

GRAC: Self-Guided and Self-Regularized Actor-Critic

Anonymous Author(s)

Affiliation

Address

email

1 **Abstract:** Deep reinforcement learning (DRL) algorithms have successfully been
2 demonstrated on a range of challenging decision making and control tasks. One
3 dominant component of recent deep reinforcement learning algorithms is the target
4 network which mitigates the divergence when learning the Q function. However,
5 target networks can slow down the learning process due to delayed function updates.
6 Our main contribution in this work is a self-regularized TD-learning method to
7 address divergence without requiring a target network. Additionally, we propose a
8 self-guided policy improvement method by combining policy-gradient with zero-
9 order optimization to search for actions associated with higher Q-values in a broad
10 neighborhood. This makes learning more robust to local noise in the Q function
11 approximation and guides the updates of our actor network. Taken together, these
12 components define *GRAC*, a novel self-guided and self-regularized actor critic
13 algorithm. We evaluate *GRAC* on the OpenAI gym tasks, outperforming popular
14 methods such as *TD3* [1] and *SAC* [2] on four tasks and achieving competitive re-
15 sults on two environments. We also apply *GRAC* to enable a non-anthropomorphic
16 robotic hand to successfully accomplish an in-hand manipulation task in the real
17 world.

18 **Keywords:** Deep Reinforcement Learning, Q-learning

19 1 Introduction

20 Reinforcement learning (RL) studies decision-making with the goal of maximizing total discounted
21 reward when interacting with an environment. Leveraging high-capacity function approximators such
22 as neural networks, Deep reinforcement learning (DRL) algorithms have been successfully applied to
23 a range of challenging domains, from video games [3] to robotic control [4].

24 Actor-critic algorithms are among the most popular approaches in DRL, e.g. *DDPG* [5], *TRPO* [4],
25 *TD3* [1] and *SAC* [2]. These methods are based on policy iteration, which alternates between policy
26 evaluation and policy improvement [6]. Actor-critic methods jointly optimize the value function
27 (critic) and the policy (actor) as it is often impractical to run either of these to convergence [2].

28 In DRL, both the actor and critic use deep neural networks as the function approximator. However,
29 DRL is known to assign unrealistically high values to state-action pairs represented by the Q-function.
30 This is detrimental to the quality of the greedy control policy derived from Q [7]. Mnih et al. [8]
31 proposed to use a *target network* to mitigate divergence. A target network is a copy of the current Q
32 function that is held fixed to serve as a stable target within the TD error update. The parameters of the
33 target network are either infrequently copied [8] or obtained by Polyak averaging [5]. A limitation of
34 using a target network is that it can slow down learning due to delayed function updates. We propose
35 an approach that reduces the need for a target network in DRL while still ensuring stable learning
36 and good performance in high-dimensional domains. We add a self-regularization term to encourage
37 small changes to the target value while minimizing the *Temporal Difference* (TD)-error [6].

38 *Evolution Strategies* (ES) are a family of black-box optimization algorithms which are typically
39 very stable, but scale poorly in high-dimensional search spaces e.g. neural networks [9]. Policy
40 gradient-based DRL methods, unlike evolutionary search methods, can continue to sample previous

41 experiences to improve value estimation, particularly in the off-policy setting. At the same time,
42 these approaches can also be unstable and highly sensitive to hyper-parameter tuning [9]. We
43 propose a novel policy improvement method which combines both approaches to get the best of both
44 worlds. Specifically, after the actor network first outputs an initial action, we apply the *Cross Entropy*
45 *Method* (CEM) [10] to search the neighborhood of the initial action to find a second action associated
46 with a higher Q value. Then we leverage the second action in the policy improvement stage to speed
47 up the learning process.

48 Our main contribution in this work is a **self-regularized** TD-learning method to address divergence
49 without requiring a target network that may slow down learning progress. In addition, we propose a
50 **self-guided** policy improvement method which combines policy-gradients and zero-order optimiza-
51 tion. This helps to speed up learning and is robust to local noise in the Q function approximation.
52 Taken together, these components define *GRAC*, a novel self-guided and self-regularized actor critic
53 algorithm. We evaluate *GRAC* on six continuous control domains from OpenAI gym [11], where
54 *GRAC* outperforms popular methods such as *TD3* [1] and *SAC* [2] on four tasks and achieves com-
55 petitive results on two environments. We run our experiments across a large number of seeds with
56 fair evaluation metrics [12], perform extensive ablation studies, and open source both our code
57 and learning curves. We also run *GRAC* to enable a non-anthropomorphic, real robotic hand to
58 successfully rotate a cube to the target pose.

59 2 Related Work

60 The proposed algorithm incorporates key ingredients within the actor-critic method: a self-regularized
61 TD update and self-guided policy improvements based on evolution strategies. In this section, we
62 review prior work related to these ideas.

63 **Divergence in Deep Q-Learning** In Deep Q-Learning, we use a nonlinear function approximator
64 such as a neural network to approximate the Q-function that represents the value of each state-action
65 pair. Learning the Q-function in this way is known to suffer from divergence issues [13] such as
66 assigning unrealistically high values to state-action pairs [7]. For the case when the control policy is
67 greedily derived from Q [7], unrealistically high Q values are detrimental. To mitigate the divergence
68 issue, Mnih et al. [8] introduce a target network which is a copy of the estimated Q-function and
69 is held fixed to serve as a stable target for some number of steps. However, because of the delayed
70 function updates, target networks can slow down learning [14]. Durugkar and Stone [15] propose
71 Constrained Q-Learning, which uses a constraint to prevent the average target value from changing
72 after an update. Achiam et al. [16] give a simple analysis based on a linear approximation of the
73 Q function and develop a stable Deep Q-Learning algorithm for continuous control without target
74 networks. However, their proposed method requires separately calculating backward passes for each
75 state-action pair in the batch, and solving a system of equations at each timestep. Bhatt et al. [17]
76 introduce a normalization called cross-normalization which is regarded as an extension of batch
77 normalization that re-centers data for on- and off-policy transitions. Peng et al. [18] uses Monte
78 Carlo instead of TD error to update the Q network and removes the need for a target network. Our
79 proposed *GRAC* algorithm adds a self-regularization term to the TD-Learning objective to keep the
80 change of the state-action value small.

