

Lecture 12: Deterministic Policy Gradient Methods

Oliver Wallscheid



Table of contents

- 1 Deterministic gradient policy
- 2 Deep deterministic policy gradient (DDPG)
- 3 Twin delayed deep deterministic policy gradient (TD3)

Background and motivation

Recap on policy gradient so far:

- ▶ The previously discussed policy functions and the policy gradient theorem were assuming stochastic policies.
- ▶ The resulting on-policy algorithms may not provide top-class learning performance:
 - ▶ Non-guided exploration with step-by-step updates and
 - ▶ Greedy actions only in the limit (i.e., infeasible long learning).

The alternative:

- ▶ Apply a deterministic policy with separate exploration.
- ▶ Enable off-policy learning (with experience replay as a possible extension).
- ▶ Hence, we will focus on a **deterministic policy function**

$$\pi(\mathbf{x}, \boldsymbol{\theta}) = \mu(\mathbf{x}, \boldsymbol{\theta}). \quad (12.1)$$

Deterministic policy gradient (DPG) theorem

Theorem 12.1: Deterministic Policy Gradient

Given a metric $J(\boldsymbol{\theta})$ for the undiscounted episodic (11.7) or continuing tasks (11.8) and a parameterizable policy $\mu(\mathbf{x}, \boldsymbol{\theta})$ the deterministic policy gradient is

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{\mu} \left[\nabla_{\boldsymbol{\theta}} \mu(\mathbf{x}, \boldsymbol{\theta}) \nabla_{\mathbf{u}} q(\mathbf{x}, \mathbf{u}) |_{\mathbf{u}=\mu(\mathbf{x})} \right]. \quad (12.2)$$

- ▶ Again, q needs to be approximated using samples, e.g., implementing a critic via TD learning.
- ▶ It turns out that (12.2) is also (approximately) valid in the off-policy case, i.e., if the sample distribution is obtained from a behavior policy.
- ▶ Proof can be found in [D. Silver et al., *Deterministic Policy Gradient Algorithms*, International Conference on Machine Learning, 2014](#)

Exploration with a deterministic policy

- ▶ If the DPG approach is applied on-policy there is no inherent exploration.
- ▶ How to learn something?
 - ▶ The environment itself is sufficiently noisy (random impacts, measurement noise).
 - ▶ Or we have to add noise to the actions, i.e., making the approach off-policy.
 - ▶ Hence, utilizing a behavior policy is also possible.
- ▶ That additional action noise could be:
 - ▶ Simple Gaussian noise or
 - ▶ a shaped noise process like a discrete-time Ornstein-Uhlenbeck (OU) process

$$\nu_{k+1} = \lambda \nu_k + \sigma \epsilon_k$$

where ν_k is the OU noise output, $0 < \lambda < 1$ is a smoothing factor and σ is the variance scaling a standard Gaussian sequence (no mean) ϵ_k .

Algo. implementation: deterministic actor-critic

input: a differentiable deterministic policy function $\mu(\mathbf{x}, \boldsymbol{\theta})$
input: a differentiable action-value function $\hat{q}(\mathbf{x}, \mathbf{u}, \mathbf{w})$
parameter: step sizes $\{\alpha_w, \alpha_\theta\} \in \{\mathbb{R} | 0 < \alpha < 1\}$
init: parameter vectors $\mathbf{w} \in \mathbb{R}^\zeta$ and $\boldsymbol{\theta} \in \mathbb{R}^d$ arbitrarily
for $j = 1, 2, \dots$, *episodes* **do**
 initialize \mathbf{x}_0 ;
 for $k = 0, 1, \dots, T - 1$ *time steps* **do**
 $\mathbf{u}_k \leftarrow$ apply from $\mu(\mathbf{x}_k, \boldsymbol{\theta})$ w/wo noise or from behavior policy;
 observe \mathbf{x}_{k+1} and r_{k+1} ;
 choose \mathbf{u}' from $\mu(\mathbf{x}_{k+1}, \boldsymbol{\theta})$;
 $\delta \leftarrow r_{k+1} + \gamma \hat{q}(\mathbf{x}_{k+1}, \mathbf{u}', \mathbf{w}) - \hat{q}(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w})$;
 $\mathbf{w} \leftarrow \mathbf{w} + \alpha_w \delta \nabla_{\mathbf{w}} \hat{q}(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w})$;
 $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \alpha_\theta \gamma^k \nabla_{\boldsymbol{\theta}} \mu(\mathbf{x}_k, \boldsymbol{\theta}) \nabla_{\mathbf{u}} \hat{q}(\mathbf{x}_k, \mathbf{u}_k, \mathbf{w})|_{\mathbf{u}=\mu(\mathbf{x})}$;

