
Self-Guided and Self-Regularized Actor-Critic

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Deep reinforcement learning (DRL) algorithms have successfully been demon-
2 strated on a range of challenging decision making and control tasks. One dominant
3 component of recent deep reinforcement learning algorithms is the target network
4 which mitigates the divergence when learning the Q function. However, target
5 networks can slow down the learning process due to delayed function updates. An-
6 other dominant component especially in continuous domains is the policy gradient
7 method which models and optimizes the policy directly. However, when Q func-
8 tions are approximated with neural networks, their landscapes can be complex and
9 therefore mislead the local gradient. In this work, we propose a self-regularized and
10 self-guided actor-critic method. We introduce a self-regularization term within the
11 TD-error minimization and remove the need for the target network. In addition, we
12 propose a self-guided policy improvement method by combining policy-gradient
13 with zero-order optimization such as the Cross Entropy Method. It helps to search
14 for actions associated with higher Q-values in a broad neighborhood and is robust
15 to local noise in the Q function approximation. These actions help to guide the
16 updates of our actor network. We evaluate our method on the suite of OpenAI gym
17 tasks, achieving or outperforming state of the art in every environment tested.

18 1 Introduction

19 Reinforcement learning (RL) studies decision-making with the goal of maximizing total discounted
20 reward when interacting with an environment. Leveraging high-capacity function approximators such
21 as neural networks, Deep reinforcement learning (DRL) algorithms have been successfully applied to
22 a range of challenging domains, from video games [12] to robotic control [16].

23 Actor-critic algorithms are among the most popular approaches in DRL, e.g. *DDPG* [11], *TRPO* [16],
24 *TD3* [6] and *SAC* [7]. These methods are based on policy iteration, which alternates between policy
25 evaluation and policy improvement [17]. Actor-critic methods jointly optimize the value function
26 (critic) and the policy (actor) as it is often impractical to run either of these to convergence [7].

27 In DRL, both the actor and critic use deep neural networks as the function approximator. However,
28 DRL is known to assign unrealistically high values to state-action pairs represented by the Q-function.
29 This is detrimental to the quality of the greedy control policy derived from Q [21]. Mnih *et al.* [13]
30 proposed to use a *target network* to mitigate divergence. A target network is a copy of the current Q
31 function that is held fixed to serve as a stable target within the TD error update. The parameters of the
32 target network are either infrequently copied [13] or obtained by Polyak averaging [11]. A limitation
33 of using a target network is that it can slow down learning due to delayed function updates. We propose
34 an approach that reduces the need for a target network in DRL while still ensuring stable learning
35 and good performance in high-dimensional domains. We add a self-regularization term to encourage
36 small changes to the target value while minimizing the Temporal Difference (TD)-error [17].

37 Evolution Strategies (ES) are a family of black-box optimization algorithms which are typically very
38 stable, but scale poorly in high-dimensional search spaces, (e.g. neural networks) [14]. Gradient-

39 based DRL methods are often sample efficient, particularly in the off-policy setting when, unlike
40 evolutionary search methods, they can continue to sample previous experiences to improve value
41 estimation. But these approaches can also be unstable and highly sensitive to hyper-parameter
42 tuning [14]. We propose a novel policy improvement method which combines both approaches to
43 get the best of both worlds. Specifically, after the actor network first outputs an initial action, we
44 apply the *Cross Entropy Method* (CEM) [15] to search the neighborhood of the initial action to find a
45 second action associated with a higher Q value. Then we leverage the second action in the policy
46 improvement stage to speed up the learning process.

47 To mitigate the overestimation issue in Q learning [18], Fujimoto *et al.* [6] proposed Clipped Double
48 Q-Learning in which the authors learn two Q-functions and use the smaller one to form the targets
49 in the TD-Learning process. This method may suffer from under-estimation. In practice, we also
50 observe that the discrepancy between the two Q-functions can increase dramatically which hinders
51 the learning process. We propose Max-min Double Q-Learning to address this discrepancy. Our
52 method also provides a better approximation of the Bellman optimality operator [17].

53 We propose a novel self-Guided and self-Regularized Actor Critic (*GRAC*) algorithm. *GRAC* uses
54 **self-regularized TD-Learning** removing the need for a target network and utilizes a novel **policy**
55 **improvement method which combines policy-gradients and zero-order optimization** to speed
56 up learning. Following Clipped Double Q-Learning, we propose **Max-min Double Q-learning** to
57 address underestimation and the discrepancy between the two Q functions. We evaluate *GRAC* on
58 six continuous control domains from OpenAI gym [3], where we achieve or outperform state of the
59 art result in every environment tested. We run our experiments across a large number of seeds with
60 fair evaluation metrics [4], perform extensive ablation studies, and open source both our code and
61 learning curves.

62 2 Related Work

63 The proposed algorithm incorporates three key ingredients within the actor-critic method: a self-
64 regularized TD update, self-guided policy improvements based on evolution strategies, and Max-min
65 double Q-Learning. In this section, we review prior work related to these ideas.