81 **Evolution Strategies in Deep Reinforcement Learning** Evolution Strategies are typically stable
82 but suffer from scaling to high-dimensional search spaces. Policy gradient-based deep RL methods,
83 such as *DDPG* [5] can continue to reuse previous experience to improve value estimations, particularly
84 in the off-policy setting, but can be unstable and highly sensitive to hyper-parameter tuning [9].
85 Researchers have proposed to combine these approaches to get the best of both worlds. Pourchot
86 and Sigaud [9] proposed *CEM-RL* to combine CEM with either *DDPG* [5] or *TD3* [1]. However,
87 *CEM-RL* applies *CEM* within the actor parameter space which is extremely high-dimensional, making
88 the search not efficient. Kalashnikov et al. [19] introduce *QT-Opt*, which leverages *CEM* to search
89 the landscape of the Q function, and enables Q-Learning in continuous action spaces without using
90 an actor. Based on *QT-Opt*, Simmons-Edler et al. [20] leverage CEM to search the landscape the
91 Q-function but propose to initialize CEM with a Gaussian prior covering the action space, independent
92 of observations. However, *CEM* does not scale well to high-dimensional action spaces [21], such as
93 in the Humanoid task. We first let the actor network output an initial Gaussian action distribution
94 conditioned on current state. Then we use CEM to search for an action with a higher Q value than the

95 Q value of the Gaussian mean. Starting from the predicted distribution, we show that *GRAC* speeds
 96 up the learning process compared to popular actor-critic methods.

97 3 Preliminaries

98 In this section, we define the notation used in subsequent sections. Consider a *Markov Decision*
 99 *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
 100 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
 101 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
 102 factor. At each discrete time step t , with a given state $s_t \in \mathcal{S}$, the agent selects an action $a_t \in \mathcal{A}$,
 103 receiving a reward r and the new state s_{t+1} of the environment.

104 Let π denote the policy which maps a state to a probability distribution over the actions, $\pi : \mathcal{S} \rightarrow$
 105 $\mathcal{P}(\mathcal{A})$. The return from a state is defined as the sum of discounted reward $R_t = \sum_{i=t} \gamma^{i-t} r(s_i, a_i)$.
 106 In reinforcement learning, the objective is to find the optimal policy π^* , with parameters ϕ , which max-
 107 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 $\rho_\pi(s_t, a_t)$
 108 denote the state and state-action marginals of the trajectory distribution induced by the policy $\pi(a_t | s_t)$.

109 We use the following standard definitions of the state-action value function Q_π . It describes the
 110 expected discounted reward after taking an action a_t in state s_t and thereafter following policy π :

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

111 In this work we use *CEM* to find optimal actions with maximum Q values. *CEM* is a randomized zero-
 112 order optimization algorithm. To find the action a that maximizes $Q(s, a)$, *CEM* is initialized with a
 113 parameterized distribution over a , $P(a; \psi)$. Then it iterates between the following two steps [22]:
 114 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
 115 decreasing Q values. Then keep the first K actions, and solve for updated parameters ψ' :

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

116 In the following, let $CEM(Q(s, \cdot), \pi(\cdot | s))$ denote the action found by *CEM* to maximize $Q(s, \cdot)$,
 117 when *CEM* is initialized with the distribution predicted by the policy.

118 4 Technical Approach

119 4.1 Self-Regularized TD Learning

120 Reinforcement learning is prone to instability and divergence when a nonlinear function approximator
 121 such as a neural network is used to represent the Q function [13]. Mnih et al. [8] identified several
 122 reasons for this. One is the correlation between the current action-values and the target value. Updates
 123 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
 124 updates also increase the target value y_t which may lead to oscillations or the divergence of the policy.

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

$$\mathcal{L}(\theta_i) = \mathbb{E}_{(s_t, a_t) \sim \mathcal{B}} \|(Q(s_t, a_t; \theta_i) - y_i)\|^2 \quad (1)$$

127 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
 128 based on the current Q network parameters θ_i . $a_{t+1}^* = \operatorname{argmax}_a Q(s_{t+1}, a)$. If we update the
 129 parameter θ_{i+1} to reduce the loss $\mathcal{L}(\theta_i)$, it changes both $Q(s_t, a_t; \theta_{i+1})$ and $Q(s_{t+1}, a_{t+1}^*; \theta_{i+1})$.
 130 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
 131 the next iteration will also increase leading to an explosion of the Q function. We demonstrated this
 132 behavior in an ablation experiment with results in Fig. 2. We also show how maintaining a separate
 133 target network [8] with frozen parameters θ^- to compute $y_{i+1} = r_t + \gamma Q(s_{t+1}, a_{t+1}^*; \theta^-)$ delays the
 134 update of the target and therefore leads to more stable learning of the Q function. However, delaying
 135 the function updates also comes with the price of slowing down the learning process.

