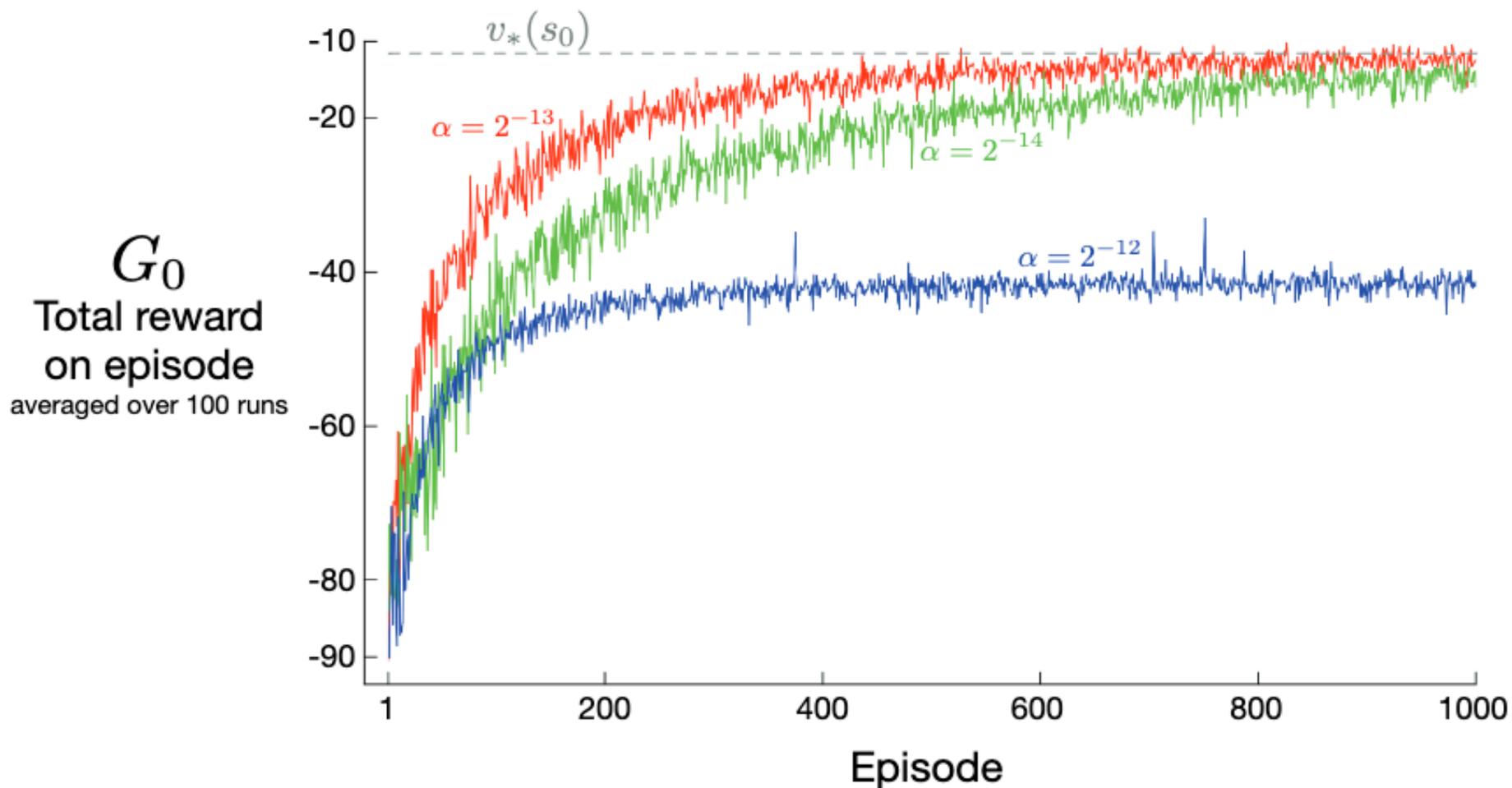
 Loop for each step of the episode $t = 0, 1, \dots, T - 1$:

$$\begin{aligned} G &\leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k && (G_t) \\ \boldsymbol{\theta} &\leftarrow \boldsymbol{\theta} + \alpha \gamma^t G \nabla \ln \pi(A_t | S_t, \boldsymbol{\theta}) \end{aligned}$$

Example: REINFORCE



Example: REINFORCE



Important "Trick" / Identity

$$\sum_a b(s) \nabla_{\boldsymbol{\theta}} \pi(a|s, \boldsymbol{\theta}) = b(s) \nabla_{\boldsymbol{\theta}} \sum_a \pi(a|s, \boldsymbol{\theta}) = b(s) \nabla_{\boldsymbol{\theta}} 1 = 0 \quad \forall s \in \mathcal{S}$$

Important "Trick" / Identity

$$\sum_a b(s) \nabla_{\boldsymbol{\theta}} \pi(a|s, \boldsymbol{\theta}) = b(s) \nabla_{\boldsymbol{\theta}} \sum_a \pi(a|s, \boldsymbol{\theta}) = b(s) \nabla_{\boldsymbol{\theta}} 1 = 0 \quad \forall s \in \mathcal{S}$$

Or written in a different way:

$$\mathbb{E} (b(s) \nabla_{\boldsymbol{\theta}} \log(\pi(a|s, \boldsymbol{\theta}))) = \sum_{s,a} b(s) p(s) \pi(a|s, \boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}} \log(\pi(a|s, \boldsymbol{\theta}))$$

Important "Trick" / Identity

$$\sum_a b(s) \nabla_{\theta} \pi(a|s, \theta) = b(s) \nabla_{\theta} \sum_a \pi(a|s, \theta) = b(s) \nabla_{\theta} 1 = 0 \quad \forall s \in \mathcal{S}$$

