NPFL139, Lecture 8



Continuous Action Space: DDPG, TD3, SAC

Milan Straka

🖬 April 09, 2025





EUROPEAN UNION European Structural and Investment Fund Operational Programme Research, Development and Education Charles University in Prague Faculty of Mathematics and Physics Institute of Formal and Applied Linguistics



unless otherwise stated



Continuous Action Space

NPFL139, Lecture 8

ContinuousActionSpace

ace DPG

DDPG

MuJoCo TD3

SAC SPE

SPI

SAC Algorithm

Continuous Action Space

Until now, the actions were discrete. However, many environments naturally accept actions from continuous space. We now consider actions which come from range [a, b] for $a, b \in \mathbb{R}$, or more generally from a Cartesian product of several such ranges:

 $\prod [a_i, b_i].$

A simple way how to parametrize the action distribution is to choose them from the normal distribution. Given mean μ and variance σ^2 , probability density function of $\mathcal{N}(\mu, \sigma^2)$ is





SPI

ContinuousActionSpace

DDPG

DPG

MuJoCo TD3

SAC SPE



Continuous Action Space in Gradient Methods

Ú F_AL

Utilizing continuous action spaces in gradient-based methods is straightforward. Instead of the softmax distribution, we suitably parametrize the action value, usually using the normal distribution.

Considering only one real-valued action, we therefore have

$$\pi(a|s;oldsymbol{ heta}) \stackrel{ ext{def}}{=} P\Big(a \sim \mathcal{N}ig(\mu(s;oldsymbol{ heta}), \sigma(s;oldsymbol{ heta})^2ig)\Big),$$

where $\mu(s; \theta)$ and $\sigma(s; \theta)$ are function approximation of mean and standard deviation of the action distribution.

The mean and standard deviation are usually computed from the shared representation, with

- the mean being computed as a usual regression (i.e., one output neuron without activation);
- the standard deviation (which must be positive) being computed again as a single neuron, but with either exp or softplus, where $ext{softplus}(x) \stackrel{\text{\tiny def}}{=} \log(1 + e^x)$.

Continuous Action Space in Gradient Methods

Ú F_AL

During training, we compute $\mu(s; \theta)$ and $\sigma(s; \theta)$ and then sample the action value (clipping it to [a, b] if required). To compute the loss, we utilize the probability density function of the normal distribution (and usually also add the entropy penalty).

```
mus = torch.nn.Linear(..., actions)(hidden_layer)
sds = torch.nn.Linear(..., actions)(hidden_layer)
sds = torch.exp(sds)  # or sds = torch.nn.softplus(sds)
```

action_dist = torch.distributions.Normal(mus, sds)

ContinuousActionSpace

SPE

SPI

Continuous Action Space

When the action consists of several real values, i.e., the action is a suitable subregion of \mathbb{R}^n for n > 1, we can:

- either use multivariate Gaussian distribution;
- or factorize the probability into a product of univariate normal distributions.
 - This is the most commonly used approach; we then consider the action to be composed of several independent **action components**.

If modeling the action distribution using a unimodal normal distribution is insufficient, a mixture of normal distributions (or mixture of logistic) can be used, capable of representing also multimodal distributions.

DDPG

SAC

SPE

SPI



Deterministic Policy Gradient Theorem

NPFL139, Lecture 8

ContinuousActionSpace

DPG DDPG

MuJoCo TD3

SPE SPI

SAC

SAC Algorithm

Deterministic Policy Gradient Theorem

Combining continuous actions and Deep Q Networks is not straightforward. In order to do so, we need a different variant of the policy gradient theorem.

Recall that in policy gradient theorem,

$$abla_{oldsymbol{ heta}} J(oldsymbol{ heta}) \propto \mathbb{E}_{s \sim \mu} \Big[\sum_{a \in \mathcal{A}} q_{\pi}(s,a)
abla_{oldsymbol{ heta}} \pi(a|s;oldsymbol{ heta}) \Big].$$

Deterministic Policy Gradient Theorem

Assume that the policy $\pi(s; \theta)$ is deterministic and computes an action $a \in \mathbb{R}$. Further, assume the reward r(s, a) is actually a deterministic function of the given state-action pair. Then, under several assumptions about continuousness, the following holds:

$$abla_{oldsymbol{ heta}} J(oldsymbol{ heta}) \propto \mathbb{E}_{s \sim \mu} \Big[
abla_{oldsymbol{ heta}} \pi(s;oldsymbol{ heta})
abla_a q_{\pi}(s,a) ig|_{a=\pi(s;oldsymbol{ heta})} \Big].$$

The theorem was first proven in the paper Deterministic Policy Gradient Algorithms by David Silver et al in 2014.

NPFL139, Lecture 8

MuJoCo TD3

SAC

SPE

SPI



Deterministic Policy Gradient Theorem – Proof

Ú F_ÅL

The proof is very similar to the original (stochastic) policy gradient theorem. However, we will be exchanging derivatives and integrals, for which we need several assumptions:

- we assume that $h(s), p(s'|s, a), \nabla_a p(s'|s, a), r(s, a), \nabla_a r(s, a), \pi(s; \theta), \nabla_{\theta} \pi(s; \theta)$ are continuous in all parameters and variables;
- we further assume that $h(s), p(s'|s, a), \nabla_a p(s'|s, a), r(s, a), \nabla_a r(s, a)$ are bounded.

Details (which assumptions are required and when) can be found in Appendix B of the paper *Deterministic Policy Gradient Algorithms: Supplementary Material* by David Silver et al.

