

Deep Deterministic Policy Gradients: Components and Extensions

Yannik Frisch · Tabea Wilke ·
Maximilian Gehrke

Received: date / Accepted: date

Abstract TODO

Keywords DDPG · DQN · DPG

1 Introduction

The field of Reinforcement Learning deals with solving problems that are accessible through the interaction of an agent with its environment. Such problems can be defined as Markov Decision Processes [x], which consist of a tuple (S, A, R, P, γ) , where S is the state-distribution, A is the action-distribution, $R : S, A \rightarrow r$ is the reward function mapping states and actions to a scalar reward r , $P : S, A \rightarrow S$ is the state transition function mapping states and actions to states, and γ is the discount factor used to make an agent more or less farsighted.

For many applications the environment details, i.e. R and P , are not available. This requires the use of so called model-free algorithms. An early one was Q-Learning [x], which updates an internal representation of the action-value function $Q(s, a)$ by the temporal-difference error [EDIT: EXPLAIN BETTER]. The internal representation of this function is not tractable for large state-action spaces. This problem is addressed by value function methods, e.g. the DQN-Algorithm (Mnih et al. 2013), which is an adaption to Q-Learning, where the action-value function is approximated with deep neural networks. Instead of approximating the value-function, one could also approximate the

F. Author
first address
Tel.: +123-45-678910
Fax: +123-45-678910
E-mail: fauthor@example.com

S. Author
second address

policy $\pi(s|a)$ directly. [EDIT: BRIEFLY DESCRIBE DPG]

Actor critic methods (Konda and Tsitsiklis 2000) approximate the value-function as well as the policy. Finally, the Deep Deterministic Policy Gradient (DDPG) approach (Lillicrap et al. 2015) combines the above mentioned methods to an actor critic algorithm using neural network function approximation for the policy and the action-value function, which learns a deterministic policy.

We will give more detailed insights into the algorithms and how DDPG evolved from them in the next section, before we describe some possible extensions to it in section 4.4.

2 Preliminaries

The general goal of an reinforcement learning algorithm is to find an optimal behavior policy $\pi(a|s)$, or $\pi(s)$ in the deterministic case, which maximizes the expected total reward an agent collects while following it. An optimal policy can be defined by

$$\pi^*(a|s) = \max_a Q^*(s, a) = E \left[\sum_{t=0}^T \gamma^t r_t \right]$$

Where t is the current time-step and T the final time-step ending an episode. This policy is just greedily choosing the action maximizing the optimal action-value function $Q^*(s, a)$. By definition this optimal value function yields the bellman equation (Sutton and Barto 2018) and can be reinterpreted as maximizing the current reward and the discounted action-value of the resulting state. In formula this gives:

$$Q^*(s, a) = E_{s' \sim P(s, a), r \sim R(s, a)} \left[r + \gamma \max_{a'} Q^*(s', a') | s, a \right]$$

2.1 Q-Learning

The Q-Learning approach [x] calculates the optimal value function by ...

2.2 Deep Q-Learning

The Deep Q-Network approach (DQN) (Mnih et al. 2013) combines the approximation power of neural networks with the traditional Q-learning. The algorithm is an off-policy, model-free approach and is able to find a close to optimal action-value function for many cases [x] and from this a close to optimal policy. For approximating the action-value function $Q(s, a|\theta) \approx Q(s, a)$

the approach uses a deep neural network with parameters θ , called the Q-Network. The Q-Network can be trained by sequentially minimizing the loss function $L_i(\theta)$, depending on the parameters

$$L^{(i)}(\theta) = E_{s, s' \sim \rho^\pi, a \sim \pi^\theta} \left[\left(r + \gamma \max_{a'} Q(s', a' | \theta^{(i-1)}) - Q(s, a | \theta^{(i)}) \right)^2 \right]$$

This loss function is similar to the classical temporal-difference loss used in Q-Learning, but with approximated action-value functions instead of lookup-tables. Derivating this loss w.r.t. the approximation's weights gives:

$$\nabla_{\theta^{(i)}} L^{(i)}(\theta) = E_{s, s' \sim \rho^\pi, a \sim \pi^\theta} \left[\left(r + \gamma \max_{a'} Q(s', a' | \theta^{(i-1)}) - Q(s, a | \theta^{(i)}) \right) \nabla_{\theta^{(i)}} Q(s, a | \theta^{(i)}) \right]$$

The expectation can be approximated by sampling from an environment and this gradient can be used to optimize the loss function by using stochastic gradient descent.