Or written in a different way:

$$\begin{aligned} \mathbb{E} (b(s) \nabla_{\theta} \log(\pi(a|s, \theta))) &= \sum_{s,a} b(s) p(s) \pi(a|s, \theta) \nabla_{\theta} \log(\pi(a|s, \theta)) \\ &= \sum_{s,a} b(s) p(s) \pi(a|s, \theta) \frac{\nabla_{\theta} \pi(a|s, \theta)}{\pi(a|s, \theta)} \end{aligned}$$

Important "Trick" / Identity

$$\sum_a b(s) \nabla_{\theta} \pi(a|s, \theta) = b(s) \nabla_{\theta} \sum_a \pi(a|s, \theta) = b(s) \nabla_{\theta} 1 = 0 \quad \forall s \in \mathcal{S}$$

Or written in a different way:

$$\begin{aligned} \mathbb{E} (b(s) \nabla_{\theta} \log(\pi(a|s, \theta))) &= \sum_{s,a} b(s) p(s) \pi(a|s, \theta) \nabla_{\theta} \log(\pi(a|s, \theta)) \\ &= \sum_{s,a} b(s) p(s) \pi(a|s, \theta) \frac{\nabla_{\theta} \pi(a|s, \theta)}{\pi(a|s, \theta)} \\ &= \sum_{s,a} b(s) p(s) \nabla_{\theta} \pi(a|s, \theta) \end{aligned}$$

Important "Trick" / Identity

$$\sum_a b(s) \nabla_{\theta} \pi(a|s, \theta) = b(s) \nabla_{\theta} \sum_a \pi(a|s, \theta) = b(s) \nabla_{\theta} 1 = 0 \quad \forall s \in \mathcal{S}$$

Or written in a different way:

$$\begin{aligned} \mathbb{E} (b(s) \nabla_{\theta} \log(\pi(a|s, \theta))) &= \sum_{s,a} b(s) p(s) \pi(a|s, \theta) \nabla_{\theta} \log(\pi(a|s, \theta)) \\ &= \sum_{s,a} b(s) p(s) \pi(a|s, \theta) \frac{\nabla_{\theta} \pi(a|s, \theta)}{\pi(a|s, \theta)} \\ &= \sum_{s,a} b(s) p(s) \nabla_{\theta} \pi(a|s, \theta) \\ &= 0 \end{aligned}$$

Improvements to REINFORCE

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^T (\gamma^t G_t) \nabla_{\theta} \log \pi(A_t | S_t) \right]$$

- Can we use our "trick" $\mathbb{E} (b(s) \nabla_{\theta} \log(\pi(a | s, \theta))) = 0$ to improve REINFORCE?

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E} \left[\sum_{t=0}^T \gamma^t (G_t - \bar{G}) \nabla_{\theta} \log(\pi) \right]$$

Improvements to REINFORCE

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^T (\gamma^t G_t) \nabla_{\theta} \log \pi(A_t | S_t) \right]$$

- Can we use our "trick" $\mathbb{E} (b(s) \nabla_{\theta} \log(\pi(a | s, \theta))) = 0$ to improve REINFORCE?

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E} \left[\sum_{t=0}^T \gamma^t (G_t - \bar{G}) \nabla_{\theta} \log(\pi) \right]$$

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E} \left[\sum_{t=0}^T \gamma^t (q_{\pi}(S_t, A_t) - v_{\pi}(S_t)) \nabla_{\theta} \log(\pi) \right]$$

Improvements to REINFORCE

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^T (\gamma^t G_t) \nabla_{\theta} \log \pi(A_t | S_t) \right]$$

- Can we use our "trick" $\mathbb{E} (b(s) \nabla_{\theta} \log(\pi(a | s, \theta))) = 0$ to improve REINFORCE?

Reducing Variance:

X, Y are two random variables.

$$\bar{X} = \mathbb{E}(X) = 0$$

$$\bar{Y} = \mathbb{E}(Y) \neq 0$$

Using samples of X, Y , I want to estimate: $\mathbb{E}(YX) \equiv J$

Reducing Variance:

X, Y are two random variables.

$$\bar{X} = \mathbb{E}(X) = 0$$

$$\bar{Y} = \mathbb{E}(Y) \neq 0$$

Using samples of X, Y , I want to estimate: $\mathbb{E}(YX) \equiv J$

$$J \approx \frac{1}{N} \sum_i^N Y_i X_i$$

Reducing Variance:

X, Y are two random variables.

$$\bar{X} = \mathbb{E}(X) = 0$$

$$\bar{Y} = \mathbb{E}(Y) \neq 0$$

Using samples of X, Y , I want to estimate: $\mathbb{E}(YX) \equiv J$

$$J \approx \frac{1}{N} \sum_i^N Y_i X_i$$

Can I do it with less variance??

Reducing Variance:

X, Y are two random variables.

$$\bar{X} = \mathbb{E}(X) = 0$$

$$\bar{Y} = \mathbb{E}(Y) \neq 0$$

Using samples of X, Y , I want to estimate: $\mathbb{E}(YX) \equiv J$

$$\mathbb{E}(YX) = \mathbb{E}[(Y - \bar{Y})X + \bar{Y}X] = \mathbb{E}[(Y - \bar{Y})X] + \bar{Y}\mathbb{E}[X]$$

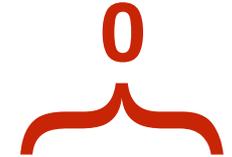
Reducing Variance:

X, Y are two random variables.

$$\bar{X} = \mathbb{E}(X) = 0$$

$$\bar{Y} = \mathbb{E}(Y) \neq 0$$

Using samples of X, Y , I want to estimate: $\mathbb{E}(YX) \equiv J$

$$\mathbb{E}(YX) = \mathbb{E}[(Y - \bar{Y})X + \bar{Y}X] = \mathbb{E}[(Y - \bar{Y})X] + \bar{Y}\mathbb{E}[X]$$