Algo. 12.1: Deterministic actor-critic for episodic tasks using SARSA(0) targets applicable on- and off-policy (output: parameter vector $\boldsymbol{\theta}^*$ for $\mu^*(\mathbf{x}, \boldsymbol{\theta}^*)$) and \mathbf{w}^* for $\hat{q}^*(\mathbf{x}, \mathbf{u}, \mathbf{w}^*)$)

Exemplary comparison to stochastic policy gradient

- ▶ DPG-based approach uses compatible function approximation, i.e., suitable linear \hat{q} estimation. A fixed Gaussian behavior policy is applied for exploration.
- ▶ SAC uses a Gaussian policy with linear function approximation.

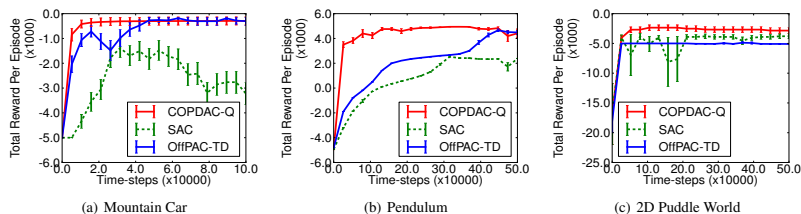


Fig. 12.1: Comparison of stochastic on-policy actor-critic (SAC), stochastic off-policy actor-critic (OffPAC), and deterministic off-policy actor-critic (COPDAC) on continuous-action reinforcement learning (source: D. Silver et al., *Deterministic Policy Gradient Algorithms*, International Conference on Machine Learning, 2014)

Table of contents

- 1 Deterministic gradient policy
- 2 Deep deterministic policy gradient (DDPG)
- 3 Twin delayed deep deterministic policy gradient (TD3)

- ▶ The upcoming **deep deterministic policy gradient (DDPG)** algorithm was very much inspired by the successes of DQNs (cf. Algo. 10.6 and landmark [paper by Mnih et al.](#)) on discrete action spaces.
- ▶ However, **DQNs are not directly applicable to (quasi-)continuous action spaces.**
- ▶ Recall the incremental Q -learning equation using function approximation

$$\mathbf{w} \leftarrow \mathbf{w} + \alpha \left[r + \gamma \max_u \hat{q}(\mathbf{x}', u, \mathbf{w}) - \hat{q}(\mathbf{x}, u, \mathbf{w}) \right] \nabla_{\mathbf{w}} \hat{q}(\mathbf{x}, u, \mathbf{w}).$$

- ▶ For every policy inference and updating step we need to find $\max_u \hat{q}(\mathbf{x}', u, \mathbf{w})$.
- ▶ If $u \in \mathcal{U} \subset \mathbb{Z}$ (i.e., using integer-encoded actions) is a sufficiently small discrete set, that is straightforward by an exhaustive search.
- ▶ In contrast, if $\mathbf{u} \in \mathcal{U} \subset \mathbb{R}^m$ is a (quasi-)continuous variable solving $\max_{\mathbf{u}} \hat{q}(\mathbf{x}', \mathbf{u}, \mathbf{w})$ requires an own **optimization routine** which is computationally expensive if we use nonlinear function approximation.

