

Appendix A. Appendix

In this appendix we provide missing proofs and implementation details. Specifically, we present Theorem 1 for CAC convergence in Section A.1, Proposition 4 for calculating CAC actor gradient expression in Section A.2 and implementation details in Section A.3.

A.1. Proof for the CAC convergence

We prove the convergence of CAC using policy iteration style argument. Similar proofs have also been used in (Haarnoja et al., 2018, Theorem 4). The following lemmas establish the convergence of the policy evaluation and policy improvement of CAC, respectively.

Lemma 5 (CAC Policy Evaluation) *Given the current policy π and a baseline policy $\tilde{\pi}$, the policy evaluation step of CAC is formulated as:*

$$Q_{\pi_{k+1}} \leftarrow \mathcal{R}(s, a) + \gamma (\mathbf{P}\mathbb{E}^\pi [Q_{\pi_k}(s, a)] + \mathcal{I}_{\tilde{\pi}}^\pi(s)), \quad (14)$$

Consider an initial Q value $Q_0 : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ with $|\mathcal{A}| < \infty$. With the repeated application of Eq. (14), the sequence Q_k converges to the following entropy-regularized Q -value $Q_{\tilde{\pi}}^\pi$ as $k \rightarrow \infty$.

$$Q_{\tilde{\pi}}^\pi(s, a) := \mathbb{E}^{d^\pi} \left[\sum_{t=0}^{\infty} \gamma^t (r_t + \mathcal{I}_{\tilde{\pi}}^\pi(s_{t+1})) \mid s_0 = s, a_0 = a \right]. \quad (15)$$

Proof Define the entropy augmented reward as $\mathcal{R}_{\tilde{\pi}}^\pi(s, a) \triangleq \mathcal{R}(s, a) + \mathcal{I}_{\tilde{\pi}}^\pi(s)$ and rewrite the update rule as:

$$Q_{\pi_{k+1}} \leftarrow \mathcal{R}_{\tilde{\pi}}^\pi(s, a) + \gamma \mathbf{P}\mathbb{E}^\pi [Q_{\pi_k}(s, a)]. \quad (16)$$

With the assumption $|\mathcal{A}| < \infty$ for bounded reward, we can apply the standard convergence results for policy evaluation (Sutton and Barto, 2018). \blacksquare

Lemma 6 (CAC Policy Improvement) *Given the current policy π , a baseline policy $\tilde{\pi}$ and the updated policy π_{new} . CAC has the following policy update:*

$$\begin{aligned} \pi_{new} &= (1 - \zeta)\pi + \zeta\hat{\pi}, \\ \text{where } \hat{\pi}(a \mid s) &= \frac{\tilde{\pi}^\alpha(a \mid s) \exp(\beta Q_{\tilde{\pi}}^\pi(s, a))}{Z(s)}, \end{aligned} \quad (17)$$

with $\zeta \in [0, 1]$. Then, $Q_{\tilde{\pi}}^{\pi_{new}}(s, a) \geq Q_{\tilde{\pi}}^\pi(s, a)$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$ with $|\mathcal{A}| < \infty$.

Proof Consider a function $f : \zeta \rightarrow \mathbb{R}$ with $\zeta \in [0, 1]$:

$$f(\zeta) = \mathbb{E}^{a \sim \pi_{new}} [Q_{\tilde{\pi}}^\pi(s, a)] + \mathcal{I}_{\tilde{\pi}}^{\pi_{new}}(s). \quad (18)$$

From the definition of $\hat{\pi} = \arg \max_{\pi} \mathbb{E}^{a \sim \pi} [Q_{\pi}^{\pi}(s, a)] + \mathcal{I}_{\pi}^{\pi}(s)$, $f(\zeta)$ takes the maximum value when $\zeta = 1$. The first and the second derivative of $f(\zeta)$ w.r.t. ζ are:

$$f'(\zeta) = \sum_a (\hat{\pi}(a|s) - \pi(a|s)) (Q_{\hat{\pi}}^{\pi}(s, a) + \tau \log \hat{\pi}(a|s)) - (\sigma + \tau) \log((1 - \zeta)\pi(a|s) + \zeta\hat{\pi}(a|s)). \quad (19)$$

$$f''(\zeta) = -(\sigma + \tau) \sum_a \frac{(\hat{\pi}(a|s) - \pi(a|s))^2}{(1 - \zeta)\pi(a|s) + \zeta\hat{\pi}(a|s)} \leq 0. \quad (20)$$