SAC

SPE

SPI

Deterministic Policy Gradient Theorem – Proof

$$\begin{aligned} \nabla_{\boldsymbol{\theta}} v_{\pi}(s) &= \nabla_{\boldsymbol{\theta}} q_{\pi}(s, \pi(s; \boldsymbol{\theta})) \\ &= \nabla_{\boldsymbol{\theta}} \left(r\left(s, \pi(s; \boldsymbol{\theta})\right) + \int_{s'} \gamma p\left(s'|s, \pi(s; \boldsymbol{\theta})\right) \left(v_{\pi}(s')\right) \mathrm{d}s' \right) \\ &= \nabla_{\boldsymbol{\theta}} \pi(s; \boldsymbol{\theta}) \nabla_{a} r(s, a) \big|_{a=\pi(s; \boldsymbol{\theta})} + \nabla_{\boldsymbol{\theta}} \int_{s'} \gamma p\left(s'|s, \pi(s; \boldsymbol{\theta})\right) v_{\pi}(s') \mathrm{d}s' \\ &= \nabla_{\boldsymbol{\theta}} \pi(s; \boldsymbol{\theta}) \nabla_{a} \left(r(s, a) + \int_{s'} \gamma p\left(s'|s, a\right) v_{\pi}(s') \mathrm{d}s' \right) \Big|_{a=\pi(s; \boldsymbol{\theta})} \\ &+ \int_{s'} \gamma p\left(s'|s, \pi(s; \boldsymbol{\theta})\right) \nabla_{\boldsymbol{\theta}} v_{\pi}(s') \mathrm{d}s' \end{aligned}$$

We finish the proof as in the gradient theorem by continually expanding $\nabla_{\theta} v_{\pi}(s')$, getting $abla_{m{ heta}} v_{\pi}(s) = \int_{s'} \sum_{k=0}^{\infty} \gamma^k P(s o s' \text{ in } k \text{ steps } |\pi) \left[
abla_{m{ heta}} \pi(s'; m{ heta})
abla_a q_{\pi}(s', a) \Big|_{a=\pi(s'; m{ heta})} \right] \mathrm{d}s' \text{ and}$ then $abla_{m{ heta}} J(m{ heta}) = \mathbb{E}_{s \sim h}
abla_{m{ heta}} v_{\pi}(s) \propto \mathbb{E}_{s \sim \mu} ig[
abla_{m{ heta}} \pi(s; m{ heta})
abla_{a} q_{\pi}(s, a) ig|_{a = \pi(s; m{ heta})} ig].$

DDPG ContinuousActionSpace DPG MuJoCo

TD3

SAC Algorithm



NPFL139, Lecture 8

ContinuousActionSpace

DPG DDPG

MuJoCo TD3

SAC SPE

SPI

SAC Algorithm

Ú_F≩L

Note that the formulation of deterministic policy gradient theorem allows an off-policy algorithm, because the loss functions no longer depends on actions sampled from the behaviour policy (similarly to how expected Sarsa is also an off-policy algorithm).

We therefore train function approximation for both $\pi(s; \theta)$ and $q(s, a; \theta)$, training $q(s, a; \theta)$ using a deterministic variant of the Bellman equation

$$q(S_t,A_t;oldsymbol{ heta}) = \mathbb{E}_{S_{t+1}}ig[r(S_t,A_t) + \gamma q(S_{t+1},\pi(S_{t+1};oldsymbol{ heta}))ig]$$

and $\pi(s; \theta)$ according to the deterministic policy gradient theorem.

The algorithm was first described in the paper Continuous Control with Deep Reinforcement Learning by Timothy P. Lillicrap et al. (2015).

The authors utilize a replay buffer, a target network (updated by exponential moving average with $\tau = 0.001$), batch normalization for CNNs, and perform exploration by adding a Ornstein-Uhlenbeck noise to the predicted actions. Training is performed by Adam with learning rates of 1e-4 and 1e-3 for the policy and the critic networks, respectively.



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 Rfor 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 RSample a random minbatch of N transitions (s_i, a_i, r_i, s_{i+1}) from RSet $y_i = r_i + \gamma Q'(s_{i+1}, \mu'(s_{i+1}|\theta^{\mu'})|\theta^{Q'})$ 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:

$$\nabla_{\theta^{\mu}} J \approx \frac{1}{N} \sum_{i} \nabla_{a} Q(s, a | \theta^{Q}) |_{s=s_{i}, a=\mu(s_{i})} \nabla_{\theta^{\mu}} \mu(s | \theta^{\mu}) |_{s=s_{i}, a=\mu(s_{i$$

Update the target networks:

$$\begin{split} \theta^{Q'} &\leftarrow \tau \theta^Q + (1-\tau) \theta^{Q'} \\ \theta^{\mu'} &\leftarrow \tau \theta^\mu + (1-\tau) \theta^{\mu'} \end{split}$$

MuJoCo

DDPG

DPG

end for end for

ContinuousActionSpace

Algorithm 1 of "Continuous Control with Deep Reinforcement Learning" by Timothy P. Lillicrap et al.

NPFL139, Lecture 8

TD3 SAC SPE

SPI



Figure 2: Performance curves for a selection of domains using variants of DPG: original DPG algorithm (minibatch NFQCA) with batch normalization (light grey), with target network (dark grey), with target networks and batch normalization (green), with target networks from pixel-only inputs (blue). Target networks are crucial.

Figure 3 of "Continuous Control with Deep Reinforcement Learning" by Timothy P. Lillicrap et al.



Results using low-dimensional (*lowd*) version of the environment, pixel representation (*pix*) and DPG reference (*cntrl*).

The architecture in the *lowd* case consists of two hidden layers with 400 and 300 units and ReLU activation, in both the actor and the critic. The actor additionally uses tanh activation to bound the action in a given range.