The deterministic policy trick

- ▶ When using a greedy, deterministic policy $\pi(\mathbf{x}, \boldsymbol{\theta}) = \boldsymbol{\mu}(\mathbf{x}, \boldsymbol{\theta})$ we can utilize it to approximate

$$\max_{\mathbf{u}} \hat{q}(\mathbf{x}', \mathbf{u}, \mathbf{w}) \approx \hat{q}(\mathbf{x}', \boldsymbol{\mu}(\mathbf{x}', \boldsymbol{\theta}), \mathbf{w}). \quad (12.3)$$

- ▶ Hence, we can obtain explicit Q -learning targets for continuous actions when using a deterministic policy.
- ▶ For improving the policy we reuse the deterministic policy gradient theorem in an off-policy fashion

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_b [\nabla_{\boldsymbol{\theta}} \boldsymbol{\mu}(\mathbf{X}, \boldsymbol{\theta}) \nabla_{\mathbf{u}} q(\mathbf{X}, \mathbf{U}) | \mathbf{U} = \boldsymbol{\mu}(\mathbf{X}, \boldsymbol{\theta})] \quad (12.4)$$

given a behavior policy $b(\mathbf{u}|\mathbf{x})$.

DDPG \approx DQN + DPG

- ▶ Hence, we can consider the DDPG approach as a combination of DQN + DPG rendering it an **actor-critic off-policy approach for continuous state and action spaces**.
- ▶ Similarly to DQN we will introduce **several 'tweaks'** to stabilize and improve the DDPG learning process.

Tweak #1: experience replay buffer

- ▶ We store $\langle \mathbf{x}, \mathbf{u}, r, \mathbf{x}' \rangle$ in \mathcal{D} after each transition step.
- ▶ The replay buffer \mathcal{D} is of limited capacity, i.e., it discards the oldest data sample when updating once it is full (ring memory).
- ▶ This allows us to improve the Q -learning critic minimizing the mean-squared Bellman error (MSBE):

$$\mathcal{L}(\mathbf{w}) = \left[\left(r + \gamma q(\mathbf{x}', \boldsymbol{\mu}(\mathbf{x}', \boldsymbol{\theta}), \mathbf{w}) \right) - q(\mathbf{x}, \mathbf{u}, \mathbf{w}) \right]_{\mathcal{D}}^2. \quad (12.5)$$

Additional DDPG tweaks (1)

Tweak #2: target networks

- ▶ Similar to DQN we introduce a (delayed) target network to estimate the Q -learning target

$$r + \gamma q(\mathbf{x}', \boldsymbol{\mu}(\mathbf{x}', \boldsymbol{\theta}), \mathbf{w})$$

since it depends on the same parameters \mathbf{w} which we want to update.

- ▶ Hence, the target network's purpose is to mimic the generation of i.i.d. data as the ground truth to minimize (12.5).
- ▶ Since the policy parameters $\boldsymbol{\theta}$ are also part of the target calculation it turns out that an additional policy target network is also beneficial to stabilize the Q -learning.
- ▶ In contrast to the classical DQN implementation, the original DDPG algorithm does not perform periodically hard target network updates but continuous ones using a low-pass filter characteristic

$$\mathbf{w}^- \leftarrow (1 - \tau)\mathbf{w}^- + \tau\mathbf{w}, \quad \boldsymbol{\theta}^- \leftarrow (1 - \tau)\boldsymbol{\theta}^- + \tau\boldsymbol{\theta} \quad (12.6)$$

with τ representing the equivalent filter constant (hyperparameter).

Additional DDPG tweaks (2)

Tweak #3: mini-batch sampling

- ▶ Given a sufficiently filled memory \mathcal{D} and the target networks parametrized by w^- and θ^- we draw uniformly distributed mini-batch samples \mathcal{D}_b from \mathcal{D} .
- ▶ The actual Q -learning is then based on the loss

$$\mathcal{L}(w) = \left[(r + \gamma q(x', \mu(x', \theta^-), w^-)) - q(x, u, w) \right]_{\mathcal{D}_b}^2. \quad (12.7)$$

Tweak #4: batch normalization

- ▶ Minimizing (12.7) is a supervised learning step within the DDPG.
- ▶ The [original DDPG paper by Lillicrap et al.](#) back in 2015/16 suggested to use batch normalization, i.e., re-centering and re-scaling the inputs of each layer in an ANN.
- ▶ This idea of batch normalization was presented at that time shortly before by Ioffe and Szegedy (cf. [original paper](#)).
- ▶ Today's perspective: stick to the current state-of-the-art supervised ML algorithms for top-class Q -learning stability and speed (which are normally well-covered in popular supervised ML toolboxes).