Thus, the function f is concave in $\zeta \in [0, 1]$. Since $f(\zeta)$ takes the maximum value with $\zeta = 1$, $f(\zeta)$ is monotonically increasing in ζ and $f(0) \leq f(\zeta)$.

Therefore, the following inequality about the entropy-regularized V-value $V_{\hat{\pi}}^{\pi}(s)$ holds:

$$\begin{aligned} V_{\hat{\pi}}^{\pi}(s) &= \mathbb{E}^{a \sim \pi} [Q_{\hat{\pi}}^{\pi}(s, a)] + \mathcal{I}_{\hat{\pi}}^{\pi}(s) \\ &\leq \mathbb{E}^{a \sim \pi_{\text{new}}} [Q_{\hat{\pi}}^{\pi}(s, a)] + \mathcal{I}_{\hat{\pi}}^{\pi_{\text{new}}}(s). \end{aligned} \quad (21)$$

By repeatedly applying Eq. (21), we obtain the following inequalities:

$$\begin{aligned} Q_{\hat{\pi}}^{\pi}(s, a) &= \mathcal{R}(s, a) + \gamma \mathbf{P} [Q_{\hat{\pi}}^{\pi}(s, a) + \mathcal{I}_{\hat{\pi}}^{\pi}(s)] \\ &\leq \mathcal{R}(s, a) + \gamma \mathbf{P} [\mathbb{E}^{\pi_{\text{new}}} [Q_{\hat{\pi}}^{\pi}(s, a)] + \mathcal{I}_{\hat{\pi}}^{\pi_{\text{new}}}(s)] \\ &\vdots \\ &\leq Q_{\hat{\pi}}^{\pi_{\text{new}}}(s, a). \end{aligned} \quad (22)$$

Convergence to $Q_{\hat{\pi}}^{\pi_{\text{new}}}$ follows from Lemma 5. ■

Combining the policy evaluation and policy improvement, we are now ready to prove Theorem 1.

Theorem 1 *Repeated application of CAC Eq. (3) on any initial policy π will make it converges to the entropy regularized optimal policy $\pi^*(a|s) = \frac{\exp(\frac{1}{\kappa} Q^*(s, a))}{\int_{a \in \mathcal{A}} \exp(\frac{1}{\kappa} Q^*(s, a))}$.*

Proof According to Lemma 5 and Lemma 6, the entropy-regularized Q-value at k -th update satisfies $Q_{\pi_{k-1}}^{\pi_k}(s, a) \geq Q_{\pi_{k-1}}^{\pi_{k-1}}(s, a)$. Given bounded reward, $Q_{\pi_k}^{\pi_k}$ is also bounded from above and the sequence converges to a unique π^* . Note that when reaching the optimum the KL regularization term becomes 0. Hence, using the same iterative argument as in the proof of Lemma 6, we get $Q_{\pi^*}^{\pi^*}(s, a) > Q_{\pi}^{\pi}(s, a)$ for all $(s, a) \in \mathcal{S} \times \mathcal{A}$ and any π . By Ziebart (2010); Haarnoja et al. (2018), the optimal policy is entropy-regularized and hence has the softmax form $\pi^*(a|s) = \frac{\exp(\frac{1}{\kappa} Q^*(s, a))}{\int_{a \in \mathcal{A}} \exp(\frac{1}{\kappa} Q^*(s, a))}$ (recall from Eq. (3) that κ is the weight coefficient of entropy). The convergence of general interpolated policy to the optimal policy follows the argument of Scherrer and Geist (2014). ■