66 **Divergence in Deep Q-Learning** In Deep Q-Learning, we use a nonlinear function approximator
67 such as a neural network to approximate the Q-function that represents the value of each state-action
68 pair. Learning the Q-function in this way is known to suffer from divergence issues [20] such as
69 assigning unrealistically high values to state-action pairs [21]. This is detrimental to the quality of
70 the greedy control policy derived from Q [21]. To mitigate the divergence issue, Mnih *et al.* [13]
71 introduce a target network which is a copy of the estimated Q-function and is held fixed to serve as a
72 stable target for some number of steps. However, target networks can slow down the learning process
73 due to delayed function updates [10]. Durugkar *et al.* [5] propose Constrained Q-Learning, which
74 uses a constraint to prevent the average target value from changing after an update. Achiam *et al.* [1]
75 give a simple analysis based on a linear approximation of the Q function and develop a stable Deep Q-
76 Learning algorithm for continuous control without target networks. However, their proposed method
77 requires separately calculating backward passes for each state-action pair in the batch, and solving a
78 system of equations at each timestep. The proposed *GRAC* algorithm adds a self-regularization term
79 to the TD-Learning objective to keep the change of the state-action value small.

80 **Evolution Strategies in Deep Reinforcement Learning** *Evolution Strategies* (ES) are a family
81 of black-box optimization algorithms which are typically very stable, but scale poorly in high-
82 dimensional search spaces [23]. Gradient-based deep RL methods, such as *DDPG* [11], are often
83 sample efficient, particularly in the off-policy setting. These off-policy methods can continue to
84 reuse previous experience to improve value estimations but can be unstable and highly sensitive to
85 hyper-parameter tuning [14]. Researchers have proposed to combine these approaches to get the best
86 of both worlds. Pourchot *et al.* [14] proposed *CEM-RL* to combine CEM with either *DDPG* [11]
87 or *TD3* [6]. However, *CEM-RL* applies *CEM* within the actor parameter space which is extremely
88 high-dimensional, making the search not efficient. Kalashnikov *et al.* [9] introduce *QT-Opt*, which
89 leverages *CEM* to search the landscape of the Q function, and enables Q-Learning in continuous
90 action spaces without using an actor. However, as shown in [23], *CEM* does not scale well to
91 high-dimensional action spaces, such as in the Humanoid task we used in this paper. A near-optimal
92 actor is needed to initialize the *CEM* process in such tasks. Different from *QT-Opt*, we adopt the

93 actor-critic framework and leverage *CEM* in both Q-Learning and policy improvement. *GRAC* speeds
 94 up the learning process compared to popular actor-critic methods.

95 **Double-Q Learning** Using function approximation, Q-learning [22] is known to suffer from
 96 overestimation [18]. To mitigate this problem, Hasselt *et al.* [8] proposed *Double Q-learning* which
 97 uses two Q functions with independent sets of weights. *TD3* [6] proposed *Clipped Double Q-learning*
 98 to learn two Q-functions and uses the smaller of the two to form the targets in the TD-Learning
 99 process. However, *TD3* [6] may lead to underestimation. Besides, the actor network in *TD3* [6] is
 100 trained to select the action to maximize the first Q function throughout the training process which
 101 may make it very different from the second Q-function. A large discrepancy results in large TD-errors
 102 which in turn results in large gradients during the update of the actor and critic networks. This makes
 103 instability of the learning process more likely. We propose *Max-min Double Q-Learning* to balance
 104 the differences between the two Q functions and provide a better approximation of the Bellman
 105 optimality operator [17].

106 3 Preliminaries

107 In this section, we define the notation used in subsequent sections. Consider a *Markov Decision*
 108 *Process* (MDP), defined by the tuple $(\mathcal{S}, \mathcal{A}, \mathcal{P}, r, \rho_0, \gamma)$, where \mathcal{S} is a finite set of states, \mathcal{A} is a finite
 109 set of actions, $\mathcal{P} : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ is the transition probability distribution, $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the
 110 reward function, $\rho_0 : \mathcal{S} \rightarrow \mathbb{R}$ is the distribution of the initial state s_0 , and $\gamma \in (0, 1)$ is the discount
 111 factor. At each discrete time step t , with a given state $s_t \in \mathcal{S}$, the agent selects an action $a_t \in \mathcal{A}$,
 112 receiving a reward r and the new state s_{t+1} of the environment.

113 Let π denote the policy which maps a state to a probability distribution over the actions, $\pi : \mathcal{S} \rightarrow$
 114 $\mathcal{P}(\mathcal{A})$. The return from a state is defined as the sum of discounted reward $R_t = \sum_{i=t} \gamma^{i-t} r(s_i, a_i)$.
 115 In reinforcement learning, the objective is to find the optimal policy π^* , with parameters ϕ , which max-
 116 imizes the expected return $J(\phi) = \sum_t \mathbb{E}_{(s_t, a_t) \sim \rho_{\pi}(s_t, a_t)} [\gamma^t r(s_t, a_t)]$ where $\rho_{\pi}(s_t)$ and $\rho_{\pi}(s_t, a_t)$
 117 denote the state and state-action marginals of the trajectory distribution induced by the policy $\pi(a_t | s_t)$.

118 We use the following standard definitions of the state-action value function Q_{π} . It describes the
 119 expected discounted reward after taking an action a_t in state s_t and thereafter following policy π :

$$Q_{\pi}(s_t, a_t) = \mathbb{E}_{\pi}[R_t | s_t, a_t]. \quad (1)$$

120 In this work we use *CEM* to find optimal actions with maximum Q values. *CEM* is a randomized
 121 zero-order optimization algorithm. To find the action a that maximizes $Q(s, a)$, *CEM* is initialized
 122 with a parameterized distribution over a , $P(a; \psi)$. Then it iterates between the following two steps [2]:
 123 First generate $a_1, \dots, a_N \sim P(s; \psi)$. Retrieve their Q values $Q(s, a_i)$ and sort the actions to have
 124 decreasing Q values. Then keep the first K actions, and solve for an updated parameters ψ' :

$$\psi' = \operatorname{argmax}_{\psi} \frac{1}{K} \sum_{i=1}^K \log(P(a_i; \psi))$$

125 In the following sections, we denote $CEM(Q(s, \cdot), \pi(\cdot | s))$ as the action found by *CEM* to maximize
 126 $Q(s, \cdot)$, when *CEM* is initialized by the distribution predicted by the policy.