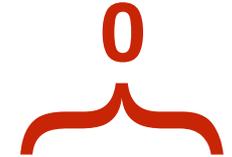
Reducing Variance:

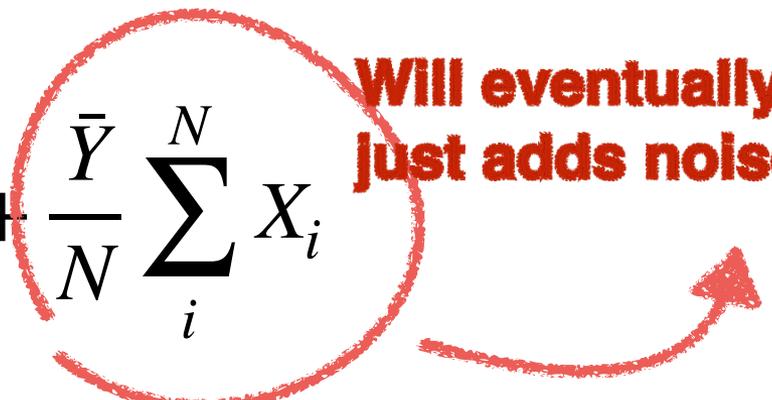
X, Y are two random variables.

$$\bar{X} = \mathbb{E}(X) = 0$$

$$\bar{Y} = \mathbb{E}(Y) \neq 0$$

Using samples of X, Y, I want to estimate: $\mathbb{E}(YX) \equiv J$

$$\mathbb{E}(YX) = \mathbb{E}[(Y - \bar{Y})X + \bar{Y}X] = \mathbb{E}[(Y - \bar{Y})X] + \bar{Y}\mathbb{E}[X]$$


$$J \approx \frac{1}{N} \sum_i (Y_i - \bar{Y})X_i + \frac{\bar{Y}}{N} \sum_i X_i$$


**Will eventually go to 0,
just adds noise!**

Improvements to REINFORCE

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^T (\gamma^t G_t) \nabla_{\theta} \log \pi(A_t | S_t) \right]$$

- Can we use our "trick" $\mathbb{E} (b(s) \nabla_{\theta} \log(\pi(a | s, \theta))) = 0$ to improve REINFORCE?

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E} \left[\sum_{t=0}^T \gamma^t (G_t - \bar{G}) \nabla_{\theta} \log(\pi) \right]$$

Improvements to REINFORCE

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^T (\gamma^t G_t) \nabla_{\theta} \log \pi(A_t | S_t) \right]$$

- Can we use our "trick" $\mathbb{E} (b(s) \nabla_{\theta} \log(\pi(a | s, \theta))) = 0$ to improve REINFORCE?

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E} \left[\sum_{t=0}^T \gamma^t (G_t - \bar{G}) \nabla_{\theta} \log(\pi) \right]$$

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E} \left[\sum_{t=0}^T \gamma^t (q_{\pi}(S_t, A_t) - v_{\pi}(S_t)) \nabla_{\theta} \log(\pi) \right]$$

Improvements to REINFORCE

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^T (\gamma^t G_t) \nabla_{\theta} \log \pi(A_t | S_t) \right]$$

- Can we use our "trick" $\mathbb{E} (b(s) \nabla_{\theta} \log(\pi(a | s, \theta))) = 0$ to improve REINFORCE?

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E} \left[\sum_{t=0}^T \gamma^t (G_t - \bar{G}) \nabla_{\theta} \log(\pi) \right]$$

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E} \left[\sum_{t=0}^T \gamma^t \underbrace{(q_{\pi}(S_t, A_t) - v_{\pi}(S_t))}_{\text{Advantage}} \nabla_{\theta} \log(\pi) \right]$$

Advantage

REINFORCE with baseline:

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

Input: a differentiable policy parameterization $\pi(a|s, \theta)$

Input: a differentiable state-value function parameterization $\hat{v}(s, \mathbf{w})$

Algorithm parameters: step sizes $\alpha^{\theta} > 0$, $\alpha^{\mathbf{w}} > 0$

Initialize policy parameter $\theta \in \mathbb{R}^{d'}$ and state-value weights $\mathbf{w} \in \mathbb{R}^d$ (e.g., to $\mathbf{0}$)

Loop forever (for each episode):

 Generate an episode $S_0, A_0, R_1, \dots, S_{T-1}, A_{T-1}, R_T$, following $\pi(\cdot|\cdot, \theta)$

 Loop for each step of the episode $t = 0, 1, \dots, T - 1$:

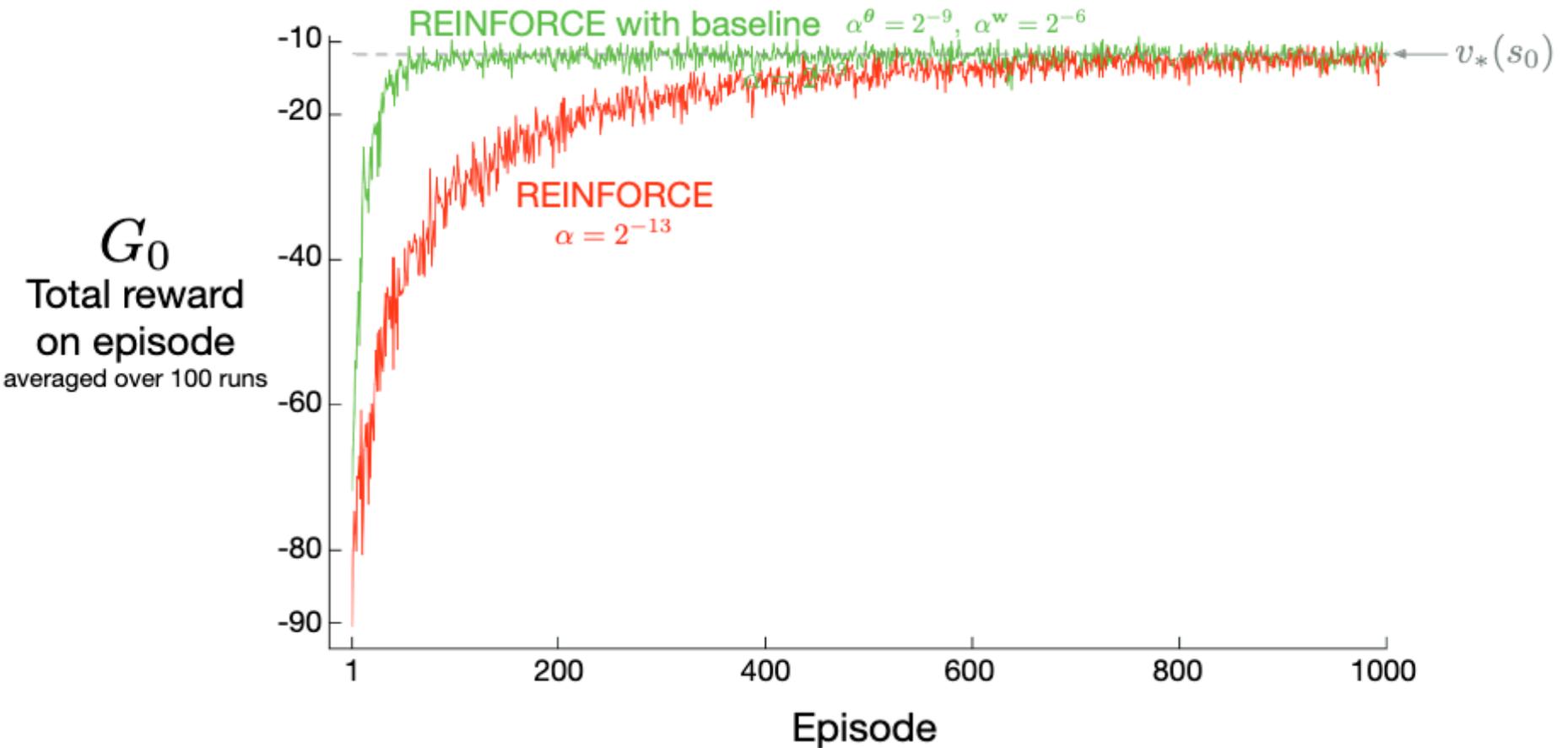
$$G \leftarrow \sum_{k=t+1}^T \gamma^{k-t-1} R_k \tag{G_t}$$

$$\delta \leftarrow G - \hat{v}(S_t, \mathbf{w})$$

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha^{\mathbf{w}} \delta \nabla \hat{v}(S_t, \mathbf{w})$$

$$\theta \leftarrow \theta + \alpha^{\theta} \gamma^t \delta \nabla \ln \pi(A_t | S_t, \theta)$$

REINFORCE with baseline:



Actor-Critic Algorithms

- ACTOR: policy π
- CRITIC: value fct V (or Q)

Actor-Critic Algorithms

- ACTOR: policy π
- CRITIC: value fct V (or Q)

Policy Gradient Theorem:

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^T (\gamma^t G_t) \nabla_{\theta} \log \pi(A_t | S_t) \right]$$

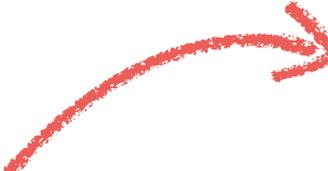
Actor-Critic Algorithms

- ACTOR: policy π
- CRITIC: value fct V (or Q)

Policy Gradient Theorem:

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^T (\gamma^t G_t) \nabla_{\theta} \log \pi(A_t | S_t) \right]$$

REINFORCE Estimates G with Monte-Carlo



Actor-Critic Algorithms

- ACTOR: policy π
- CRITIC: value fct V (or Q)

Policy Gradient Theorem:

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^T (\gamma^t G_t) \nabla_{\theta} \log \pi(A_t | S_t) \right]$$

Actor-Critic: use V and/or Q to estimate G , e.g. TD(0)

Actor-Critic 1-step TD / TD(0) estimate:

Policy Gradient Theorem:

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E} \left[\sum_{t=0}^T \gamma^t \underbrace{(q_{\pi}(S_t, A_t) - v_{\pi}(S_t))}_{\text{Advantage}} \nabla_{\theta} \log(\pi) \right]$$

$$\begin{aligned} \boldsymbol{\theta}_{t+1} &\doteq \boldsymbol{\theta}_t + \alpha \left(G_{t:t+1} - \hat{v}(S_t, \mathbf{w}) \right) \frac{\nabla \pi(A_t | S_t, \boldsymbol{\theta}_t)}{\pi(A_t | S_t, \boldsymbol{\theta}_t)} \\ &= \boldsymbol{\theta}_t + \alpha \left(R_{t+1} + \gamma \hat{v}(S_{t+1}, \mathbf{w}) - \hat{v}(S_t, \mathbf{w}) \right) \frac{\nabla \pi(A_t | S_t, \boldsymbol{\theta}_t)}{\pi(A_t | S_t, \boldsymbol{\theta}_t)} \\ &= \boldsymbol{\theta}_t + \alpha \delta_t \frac{\nabla \pi(A_t | S_t, \boldsymbol{\theta}_t)}{\pi(A_t | S_t, \boldsymbol{\theta}_t)}. \end{aligned}$$