A.2. Proof for the CAC gradient

In this subsection we derive the gradient expression for CAC. For the ease of reading we rephrase the proposition here:

Proposition 4 *Let the actor network be parametrized by weights ϕ and critic by θ . Define $\mathcal{G}Q_{\bar{\phi},\theta}$ as the greedy policy with respect to the CAC critic. The subscript $\bar{\phi}$ comes from the baseline policy introduced by KL divergence. Then the gradient of the actor update can be expressed as:*

$$\begin{aligned} & \nabla_{\phi} \mathbb{E}^{s \sim \mathcal{B}} \left[D_{\bar{\phi}}^{\phi} - \frac{\beta}{1 + \mathcal{X}} Q_{\theta}(s, a) \right], \\ & \text{where } D_{\bar{\phi}}^{\phi} = \log \pi_{\phi}(a | s) - \frac{\alpha + \mathcal{X}}{1 + \mathcal{X}} \log \pi_{\bar{\phi}}(a; s) \\ & \mathcal{X} = \frac{1 - \zeta}{\zeta} \cdot \frac{\pi_{\bar{\phi}}(a | s)}{\mathcal{G}Q_{\bar{\phi},\theta}(a | s)}. \end{aligned} \quad (23)$$

Proof Using the reparameterization trick $a = f_{\phi}(\epsilon; s_t)$ with ϵ a noise vector (Haarnoja et al., 2018), the gradient of Eq. (8c) can be expressed as:

$$\begin{aligned} \hat{\nabla}_{\phi} J_{\pi}(\phi) &= \nabla_{\phi} \log \pi_{\phi}(a | s) + \nabla_a \log \pi_{\phi}(a | s) \nabla_{\phi} f_{\phi}(\epsilon; s_t) \\ &\quad - \nabla_a \log \left((1 - \zeta) \pi_{\bar{\phi}}(a | s) + \zeta \mathcal{G}Q_{\bar{\phi},\theta}(a | s) \right) \nabla_{\phi} f_{\phi}(\epsilon; s_t). \end{aligned} \quad (24)$$

We expand the term $\nabla_a \log \left((1 - \zeta) \pi_{\bar{\phi}}(a | s) + \zeta \mathcal{G}Q_{\bar{\phi},\theta}(a | s) \right)$ by using that $\nabla_{x_i} \log \left(\sum_j \exp x_j \right) = \frac{\exp x_i}{\sum_j \exp x_j}$.

Let $\exp(C_1(a)) = (1 - \zeta) \pi_{\bar{\phi}}(a | s)$ and $\exp(C_2(a)) = \zeta \mathcal{G}Q_{\bar{\phi},\theta}(a | s)$. We have the following transformation:

$$\begin{aligned} & \nabla_a \log \left((1 - \zeta) \pi_{\bar{\phi}}(a | s) + \zeta \mathcal{G}Q_{\bar{\phi},\theta}(a | s) \right) \\ &= \nabla_a \log \left(\exp(C_1(a)) + \exp(C_2(a)) \right) \\ &= \frac{\left(D \frac{\partial}{\partial a} \pi_{\bar{\phi}}(a | s) + \alpha \frac{\partial}{\partial a} \pi_{\bar{\phi}}(a | s) + \beta \frac{\partial}{\partial a} Q_{\theta}(s, a) \right)}{1 + D}, \\ & \text{where } D = \exp(C_1(a) - C_2(a)). \end{aligned} \quad (25)$$

After replacing D with \mathcal{X} and inserting Eq. (25) into Eq. (24), we obtain Eq. (10). ■

A.3. Implementation details

This section presents implementation details of CAC with deep networks. Pseudo-code is provided in Algorithm 1.

On-policy replay buffer To make the algorithm off-policy, we approximate the on-policy samples with on-policy replay buffer \mathcal{B}_K which stores K recent samples where K is smaller than the size of the main replay buffer \mathcal{B} .