Additional DDPG tweaks (3)

Tweak #5: exploration

- ▶ Since our policy is deterministic we require an exploratory behavior policy.
- ▶ Similar to DPG the standard approach is to add noise to the greedy actions, e.g., again from an Ornstein-Uhlenbeck (OU) process

$$\mathbf{u}_k \sim \mathbf{b}(\mathbf{u}|\mathbf{x}_k) = \boldsymbol{\mu}(\mathbf{x}_k, \boldsymbol{\theta}_k) + \boldsymbol{\nu}_k, \quad \boldsymbol{\nu}_k = \lambda \boldsymbol{\nu}_{k-1} + \sigma \boldsymbol{\epsilon}_{k-1}.$$

- ▶ One might also add a schedule for λ and σ along the training procedure, e.g., starting with significant noise levels (increased exploration) while reducing it over time (focusing exploitation)¹.
- ▶ However, many other behavior policies are possible, e.g., using model or expert-based guidance.

¹Please note that this 'lambda' is not related to TD(λ), SARSA(λ), etc. Here, it is representing the stiffness of the OU noise process.

Visual summary of DDPG working principle

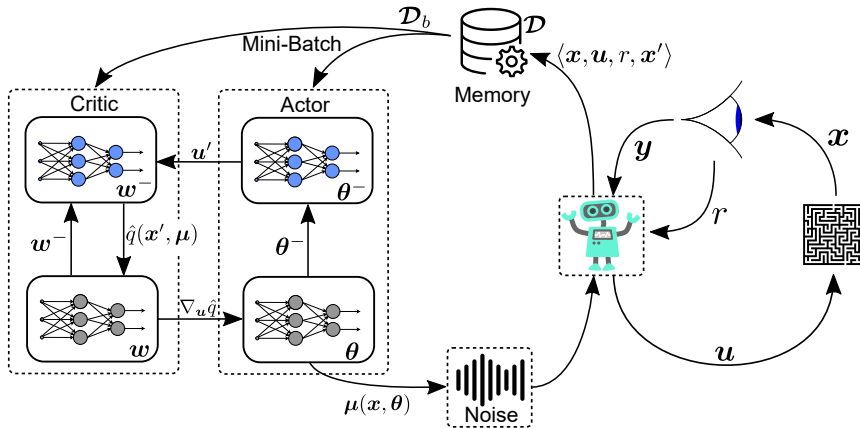


Fig. 12.2: DDPG structure from a bird's-eye perspective (derivative work of Fig. 1.1 and [wikipedia.org](https://en.wikipedia.org/wiki/Deep_Deterministic_Policy_Gradient), CC0 1.0)