Actor-Critic 1-step TD / TD(0) estimate:

Policy Gradient Theorem:

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E} \left[\sum_{t=0}^T \gamma^t \underbrace{(q_{\pi}(S_t, A_t) - v_{\pi}(S_t))}_{\text{Advantage}} \nabla_{\theta} \log(\pi) \right]$$

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

Input: a differentiable policy parameterization $\pi(a|s, \theta)$

Input: a differentiable state-value function parameterization $\hat{v}(s, \mathbf{w})$

Parameters: step sizes $\alpha^{\theta} > 0$, $\alpha^{\mathbf{w}} > 0$

Initialize policy parameter $\theta \in \mathbb{R}^{d'}$ 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):

$A \sim \pi(\cdot|S, \theta)$

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 $\hat{v}(S', \mathbf{w}) \doteq 0$)

$\mathbf{w} \leftarrow \mathbf{w} + \alpha^{\mathbf{w}} \delta \nabla \hat{v}(S, \mathbf{w})$

$\theta \leftarrow \theta + \alpha^{\theta} I \delta \nabla \ln \pi(A|S, \theta)$

$I \leftarrow \gamma I$

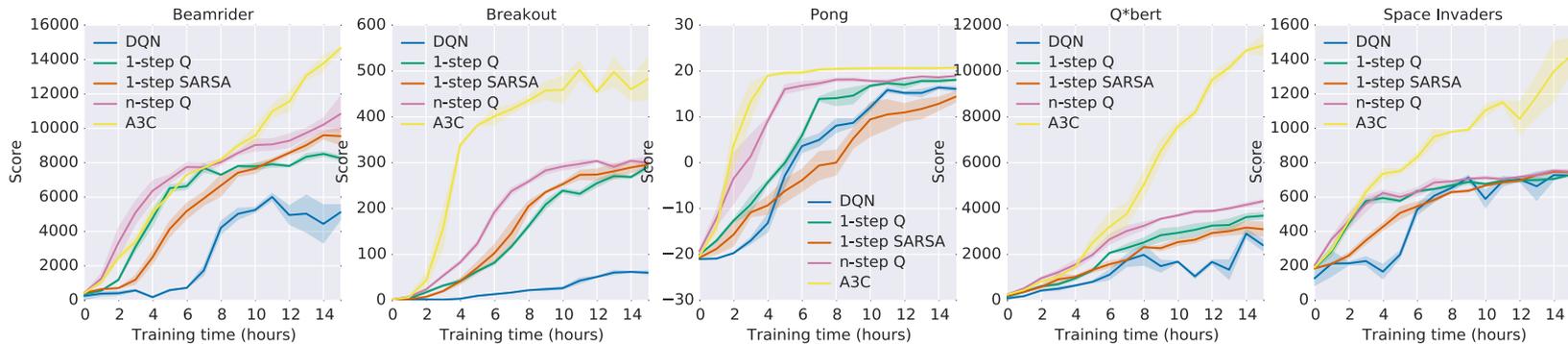
$S \leftarrow S'$

A3C: Asynchronous Advantage Actor Critic:

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E} \left[\sum_{t=0}^T \gamma^t \underbrace{(q_{\pi}(S_t, A_t) - v_{\pi}(S_t))}_{\text{Advantage}} \nabla_{\theta} \log(\pi) \right]$$

A3C: Asynchronous Advantage Actor Critic:

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E} \left[\sum_{t=0}^T \underbrace{\gamma^t (q_{\pi}(S_t, A_t) - v_{\pi}(S_t))}_{\text{Advantage}} \nabla_{\theta} \log(\pi) \right]$$