127 4 Technical Approach

128 4.1 Self-Regularized TD Learning

129 Reinforcement learning is prone to instability and divergence when a nonlinear function approximator
 130 such as a neural network is used to represent the Q function [20]. Mnih *et al.* [13] identified several
 131 reasons for this. One is the correlation between the current action-values and the target value. Updates
 132 to $Q(s_t, a_t)$ often also increase $Q(s_{t+1}, a_{t+1}^*)$ where a_{t+1}^* is the optimal next action. Hence, these
 133 updates also increase the target value y_t which may lead to oscillations or the divergence of the policy.

134 More formally, given transitions (s_t, a_t, r_t, s_{t+1}) sampled from the replay buffer distribution \mathcal{B} , the
 135 Q network can be trained by minimising the loss functions $\mathcal{L}(\theta_i)$ at iteration i :

$$\mathcal{L}(\theta_i) = \mathbb{E}_{(s_t, a_t) \sim \mathcal{B}} \|(Q(s_t, a_t; \theta_i) - y_i)\|^2 \quad (2)$$

Algorithm 1 GRAC

Initialize critic network $Q_{\theta_1}, Q_{\theta_2}$ and actor network π_ϕ with random parameters θ_1, θ_2 and ϕ
Initialize replay buffer \mathcal{B} , Set $\alpha < 1$

- 1: **for** $i = 1, \dots$ **do**
- 2: Select action $a \sim \pi_{\phi_i}(s)$ and observe reward r and new state s'
- 3: Store transition tuple (s, a, r, s') in \mathcal{B}
- 4: Sample mini-batch of N transitions (s_t, a_t, r_t, s_{t+1}) from \mathcal{B}
- 5: $\hat{a}_{t+1} \sim \pi_{\phi_i}(s_{t+1})$
- 6: $\tilde{a}_{t+1} \leftarrow \text{CEM}(Q(s_{t+1}, \cdot; \theta_2), \pi_{\phi_i}(\cdot | s_{t+1}))$
- 7: $y \leftarrow r_t + \gamma \max\{\min_{j=1,2} Q(s_{t+1}, \tilde{a}_{t+1}; \theta_j), \min_{j=1,2} Q(s_{t+1}, \hat{a}_{t+1}; \theta_j)\}$
- 8: $a^\dagger \leftarrow \arg \max_{\{\bar{a}, \hat{a}\}} \{\min_{j=1,2} Q(s_{t+1}, \tilde{a}_{t+1}; \theta_j), \min_{j=1,2} Q(s_{t+1}, \hat{a}_{t+1}; \theta_j)\}$
- 9: $y'_1, y'_2 \leftarrow Q(s_{t+1}, a^\dagger; \theta_1), Q(s_{t+1}, a^\dagger; \theta_2)$
- 10: **for** $k = 1$ to K **do**
- 11: $\mathcal{L}_k = \|y - Q(s_t, a_t; \theta_1)\|_2 + \|y - Q(s_t, a_t; \theta_2)\|_2 + \|y'_1 - Q(s_{t+1}, a^\dagger; \theta_1)\|_2 + \|y'_2 - Q(s_{t+1}, a^\dagger; \theta_2)\|_2$
- 12: $\theta_1 \leftarrow \theta_1 - \lambda \nabla_{\theta_1} \mathcal{L}_k, \theta_2 \leftarrow \theta_2 - \lambda \nabla_{\theta_2} \mathcal{L}_k$
- 13: **if** $k > 1$ and $\mathcal{L}_k < \alpha \mathcal{L}_2$ **then**
- 14: Break
- 15: **end if**
- 16: **end for**
- 17: $\hat{a}_t \sim \pi_{\phi_i}(s_t)$
- 18: $J_\pi(\phi) = \mathbb{E}_{(s_t, \hat{a}_t)}[Q(s_t, \hat{a}_t; \theta_1)]$
- 19: $\bar{a}_t \leftarrow \text{CEM}(Q(s_t, \cdot; \theta_1), \pi_{\phi_i}(\cdot | s_t))$
- 20: $\phi \leftarrow \phi - \lambda \nabla_\phi J_\pi(\phi) - \lambda \mathbb{E}_{(s_t, \hat{a}_t)}[Q(s_t, \bar{a}_t; \theta_1) - Q(s_t, \hat{a}_t; \theta_1)]_+ \nabla_\phi \log \pi(\bar{a}_t | s_t; \phi)$
- 21: **end for**

136 where for now let us assume $y_i = r_t + \gamma \max_a Q(s_{t+1}, a; \theta_i)$ to be the target for iteration i computed
137 based on the current Q network parameters θ_i . $a_{t+1}^* = \arg \max_a Q(s_{t+1}, a)$. If we update the
138 parameter θ_{i+1} to reduce the loss $\mathcal{L}(\theta_i)$, it changes both $Q(s_t, a_t; \theta_{i+1})$ and $Q(s_{t+1}, a_{t+1}^*; \theta_{i+1})$.
139 Assuming an increase in both values, then the new target value $y_{i+1} = r_t + \gamma Q(s_{t+1}, a_{t+1}^*; \theta_{i+1})$ for
140 the next iteration will also increase leading to an explosion of the Q function. We demonstrated this
141 behavior in an ablation experiment with results in Fig. 2. We also show how maintaining a separate
142 target network [13] with frozen parameters θ^- to compute $y_{i+1} = r_t + \gamma Q(s_{t+1}, a_{t+1}^*; \theta^-)$ delays
143 the update of the target and therefore leads to more stable learning of the Q function. However,
144 delaying the function updates also comes with the price of slowing down the learning process.