136 We propose a self-Regularized TD-learning approach to minimize the TD-error while also keeping
 137 the change of $Q(s_{t+1}, a_{t+1}^*)$ small. This regularization mitigates the divergence issue [13], and

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 $a \sim \pi_{\phi_i}(s)$, 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: $y \leftarrow r_t + \gamma \max_{j=1,2} \{ \min_{j=1,2} Q(s_{t+1}, \tilde{a}_{t+1}; \theta_j), \min_{j=1,2} Q(s_{t+1}, \hat{a}_{t+1}; \theta_j) \}$
- 8: $a^\dagger \leftarrow \arg \max_{\{\tilde{a}, \hat{a}\}} \{ \min_{j=1,2} Q(s_{t+1}, \tilde{a}_{t+1}; \theta_j), \min_{j=1,2} Q(s_{t+1}, \hat{a}_{t+1}; \theta_j) \}$
- 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: $\mathcal{L}_k = \|y - Q(s_t, a_t; \theta_1)\|_2^2 + \|y - Q(s_t, a_t; \theta_2)\|_2^2 + \|y'_1 - Q(s_{t+1}, a^\dagger; \theta_1)\|_2^2 + \|y'_2 - Q(s_{t+1}, a^\dagger; \theta_2)\|_2^2$
- 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** $\mathcal{L}_k < \alpha \mathcal{L}_1$ **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**

138 no longer requires a target network that would otherwise slow down the learning process. Let
139 $y'_i = Q(s_{t+1}, a_{t+1}^*; \theta_i)$, and $y_i = r_t + \gamma y'_i$. We define the learning objective as

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

140 where the first term is the original TD-Learning objective and the second term is the regularization
141 term penalizing large updates to $Q(s_{t+1}, a_{t+1}^*)$. Note that when the current Q network updates its
142 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
143 which is different from the approach of keeping a frozen target network for a few iterations. We will
144 demonstrate in our experiments that this self-regularized TD-Learning approach removes the delays
145 in the update of the target value thereby achieves faster and stable learning.

146 4.2 Self-Guided Policy Improvement with Evolution Strategies

147 The policy, known as the actor, can be updated through a combination of two parts. The first part,
148 which we call Q-loss policy update, improves the policy through local gradients of the current Q
149 function, while the second part, which we call *CEM* policy update, finds a high-value action via
150 *CEM* in a broader neighborhood of the Q function landscape and updates the action distribution to be
151 around this high-value action. We formally describe the two parts below.

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

$$J_\pi(\phi) = \mathbb{E}_{s_t \sim \mathcal{B}, \hat{a}_t \sim \pi} [Q(s_t, \hat{a}_t)], \quad (3)$$

153 where \hat{a}_t is sampled from the current policy $\pi(\cdot | s_t)$. The gradient is taken through the reparameteri-
154 zation trick. We reparameterize the policy using a neural network transformation as described by
155 Haarnoja et al. [2],

$$\hat{a}_t = f_\phi(\epsilon_t | s_t) \quad (4)$$

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

$$\nabla J_\pi(\phi) = \mathbb{E}_{s_t \sim \mathcal{B}, \epsilon_t \sim \mathcal{N}} \left[\frac{\partial Q(s_t, f_\phi(\epsilon_t | s_t))}{\partial f} \frac{\partial f_\phi(\epsilon_t | s_t)}{\partial \phi} \right] \quad (5)$$

158 For the CEM policy update, given a minibatch of states s_t , we first find a high-value action \bar{a}_t for
 159 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
 160 is updated to increase the probability of this high-value action. The guided update on the parameter ϕ
 161 of π at iteration i is

$$\mathbb{E}_{s_t \sim \mathcal{B}, \hat{a}_t \sim \pi} [Q(s_t, \bar{a}_t) - Q(s_t, \hat{a}_t)]_+ \nabla_\phi \log \pi_i(\bar{a}_t | s_t). \quad (6)$$

162 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
 163 baseline $V(s_t)$:

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

164 In our implementation, we only perform an update if the improvement on the Q function, $Q(s_t, \bar{a}_t) -$
 165 $Q(s_t, \hat{a}_t)$, is non-negative, to guard against the occasional cases where *CEM* fails to find a better
 166 action.

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

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

168 where λ is the step size.

169 Let Q^π be the state-action value function induced by the current policy. We can prove that if the Q
 170 function has converged to Q^π then both the Q-loss policy update and the *CEM* policy update will be
 171 guaranteed to improve the current policy. We formalize this result in Theorem 1 and Theorem 2, and
 172 prove them in Appendix.

173 **Theorem 1. Q-loss Policy Improvement** Starting from the current policy π , we maximize the objec-
 174 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} .
 175 Then the induced Q function, $Q^{\pi_{new}}$, satisfies $\forall (s, a), Q^{\pi_{new}}(s, a) \geq Q^\pi(s, a)$.

176 **Theorem 2. CEM Policy Improvement** Assuming the CEM process is able to find the optimal
 177 action of the state-action value function, $a^*(s) = \arg \max_a Q^\pi(s, a)$, where Q^π is the Q function
 178 induced by the current policy π . By iteratively applying the update $\mathbb{E}_{(s,a) \sim \rho_\pi(s,a)} [Q(s, a^*) -$
 179 $Q(s, a)]_+ \nabla \log \pi(a^* | s)$, the policy converges to π_{new} . Then $Q^{\pi_{new}}$ satisfies $\forall (s, a), Q^{\pi_{new}}(s, a) \geq$
 180 $Q^\pi(s, a)$.