In the case of pixel representation, 3 convolution layers with 32 channels and ReLU activation are used (no pooling), followed by two fully-connected ReLUactivated layers with 200 units each.

environment	$R_{av,lowd}$	$R_{best,lowd}$	$R_{av,pix}$	$R_{best,pix}$	$R_{av,cntrl}$	$R_{best,cntrl}$
blockworld1	1.156	1.511	0.466	1.299	-0.080	1.260
blockworld3da	0.340	0.705	0.889	2.225	-0.139	0.658
canada	0.303	1.735	0.176	0.688	0.125	1.157
canada2d	0.400	0.978	-0.285	0.119	-0.045	0.701
cart	0.938	1.336	1.096	1.258	0.343	1.216
cartpole	0.844	1.115	0.482	1.138	0.244	0.755
cartpoleBalance	0.951	1.000	0.335	0.996	-0.468	0.528
cartpoleParallelDouble	0.549	0.900	0.188	0.323	0.197	0.572
cartpoleSerialDouble	0.272	0.719	0.195	0.642	0.143	0.701
cartpoleSerialTriple	0.736	0.946	0.412	0.427	0.583	0.942
cheetah	0.903	1.206	0.457	0.792	-0.008	0.425
fixedReacher	0.849	1.021	0.693	0.981	0.259	0.927
fixedReacherDouble	0.924	0.996	0.872	0.943	0.290	0.995
fixedReacherSingle	0.954	1.000	0.827	0.995	0.620	0.999
gripper	0.655	0.972	0.406	0.790	0.461	0.816
gripperRandom	0.618	0.937	0.082	0.791	0.557	0.808
hardCheetah	1.311	1.990	1.204	1.431	-0.031	1.411
hopper	0.676	0.936	0.112	0.924	0.078	0.917
hyq	0.416	0.722	0.234	0.672	0.198	0.618
movingGripper	0.474	0.936	0.480	0.644	0.416	0.805
pendulum	0.946	1.021	0.663	1.055	0.099	0.951
reacher	0.720	0.987	0.194	0.878	0.231	0.953
reacher3daFixedTarget	0.585	0.943	0.453	0.922	0.204	0.631
reacher3daRandomTarget	0.467	0.739	0.374	0.735	-0.046	0.158
reacherSingle	0.981	1.102	1.000	1.083	1.010	1.083
walker2d	0.705	1.573	0.944	1.476	0.393	1.397
torcs	-393.385	1840.036	-401.911	1876.284	-911.034	1961.600

Table 1 of "Continuous Control with Deep Reinforcement Learning" by Timothy P. Lillicrap et al.

DPG DDPG

MuJoCo TD3

SPE SPI

SAC

Ornstein-Uhlenbeck Exploration

Ú F_AL

16/52

While it is natural to use Gaussian noise for the exploration policy, the authors claim that temporally-correlated noise is more effective for physical control problems with inertia.

They therefore generate noise using Ornstein-Uhlenbeck process, by computing

$$n_t \leftarrow n_{t-1} + heta \cdot (\mu - n_{t-1}) + arepsilon \sim \mathcal{N}(0,\sigma^2),$$

utilizing hyperparameter values heta=0.15 and $\sigma=0.2$.

Ornstein-Uhlenbeck Exploration Visualization



https://upload.wikimedia.org/wikipedia/commons/7/79/Wiener-process-5traces.svg



SPL

- On the left, there is a continuous *Wiener process* (a "brownian path"), corresponding to $\theta = 0$ and $\sigma = 1$.
- On the right, there is Ornstein-Uhlenbeck process example with $heta=\sigma=1$ and $\mu=0$.

The gray are corresponds to the standard deviation of x (n in our notation).

TD3

SAC

SPE

SAC Algorithm



MuJoCo

NPFL139, Lecture 8

ContinuousActionSpace

ce DPG DDPG

MuJoCo TD3

SAC SPE

SPI

SAC Algorithm

MuJoCo





Figure 4. Example MuJoCo environments (a) HalfCheetah-v1, (b) Hopper-v1, (c) Walker2d-v1, (d) Ant-v1.

Figure 4 of "Addressing Function Approximation Error in Actor-Critic Methods" by Scott Fujimoto et al.

SPE

SPL

See the Gymnasium documentation of the <u>HalfCheetah</u>, <u>Hopper</u>, <u>Walker2D</u>, <u>Ant</u>, <u>Humanoid</u> environments for a detailed description of observation spaces and action spaces.

NPFL139, Lecture 8

TD3

SAC

SAC Algorithm



Twin Delayed Deep Deterministic Policy Gradient (TD3)

NPFL139, Lecture 8

ContinuousActionSpace

ace DPG

DDPG MuJoCo

oCo TD3

SAC

SPE SPI

SAC Algorithm

Twin Delayed Deep Deterministic Policy Gradient



The paper Addressing Function Approximation Error in Actor-Critic Methods by Scott Fujimoto et al. from February 2018 proposes improvements to DDPG which

- decrease maximization bias by training two critics and choosing the minimum of their predictions;
- introduce several variance-lowering optimizations:
 - delayed policy updates;
 - $^{\circ}\,$ target policy smoothing.

In 2022, together with the SAC algorithm, the TD3 algorithm has been one of the best algorithms for off-policy continuous-actions RL training.

DPG DDPG

MuJoCo TD3

SPE

SPI

TD3 – Maximization Bias

Ú_F≩L

Similarly to Q-learning, the DDPG algorithm suffers from maximization bias. In Q-learning, the maximization bias was caused by the explicit max operator. For DDPG methods, it can be caused by the gradient descent itself. Let θ_{approx} be the parameters maximizing the q_{θ} and let θ_{true} be the hypothetical parameters which maximise true q_{π} , and let π_{approx} and π_{true} denote the corresponding policies.