Furthermore, a replay buffer is used which stores samples of the environment. This allows random mini-batch sampling, which decorrelates the samples and is proven to improve the data efficiency [x]. The mini-batch sampling also enables the use of improved derivatives of vanilla stochastic gradient descent, e.g. *RPROP* as in the *Neural Fitted Q-Learning* approach (Riedmiller 2005) or *ADAM Update* (Kingma and Ba 2014).

There are different ways of estimating the expected Q-values. Either with a target network with the same structure as the network for the action-value function or the normal network. If a target network is used, the target weights need to be updated after some training steps (Mnih et al. 2015).

A pseudo-code for the DQN approach can be found in algorithm 1. The DQN approach was able to significantly outperform earlier learning methods despite incorporating almost no prior knowledge about the inputs (Mnih et al. 2013). However, it is limited by the disability to cope with continuous and high-dimensional action spaces due to the max operator in the action selection (Lillicrap et al. 2015). This limitations can be addressed by combining the approach with the Deterministic Policy Gradient, which is described in the following section.

2.3 Deterministic Policy Gradient

Most problems in reinforcement learning consist of a continuous action space which makes it very difficult to greedily choose the best action given a policy, due to the max operator. From a stochastic point of view, the policy is a probability distribution $a \sim \pi(a|s)$ over all actions. In order to calculate the gradient of a parameterized policy $\pi(a, s|\theta)$ over the total reward w.r.t. the weights, one needs to solve an integral over all actions and states, which becomes intractable for large state-action spaces. From a deterministic view the policy is a discrete mapping from states to actions $a = \pi(s)$ and thus only one integration over the state space is sufficient.

Algorithm 1 Deep Q-Learning (DQN)

Initialize: Replay buffer D with high capacity
Initialize: Neural network for action-value function Q with random weights θ
Initialize: Neural network for target action-value function \hat{Q} with weights $\theta^- = \theta$

for episode 1 **to** M **do**
 reset environment to state s_1
 for $t = 1$ **to** T **do**
 if random $i \leq \epsilon$ **then**
 random action a_t
 else
 $a_t = \operatorname{argmax}_a Q(s_t, a|\theta)$
 end if
 execute $a_t \rightarrow$ reward r_t and next state s_{t+1}
 save (s_t, a_t, r_t, s_{t+1}) in D
 sample mini-batch $(s_i, a_i, r_i, s_{i+1})_k$ of size k from D
 $q_i = \begin{cases} r_i & \text{if episode terminates at step } i+1 \\ r_i + \gamma \max_{a'} \hat{Q}(s_{i+1}, a'|\theta^-) & \text{else} \end{cases}$
 perform gradient descent on $(q_i - Q(s_i, a_i|\theta))_\theta^2$
 every C steps update $\hat{Q} = Q$
 end for
end for

The *policy gradient theorem* (Sutton and Barto 2018) gives the update rule for a parameterized policy, optimizing the loss function:

$$\nabla_\theta J(\theta) = E_{s \sim \rho^\pi, a \sim \pi_\theta} [\nabla_\theta \log \pi_\theta(a|s) Q^\pi(s, a)]$$

From this, a deterministic approach is derived in [EDIT: by?] [Silver et al. (2014)], which gives the update rule for a parameterized deterministic policy function $\pi(s|\theta)$. Rather than trying to maximize the action-value function $Q(s, a)$ globally by greedy improvements of the policy, the authors move the policy in the direction of the gradient of $Q(s, a)$:

$$\nabla_{\theta^\pi} J \approx E_{s \sim \rho^\pi} [\nabla_{\theta^\pi} Q(s, a|\theta^\pi)]$$