145 We propose a self-Regularized TD-learning approach to minimize the TD-error while also keeping
146 the change of $Q(s_{t+1}, a_{t+1}^*)$ small. This regularization mitigates the divergence issue [20], and
147 no longer requires a target network that would otherwise slow down the learning process. Let
148 $y'_i = Q(s_{t+1}, a_{t+1}^*; \theta_i)$, and $y_i = r_t + \gamma y'_i$. We define the learning objective as

$$\min_{\theta} \|Q(s_t, a_t; \theta) - y_i\|^2 + \|Q(s_{t+1}, a_{t+1}^*; \theta) - y'_i\|^2 \quad (3)$$

149 where the first term is the original TD-Learning objective and the second term is the regularization
150 term penalizing large updates to $Q(s_{t+1}, a_{t+1}^*)$. Note that when the current Q network updates its
151 parameters θ , both $Q(s_t, a_t)$ and $Q(s_{t+1}, a_{t+1}^*)$ change. Hence, the target value y_i will also change
152 which is different from the approach of keeping a frozen target network for a few iterations. We will
153 demonstrate in our experiments that this self-regularized TD-Learning approach removes the delays
154 in the update of the target value thereby achieves faster and stable learning.

155 4.2 Self-Guided Policy Improvement with Evolution Strategies

156 The policy, known as the actor, can be updated through a combination of two parts. The first part,
157 which we call Q-loss policy update, improves the policy through local gradients of the current Q
158 function, while the second part, which we call CEM policy update, finds a high-value action via
159 CEM in a broader neighborhood of the Q function landscape, and update the action distribution to
160 concentrate towards this high-value action. We describe the two parts formally below.

161 Given states s_t sampled from the replay buffer, the Q-loss policy update maximizes the objective

$$J_\pi(\phi) = \mathbb{E}_{s_t \sim \mathcal{B}, \hat{a}_t \sim \pi} [Q(s_t, \hat{a}_t)], \quad (4)$$

162 where \hat{a}_t is sampled from the current policy $\pi(\cdot|s_t)$. The gradient is taken through the reparameter-
163 ization trick. We reparameterize the policy using a neural network transformation as described in
164 Haarnoja *et al.* [7],

$$\hat{a}_t = f_\phi(\epsilon_t|s_t) \quad (5)$$

165 where ϵ_t is an input noise vector, sampled from a fixed distribution, such as a standard multivariate
166 Normal distribution. Then the gradient of $J_\pi(\phi)$ is:

$$\nabla J_\pi(\phi) = \mathbb{E}_{s_t \sim \mathcal{B}, \epsilon_t \sim \mathcal{N}} \left[\frac{\partial Q(s_t, f_\phi(\epsilon_t|s_t))}{\partial f} \frac{\partial f_\phi(\epsilon_t|s_t)}{\partial \phi} \right] \quad (6)$$

167 For the CEM policy update, given a minibatch of states s_t , we first find a high-value action \bar{a}_t for
168 each state by running *CEM* on the current Q function, $\bar{a}_t = \text{CEM}(Q(s_t, \cdot), \pi(\cdot|s_t))$. Then the policy
169 is updated to increase the probability of this high-value action. The guided update on the parameter ϕ
170 of π at iteration i is

$$\mathbb{E}_{s_t \sim \mathcal{B}, \hat{a}_t \sim \pi} [Q(s_t, \bar{a}_t) - Q(s_t, \hat{a}_t)]_+ \nabla_\phi \log \pi_i(\bar{a}_t|s_t). \quad (7)$$

171 We used $Q(s_t, \hat{a}_t)$ as a baseline term, since its expectation over actions \hat{a}_t will give us the normal
172 baseline $V(s_t)$:

$$\mathbb{E}_{s_t \sim \mathcal{B}} [Q(s_t, \bar{a}_t) - V(s_t)] \nabla_\phi \log \pi_i(\bar{a}_t|s_t) \quad (8)$$

173 In our implementation, we only perform an update if the improvement on the Q function, $Q(s_t, \bar{a}_t) -$
174 $Q(s_t, \hat{a}_t)$, is non-negative, to guard against the occasional cases where *CEM* fails to find a better
175 action.

176 Combining both parts of updates, the final update rule on the parameter ϕ_i of policy π_i is

$$\phi_{i+1} = \phi_i - \lambda \nabla_\phi J_{\pi_i}(\phi_i) - \lambda \mathbb{E}_{s_t \sim \mathcal{B}, \hat{a}_t \sim \pi_i} [Q(s_t, \bar{a}_t) - Q(s_t, \hat{a}_t)]_+ \nabla_\phi \log \pi_i(\bar{a}_t|s_t)$$

177 where λ is the step size.