Because the gradient direction is a local maximizer, for sufficiently small $lpha < arepsilon_1$ we have

$$\mathbb{E}ig[q_{oldsymbol{ heta}}(s,\pi_{approx})ig] \geq \mathbb{E}ig[q_{oldsymbol{ heta}}(s,\pi_{true})ig].$$

However, for real q_{π} and for sufficiently small $lpha < arepsilon_2$, it holds that

$$\mathbb{E}ig[q_{\pi}(s,\pi_{true})ig] \geq \mathbb{E}ig[q_{\pi}(s,\pi_{approx})ig].$$

Therefore, if $\mathbb{E}[q_{\theta}(s, \pi_{true})] \geq \mathbb{E}[q_{\pi}(s, \pi_{true})]$, for $\alpha < \min(\varepsilon_1, \varepsilon_2)$

$$\mathbb{E}ig[q_{oldsymbol{ heta}}(s,\pi_{approx})ig] \geq \mathbb{E}ig[q_{\pi}(s,\pi_{approx})ig].$$

NPFL139, Lecture 8

ContinuousActionSpace DPG

DDPG MuJoCo

DCO TD3 SAC

SPE

SPI

SAC Algorithm

TD3 – Maximization Bias



Analogously to Double DQN we could compute the learning targets using the current policy and the target critic, i.e., $r + \gamma q_{\theta'}(s', \pi_{\varphi(s')})$ (instead of using the target policy and the target critic as in DDPG), obtaining DDQN-AC algorithm. However, the authors found out that the policy changes too slowly and the target and current networks are too similar.

Using the original Double Q-learning, two pairs of actors and critics could be used, with the learning targets computed by the opposite critic, i.e., $r + \gamma q_{\theta_2}(s', \pi_{\varphi_1}(s'))$ for updating q_{θ_1} . The resulting DQ-AC algorithm is slightly better, but still suffering from overestimation.

NPFL139, Lecture 8

TD3 – Algorithm



The authors instead suggest to employ two critics and one actor. The actor is trained using one of the critics, and both critics are trained using the same target computed using the *minimum* value of both critics as

$$r+\gamma\min_{i=1,2}q_{oldsymbol{ heta}_i}(s',\pi_{oldsymbol{arphi}'}(s')).$$

The resulting algorithm is called CDQ – Clipped Double Q-learning.

Furthermore, the authors suggest two additional improvements for variance reduction.

- For obtaining higher quality target values, the authors propose to train the critics more often. Therefore, critics are updated each step, but the actor and the target networks are updated only every d-th step (d = 2 is used in the paper).
- To explicitly model that similar actions should lead to similar results, a small random noise is added to the performed actions when computing the target value:

$$r+\gamma \min_{i=1,2} q_{oldsymbol{ heta}_i}(s',\pi_{oldsymbol{arphi}'}(s')+arepsilon) ~~ ext{ for }~~arepsilon\sim ext{clip}(\mathcal{N}(0,\sigma),-c,c), ext{ with } \sigma=0.2, c=0.5.$$

TD3 – Algorithm



Algorithm 1 TD3

Initialize critic networks $Q_{\theta_1}, Q_{\theta_2}$, and actor network π_{ϕ} with random parameters θ_1, θ_2, ϕ Initialize target networks $\theta'_1 \leftarrow \theta_1, \theta'_2 \leftarrow \theta_2, \phi' \leftarrow \phi$ Initialize replay buffer \mathcal{B} for t = 1 to T do Select action with exploration noise $a \sim \pi_{\phi}(s) + \epsilon$, $\epsilon \sim \mathcal{N}(0, \sigma)$ and observe reward r and new state s'Store transition tuple (s, a, r, s') in \mathcal{B}

Sample mini-batch of N transitions (s, a, r, s') from \mathcal{B} $\tilde{a} \leftarrow \pi_{\phi'}(s') + \epsilon$, $\epsilon \sim \operatorname{clip}(\mathcal{N}(0, \tilde{\sigma}), -c, c)$ $y \leftarrow r + \gamma \min_{i=1,2} Q_{\theta'_i}(s', \tilde{a})$ Update critics $\theta_i \leftarrow \operatorname{argmin}_{\theta_i} N^{-1} \sum (y - Q_{\theta_i}(s, a))^2$ if $t \mod d$ then Update ϕ by the deterministic policy gradient: $\nabla_{\phi} J(\phi) = N^{-1} \sum \nabla_a Q_{\theta_1}(s, a)|_{a = \pi_{\phi}(s)} \nabla_{\phi} \pi_{\phi}(s)$

Update target networks:

$$\theta'_i \leftarrow \tau \theta_i + (1 - \tau) \theta'_i$$

 $\phi' \leftarrow \tau \phi + (1 - \tau) \phi'$
d if

end if end for

Algorithm 1 of "Addressing Function Approximation Error in Actor-Critic Methods" by Scott Fujimoto et al.

NPFL139, Lecture 8

ContinuousActionSpace DPG DDPG MuJoCo TD3 SAC

SAC Algorithm

SPE

SPL

TD3 – **Hyperparameters**

ÚF	AL

Hyper-parameter	Ours	DDPG
Critic Learning Rate	10^{-3}	10^{-3}
Critic Regularization	None	$10^{-2} \cdot \theta ^2$
Actor Learning Rate	10^{-3}	10^{-4}
Actor Regularization	None	None
Optimizer	Adam	Adam
Target Update Rate (τ)	$5 \cdot 10^{-3}$	10^{-3}
Batch Size	100	64
Iterations per time step	1	1
Discount Factor	0.99	0.99
Reward Scaling	1.0	1.0
Normalized Observations	False	True
Gradient Clipping	False	False
Exploration Policy	$\mathcal{N}(0, 0.1)$	OU, $\theta = 0.15, \mu = 0, \sigma = 0.2$

Table 3 of "Addressing Function Approximation Error in Actor-Critic Methods" by Scott Fujimoto et al.