181 4.3 Max-min Double Q-Learning

182 Q-learning [23] is known to suffer from overestimation [24]. Hasselt [25] proposed Double-Q
 183 learning which uses two Q functions with independent sets of weights to mitigate the overestimation
 184 problem. Fujimoto et al. [1] proposed Clipped Double Q-learning with two Q functions denoted
 185 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
 186 Double Q-learning uses the minimum between the two estimates of the Q functions when calculating
 187 the target value in TD-error [6]:

$$y = r_t + \gamma \min_{j=1,2} Q(s_{t+1}, \hat{a}_{t+1}; \theta_j) \quad (8)$$

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

189 Fujimoto et al. [1] mentioned that such an update rule may induce an underestimation bias. In
 190 addition, $\hat{a}_{t+1} = \pi_\phi(s_{t+1})$ is the prediction of the actor network. The actor network's parameter
 191 ϕ is optimized according to the gradients of Q_1 . In other words, \hat{a}_{t+1} tends to be selected according
 192 to the Q_1 network which consistently increases the discrepancy between the two Q-functions. In
 193 practice, we observe that the discrepancy between the two estimates of the Q-function, $|Q_1 - Q_2|$,
 194 can increase dramatically leading to an unstable learning process.

195 We found that a technique we call *Max-min Double Q-Learning* reduces the discrepancy between
 196 the Q-functions. We first select \hat{a}_{t+1} according to the actor network $\pi_\phi(s_{t+1})$. Then we run *CEM*

197 to search the landscape of Q_2 within a broad neighborhood of \hat{a}_{t+1} to return a second action \tilde{a}_{t+1} .
 198 Note that *CEM* selects an action \tilde{a}_{t+1} that maximises Q_2 while the actor network selects an action
 199 \hat{a}_{t+1} that maximises Q_1 . We gather four different Q-values: $Q(s_{t+1}, \hat{a}_{t+1}; \theta_1)$, $Q(s_{t+1}, \hat{a}_{t+1}; \theta_2)$,
 200 $Q(s_{t+1}, \tilde{a}_{t+1}; \theta_1)$, and $Q(s_{t+1}, \tilde{a}_{t+1}; \theta_2)$. We then run a max-min operation to compute the target
 201 value that cancels the biases induced by \hat{a}_{t+1} and \tilde{a}_{t+1} .

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

202 The inner min-operation $\min_{j=1,2} Q(s_{t+1}, \hat{a}_{t+1}; \theta_j)$ is adopted from Eq. 8 and mitigates overestima-
 203 tion [24]. The outer max operation helps to reduce the difference between Q_1 and Q_2 . In addition, the
 204 max operation provides a better approximation of the Bellman optimality operator [6]. We visualize
 205 Q_1 and Q_2 during the learning process in the supplementary material. We formalize the convergence
 206 of the proposed Max-min Double Q-Learning approach in the finite MDP setting and prove this
 207 theorem in the Appendix.

208 5 Experiments

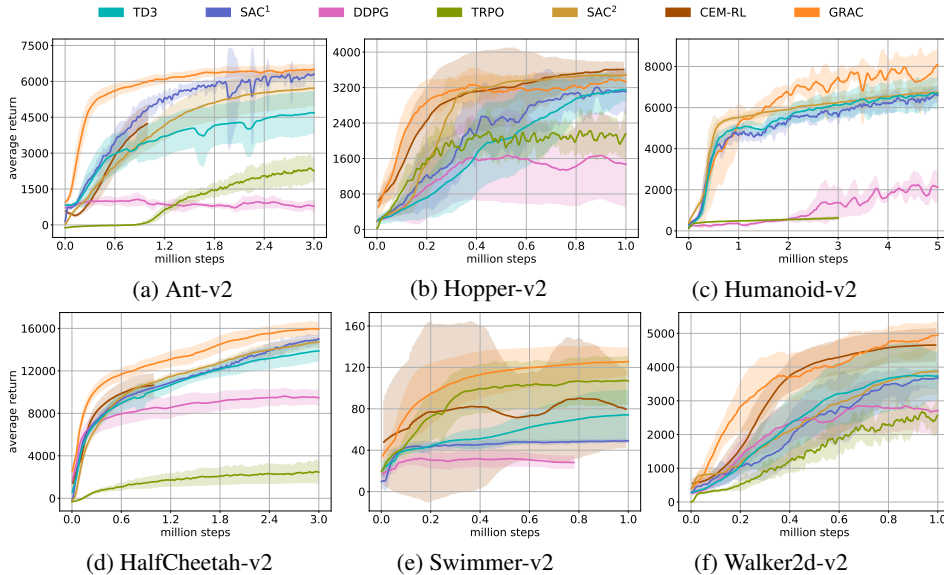


Figure 1: Learning curves for the OpenAI gym continuous control tasks. For each task, we train 10 instances of each algorithm, using 10 different seeds. Evaluations are performed every 5000 interactions with the environment. Each evaluation reports the return (total reward), averaged over 10 episodes. 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 of the average return over 10 seeds. All curves are smoothed with window size 10 for visual clarity (we graph the same learning curves without smoothing in the appendix.) *GRAC* (orange) learns faster than other methods across all tasks. *GRAC* achieves comparable result to the popular methods on the Hopper-v2 task and the Ant-v2 task and outperforms prior methods on the other four tasks including the complex high-dimensional Humanoid-v2.