Advantage estimation While it is possible to simply use the entropy-regularized advantage function $A_{\mathcal{I}}(s, a) = Q_{\mathcal{I}}(s, a) - V_{\mathcal{I}}(s)$ for computing ζ , we are interested in studying the guidance of ζ when no entropy is involved since it might provide a more informative gradient improving direction. This corresponds to the case of (Kakade and Langford, 2002; Pirota et al., 2013). To this end, we train another Q-network Q_{ω} by solving:

$$\begin{aligned} \omega &\leftarrow \arg \min \mathbb{E}^{\mathcal{B}} \left[(Q_{\omega}(s, a) - y)^2 \right], \\ \text{where } y &= r + \gamma \left(\mathbb{E}^{a \sim \pi_{\phi}(\cdot | s')} [Q_{\bar{\omega}}(s', a)] \right), \end{aligned} \tag{26}$$

where $\bar{\omega}$ is the target network. Then we approximate the advantage as $A^{\pi}(s, a) = Q_{\omega}(s, a) - \mathbb{E}^{\pi_{\phi}}[Q_{\omega}(s, a)]$. While the advantage estimation is expected to be further improved with the recent generalized advantage estimation, we found that the above simple implementation is sufficient to stabilize the learning.

Target smoothing For the target Q-networks $Q_{\bar{\theta}}$ and $Q_{\bar{\omega}}$, we update the parameters using the moving average (Haarnoja et al., 2018):

$$\begin{aligned} \bar{\theta} &\leftarrow \nu_{\bar{\theta}} \theta + (1 - \nu_{\bar{\theta}}) \bar{\theta}, \\ \bar{\omega} &\leftarrow \nu_{\bar{\omega}} \omega + (1 - \nu_{\bar{\omega}}) \bar{\omega}, \end{aligned} \tag{27}$$

where $\nu_{\bar{\theta}}$ and $\nu_{\bar{\omega}}$ are the target smoothing coefficients. In the mixing step, we use the previous policy $\pi_{\bar{\phi}}$ rather than the current policy π_{ϕ} to stabilize the training:

$$\mathbb{E}^{\mathcal{B}} [D_{KL}(\pi_{\phi} \| (1 - \zeta)\pi_{\bar{\phi}} + \zeta\hat{\pi})].$$

Thus, the target policy $\pi_{\bar{\phi}}$ corresponds to the monotonically improved policy in the CPI algorithm that is not updated when performance oscillation happens. To reflect this fact, we update the weight of the target policy network as:

$$\bar{\phi} \leftarrow \zeta \nu_{\bar{\phi}} \theta + (1 - \zeta \nu_{\bar{\phi}}) \bar{\phi}, \tag{28}$$

where $\nu_{\bar{\phi}}$ is the target smoothing coefficient.

Normalization factor estimation Since CAC algorithm requires the density of the reference policy $\hat{\pi}(a | s) = \pi_{\bar{\phi}}^{\alpha}(a | s) \exp(\beta Q_{\theta}(s, a))(Z(s))^{-1}$, we need to estimate the normalization factor $Z(s)$.

A simple approach to estimate $Z(s)$ is by Monte-Carlo sampling with some distribution q that is easier to sample from:

$$Z(s) = \mathbb{E}^q \left[\frac{\pi_{\bar{\phi}}^{\alpha}(a | s)^{\alpha} \exp(\beta Q_{\theta}(s, a))}{q(a | s)} \right]. \tag{29}$$

The closer $q(\cdot | s)$ and the reference policy $\hat{\pi}(\cdot | s)$ are, the better the accuracy of the $Z(s)$ approximation.