In TD3, the actor and the critic also use two fully-connected ReLU-activated layers with 400 and 300 units, respectively. The actor actions are bounded in a given range using a suitably scaled tanh activation.

In TD3, the authors state that they also tried the Ornstein-Uhlenbeck noise, but it provided no benefit compared to $\mathcal{N}(0, 0.1)$.

26/52

SAC Algorithm

SPL

TD3 – **Results**



Figure 5 of "Addressing Function Approximation Error in Actor-Critic Methods" by Scott Fujimoto et al.

Environment	TD3	DDPG	Our DDPG	PPO	TRPO	ACKTR	SAC
HalfCheetah	$\textbf{9636.95} \pm \textbf{859.065}$	3305.60	8577.29	1795.43	-15.57	1450.46	2347.19
Hopper	$\textbf{3564.07} \pm \textbf{114.74}$	2020.46	1860.02	2164.70	2471.30	2428.39	2996.66
Walker2d	$\textbf{4682.82} \pm \textbf{539.64}$	1843.85	3098.11	3317.69	2321.47	1216.70	1283.67
Ant	$\textbf{4372.44} \pm \textbf{1000.33}$	1005.30	888.77	1083.20	-75.85	1821.94	655.35
Reacher	$\textbf{-3.60}\pm\textbf{0.56}$	-6.51	-4.01	-6.18	-111.43	-4.26	-4.44
InvPendulum	$\textbf{1000.00} \pm \textbf{0.00}$	1000.00	1000.00	1000.00	985.40	1000.00	1000.00
InvDoublePendulum	$\textbf{9337.47} \pm \textbf{14.96}$	9355.52	8369.95	8977.94	205.85	9081.92	8487.15

Table 1 of "Addressing Function Approximation Error in Actor-Critic Methods" by Scott Fujimoto et al.

SAC

SPE

SPI

Co TD3

TD3 – **Ablations**





The AHE is the authors' reimplementation of DDPG using updated architecture, hyperparameters, and exploration. TPS is Target Policy Smoothing, DP is Delayed Policy update, and CDQ is Clipped Double Q-learning.

NPFL139, Lecture 8

ContinuousActionSpace DPG DDPG MuJoCo TD3 SAC SPE SPI SAC Algorithm

TD3 – **Ablations**



Method	HCheetah	Hopper	Walker2d	Ant
TD3	9532.99	3304.75	4565.24	4185.06
DDPG	3162.50	1731.94	1520.90	816.35
AHE	8401.02	1061.77	2362.13	564.07
AHE + DP	7588.64	1465.11	2459.53	896.13
AHE + TPS	9023.40	907.56	2961.36	872.17
AHE + CDQ	6470.20	1134.14	3979.21	3818.71
TD3 - DP	9590.65	2407.42	4695.50	3754.26
TD3 - TPS	8987.69	2392.59	4033.67	4155.24
TD3 - CDQ	9792.80	1837.32	2579.39	849.75
DQ-AC	9433.87	1773.71	3100.45	2445.97
DDQN-AC	10306.90	2155.75	3116.81	1092.18

Table 2 of "Addressing Function Approximation Error in Actor-Critic Methods" by Scott Fujimoto et al.

NPFL139, Lecture 8

ContinuousActionSpace

DPG DDPG

MuJoCo TD3

SAC

SPE

SPI SA



NPFL139, Lecture 8

ContinuousActionSpace

e DPG DDPG

MuJoCo TD3

SAC SPE

SPI

SAC Algorithm



The paper *Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor* by Tuomas Haarnoja et al. from Jan 2018 introduces a different off-policy algorithm for continuous action space.

It was followed by a continuation paper *Soft Actor-Critic Algorithms and Applications* in Dec 2018.

The general idea is to introduce entropy directly in the value function we want to maximize, instead of just ad-hoc adding the entropy penalty. Such an approach is an instance of *regularized policy optimization*.

DPG DDPG

MuJoCo TD3

SAC SPE

SPI

SAC Algorithm

Soft Actor Critic Objective

32/52

Until now, our goal was to optimize

 $\mathbb{E}_{\pi}[G_0].$

Assume the rewards are deterministic and that μ_{π} is on-policy distribution of a policy π .

In the soft actor-critic, the authors consider infinite-horizon MDPs and propose to optimize the maximum entropy objective

$$egin{aligned} \pi_* &= rg\max_{\pi} \mathbb{E}_{s \sim \mu_{\pi}} \Big[\mathbb{E}_{a \sim \pi(s)} ig[r(s,a) ig] + lpha H(\pi(\cdot|s)) \Big] \ &= rg\max_{\pi} \mathbb{E}_{s \sim \mu_{\pi}, a \sim \pi(s)} ig[r(s,a) - lpha \log \pi(a|s) ig]. \end{aligned}$$

Note that the value of α is dependent on the magnitude of returns and that for a fixed policy, the entropy penalty can be "hidden" in the reward.

NPFL139, Lecture 8

Soft Actor Critic Objective



To maximize the regularized objective, we define the following augmented reward:

$$r_{\pi}(s,a) \stackrel{ ext{\tiny def}}{=} r(s,a) + \mathbb{E}_{s' \sim p(s,a)}ig[lpha H(\pi(\cdot|s'))ig].$$

From now on, we consider **soft action-value** function corresponding to this augmented reward.