209 5.1 Comparative Evaluation

210 We present *GRAC*, a self-guided and self-regularized actor-critic algorithm as summarized in Algo-
 211 rithm 1. To evaluate *GRAC*, we measure its performance on the suite of MuJoCo continuous control
 212 tasks [26], interfaced through OpenAI Gym [11]. We compare our method with *DDPG* [5], *TD3* [1],
 213 *CEM-RL* [9], *TRPO* [4], *SAC*¹ [2], and *SAC*² [27]. We use the source code released by the original
 214 authors and adopt the same hyperparameters reported in the original papers and the number of training
 215 steps according to *SAC*¹ [2]. For *CEM-RL* [9] and *SAC*² [27], we use the results provided by their

216 corresponding authors. SAC^2 does not contain results on Swimmer-v2. $CEM-RL$ does not contain
 217 results on Humanoid-v2 and only runs one millions step on every tested task. Hyperparameters
 218 for all experiments are in the Appendix. Results are shown in Figure 1. $GRAC$ outperforms or is
 219 comparable to all other algorithms in both final performance and learning speed across all tasks.
 220 On complex tasks with high state and action dimensions such as Humanoid-v2, $GRAC$ outperforms
 221 all other algorithms by a large margin. Both $CEM-RL$ and $GRAC$ leverage CEM. $CEM-RL$ applies
 222 CEM within the actor parameter space which is extremely high-dimensional while $GRAC$ utilizes
 223 CEM within the action space of the Q function. $GRAC$ outperform $CEM-RL$ on tasks such as Ant-v2,
 224 HalfCheetah-v2, Swimmer-v2 by a large margin.

225 5.2 Ablation Study

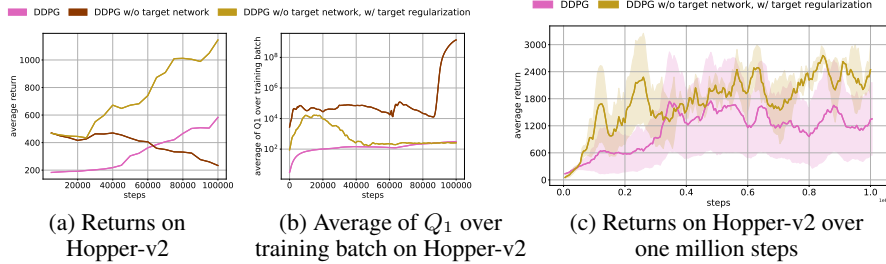


Figure 2: Learning curves and average Q_1 values (y'_1 in Alg. 1) on Hopper-v2. $DDPG$ w/o target network quickly diverges as seen by the unrealistically high Q values. $DDPG$ is stable but progresses slowly. If we remove the target network and add the proposed target regularization, we both maintain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

226 In this section, we present ablation studies to under-
 227 stand the contribution of each proposed component: Self-
 228 Regularized TD-Learning (Section 4.1) and Self-Guided
 229 Policy Improvement (Section 4.2). We present our results
 230 in Fig. 4 in which we compare the performance of $GRAC$
 231 with alternatives, each removing one component from
 232 $GRAC$. Additional learning curves can be found in the Appen-
 233 dix. We also run experiments to examine how sensitive
 234 $GRAC$ is to some hyperparameters such as α and K listed
 235 in Alg. 1, and the results can be found in the Appendix.

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

243 In Fig. 2, we visualize the average returns, and average Q_1 values over training batches (y'_1 in Alg.1).
 244 The Q_1 values of $DDPG$ w/o target network changes dramatically which leads to poor average returns.
 245 $DDPG$ maintains stable Q values but makes slow progress. Our proposed $DDPG$ w/o target network
 246 w/ target regularization maintains stable Q values. In addition, we compare the average returns of
 247 $DDPG$ w/o target network, w/ target regularization and $DDPG$ within one million steps over four
 248 random seeds. $DDPG$ w/o target network, w/ target regularization outperforms $DDPG$ by large
 249 margins in five out of six Mujoco tasks (Fig. 2 shows results on Hopper-v2. The remaining results
 250 can be found in the Appendix). We also apply self-regularized TD-learning to DQN called DQN w/o
 251 target network w/ target regularization on the Atari Breakout environment and it outperforms DQN
 252 by 25%. The learning curve over 50 million steps is shown in Fig.3. These results on two different
 253 Q-learning methods demonstrate the effectiveness of our method and its potentials to be applied to a
 254 wide range of DRL methods.

255 **Policy Improvement with Evolution Strategies** The $GRAC$ actor network uses a combination of
 256 two actor loss functions, denoted as $QLoss$ and $CEMLoss$. $QLoss$ refers to the unbiased gradient

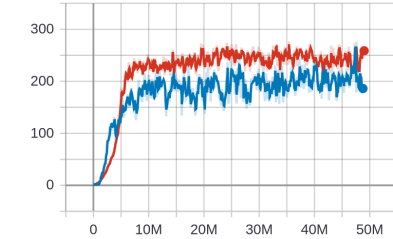


Figure 3: Average returns for the BreakoutNoFrameskip-v4 environment on OpenAI gym. Blue is DQN , red is DQN w/o target network w/ target regularization.

257 estimators which extend the *DDPG*-style policy gradients [5] to stochastic policies. *CEMLoss*
 258 represents the policy improvement guided by the action found with the zero-order optimization
 259 method CEM. We run another two ablation experiments on all six control tasks and compare it
 260 with our original policy training method denoted as *GRAC*. As seen in Fig. 4, in general *GRAC*
 261 achieves a better performance compared to either using *CEMLoss* or *QLoss*. The significance of
 262 the improvements varies over the six control tasks. For example, *CEMLoss* plays a dominant role
 263 in Swimmer while *QLoss* has a major effect in HalfCheetah. This suggests that the *CEMLoss* and
 264 *QLoss* are complementary.