Algo. implementation: DDPG

```
input: diff. deterministic policy function  $\mu(x, \theta)$  and action-value function  $\hat{q}(x, u, w)$   
parameter: step sizes and filter constant  $\{\alpha_w, \alpha_\theta, \tau\} \in \{\mathbb{R} | 0 < \alpha, \tau < 1\}$   
init: weights  $w = w^- \in \mathbb{R}^\zeta$  and  $\theta = \theta^- \in \mathbb{R}^d$  arbitrarily, memory  $\mathcal{D}$   
for  $j = 1, 2, \dots$ , episodes do  
    initialize  $x_0$ ;  
    for  $k = 0, 1, \dots, T - 1$  time steps do  
         $u_k \leftarrow$  apply from  $\mu(x_k, \theta)$  w/o noise or from behavior policy;  
        observe  $x_{k+1}$  and  $r_{k+1}$ ;  
        store tuple  $\langle x_k, u_k, r_{k+1}, x_{k+1} \rangle$  in  $\mathcal{D}$ ;  
        sample mini-batch  $\mathcal{D}_b$  from  $\mathcal{D}$  (after initial memory warmup);  
        for  $i = 1, \dots, b$  samples do calculate  $Q$ -targets  
            if  $x_{i+1}$  is terminal then  $y_i = r_{i+1}$ ;  
            else  $y_i = r_{i+1} + \gamma \hat{q}(x_{i+1}, \mu(x_{i+1}, \theta^-), w^-)$ ;  
        fit  $w$  on loss  $\mathcal{L}(w) = [y - \hat{q}(x, u, w)]_{\mathcal{D}_b}^2$  with step size  $\alpha_w$ ;  
         $\theta \leftarrow \theta + \alpha_\theta [\nabla_\theta \mu(x, \theta) \nabla_u \hat{q}(x, u, w)|_{u=\mu_\theta(x)}]_{\mathcal{D}_b}$ ;  
        Update target net.  $w^- \leftarrow (1 - \tau)w^- + \tau w$ ,  $\theta^- \leftarrow (1 - \tau)\theta^- + \tau \theta$ ;
```

Algo. 12.2: Deep deterministic policy gradient (output: parameter vectors θ^* for $\mu^*(x, \theta^*)$)

Table of contents

- 1 Deterministic gradient policy
- 2 Deep deterministic policy gradient (DDPG)
- 3 Twin delayed deep deterministic policy gradient (TD3)

Overestimation bias

- ▶ For Q -learning in the tabular case we have already discussed the **maximization bias** (cf. Fig. 5.13) issue.
- ▶ Recap: Due to the greedy policy targets, \hat{q} was overestimated when calculated using sampled values of stochastic MDPs.
- ▶ Additional problem when applying function approximation: the estimator itself introduces additional variance during the learning process which represents another source of the maximization bias problem.

This issue is already known in the DQN context (cf. Algo. 10.6). Similar to the tabular case, **double DQN** introduces a second Q -network counteracting the overestimation issue (cf. [paper by van Hasselt et al.](#)).

However, we did not address this possible problem in an actor-critic context using function approximation (e.g., DDPG).

Overestimation bias in actor-critic approaches (1)

- ▶ It turns out that the overestimation bias is also an issue for actor-critic methods¹.
- ▶ Consider an actor-critic policy with the current policy parameters θ .
- ▶ Let $\tilde{\theta}$ define the parameters from the actor update induced by the maximization of the approximate critic $\hat{q}_w(\mathbf{x}, \mathbf{u})$.
- ▶ Let θ^* be the parameters from the hypothetical actor update w.r.t. the true underlying value function $q^\pi(\mathbf{x}, \mathbf{u})$.
- ▶ Then, we perform the policy update

$$\begin{aligned}\tilde{\theta} &= \theta + \frac{\alpha}{Z_1} \mathbb{E}_\pi [\nabla_\theta \pi_\theta(\mathbf{X}) \nabla_{\mathbf{u}} \hat{q}_w(\mathbf{X}, \mathbf{U}) | \mathbf{U} = \pi_\theta(\mathbf{X})], \\ \theta^* &= \theta + \frac{\alpha}{Z_2} \mathbb{E}_\pi [\nabla_\theta \pi_\theta(\mathbf{X}) \nabla_{\mathbf{u}} q^\pi(\mathbf{X}, \mathbf{U}) | \mathbf{U} = \pi_\theta(\mathbf{X})],\end{aligned}\tag{12.8}$$

where Z_1 and Z_2 normalize the gradient such that $Z^{-1} \|\mathbb{E}[\cdot]\| = 1$.