DPG DDPG

MuJoCo TD3

SAC

SPE

SPI

SAC Algorithm



Soft Policy Evaluation

NPFL139, Lecture 8

ContinuousActionSpace

ce DPG

DDPG MuJoCo

SPE

SPI

SAC

TD3

SAC Algorithm

Soft Policy Evaluation



Our goal is now to derive **soft policy iteration**, an analogue of policy iteration algorithm. We start by considering soft policy evaluation. Let a modified Bellman backup operator \mathcal{T}_{π} be defined as

$$\mathcal{T}_{\pi}q(s,a) \stackrel{ ext{\tiny def}}{=} r(s,a) + \gamma \mathbb{E}_{s' \sim p(s,a)}ig[v(s')ig],$$

where the **soft (state-)value** function v(s) is defined as

$$v(s) = \mathbb{E}_{a \sim \pi}ig[q(s,a)ig] + lpha H(\pi(\cdot|s)) = \mathbb{E}_{a \sim \pi}ig[q(s,a) - lpha \log \pi(a|s)ig].$$

This modified Bellman backup operator corresponds to the usual one for the augmented rewards $r_{\pi}(s, a)$, and therefore the repeated application $\mathcal{T}_{\pi}^{k}q$ converges to q_{π} according to the original proof.



Soft Policy Improvement

NPFL139, Lecture 8

ContinuousActionSpace

e DPG DDPG

MuJoCo TD3

SAC SPE

SPI

SAC Algorithm

Soft Policy Improvement

While the soft policy evaluation was a straightforward modification of the original policy evaluation, the soft policy improvement is quite different.

Assume we have a policy π , its action-value function q_{π} from the soft policy evaluation, and we want to improve the policy. Furthermore, we should select the improved policy from a family of parametrized distributions Π .

We define the improved policy π' as

$$\pi'(\cdot|s) \stackrel{ ext{\tiny def}}{=} rgmin_{ar{\pi}\in\Pi} J_{\pi}(ar{\pi}) \stackrel{ ext{\tiny def}}{=} rgmin_{ar{\pi}\in\Pi} D_{ ext{KL}}igg(ar{\pi}(\cdot|s) \Bigg\| rac{\exp\left(rac{1}{lpha} q_{\pi}(s,\cdot)
ight)}{z_{\pi}(s)}igg),$$

where $z_{\pi}(s)$ is the partition function (i.e., normalization factor such that the right-hand side is a distribution), which does not depend on the new policy and thus can be ignored.

Soft Policy Improvement



We now prove that $q_{\pi'}(s,a) \geq q_{\pi}(s,a)$ for any state s and action a.

We start by noting that $J_{\pi}(\pi') \leq J_{\pi}(\pi)$, because we can always choose π as the improved policy. Therefore,

$$\mathbb{E}_{a \sim \pi'}ig[lpha \log \pi'(a|s) - q_{\pi}(s,a) + lpha \log z_{\pi}(s)ig] \leq \mathbb{E}_{a \sim \pi}ig[lpha \log \pi(a|s) - q_{\pi}(s,a) + lpha \log z_{\pi}(s)ig],$$

which results in

$$\mathbb{E}_{a \sim \pi'}ig[q_{\pi}(s,a) - lpha \log \pi'(a|s)ig] \geq v_{\pi}(s).$$

We now finish the proof analogously to the original one:

. . .

$$egin{aligned} q_\pi(s,a) &= r(s,a) + \gamma \mathbb{E}_{s'}[v_\pi(s')] \ &\leq r(s,a) + \gamma \mathbb{E}_{s'}[\mathbb{E}_{a' \sim \pi'}[q_\pi(s',a') - lpha \log \pi'(a'|s')] \end{aligned}$$

$$\leq q_{\pi'}(s,a).$$

NPFL139, Lecture 8

ContinuousActionSpace DPG

DDPG

MuJoCo TD3 SAC

SAC Algorithm

SPI

SPE

Soft Policy Iteration

Ú_F≩L

The soft policy iteration algorithm alternates between the soft policy evaluation and soft policy improvement steps.

The repeated application of these two steps produce better and better policies. In other words, we get a monotonically increasing sequence of soft action-value functions.

If the soft action-value function is bounded (the paper assumes a bounded reward and a finite number of actions to bound the entropy), the repeated application converges to some q_* , from which we get a π_* using the soft policy improvement step.

It remains to show that the π_* is indeed the optimal policy fulfilling $q_{\pi_*}(s, a) \ge q_{\pi}(s, a)$. However, this follows from the fact that at convergence, $J_{\pi_*}(\pi_*) \le J_{\pi_*}(\pi)$, and following the

same reasoning as in the proof of the soft policy improvement, we obtain the required $q_{\pi_*}(s,a) \geq q_{\pi}(s,a)$.

NPFL139, Lecture 8

Soft Policy Improvement Derivation

The following derivation is not in the original paper, but it is my understanding of how the softmax of the action-value function arises. For simplicity, we assume finite number of actions, but the same approach can be generalized to continuous actions.