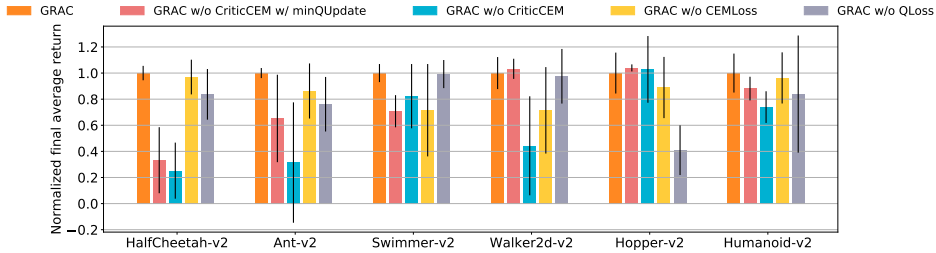


Figure 4: Final average returns, normalized w.r.t *GRAC* for all tasks. For each task, we train each ablation setting with 4 seeds, and average the last 10 evaluations of each seed (40 evaluations in total). The black lines represent one standard deviation. Actor updates without CEMLoss (*GRAC w/o CEMLoss*) 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 challenging tasks with higher-dimensional state and action spaces (Humanoid-v2, Ant-v2, HalfCheetah-v2). This suggests that individual components of *GRAC* complement each other.

265 5.3 In Hand Manipulation

266 We evaluate our approach on the problem of in-hand manipulation
 267 which remains unsolved due the high dimensionality of the prob-
 268 lem and the complexity of multi-contact control [28]. Having this
 269 capability would allow robots to perform sophisticated tasks re-
 270 quiring repositioning and reorienting of grasped objects. We apply
 271 *GRAC* to a non-anthropomorphic robotic hand with 9 degrees of
 272 freedom [29]. We train the robotic hand in simulation to rotate a
 273 cube by 50 degrees around the gravity direction. The input states
 274 are the current, previous, initial, and target object position and
 275 orientation, and all nine gripper joint positions at the current time
 276 step. The actions are nine joint positions for the gripper at the next
 277 time step. *GRAC* learns a policy successfully accomplishing the
 278 task within 500k iterations. We test the learned policy on the real
 279 hand. Videos are included in the supplementary material.

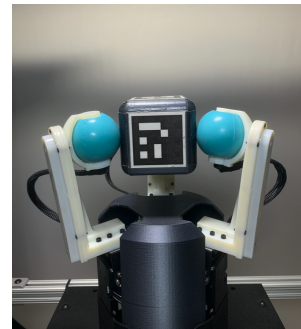


Figure 5: Real-world Experiment Setup

280 6 Conclusion

281 Leveraging neural networks as function approximators, DRL has been successfully demonstrated
 282 on a range of decision-making and control tasks. However, the nonlinear function approximators
 283 also introduce issues such as divergence and overestimation. As our main contribution, we proposed
 284 a self-regularized TD-learning method to address divergence without requiring a target network
 285 that may slow down learning progress. The proposed method is agnostic to the specific Q-learning
 286 method and can be added to any of them. In addition, we propose self-guided policy improvement by
 287 combining policy-gradient with zero-order optimization such as the Cross Entropy Method. This
 288 helps to search for actions associated with higher Q-values in a broad neighborhood and is robust to
 289 local noise in the Q function approximation. Taken together, these components define *GRAC*, a novel
 290 self-guided and self-regularized actor critic algorithm. We evaluate *GRAC* on the OpenAI gym tasks,
 291 outperforming popular methods such as *TD3* [1] and *SAC* [2] on four tasks and achieving competitive
 292 results on two environments. We also run *GRAC* to enable a non-anthropomorphic robotic hand to
 293 successfully accomplish an in-hand manipulation task in the real world.

294 **References**

- 295 [1] S. Fujimoto, H. van Hoof, and D. Meger. Addressing function approximation error in actor-critic
296 methods. *arXiv preprint arXiv:1802.09477*, 2018.
- 297 [2] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy
298 deep reinforcement learning with a stochastic actor. *arXiv preprint arXiv:1801.01290*, 2018.
- 299 [3] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller.
300 Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- 301 [4] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz. Trust region policy optimization.
302 In *International conference on machine learning*, pages 1889–1897, 2015.
- 303 [5] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra.
304 Continuous control with deep reinforcement learning. In *International Conference on Learning
305 Representations*, 2016.
- 306 [6] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. A Bradford Book,
307 Cambridge, MA, USA, 2018. ISBN 0262039249.
- 308 [7] H. Van Hasselt, Y. Doron, F. Strub, M. Hessel, N. Sonnerat, and J. Modayil. Deep reinforcement
309 learning and the deadly triad. *arXiv preprint arXiv:1812.02648*, 2018.
- 310 [8] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves,
311 M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al. Human-level control through deep rein-
312 forcement learning. *Nature*, 518(7540):529, 2015.
- 313 [9] Pourchot and Sigaud. CEM-RL: Combining evolutionary and gradient-based methods for
314 policy search. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=BkeU5j0ctQ>.
315
- 316 [10] R. Y. Rubinstein. The cross-entropy method for combinatorial and continuous optimization.
317 *Methodology and Computing in Applied Probability*, 1999.
- 318 [11] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba.
319 Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 320 [12] Y. Duan, X. Chen, R. Houthoof, J. Schulman, and P. Abbeel. Benchmarking deep reinforcement
321 learning for continuous control. In *International Conference on Machine Learning*, pages 1329–
322 1338, 2016.
- 323 [13] J. N. Tsitsiklis and B. Van Roy. Analysis of temporal-difference learning with function
324 approximation. In *Advances in neural information processing systems*, pages 1075–1081, 1997.
- 325 [14] S. Kim, K. Asadi, M. Littman, and G. Konidaris. Deepmellow: Removing the need for a target
326 network in deep q-learning. In *Proceedings of the Twenty-Eighth International Joint Conference
327 on Artificial Intelligence, IJCAI-19*, pages 2733–2739. International Joint Conferences on
328 Artificial Intelligence Organization, 7 2019. doi:10.24963/ijcai.2019/379. URL <https://doi.org/10.24963/ijcai.2019/379>.
329
- 330 [15] I. Durugkar and P. Stone. Td learning with constrained gradients. In *Proceedings of the Deep
331 Reinforcement Learning Symposium, NIPS 2017*, Long Beach, CA, USA, December 2017. URL
332 <http://www.cs.utexas.edu/users/ai-lab?NIPS17-ishand>.
- 333 [16] J. Achiam, E. Knight, and P. Abbeel. Towards characterizing divergence in deep q-learning.
334 *arXiv preprint arXiv:1903.08894*, 2019.
- 335 [17] A. Bhatt, M. Argus, A. Amiranashvili, and T. Brox. Crossnorm: Normalization for off-policy td
336 reinforcement learning. *arXiv preprint arXiv:1902.05605*, 2019.
- 337 [18] X. B. Peng, A. Kumar, G. Zhang, and S. Levine. Advantage-weighted regression: Simple and
338 scalable off-policy reinforcement learning. *arXiv preprint arXiv:1910.00177*, 2019.