Applying the chain rule to this equation gives the *deterministic policy gradient (DPG) theorem*:

$$\nabla_{\theta^\pi} J \approx E_{s \sim \rho^\pi} [\nabla_a Q(s, a|\theta^\pi)|_{a=\pi(s|\theta^\pi)} \nabla_{\theta^\pi} \pi(s|\theta^\pi)]$$

where the expectation can again be approximated by sampling from an environment. Only using deterministic action outputs will vanish the algorithms exploration, so one needs to make sure there still is exploration [EDIT: FOR-MULATION]. This is realized by using an off-policy approach which follows a stochastic policy while learning a deterministic policy. The authors also introduce the notion of *compatible function approximation*. [EDIT: DESCRIBE] Using these to estimate the gradient, an unbiased approximation is guaranteed.

The following section describes a typical structure of how to use deep neural networks for function approximation in reinforcement learning. Together with this section this led to the algorithm described in chapter 3.

2.4 Actor-Critic Methods

A lot of recent success in reinforcement learning is based on *Actor-Critic* methods [Konda and Tsitsiklis (2000)]. In contrast to value-function or policy-gradient methods, they parameterize both, the value function $Q(s, a) \approx \hat{Q}(s, a|\theta^Q)$, also known as the *Critic*, and the policy $\pi(s|a) \approx \hat{\pi}(s|a, \theta^\pi)$. To get an intuition about these methods figure 1 illustrates the update-cycle:

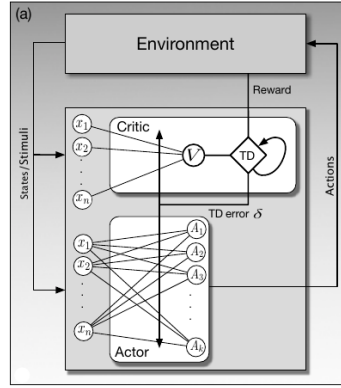


Fig. 1 Intuition about actor-critic methods (figure from Sutton and Barto (2018))

While the actor learns how to choose the right action and is responsible to update his policy, the critic has to learn and update the parameters of the state-value function. The actors' and the critics' parameters can be updated following the TD-error of the critic which is computed from the observed reward and the current error of the estimated state values in every time-step. As Fig. 1 illustrates, the actor has no information about the current reward and the critic has no direct influence on the actions. A pseudo-code for an actor-critic method using TD-errors [x] is shown in algorithm 2.

3 Deep Deterministic Policy Gradient

The combination of above approaches led to the *Deep Deterministic Policy Gradient (DDPG)* approach [Lillicrap et al. (2015)], which is a model-free and off-policy algorithm. It can be grouped into the class actor-critic methods and uses a deterministic target policy and deep Q-Learning. Both, the actor and the critic, are realized by deep neural networks. The pseudo-code for DDPG can be found in 3.

It consists of a parameterized deterministic policy, the actor, $\pi(s|\theta^\pi)$ and a parameterized action-value function $Q(s, a|\theta^Q)$, the critic. The critic is updated using the *Bellman Equation* with a TD-error similar in Q-Learning [Watkins and Dayan (1992)] [EDIT: EQUATION?] and the actor is updated using the

Algorithm 2 Episodic One-step Actor-Critic for Estimating $\pi(a|s, \theta^\pi) \approx \pi^*(a|s)$