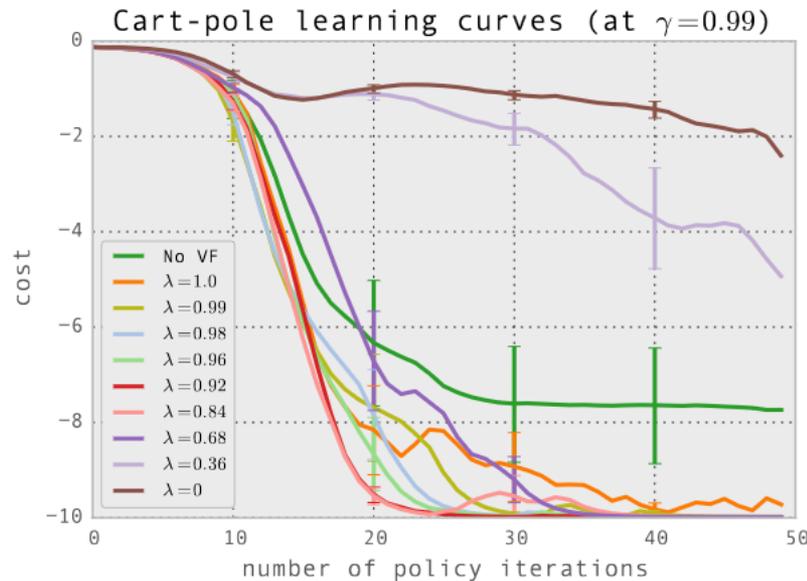


GAE: Generalized Advantage Estimation

- Use Advantage (i.e. $G - V(S)$)
- Use TD(λ) target for G

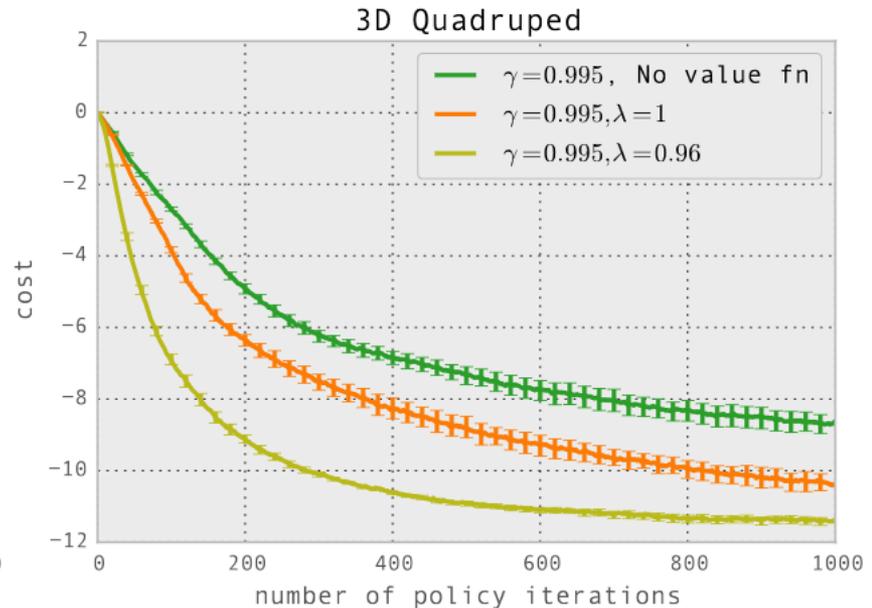
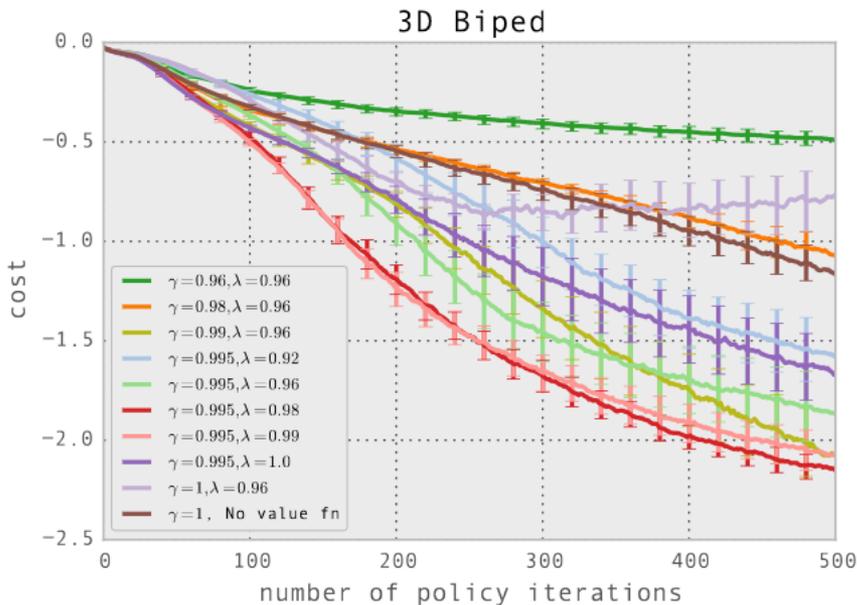
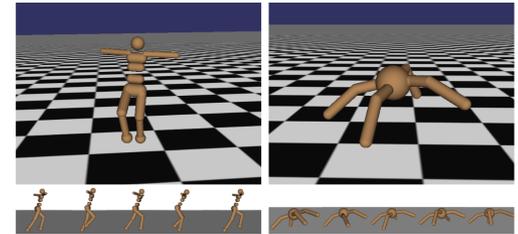
GAE: Generalized Advantage Estimation

- Use Advantage (i.e. $G - V(S)$)
- Use $TD(\lambda)$ target for G



GAE: Generalized Advantage Estimation

- Use Advantage (i.e. $G - V(S)$)
- Use $TD(\lambda)$ target for G



What about if we want a Deterministic Policy?

We can't use the Policy Gradient Theorem :

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^T (\gamma^t G_t) \nabla_{\theta} \log \pi(A_t | S_t) \right]$$

How can we estimate $\nabla_{\theta} J_{\theta}(\pi)$? When $A = \pi(S, \theta)$

What about if we want a Deterministic Policy?

We can't use the Policy Gradient Theorem :

$$\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^T (\gamma^t G_t) \nabla_{\theta} \log \pi(A_t | S_t) \right]$$

How can we estimate $\nabla_{\theta} J_{\theta}(\pi)$? When $A = \pi(S, \theta)$

$$J_{\theta}(\pi | S_0 = S) = q_{\pi}(\pi(S), S) \approx Q_{\pi}(\pi(S, \theta), S)$$

What about if we want a Deterministic Policy?

How can we estimate $\nabla_{\theta} J_{\theta}(\pi)$? When $A = \pi(S, \theta)$

$$A = (a_1, \dots, a_m), \pi = (\pi_1, \dots, \pi_m)$$

$$J_{\theta}(\pi | S_0 = S) = q_{\pi}(\pi(S), S) \approx Q_{\pi}(\pi(S, \theta), S)$$

What about if we want a Deterministic Policy?

How can we estimate $\nabla_{\theta} J_{\theta}(\pi)$? When $A = \pi(S, \theta)$

$$A = (a_1, \dots, a_m), \pi = (\pi_1, \dots, \pi_m)$$

$$J_{\theta}(\pi | S_0 = S) = q_{\pi}(\pi(S), S) \approx Q_{\pi}(\pi(S, \theta), S)$$

$$\nabla_{\theta} J_{\theta}(\pi | S_0 = S) \approx \nabla_{\theta} Q_{\pi}(\pi(S, \theta), S)$$

What about if we want a Deterministic Policy?