178 We can prove that if the Q function has converged to Q^π , the state-action value function induced by
179 the current policy, then both the Q-loss policy update and the *CEM* policy update will be guaranteed
180 to improve the current policy. We formalize this result in Theorem 1 and Theorem 2, and prove them
181 in Appendix 3.1 and 3.2.

182 **Theorem 1 Q-loss Policy Improvement** Starting from the current policy π , we maximize the objec-
183 tive $J_\pi = \mathbb{E}_{(s,a) \sim \rho_\pi(s,a)} Q^\pi(s, a)$. The maximization converges to a critical point denoted as π_{new} .
184 Then the induced Q function, $Q^{\pi_{new}}$, satisfies $\forall (s, a), Q^{\pi_{new}}(s, a) \geq Q^\pi(s, a)$.

185 **Theorem 2 CEM Policy Improvement** Assuming the CEM process is able to find the optimal
186 action of the state-action value function, $a^*(s) = \arg \max_a Q^\pi(s, a)$, where Q^π is the Q function
187 induced by the current policy π . By iteratively applying the update $\mathbb{E}_{(s,a) \sim \rho_\pi(s,a)} [Q(s, a^*) -$
188 $Q(s, a)]_+ \nabla \log \pi(a^*|s)$, the policy converges to π_{new} . Then $Q^{\pi_{new}}$ satisfies $\forall (s, a), Q^{\pi_{new}}(s, a) \geq$
189 $Q^\pi(s, a)$.

190 4.3 Max-min Double Q-Learning

191 Q-learning [22] is known to suffer from overestimation [18]. Hasselt *et al.* [8] proposed Double-Q
192 learning which uses two Q functions with independent sets of weights to mitigate the overestimation
193 problem. Fujimoto *et al.* [6] proposed Clipped Double Q-learning with two Q function denoted
194 as $Q(s, a; \theta_1)$ and $Q(s, a; \theta_2)$, or Q_1 and Q_2 in short. Given a transition (s_t, a_t, r_t, s_{t+1}) , Clipped
195 Double Q-learning uses the minimum between the two estimates of the Q functions when calculating
196 the target value in TD-error [17]:

$$y = r_t + \gamma \min_{j=1,2} Q(s_{t+1}, \hat{a}_{t+1}; \theta_j) \quad (9)$$

197 where \hat{a}_{t+1} is the predicted next action.

198 Fujimoto *et al.* [6] mentioned that such an update rule may induce an underestimation bias. In
 199 addition, $\hat{a}_{t+1} = \pi_\phi(s_{t+1})$ is the prediction of the actor network. The actor network’s parameter
 200 ϕ is optimized according to the gradients of Q_1 . In other words, \hat{a}_{t+1} tends to be selected according
 201 to the Q_1 network which consistently increases the discrepancy between the two Q-functions. In
 202 practice, we observe that the discrepancy between the two estimates of the Q-function, $|Q_1 - Q_2|$,
 203 can increase dramatically leading to an unstable learning process. An example is shown in Fig. 4
 204 where $Q(s_{t+1}, \hat{a}_{t+1}; \theta_1)$ is always bigger than $Q(s_{t+1}, \hat{a}_{t+1}; \theta_2)$.

205 We introduce *Max-min Double Q-Learning* to reduce the discrepancy between the Q-functions. We
 206 first select \hat{a}_{t+1} according to the actor network $\pi_\phi(s_{t+1})$. Then we run *CEM* to search the landscape
 207 of Q_2 within a broad neighborhood of \hat{a}_{t+1} to return a second action \tilde{a}_{t+1} . Note that *CEM* selects
 208 an action \tilde{a}_{t+1} that maximises Q_2 while the actor network selects an action \hat{a}_{t+1} that maximises
 209 Q_1 . We gather four different Q-values: $Q(s_{t+1}, \hat{a}_{t+1}; \theta_1)$, $Q(s_{t+1}, \hat{a}_{t+1}; \theta_2)$, $Q(s_{t+1}, \tilde{a}_{t+1}; \theta_1)$, and
 210 $Q(s_{t+1}, \tilde{a}_{t+1}; \theta_2)$. We then run a max-min operation to compute the target value that cancels the
 211 biases induced by \hat{a}_{t+1} and \tilde{a}_{t+1} .

$$y = r_t + \gamma \max_{j=1,2} \{ \min_{j=1,2} Q(s_{t+1}, \hat{a}_{t+1}; \theta_j), \min_{j=1,2} Q(s_{t+1}, \tilde{a}_{t+1}; \theta_j) \} \quad (10)$$

212 The inner min-operation $\min_{j=1,2} Q(s_{t+1}, \hat{a}_{t+1}; \theta_j)$ is adopted from Eq. 9 and mitigates overestima-
 213 tion [18]. The outer max operation helps to reduce the difference between Q_1 and Q_2 . In addition,
 214 the max operation provides a better approximation of the Bellman optimality operator [17]. We
 215 visualize Q_1 and Q_2 during the learning process in Fig. 4. The following theorem formalizes the
 216 convergence of the proposed Max-min Double Q-Learning approach in the finite MDP setting. We
 217 prove the theorem in Appendix 3.3.

218 5 Experiments

219 5.1 Comparative Evaluation

220 We present *GRAC*, a self-guided and self-regularized actor-critic algorithm as summarized in Algo-
 221 rithm 1. To evaluate *GRAC*, we measure its performance on the suite of MuJoCo continuous control
 222 tasks [19], interfaced through OpenAI Gym [3]. We compare our method with *DDPG* [11], *TD3* [6],
 223 *TRPO* [16], and *SAC* [7]. We use the source code released by the original authors and adopt the
 224 same hyperparameters reported in the original papers. Hyperparameters for all experiments are in
 225 Appendix 2.1. Results are shown in Figure 1. *GRAC* outperforms or performs comparably to all other
 226 algorithms in both final performance and learning speed across all tasks.