Initialize: Differentiable policy parameterization $\pi(a|s, \theta^\pi)$
Initialize: Differentiable action-value function parameterization $Q(s, a|\theta^Q)$
Initialize: Random initial weights θ^π and θ^Q
Initialize: Step size parameters $\alpha^Q > 0$ and $\alpha^\pi > 0$
Initialize: Discount factor γ

```

for episode 1 to  $M$  do
  Get initial state  $s$ 
   $i \leftarrow 1$ 
  for time-step 1 to  $T$  do
    Draw action from actor:  $a \sim \pi(s|a, \theta^\pi)$ 
    Do action  $a$ , observe reward  $r$  and successor state  $s'$ 
    Calculate the TD-error:
       $\delta \leftarrow r + \gamma \max_{a'} Q(s', a'|\theta^Q) - Q(s, a|\theta^Q)$ 
    Update the weights:
       $\theta^Q \leftarrow \theta^Q + \alpha^Q \delta \nabla_{\theta^Q} Q(s, a|\theta^Q)$ 
       $\theta^\pi \leftarrow \theta^\pi + \alpha^\pi \delta \nabla_{\theta^\pi} \log \pi(a|s, \theta^\pi)$ 
    Update:
       $i \leftarrow \gamma i$ 
       $s \leftarrow s'$ 
  end for
end for

```

DPG theorem [EDIT: LINK? EQUATIONS WITH NUMBERS?].

The use of neural networks to parameterize the above functions means that the convergence guarantees do not hold anymore. Therefore the Actor-Critic DPG approach is combined with recent successes from DQN.

To ensure independently and identically distributed data, the authors use a replay buffer and sample random mini-batches from it. This again decorrelates the samples and allows the efficient use of hardware optimization, e.g. the ADAM update [Kingma and Ba (2014)].

To address instability issues from applying deep neural network approximation to Q-Learning they also use *target networks* which are copies of the actor $\pi'(s|\theta^{\pi'})$ and the critic $Q'(s, a|\theta^{Q'})$. These target-networks track the learned networks and are constrained to slow changes by using soft updates: $\theta' \leftarrow \tau \theta + (1 - \tau) \theta'$ with $\tau \ll 1$. These consistent targets might slow down the learning process but greatly improve the stability of it.

Using low dimensional feature input might give very different scales for the single states. This can lead to problematic learning for the neural networks and is addressed by using *batch normalization* which normalizes each dimension across the samples in a mini-batch.

To ensure exploration while using a deterministic policy, a noise process N is added to the action output of the actor network. This noise process can be chosen to suit the environment. The algorithm was evaluated on more than 20 simulated physical tasks using the same algorithm, network structures and hyper-parameters, including classic control problems like the cart-pole environment. Using low-dimensional feature input, it was able to find policies performing really well on most of the tasks. Their performance is competitive

Algorithm 3 Deep Deterministic Policy Gradient (DDPG)

Initialize: Replay buffer D with high capacity
Initialize: Critic network $Q(s, a|\theta^Q)$ and actor network $\pi(s|\theta^\pi)$ with random weights θ^Q and θ^π
Initialize: Initialize target networks Q' and π' with weights $\theta^{Q'} \leftarrow \theta^Q$ and $\theta^{\pi'} \leftarrow \theta^\pi$
for episode 1 **to** M **do**
 Initialize random process N for action exploration
 Reset environment to state s_1
 for $t = 1$ **to** T **do**
 Select action $a_t = \pi(s_t|\theta^\pi) + N_t$ from local actor
 Execute action a_t and observe reward r_t and next state s_{t+1}
 Save (s_t, a_t, r_t, s_{t+1}) in replay buffer D
 Sample mini-batch $(s_i, a_i, r_i, s_{i+1})_k$ with size k from D
 Set TD-target from target networks:
 $y_i = r_i + \gamma Q'(s_{i+1}, \pi'(s_{i+1}|\theta^{\pi'}))|\theta^{Q'}$
 Update the critic by minimizing the loss:
 $L = \frac{1}{N} \sum_i (y_i - Q(s_i, a_i|\theta^Q))^2$
 Update the actor using the sampled policy gradient:
 $\nabla_{\theta^\pi} J \approx \frac{1}{N} \sum_i \nabla_a Q(s, a|\theta^Q)|_{s=s_i, a=\pi(s_i)} \nabla_{\theta^\pi} \pi(s|\theta^\pi)|_{s=s_i}$
 Update the target networks:
 $\theta^{Q'} \leftarrow \tau \theta^Q + (1 - \tau) \theta^{Q'}$
 $\theta^{\pi'} \leftarrow \tau \theta^\pi + (1 - \tau) \theta^{\pi'}$
 end for
end for

with these found by a controller with full access to the environment. The algorithm is even able to find good policies using high dimensional pixel input. For simple tasks this turned out to be as fast as using low dimensional state features.