How can we estimate $\nabla_{\theta} J_{\theta}(\pi)$? When $A = \pi(S, \theta)$

$$A = (a_1, \dots, a_m), \pi = (\pi_1, \dots, \pi_m)$$

$$J_{\theta}(\pi | S_0 = S) = q_{\pi}(\pi(S), S) \approx Q_{\pi}(\pi(S, \theta), S)$$

$$\begin{aligned} \nabla_{\theta} J_{\theta}(\pi | S_0 = S) &\approx \nabla_{\theta} Q_{\pi}(\pi(S, \theta), S) \\ &= \sum_i^m \frac{\partial Q_{\pi}(A = \pi(S, \theta), S)}{\partial a_i} \nabla_{\theta} \pi_i(S, \theta) \end{aligned}$$

What about if we want a Deterministic Policy?

How can we estimate $\nabla_{\theta} J_{\theta}(\pi)$? When $A = \pi(S, \theta)$

$$A = (a_1, \dots, a_m), \pi = (\pi_1, \dots, \pi_m)$$

$$J_{\theta}(\pi | S_0 = S) = q_{\pi}(\pi(S), S) \approx Q_{\pi}(\pi(S, \theta), S)$$

$$\begin{aligned} \nabla_{\theta} J_{\theta}(\pi | S_0 = S) &\approx \nabla_{\theta} Q_{\pi}(\pi(S, \theta), S) \\ &= \sum_i^m \frac{\partial Q_{\pi}(A = \pi(S, \theta), S)}{\partial a_i} \nabla_{\theta} \pi_i(S, \theta) \\ &= \nabla_A Q_{\pi}(A = \pi(S, \theta), S) \nabla_{\theta} \pi(S, \theta) \end{aligned}$$

Deterministic Policy Gradient:

How can we estimate $\nabla_{\theta} J_{\theta}(\pi)$? When $A = \pi(S, \theta)$

$$A = (a_1, \dots, a_m), \pi = (\pi_1, \dots, \pi_m)$$

$$\begin{aligned} \nabla_{\theta} J_{\theta}(\pi | S_0 = S) &\approx \sum_i^m \frac{\partial Q_{\pi}(A = \pi(S, \theta), S)}{\partial a_i} \nabla_{\theta} \pi_i(S, \theta) \\ &= \nabla_A Q_{\pi}(A = \pi(S, \theta), S) \nabla_{\theta} \pi(S, \theta) \end{aligned}$$

Deterministic Policy Gradient (on Continuous Control Tasks):

Deterministic Policy Gradient Algorithms

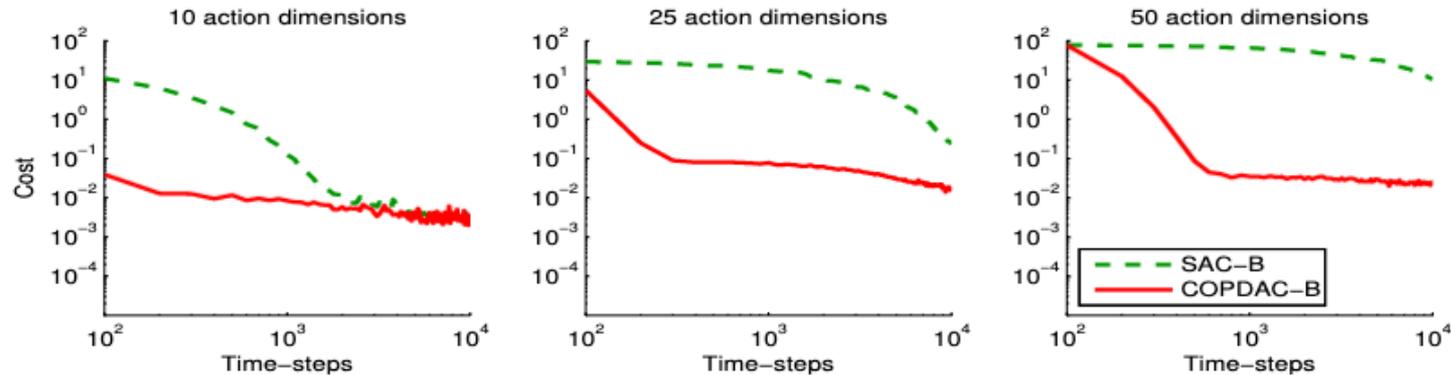


Figure 1. Comparison of stochastic actor-critic (SAC-B) and deterministic actor-critic (COPDAC-B) on the continuous bandit task.

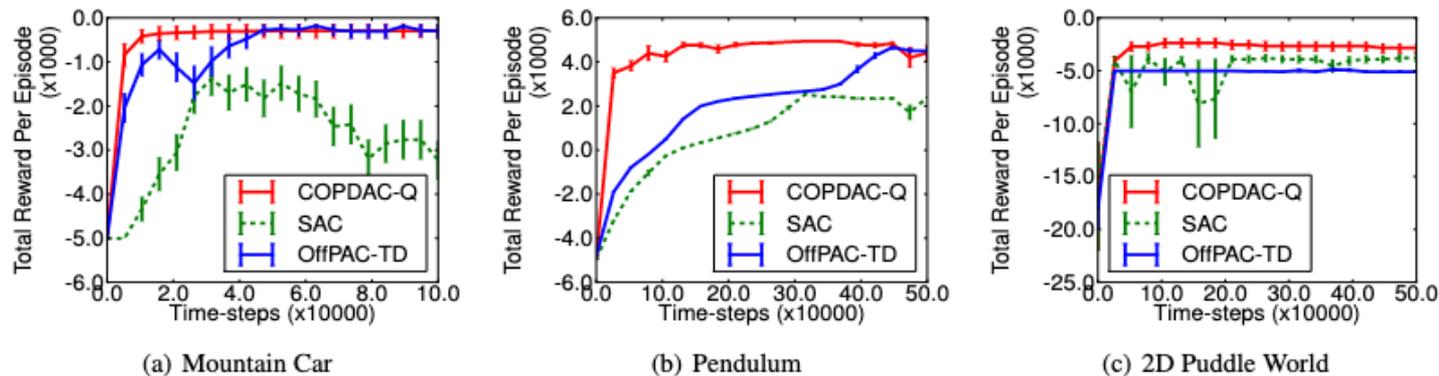


Figure 2. Comparison of stochastic on-policy actor-critic (SAC), stochastic off-policy actor-critic (OffPAC), and deterministic off-policy actor-critic (COPDAC) on continuous-action reinforcement learning. Each point is the average test performance of the mean policy.