Assuming we have a policy π and its action-value function q_{π} , we usually improve the policy using

$$egin{aligned}
u(\cdot|s) &= rg\max_{
u} \mathbb{E}_{a \sim
u(\cdot|s)}ig[q_{\pi}(s,a)ig] \ &= rg\max_{
u} \sum_{a} q_{\pi}(s,a)
u(a|s) \ &= rg\max_{
u} oldsymbol{q}_{\pi}(s,\cdot)^T oldsymbol{
u}(\cdot|s), \end{aligned}$$

which results in a greedy improvement with the form of

$$u(s) = rg\max_a q_\pi(s,a).$$

NPFL139, Lecture 8

MuJoCo 7

DDPG

Co TD3 SAC

SAC Algorithm

SPI

SPE



Soft Policy Improvement Derivation



Now consider instead the regularized objective

$$egin{split}
u(\cdot|s) &= rg\max_{
u} ig(\mathbb{E}_{a \sim
u(\cdot|s)}ig[q_{\pi}(s,a)ig] + lpha H(
u(\cdot|s))ig) \ &= rg\max_{
u} ig(\mathbb{E}_{a \sim
u}ig[q_{\pi}(s,a) - lpha \log
u(a|s)ig]ig) \end{split}$$

To maximize it for a given s, we form a Lagrangian

$$\mathcal{L} = \Big(\sum_a
u(a|s)ig(q_\pi(s,a) - lpha \log
u(a|s)ig)ig) - \lambda \Big(1 - \sum_a
u(a|s)\Big).$$

The derivative with respect to u(a|s) is

$$rac{\partial \mathcal{L}}{\partial
u(a|s)} = q_\pi(s,a) - lpha \log
u(a|s) - lpha + \lambda.$$

Setting it to zero, we get $lpha \log
u(a|s) = q_\pi(s,a) + \lambda - lpha$, resulting in $u(a|s) \propto e^{rac{1}{lpha} q_\pi(s,a)}$.

ContinuousActionSpace DPG DDPG MuJoCo TD3 SAC

41/52

SPI

SPE



Soft Actor Critic Algorithm

NPFL139, Lecture 8

ContinuousActionSpace

ce DPG

DDPG

MuJoCo TD3

SPE SPI

SAC

SAC Algorithm

Soft Actor Critic Algorithm

Ú F_AL

Our soft actor critic will be an off-policy algorithm with continuous action space. The model consist of two critics q_{θ_1} and q_{θ_2} , two target critics $q_{\bar{\theta}_1}$ and $q_{\bar{\theta}_2}$, and a single actor π_{φ} .

The authors state that

- with a single critic, all the described experiments still converge;
- they adopted the two critics from the TD3 paper;
- using two critics "significantly speed up training".

DPG DDPG

MuJoCo TD3

SAC

SPE

SPI

SAC Algorithm

Soft Actor Critic – Critic Training

To train the critic, we use the modified Bellman backup operator, resulting in the loss

$$J_q(oldsymbol{ heta}_i) = \mathbb{E}_{s \sim \mu_\pi, a \sim \pi_arphi(s)} \Big[ig(q_{oldsymbol{ heta}_i}(s,a) - ig(r(s,a) + \gamma \mathbb{E}_{s' \sim p(s,a)}[v_{\min}(s')] ig) ig)^2 \Big],$$

where

$$v_{\min}(s) = \mathbb{E}_{a \sim \pi_{oldsymbol{arphi}}(s)} \Big[\min_i ig(q_{oldsymbol{ar{ heta}}_i}(s,a) ig) - lpha \log \pi_{oldsymbol{arphi}}(a|s) \Big].$$

The target critics are updated using exponential moving averages with momentum τ .

NPFL139, Lecture 8 ContinuousActionSpace DPG DDPG MuJoCo TD3 SAC SPE SPI SAC Algorithm

Soft Actor Critic – Actor Training

The actor is updated by directly minimizing the KL divergence, resulting in the loss

$$J_{\pi}(oldsymbol{arphi}) = \mathbb{E}_{s \sim \mu_{\pi}, a \sim \pi_{oldsymbol{arphi}}(s)} \Big[lpha \logig(\pi_{oldsymbol{arphi}}(a,s)ig) - \min_iig(q_{oldsymbol{ heta}_i}(s,a)ig) \Big].$$

Given that our critics are differentiable, in order to be able to compute the gradient $\nabla_{\varphi} q_{\theta_i}(s, a)$, we only need to reprametrize the policy as

$$a=f_{oldsymbol{arphi}}(s,arepsilon).$$

Specifically, we sample $arepsilon \sim \mathcal{N}(0,1)$ and let $f_{m{arphi}}$ produce an unbounded Gaussian distribution $\mathcal{N}(\mu(s; oldsymbol{arphi}), \sigma(s; oldsymbol{arphi}))$, or a diagonal one if the actions are vectors, with the sampled action $a=\mu(s;oldsymbol{arphi})+arepsilon\sigma(s;oldsymbol{arphi}).$

Together, we obtain

$$J_{\pi}(oldsymbol{arphi}) = \mathbb{E}_{s \sim \mu_{\pi}, arepsilon \sim \mathcal{N}(0,1)} \Big[lpha \log ig(\pi_{oldsymbol{arphi}}(s,arepsilon), s) ig) - \min_i ig(q_{oldsymbol{ heta}_i}(s, f_{oldsymbol{arphi}}(s, arepsilon)) ig) \Big].$$

SPE

SPI

Soft Actor Critic – Bounding Actions

In practice, the actions need to be bounded.

The authors propose to apply an invertible squashing function tanh on the unbounded Gaussian distribution.

Consider that our policy produces an unbounded action $\pi(u|s)$. To define a distribution $\overline{\pi}(a|s)$ with $a = \tanh(u)$, we need to employ the change of variables, resulting in

$$ar{\pi}(a|s) = \pi(u|s)igg(rac{\partial a}{\partial u}igg)^{-1} = \pi(u|s)igg(rac{\partial anh(u)}{\partial u}igg)^{-1}.$$

Therefore, the log-likelihood has quite a simple form

ContinuousActionSpace

$$\log ar{\pi}(a|s) = \log \pi(u|s) - \logig(1- anh^2(u)ig).$$

NPFL139, Lecture 8

SAC Algorithm

SPL

Soft Actor Critic – Automatic Entropy Adjustment

One of the most important hyperparameters is the entropy penalty α .