- 339 [19] D. Kalashnikov, A. Irpan, P. Pastor, J. Ibarz, A. Herzog, E. Jang, D. Quillen, E. Holly, M. Kalakrishnan, V. Vanhoucke, et al. Scalable deep reinforcement learning for vision-based robotic manipulation. In *Conference on Robot Learning*, pages 651–673, 2018.
- 340
- 341
- 342 [20] R. Simmons-Edler, B. Eisner, E. Mitchell, S. Seung, and D. Lee. Q-learning for continuous actions with cross-entropy guided policies. *arXiv preprint arXiv:1903.10605*, 2019.
- 343
- 344 [21] M. Yan, A. Li, M. Kalakrishnan, and P. Pastor. Learning probabilistic multi-modal actor models for vision-based robotic grasping. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 4804–4810. IEEE, 2019.
- 345
- 346
- 347 [22] Z. Botev, D. Kroese, R. Rubinstein, and P. L’Ecuyer. *The Cross-Entropy Method for Optimization*, volume 31, pages 35–59. 12 2013.
- 348
- 349 [23] C. J. C. H. Watkins. *Learning from Delayed Rewards*. PhD thesis, King’s College, Oxford, 1989.
- 350
- 351 [24] S. Thrun and A. Schwartz. Issues in using function approximation for reinforcement learning. In M. Mozer, P. Smolensky, D. Touretzky, J. Elman, and A. Weigend, editors, *Proceedings of the 1993 Connectionist Models Summer School*, pages 255–263. Lawrence Erlbaum, 1993. URL http://www.ri.cmu.edu/pub_files/pub1/thrun_sebastian_1993_1/thrun_sebastian_1993_1.pdf.
- 352
- 353
- 354
- 355
- 356 [25] H. V. Hasselt. Double q-learning. In *Advances in neural information processing systems*, pages 2613–2621, 2010.
- 357
- 358 [26] E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–5033. IEEE, 2012.
- 359
- 360
- 361 [27] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta, P. Abbeel, et al. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*, 2018.
- 362
- 363
- 364 [28] O. M. Andrychowicz, B. Baker, M. Chociej, R. Józefowicz, B. McGrew, J. Pachocki, A. Petron, M. Plappert, G. Powell, A. Ray, J. Schneider, S. Sidor, J. Tobin, P. Welinder, L. Weng, and W. Zaremba. Learning dexterous in-hand manipulation. *The International Journal of Robotics Research*, 39(1):3–20, 2020. doi:10.1177/0278364919887447.
- 365
- 366
- 367
- 368 [29] S. Yuan, L. Shao, C. L. Yako, A. Gruebele, and J. K. Salisbury. Design and control of roller grasper v2 for in-hand manipulation. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2020.
- 369
- 370

371 References

- 372 [1] S. Fujimoto, H. van Hoof, and D. Meger. Addressing function approximation error in actor-critic methods. *arXiv preprint arXiv:1802.09477*, 2018.
- 373
- 374 [2] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. *arXiv preprint arXiv:1801.01290*, 2018.
- 375
- 376 [3] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*, 2013.
- 377
- 378 [4] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz. Trust region policy optimization. In *International conference on machine learning*, pages 1889–1897, 2015.
- 379
- 380 [5] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra. Continuous control with deep reinforcement learning. In *International Conference on Learning Representations*, 2016.
- 381
- 382
- 383 [6] R. S. Sutton and A. G. Barto. *Reinforcement Learning: An Introduction*. A Bradford Book, Cambridge, MA, USA, 2018. ISBN 0262039249.
- 384