¹Source: S. Fujimoto et al., *Addressing Function Approximation Error in Actor-Critic Methods*, <https://arxiv.org/abs/1802.09477>, 2018

Overestimation bias in actor-critic approaches (2)

- ▶ Lets denote $\tilde{\pi}$ and π^* as the policies with updated parameters $\tilde{\theta}$ and θ^* respectively.
- ▶ As the gradient direction is a local maximizer, there exists ϵ_1 sufficiently small such that if $\alpha \leq \epsilon_1$ then the *approximate* value of $\tilde{\pi}$ will be bounded below by the *approximate* value of π^* :

$$\mathbb{E} [\hat{q}_w(\mathbf{X}, \tilde{\pi}(\mathbf{X}))] \geq \mathbb{E} [\hat{q}_w(\mathbf{X}, \pi^*(\mathbf{X}))]. \quad (12.9)$$

- ▶ Conversely, there exists ϵ_2 sufficiently small such that if $\alpha \leq \epsilon_2$ then the *true* value of $\tilde{\pi}$ will be bounded above by the *true* value of π^* :

$$\mathbb{E} [q^\pi(\mathbf{X}, \pi^*(\mathbf{X}))] \geq \mathbb{E} [q^\pi(\mathbf{X}, \tilde{\pi}(\mathbf{X}))]. \quad (12.10)$$

- ▶ In other words: if the approximate and true critics differ from each other, the according policy gradient updates cannot lead to better policy updates of the respective other framework.

Overestimation bias in actor-critic approaches (3)

- ▶ If the expected, estimated action value will be at least as large as the *true* action value w.r.t. θ^*

$$\mathbb{E} [\hat{q}_w(\mathbf{X}, \pi^*(\mathbf{X}))] \geq \mathbb{E} [q^\pi(\mathbf{X}, \pi^*(\mathbf{X}))], \quad (12.11)$$

then (12.9) and (12.10) imply

$$\mathbb{E} [\hat{q}_w(\mathbf{X}, \tilde{\pi}(\mathbf{X}))] \geq \mathbb{E} [q^\pi(\mathbf{X}, \tilde{\pi}(\mathbf{X}))] \quad (12.12)$$

with a sufficiently small $\alpha < \min\{\epsilon_1, \epsilon_2\}$.

- ▶ Hence, the **maximization bias is also present in actor-critic** updates.
- ▶ It can add up over several estimation updates and, therefore, may lead to suboptimal policy updates.
- ▶ A proof for unnormalized gradients can be also found in S. Fujimoto et al., *Addressing Function Approximation Error in Actor-Critic Methods*, 2018.

Overestimation example for DDPG

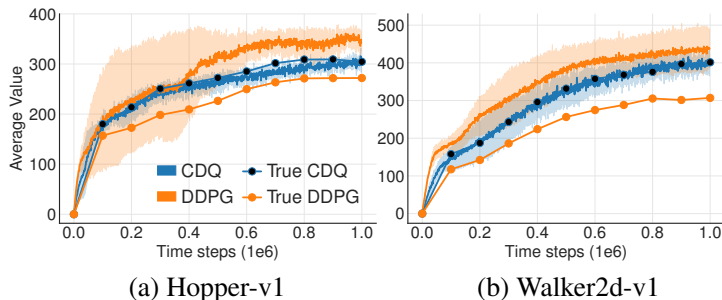


Fig. 12.3: Comparison of true and estimated values averaged over 10000 states in two robotic examples from [OpenAI Gym](#). Estimated values originate from the approximate DDPG critic while the true values are based on the average discounted return over 1000 episodes following the current policy, starting from states sampled from the replay buffer (source: S. Fujimoto et al., *Addressing Function Approximation Error in Actor-Critic Methods*, 2018).