227 5.2 Ablation Study

228 In this section, we present ablation studies to understand the contribution of each proposed component:
 229 Self-Regularized TD-Learning (Section 4.1), Self-Guided Policy Improvement (Section 4.2), and
 230 Max-min Double Q-Learning (Section 4.3). We present our results in Fig. 3 in which we compare
 231 the performance of *GRAC* with alternatives, each removing one component from *GRAC*. Additional
 232 learning curves can be found in Appendix 2.2. We also run experiments to examine how sensitive
 233 *GRAC* is to some hyperparameters such as α and K listed in Alg. 1, and the results can be found in
 234 Appendix 2.4.

235 **Self-Regularized TD Learning** To verify the effectiveness of the proposed self-regularized TD-
 236 learning method, we apply our method to *DDPG* (*DDPG w/o target network w/ target regularization*).
 237 We compare against two baselines: the original *DDPG* and *DDPG* without target networks for both
 238 actor and critic (*DDPG w/o target network*). We choose *DDPG*, because it does not have additional
 239 components such as Double Q-Learning, which may complicate the analysis of this comparison.

240 In Fig. 2, we visualize the average returns, and average Q_1 values over training batches (y'_1 in Alg.1).
 241 The Q_1 values of *DDPG w/o target network* changes dramatically which leads to poor average returns.
 242 *DDPG* maintains stable Q values but makes slow progress. Our proposed *DDPG w/o target network*
 243 *w/ target regularization* maintains stable Q values and learns considerably faster. This demonstrates
 244 the effectiveness of our method and its potentials to be applied to a wide range of DRL methods. Due
 245 to page limit, we only include results on Hopper-v2. The results on other tasks are in Appendix 2.3.
 246 All tasks exhibit a similar phenomenon.

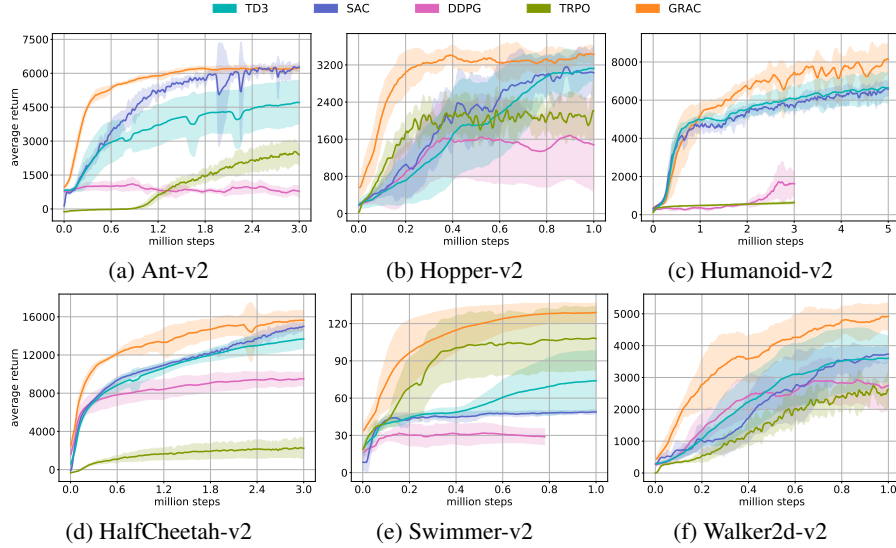


Figure 1: Learning curves for the OpenAI gym continuous control tasks. For each task, we train 8 instances of each algorithm, using 8 different seeds. Evaluations are performed every 5000 interactions with the environment. Each evaluation reports the return (total reward), averaged over 10 episodes. For each training seed, we use a different seed for evaluation, which results in different start states. The solid curves and shaded regions represent the mean and standard deviation, respectively, of the average return over 8 seeds. All curves are smoothed with window size 10 for visual clarity. *GRAC* (orange) learns faster than other methods across all tasks. *GRAC* achieves comparable result to the state-of-the-art methods on the Ant-v2 task and outperforms prior methods on the other five tasks including the complex high-dimensional Humanoid-v2.

247 **Policy Improvement with Evolution Strategies** The *GRAC* actor network uses a combination of
 248 two actor loss functions, denoted as $QLoss$ and $CEMLoss$. $QLoss$ refers to the unbiased gradient
 249 estimators which extends the DDPG-style policy gradients [11] to stochastic policies. $CEMLoss$
 250 represents the policy improvement guided by the action found with the zero-order optimization method
 251 CEM. We run another two ablation experiments on all six control tasks and compare it with our
 252 original policy training method denoted as *GRAC*. As seen in Fig.3, in general *GRAC* achieves a better
 253 performance compared to either using $CEMLoss$ or $QLoss$. The significance of the improvements
 254 varies within the six control tasks. For example, $CEMLoss$ plays a dominant role in Swimmer while
 255 $QLoss$ has a major effect in HalfCheetah. It suggests that $CEMLoss$ and $QLoss$ are complementary.