The most challenging issues of the approach are his poor sample efficiency and some instabilities. We present some possible extensions to DDPG in the next chapter which might improve on these issues.

TODO:

- discretizing the action space often suffers from the curse of dimensionality
- naive extension of DPG with nns turns out to be unstable for challenging problems

4 Improvements for DDPG

Despite it's good performance on many simulated tasks there is still some room to improve the DDPG algorithm. We show some possible extensions for it in this section.

4.1 Using importance sampling to sample from the replay-buffer

In practice the algorithm is limited by the maximum storage size N of the replay-buffer D . Overwriting older samples by current ones does nowhere differentiate between more or less important experiences, because uniform random samples does weight all experiences equally. One could use a technique similar to *prioritized sweeping* [Moore and Atkeson (1993)] which uses *importance sampling* [Glynn and Iglehart (1989)] to prefer transitions which are more important over ones that have less value for the training process.

4.2 Using Action Noise in Parameter Space

Instead of adding noise to the action space to ensure exploration, one could add adaptive noise directly to the parameters of the neural network [Plappert et al. (2017)]. This would add some randomness into the parameters of the agent and therefore into the decision it makes, while still always fully depending on it's current observation about it's environment. This parameter noise makes an agent's exploration more consistent and results in a more effective exploration, increased performance and smoother behavior.

4.3 Evolutionary Approaches

One can consider an even more extreme case of the above mentioned extension, which would be the use of *Evolutionary Strategies* to approximate the gradient of our objective function [Salimans et al. (2017)]. This does not require back-propagation at all and is competitive with state of the art RL.

4.4 Improvements of the Deep Neural Network Architectures

TODO

5 Conclusion

References

- Glynn PW, Iglehart DL (1989) Importance sampling for stochastic simulations. *Management Science* 35(11):1367–1392
- Kingma DP, Ba J (2014) Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*
- Konda VR, Tsitsiklis JN (2000) Actor-critic algorithms. In: *Advances in neural information processing systems*, pp 1008–1014
- Lillicrap TP, Hunt JJ, Pritzel A, Heess N, Erez T, Tassa Y, Silver D, Wierstra D (2015) Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*
- Mnih V, Kavukcuoglu K, Silver D, Graves A, Antonoglou I, Wierstra D, Riedmiller M (2013) Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*

- Mnih V, Kavukcuoglu K, Silver D, Rusu AA, Veness J, Bellemare MG, Graves A, Riedmiller M, Fidjeland AK, Ostrovski G, et al. (2015) Human-level control through deep reinforcement learning. *Nature* 518(7540):529
- Moore AW, Atkeson CG (1993) Prioritized sweeping: Reinforcement learning with less data and less time. *Machine learning* 13(1):103–130
- Plappert M, Houthoofd R, Dhariwal P, Sidor S, Chen RY, Chen X, Asfour T, Abbeel P, Andrychowicz M (2017) Parameter space noise for exploration. *arXiv preprint arXiv:1706.01905*
- Riedmiller M (2005) Neural fitted q iteration—first experiences with a data efficient neural reinforcement learning method. In: *European Conference on Machine Learning*, Springer, pp 317–328
- Salimans T, Ho J, Chen X, Sidor S, Sutskever I (2017) Evolution strategies as a scalable alternative to reinforcement learning. *arXiv preprint arXiv:1703.03864*
- Silver D, Lever G, Heess N, Degris T, Wierstra D, Riedmiller M (2014) Deterministic policy gradient algorithms. In: *ICML*
- Sutton RS, Barto AG (2018) *Reinforcement learning: An introduction*. MIT press
- Watkins CJ, Dayan P (1992) Q-learning. *Machine learning* 8(3-4):279–292