- 385 [7] H. Van Hasselt, Y. Doron, F. Strub, M. Hessel, N. Sonnerat, and J. Modayil. Deep reinforcement
386 learning and the deadly triad. *arXiv preprint arXiv:1812.02648*, 2018.
- 387 [8] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves,
388 M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al. Human-level control through deep rein-
389 forcement learning. *Nature*, 518(7540):529, 2015.
- 390 [9] Pourchot and Sigaud. CEM-RL: Combining evolutionary and gradient-based methods for
391 policy search. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=BkeU5j0ctQ>.
392
- 393 [10] R. Y. Rubinstein. The cross-entropy method for combinatorial and continuous optimization.
394 *Methodology and Computing in Applied Probability*, 1999.
- 395 [11] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba.
396 Openai gym. *arXiv preprint arXiv:1606.01540*, 2016.
- 397 [12] Y. Duan, X. Chen, R. Houthoofd, J. Schulman, and P. Abbeel. Benchmarking deep reinforcement
398 learning for continuous control. In *International Conference on Machine Learning*, pages 1329–
399 1338, 2016.
- 400 [13] J. N. Tsitsiklis and B. Van Roy. Analysis of temporal-difference learning with function
401 approximation. In *Advances in neural information processing systems*, pages 1075–1081, 1997.
- 402 [14] S. Kim, K. Asadi, M. Littman, and G. Konidaris. Deepmellow: Removing the need for a target
403 network in deep q-learning. In *Proceedings of the Twenty-Eighth International Joint Conference*
404 *on Artificial Intelligence, IJCAI-19*, pages 2733–2739. International Joint Conferences on
405 Artificial Intelligence Organization, 7 2019. doi:10.24963/ijcai.2019/379. URL <https://doi.org/10.24963/ijcai.2019/379>.
406
- 407 [15] I. Durugkar and P. Stone. Td learning with constrained gradients. In *Proceedings of the Deep*
408 *Reinforcement Learning Symposium, NIPS 2017*, Long Beach, CA, USA, December 2017. URL
409 <http://www.cs.utexas.edu/users/ai-lab?NIPS17-ishand>.
- 410 [16] J. Achiam, E. Knight, and P. Abbeel. Towards characterizing divergence in deep q-learning.
411 *arXiv preprint arXiv:1903.08894*, 2019.
- 412 [17] A. Bhatt, M. Argus, A. Amiranashvili, and T. Brox. Crossnorm: Normalization for off-policy td
413 reinforcement learning. *arXiv preprint arXiv:1902.05605*, 2019.
- 414 [18] X. B. Peng, A. Kumar, G. Zhang, and S. Levine. Advantage-weighted regression: Simple and
415 scalable off-policy reinforcement learning. *arXiv preprint arXiv:1910.00177*, 2019.
- 416 [19] D. Kalashnikov, A. Irpan, P. Pastor, J. Ibarz, A. Herzog, E. Jang, D. Quillen, E. Holly, M. Kalakr-
417 ishnan, V. Vanhoucke, et al. Scalable deep reinforcement learning for vision-based robotic
418 manipulation. In *Conference on Robot Learning*, pages 651–673, 2018.
- 419 [20] R. Simmons-Edler, B. Eisner, E. Mitchell, S. Seung, and D. Lee. Q-learning for continuous
420 actions with cross-entropy guided policies. *arXiv preprint arXiv:1903.10605*, 2019.
- 421 [21] M. Yan, A. Li, M. Kalakrishnan, and P. Pastor. Learning probabilistic multi-modal actor
422 models for vision-based robotic grasping. In *2019 International Conference on Robotics and*
423 *Automation (ICRA)*, pages 4804–4810. IEEE, 2019.
- 424 [22] Z. Botev, D. Kroese, R. Rubinstein, and P. L’Ecuyer. *The Cross-Entropy Method for Optimiza-*
425 *tion*, volume 31, pages 35–59. 12 2013.
- 426 [23] C. J. C. H. Watkins. *Learning from Delayed Rewards*. PhD thesis, King’s College, Oxford,
427 1989.
- 428 [24] S. Thrun and A. Schwartz. Issues in using function approximation for reinforcement learn-
429 ing. In M. Mozer, P. Smolensky, D. Touretzky, J. Elman, and A. Weigend, editors, *Pro-*
430 *ceedings of the 1993 Connectionist Models Summer School*, pages 255–263. Lawrence
431 Erlbaum, 1993. URL http://www.ri.cmu.edu/pub_files/pub1/thrun_sebastian_1993_1/thrun_sebastian_1993_1.pdf.
432

- 433 [25] H. V. Hasselt. Double q-learning. In *Advances in neural information processing systems*, pages
434 2613–2621, 2010.
- 435 [26] E. Todorov, T. Erez, and Y. Tassa. Mujoco: A physics engine for model-based control. In *2012*
436 *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 5026–5033. IEEE,
437 2012.
- 438 [27] T. Haarnoja, A. Zhou, K. Hartikainen, G. Tucker, S. Ha, J. Tan, V. Kumar, H. Zhu, A. Gupta,
439 P. Abbeel, et al. Soft actor-critic algorithms and applications. *arXiv preprint arXiv:1812.05905*,
440 2018.
- 441 [28] O. M. Andrychowicz, B. Baker, M. Chociej, R. Józefowicz, B. McGrew, J. Pachocki, A. Petron,
442 M. Plappert, G. Powell, A. Ray, J. Schneider, S. Sidor, J. Tobin, P. Welinder, L. Weng, and
443 W. Zaremba. Learning dexterous in-hand manipulation. *The International Journal of Robotics*
444 *Research*, 39(1):3–20, 2020. doi:10.1177/0278364919887447.
- 445 [29] S. Yuan, L. Shao, C. L. Yako, A. Gruebele, and J. K. Salisbury. Design and control of roller
446 grasper v2 for in-hand manipulation. In *2020 IEEE/RSJ International Conference on Intelligent*
447 *Robots and Systems (IROS)*, 2020.