256 **Max-min Double Q-Learning** We additionally verify the effectiveness of the proposed Max-min
 257 Double Q-Learning method. We run an ablation experiment by replacing Max-min by Clipped
 258 Double Q-learning [6] denoted as *GRAC w/o CriticCEM*. In Fig. 4, we visualize the learning curves
 259 of the average return, Q_1 (y'_1 in Alg. 1), and $Q_1 - Q_2$ ($y'_1 - y'_2$ in Alg. 1). *GRAC* achieves high

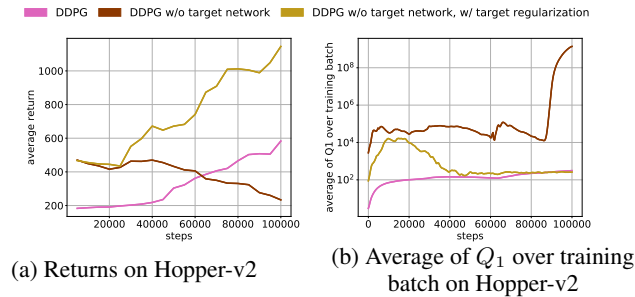


Figure 2: Learning curves and average Q_1 values (y'_1 in Alg. 1) on Hopper-v2. DDPG w/o target network quickly diverges as seen by the unrealistically high Q values. DDPG is stable but progresses slowly. If we remove the target network and add the proposed target regularization, we both maintain stability and achieve faster learning than DDPG.

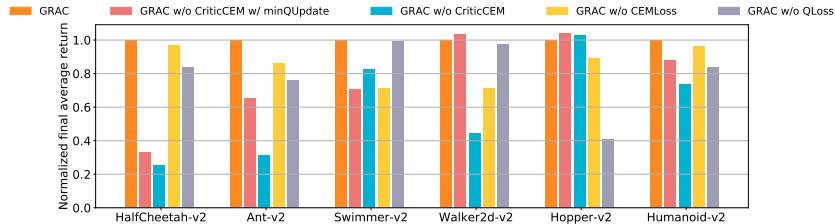


Figure 3: Final average returns, normalized w.r.t *GRAC* for all tasks. For each task, we train each ablation setting with 4 seeds, and average the last 10 evaluations of each seed (40 evaluations in total). Actor updates without CEMLoss (*GRAC w/o CEMLoss*) and actor updates w.r.t minimum of both Q networks (*GRAC w/o CriticCEM w/ minQUpdate*) achieves slightly better performance on Walker2d-v2 and Hopper-v2. *GRAC* achieves the best performance on 4 out of 6 tasks, especially on more complicated tasks with higher-dimensional state and action spaces (Humanoid-v2, Ant-v2, HalfCheetah-v2). This suggests that individual components of *GRAC* complement each other.

260 performance while maintaining a smoothly increasing Q function. Note that the difference between Q
 261 functions, $Q_1 - Q_2$, remains around zero for *GRAC*. *GRAC w/o CriticCEM* shows high variance and
 262 drastic changes in the learned Q_1 value. In addition, Q_1 and Q_2 do not always agree. Such unstable
 263 Q values result in a performance crash during the learning process. Instead of Max-min Double Q
 264 Learning, another way to address the gap between Q_1 and Q_2 is to perform actor updates on the
 265 minimum of Q_1 and Q_2 networks (as seen in SAC). Replacing Max-min Double Q Learning with
 266 this trick achieves lower performance than *GRAC* in more complicated tasks such as HalfCheetah-v2,
 267 Ant-v2, and Humanoid-v2 (See *GRAC w/o CriticCEM w/ minQUpdate* in Fig.3).

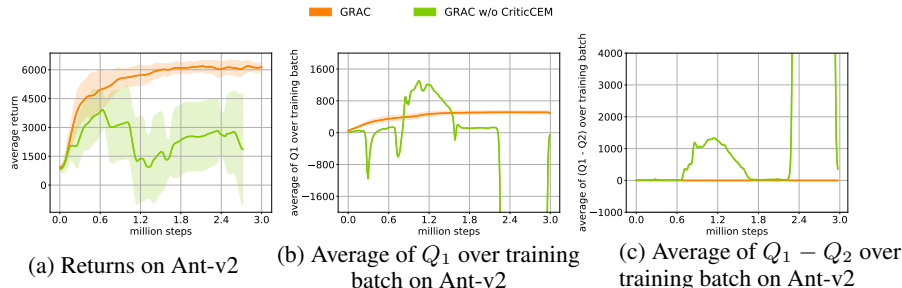


Figure 4: Learning curves (left), average Q_1 values (middle), and average of the difference between Q_1 and Q_2 (right) on Ant-v2. Average Q values are computed as minibatch average of y'_1 and y'_2 , defined in Alg. 1. *GRAC w/o CriticCEM* represents replacing Max-min Double Q-Learning with Clipped Double Q-Learning. Without Max-min double Q-Learning to balance the magnitude of Q_1 and Q_2 , Q_1 blows up significantly compared to Q_2 , leading to divergence.

268 6 Conclusion

269 Leveraging neural networks as function approximators, DRL has been successfully demonstrated on
 270 a range of decision-making and control tasks. However, the nonlinear function approximators also
 271 introduce issues such as divergence and overestimation. We proposed a self-regularized TD-learning
 272 method to address divergence without requiring a target network that may slow down learning
 273 progress. The proposed method is agnostic to the specific Q-learning method and can be added to
 274 any of them. We introduced Max-min Double Q-learning to mitigate over-estimation while reducing
 275 the discrepancy between the two Q functions and to provide a better approximation of the Bellman
 276 optimality operator. In addition, we propose self-guided policy improvement by combining policy-
 277 gradient with zero-order optimization such as the Cross Entropy Method. This helps to search for
 278 actions associated with higher Q-values in a broad neighborhood and is robust to local noise in the Q
 279 function approximation. Taken together, these three components define *GRAC*, a novel self-guided
 280 and self-regularized actor critic algorithm. We evaluate our method on the suite of OpenAI gym tasks,
 281 achieving or outperforming state of the art in every environment tested.