In the second paper, the authors presented an algorithm for automatic adjustment of its value. Instead of setting the entropy penalty α , they propose to specify target entropy value \mathcal{H} and then solve a constrained optimization problem

$$\pi_* = rg\max_{\pi} \mathbb{E}_{s \sim \mu_{\pi}, a \sim \pi(s)} ig[r(s, a) ig] ext{ such that } \mathbb{E}_{s \sim \mu_{\pi}, a \sim \pi(s)} ig[-\log \pi(a|s) ig] \geq \mathcal{H}.$$

We can then form a Lagrangian with a multiplier lpha

$$\mathbb{E}_{s\sim \mu_{\pi}, a\sim \pi(s)} \Big[r(s,a) + lpha ig(-\log \pi(a|s) - \mathcal{H} ig) \Big],$$

which should be maximized with respect to π and minimized with respect to $\alpha \geq 0$.

NPFL139, Lecture 8

ContinuousActionSpace DPG DDPG MuJoCo TD3 SAC SPE SPI SAC Algorithm

Soft Actor Critic – Automatic Entropy Adjustment

^ÚF_AL

To optimize the Lagrangian, we perform *dual gradient descent*, where we alternate between maximization with respect to π and minimization with respect to α .

While such a procedure is guaranteed to converge only under the convexity assumptions, the authors report that the dual gradient descent works in practice also with nonlinear function approximation.

To conclude, the automatic entropy adjustment is performed by introducing a final loss

$$J(lpha) = \mathbb{E}_{s \sim \mu_{\pi}, a \sim \pi(s)}ig[- lpha \log \pi(a|s) - lpha \mathcal{H}ig].$$



Algorithm 1 Soft Actor-Critic

Input: θ_1, θ_2, ϕ $\theta_1 \leftarrow \theta_1, \theta_2 \leftarrow \theta_2$ $\mathcal{D} \leftarrow \emptyset$ for each iteration do for each environment step do $\mathbf{a}_t \sim \pi_{\phi}(\mathbf{a}_t | \mathbf{s}_t)$ $\mathbf{s}_{t+1} \sim p(\mathbf{s}_{t+1}|\mathbf{s}_t,\mathbf{a}_t)$ $\mathcal{D} \leftarrow \mathcal{D} \cup \{(\mathbf{s}_t, \mathbf{a}_t, r(\mathbf{s}_t, \mathbf{a}_t), \mathbf{s}_{t+1})\}$ end for for each gradient step do $\theta_i \leftarrow \theta_i - \lambda_Q \hat{\nabla}_{\theta_i} J_Q(\theta_i) \text{ for } i \in \{1, 2\}$ $\phi \leftarrow \phi - \lambda_{\pi} \nabla_{\phi} J_{\pi}(\phi)$ $\alpha \leftarrow \alpha - \lambda \hat{\nabla}_{\alpha} J(\alpha)$ $\theta_i \leftarrow \tau \theta_i + (1 - \tau) \overline{\theta}_i \text{ for } i \in \{1, 2\}$ end for end for **Output:** θ_1, θ_2, ϕ

Initial parameters
 Initialize target network weights
 Initialize an empty replay pool

Sample action from the policy
 Sample transition from the environment
 Store the transition in the replay pool

Update the Q-function parameters
 Update policy weights
 Adjust temperature
 Update target network weights

▷ Optimized parameters

Algorithm 1 of "Soft Actor-Critic Algorithms and Applications" by Tuomas Haarnoja et al.

SPE

SPL

NPFL139, Lecture 8

MuJoCo

DDPG

TD3 SAC



Table 1: SAC Hyperparameters

Parameter	Value
optimizer	Adam (Kingma & Ba, 2015)
learning rate	$3 \cdot 10^{-4}$
discount (γ)	0.99
replay buffer size	10^{6}
number of hidden layers (all networks)	2
number of hidden units per layer	256
number of samples per minibatch	256
entropy target	$-\dim(\mathcal{A})$ (e.g., -6 for HalfCheetah-v1)
nonlinearity	ReLU
target smoothing coefficient (τ)	0.005
target update interval	1
gradient steps	1

Table 1 of "Soft Actor-Critic Algorithms and Applications" by Tuomas Haarnoja et al.

SPI

SPE

MuJoCo TD3

SAC



Figure 1: Training curves on continuous control benchmarks. Soft actor-critic (blue and yellow) performs consistently across all tasks and outperforming both on-policy and off-policy methods in the most challenging tasks.

Figure 1 of "Soft Actor-Critic Algorithms and Applications" by Tuomas Haarnoja et al.

SPI

SPE

SAC

ContinuousActionSpace DPG DDPG

MuJoCo

TD3





Figure 3. Sensitivity of soft actor-critic to selected hyperparameters on Ant-v1 task. (a) Evaluating the policy using the mean action generally results in a higher return. Note that the policy is trained to maximize also the entropy, and the mean action does not, in general, correspond the optimal action for the maximum return objective. (b) Soft actor-critic is sensitive to reward scaling since it is related to the temperature of the optimal policy. The optimal reward scale varies between environments, and should be tuned for each task separately. (c) Target value smoothing coefficient τ is used to stabilize training. Fast moving target (large τ) can result in instabilities (red), whereas slow moving target (small τ) makes training slower (blue).

Figure 3 of "Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor" by Tuomas Haarnoja et al.

NPFL139, Lecture 8 ContinuousActionSpace DPG DDPG MuJoCo TD3 SAC SPE SPI SAC Algorithm