Deterministic Policy Gradient (on Continuous Control Tasks):

Actor Critic with stochastic policy

Deterministic Policy Gradient

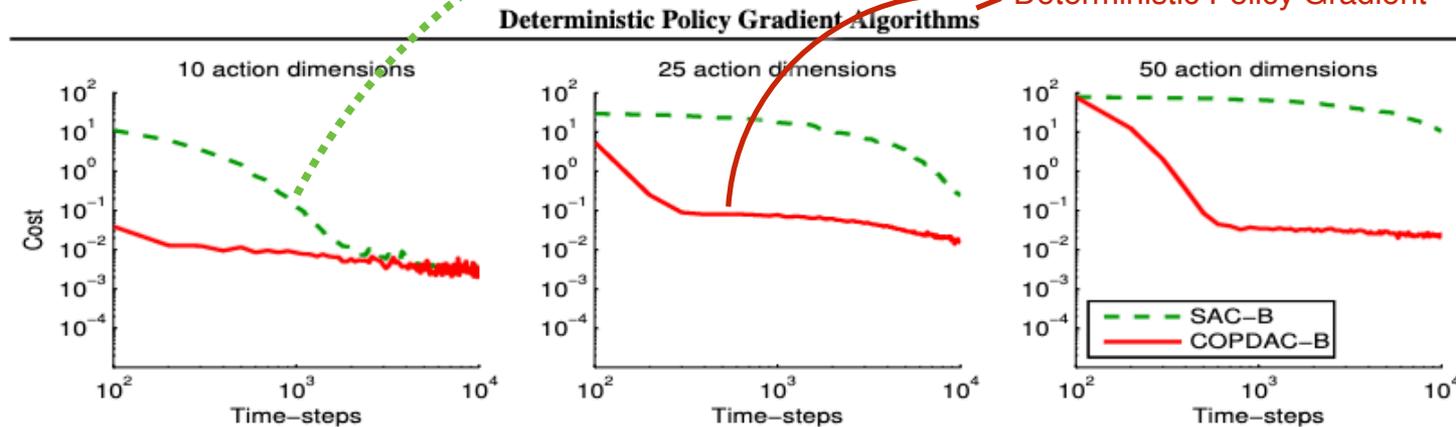


Figure 1. Comparison of stochastic actor-critic (SAC-B) and deterministic actor-critic (COPDAC-B) on the continuous bandit task.

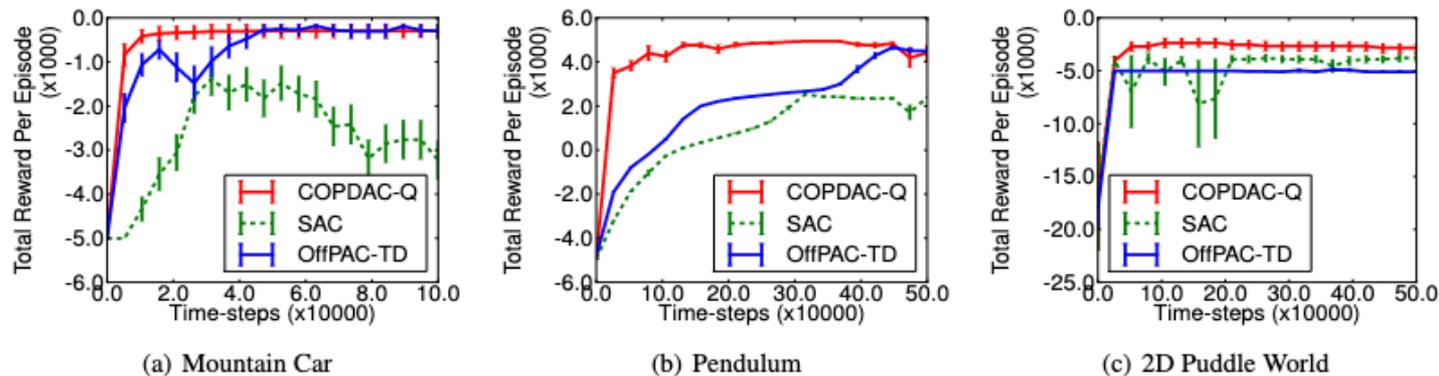


Figure 2. Comparison of stochastic on-policy actor-critic (SAC), stochastic off-policy actor-critic (OffPAC), and deterministic off-policy actor-critic (COPDAC) on continuous-action reinforcement learning. Each point is the average test performance of the mean policy.

Deep Deterministic Policy Gradient (DDPG):

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**

 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

Conclusion

- Policy Gradient Theorem: $\nabla_{\theta} J_{\theta}(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^T (\gamma^t G_t) \nabla_{\theta} \log \pi(A_t | S_t) \right]$
- REINFORCE: PGT + MC for estimate of G
- Actor-Critic: PGT + V,Q for estimate of G
- Deterministic Policy Gradient: $\nabla_{\theta} J_{\theta}(\pi | S_0 = S) \approx \nabla_{\theta} Q_{\pi}(\pi(S, \theta), S)$