282 7 Broader Impact

283 This work propose new methods to improve reinforcement learning in continuous control tasks. (1)
284 DRL takes a lot of data, thus compute, to train. Our method speeds up training, thus reduces the
285 necessary compute and the corresponding carbon footprint and energy consumption. (2) DRL's
286 successful application in robotics would have societal impacts. On the positive side, automating
287 repetitive manual labour increases productivity and thus increases the wealth of the society. On
288 the other hand, automation may lead to unemployment of workers who have previously performed
289 this manual labour. New policies are required to ensure decent income and access to professional
290 education, so the impacted workers can transition into newly created jobs.

291 References

- 292 [1] Joshua Achiam, Ethan Knight, and Pieter Abbeel. Towards characterizing divergence in deep
293 q-learning. *arXiv preprint arXiv:1903.08894*, 2019.
- 294 [2] Zdravko Botev, Dirk Kroese, Reuven Rubinstein, and Pierre L'Ecuyer. *The Cross-Entropy*
295 *Method for Optimization*, volume 31, pages 35–59. 12 2013.
- 296 [3] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang,
297 and Wojciech Zaremba. Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 298 [4] Yan Duan, Xi Chen, Rein Houthoofd, John Schulman, and Pieter Abbeel. Benchmarking
299 deep reinforcement learning for continuous control. In *International Conference on Machine*
300 *Learning*, pages 1329–1338, 2016.
- 301 [5] Ishan Durugkar and Peter Stone. Td learning with constrained gradients. In *Proceedings of*
302 *the Deep Reinforcement Learning Symposium, NIPS 2017*, Long Beach, CA, USA, December
303 2017.
- 304 [6] Scott Fujimoto, Herke van Hoof, and David Meger. Addressing function approximation error in
305 actor-critic methods. *arXiv preprint arXiv:1802.09477*, 2018.
- 306 [7] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. Soft actor-critic: Off-
307 policy maximum entropy deep reinforcement learning with a stochastic actor. *arXiv preprint*
308 *arXiv:1801.01290*, 2018.
- 309 [8] Hado V Hasselt. Double q-learning. In *Advances in neural information processing systems*,
310 pages 2613–2621, 2010.
- 311 [9] Dmitry Kalashnikov, Alex Irpan, Peter Pastor, Julian Ibarz, Alexander Herzog, Eric Jang,
312 Deirdre Quillen, Ethan Holly, Mrinal Kalakrishnan, Vincent Vanhoucke, et al. Scalable deep
313 reinforcement learning for vision-based robotic manipulation. In *Conference on Robot Learning*,
314 pages 651–673, 2018.
- 315 [10] Seungchan Kim, Kavosh Asadi, Michael Littman, and George Konidaris. Deepmellow: remov-
316 ing the need for a target network in deep q-learning.
- 317 [11] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa,
318 David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. In
319 *International Conference on Learning Representations*, 2016.
- 320 [12] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan
321 Wierstra, and Martin Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint*
322 *arXiv:1312.5602*, 2013.
- 323 [13] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G
324 Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al.
325 Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015.
- 326 [14] Pourchot and Sigaud. CEM-RL: Combining evolutionary and gradient-based methods for policy
327 search. In *International Conference on Learning Representations*, 2019.

- 328 [15] R. Y. Rubinstein. The cross-entropy method for combinatorial and continuous optimization.
329 *Methodology and Computing in Applied Probability*, 1:127—190, 1999.
- 330 [16] John Schulman, Sergey Levine, Pieter Abbeel, Michael Jordan, and Philipp Moritz. Trust region
331 policy optimization. In *International conference on machine learning*, pages 1889–1897, 2015.
- 332 [17] Richard S Sutton et al. *Introduction to reinforcement learning*, volume 135.
- 333 [18] Sebastian Thrun and Anton Schwartz. Issues in using function approximation for reinforcement
334 learning.
- 335 [19] Emanuel Todorov, Tom Erez, and Yuval Tassa. Mujoco: A physics engine for model-based
336 control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages
337 5026–5033. IEEE, 2012.
- 338 [20] John N Tsitsiklis and Benjamin Van Roy. Analysis of temporal-difference learning with function
339 approximation. In *Advances in neural information processing systems*, pages 1075–1081, 1997.
- 340 [21] Hado Van Hasselt, Yotam Doron, Florian Strub, Matteo Hessel, Nicolas Sonnerat, and Joseph
341 Modayil. Deep reinforcement learning and the deadly triad. *arXiv preprint arXiv:1812.02648*,
342 2018.
- 343 [22] C. J. C. H. Watkins. *Learning from Delayed Rewards*. PhD thesis, King’s College, Oxford,
344 1989.
- 345 [23] M. Yan, A. Li, M. Kalakrishnan, and P. Pastor. Learning probabilistic multi-modal actor
346 models for vision-based robotic grasping. In *2019 International Conference on Robotics and*
347 *Automation (ICRA)*, pages 4804–4810, 2019.