Increased variance due to accumulating TD errors

- ▶ Using function approximation, the **Bellman equation is never exactly satisfied** leaving room for some amount of **residual TD-error** $\tilde{\delta}(\mathbf{x}, \mathbf{u})$:

$$\hat{q}_{\mathbf{w}}(\mathbf{x}, \mathbf{u}) = r + \gamma \mathbb{E}_{\pi} [\hat{q}_{\mathbf{w}}(\mathbf{X}', \mathbf{U}') | \mathbf{X}' = \mathbf{x}', \mathbf{U}' = \mathbf{u}'] - \tilde{\delta}(\mathbf{x}, \mathbf{u}). \quad (12.13)$$

- ▶ Although this error might be considered small per update step, it may accumulate over future steps if biased:

$$\hat{q}_{\mathbf{w}}(\mathbf{x}, \mathbf{u}) = \mathbb{E}_{\pi} \left[\sum_{k=0}^{\infty} \gamma^k \left(R_k - \tilde{\delta}_k(\mathbf{X}, \mathbf{U}) \right) \middle| \mathbf{X} = \mathbf{x}, \mathbf{U} = \mathbf{u} \right]. \quad (12.14)$$

- ▶ Observation: the **variance of \hat{q} will be proportional to the variance of future reward and residual TD-errors**.
- ▶ If γ is large, the estimation variance might increase significantly.
- ▶ Mini-batch sampling will contribute to this variance issue.

TD3 extensions and modifications (1)

In order to reduce both the maximization bias and the learning variance, TD3 introduces mainly three measures on top of the DDPG algorithm. Hence, **TD3 is a direct successor of DDPG**.

Measure #1: clipped double Q -learning for actor-critic

- ▶ Following double Q -learning, a pair of critics $\{\hat{q}_{w_1}, \hat{q}_{w_2}\}$ is introduced.
- ▶ In contrast, the clipped target (with target networks $\{w_1^-, w_2^-\}$)

$$y = r + \gamma \min_{i=1,2} \hat{q}_{w_i^-}(x', u') \quad (12.15)$$

provides an upper-bound on the estimated action value.

- ▶ May introduce some underestimation, which is considered less critical than overestimation, since the value of underestimated actions will not be explicitly propagated through the policy update.
- ▶ The min operator will also (indirectly) favor actions leading to values with estimation errors of lower variance.

TD3 extensions and modifications (2)

Measure #2: target policy smoothing regularization

- ▶ Background: deterministic policies μ tend to overfit to narrow peaks in the action-value estimate.
- ▶ Counteraction: fit the action value of a small area around the target action (i.e., smoothing \hat{q} in the action space):

$$y = r + \gamma \hat{q}_{w^-}(\mathbf{x}', \mu_{\theta^-}(\mathbf{x}') + \epsilon). \quad (12.16)$$

- ▶ Here, $\epsilon \sim \text{clip}(\mathcal{N}(\mathbf{0}, \Sigma), -c, c)$ is a mean-free, Gaussian noise with covariance Σ , which is clipped at $\pm c$ while θ^- are the policy target network parameters.
- ▶ To satisfy possible action constraints (denoted by upper and lower box constraints $\{\underline{u}, \bar{u}\}$), we add an additional clipping:

$$\mathbf{u}' = \text{clip}(\mu_{\theta^-}(\mathbf{x}') + \epsilon, \underline{u}, \bar{u}). \quad (12.17)$$

- ▶ This modified action is then used for the target calculation (12.15).

TD3 extensions and modifications (3)

Measure #3: delayed policy updates

- ▶ Similar to DDPG, TD3 uses policy target networks θ^- and (two) critic target networks $\{w_1^-, w_2^-\}$ in order to provide (rather) fixed Q -learning targets trying to stabilize the learning of \hat{q} .
- ▶ The target networks are also continuously updated using

$$w_i^- \leftarrow (1 - \tau)w_i^- + \tau w_i, \quad \theta^- \leftarrow (1 - \tau)\theta^- + \tau \theta.$$

- ▶ However, each policy update will inherently change the (true) Q -learning target directly adding variance to the learning process (cf. Fig. 12.4 on next slide).
- ▶ Therefore, it is argued that a policy update should not follow after each Q -learning update such that the critic can adapt properly to the previous policy update.
- ▶ The original TD3 implementation suggests a policy update every second Q -learning update, however, we can consider this update rate a hyperparameter.

TD3 extensions and modifications (4)

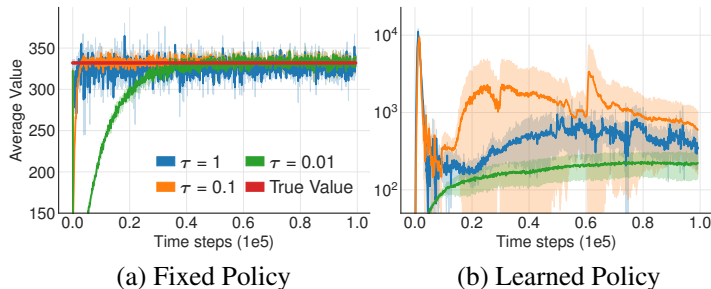


Fig. 12.4: Average estimated action value of a randomly selected state on Hopper-v1 environment from OpenAI Gym (source: S. Fujimoto et al., *Addressing Function Approximation Error in Actor-Critic Methods*, 2018).

input: diff. deterministic policy function $\mu(x, \theta)$ and action-value function $\hat{q}(x, u, w)$

parameter: step sizes and filter constant $\{\alpha_w, \alpha_\theta, \tau\} \in \{\mathbb{R} | 0 < \alpha, \tau < 1\}$, policy update rate $k_w \in \{\mathbb{N} | 1 \leq k_w\}$, target noise $\Sigma \in \mathbb{R}^{m \times m}$ and $c \in \mathbb{R}^m$

init: weights $\{w_1 = w_1^-, w_2 = w_2^-\} \in \mathbb{R}^\zeta$, $\theta = \theta^- \in \mathbb{R}^d$ arbitrarily, memory \mathcal{D}

for $j = 1, 2, \dots$, *episodes* **do**
 initialize x_0 ;
 for $k = 0, 1, \dots, T - 1$ *time steps* **do**
 $u_k \leftarrow$ apply from $\mu(x_k, \theta)$ w/wo noise or from behavior policy;
 observe x_{k+1} and r_{k+1} ;
 store tuple $\langle x_k, u_k, r_{k+1}, x_{k+1} \rangle$ in \mathcal{D} ;
 sample mini-batch \mathcal{D}_b from \mathcal{D} (after initial memory warmup);
 for $i = 1, \dots, b$ *samples* **do** calculate Q -targets
 if x_{i+1} *is terminal* **then** $y_i = r_{i+1}$;
 else
 $u' = \text{clip}(\mu_{\theta^-}(x_{i+1}) + \text{clip}(\mathcal{N}(\mathbf{0}, \Sigma), -c, c), \underline{u}, \bar{u})$;
 $y_i = r_{i+1} + \gamma \min_{l=1,2} \hat{q}(x_{i+1}, u', w_l^-)$;
 fit w_l on loss $\mathcal{L}(w_l) = [y - \hat{q}(x, u, w_l)]_{\mathcal{D}_b}^2$ with step size $\alpha_w \forall l$;
 if $k \bmod k_w = 0$ **then**
 $\theta \leftarrow \theta + \alpha_\theta [\nabla_\theta \mu(x, \theta) \nabla_u \hat{q}(x, u, w_1) |_{u=\mu_\theta(x)}]_{\mathcal{D}_b}$;
 $w_l^- \leftarrow (1 - \tau)w_l^- + \tau w_l$, $\theta^- \leftarrow (1 - \tau)\theta^- + \tau \theta$;

Algo. 12.3: Twin delayed deep deterministic policy gradient (TD3)

Summary: what you've learned today

- ▶ The deep deterministic policy gradient (DDPG) approach 'transfers' many deep Q -network (DQN) ideas to continuous action spaces.
- ▶ It mainly combines DQN + deterministic policy gradients + policy and value target networks (plus additional minor tweaks).
- ▶ However, the DDPG actor-critic suffers from value overestimation and high variance during learning. Hence, sampled policy gradients might not be optimal (pointing towards overrated action values).
- ▶ Twin delayed DDPG (TD3) adds clipped double Q -learning, delayed policy updates and target policy smoothing to counteract these issues.