Theorem 2 indicates that by choosing the current policy π as the proposal distribution, we can control the closeness of the two distributions and the accuracy of the MC approximation via changing the entropy regularization weighting coefficients. We empirically study the effectiveness of entropy regularization against the closeness and the accuracy when $\zeta \leq 1$

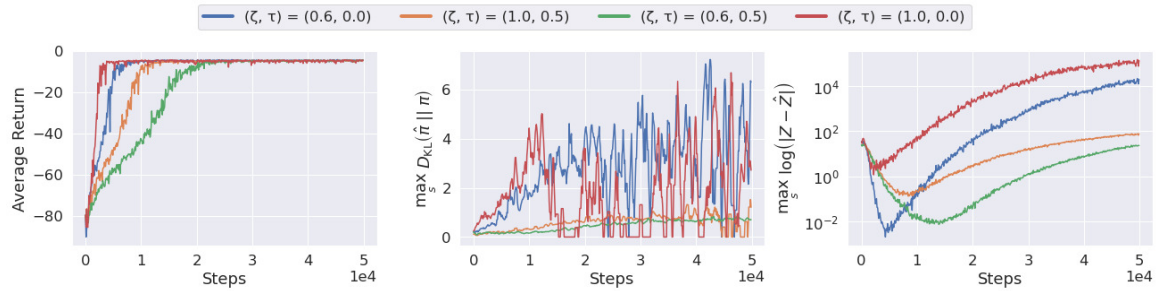


Figure 3: **Left:** The learning curves of CAC with different parameters on pendulum. **Middle:** The maximum KL divergence over all the states. **Right:** The maximum error of log-scaled $Z(s)$ approximation over all the states.

Algorithm 1 Cautious Actor-Critic

- 1: Initialize parameter vectors $\theta, \phi, \bar{\theta}, \bar{\phi}, \omega, \bar{\omega}$
 - 2: Initialize variable $\tilde{\mathbb{A}}$ and $\tilde{\mathbb{A}}^{\text{MaxDiff}}$
 - 3: **for** each iteration **do**
 - 4: Collect transitions by π_θ and add them to \mathcal{B} and \mathcal{B}_K
 - 5: **for** each gradient step **do**
 - 6: Update θ with one step of SGD using Eq. (8b)
 - 7: Update ω with one step of SGD using Eq. (26)
 - 8: Update $\tilde{\mathbb{A}}$ and $\tilde{\mathbb{A}}^{\text{MaxDiff}}$ using Eq. (12)
 - 9: Update ϕ with one step of SGD using Eq. (8c)
 - 10: Update ζ using Eq. (11)
 - 11: $\bar{\theta} \leftarrow \nu_{\bar{\theta}}\theta + (1 - \nu_{\bar{\theta}})\bar{\theta}$
 - 12: $\bar{\omega} \leftarrow \nu_{\bar{\omega}}\omega + (1 - \nu_{\bar{\omega}})\bar{\omega}$
 - 13: $\bar{\phi} \leftarrow \zeta\nu_{\bar{\phi}}\phi + (1 - \zeta\nu_{\bar{\phi}})\bar{\phi}$
 - 14: **end for**
 - 15: **end for**
-

We use the pendulum environment from Fu et al. (2019) where the dynamics are discretized so that we can compute the oracle values such as the KL divergence between the current and the reference policy. The hyperparameters used in the experiment is listed in Section A.4. Figure 3 shows how the learning behavior of CAC changes when the interpolation coefficient ζ and KL regularization weight τ vary: When KL regularization is present, the approximation quality of $Z(s)$ is improved significantly.

A.4. Hyperparameters

This section lists the hyperparameters used in the comparative evaluation Section 5.1.

Table 2: Hyperparameters of off-policy algorithms in mujoco tasks

Parameter	Value
<i>Shared</i>	
optimizer	Adam
learning rate	10^{-3}
discount factor (γ)	0.99
replay buffer size (\mathcal{B})	10^6
number of hidden layers	2
number of hidden units per layer	256
number of samples per minibatch	100
activations	ReLU
<i>TD3</i>	
Stddev for Gaussian noise	0.1
Stddev for target smoothing noise	0.2
policy delay	2
<i>SAC</i>	
entropy coefficient (κ)	0.2
$\bar{\theta}$ smoothing coefficient	0.995
<i>CAC</i>	
entropy coefficient (κ)	0.2
KL coefficient (τ)	0.1
$\bar{\theta}$ smoothing coefficient ($\nu_{\bar{\theta}}$)	0.995
$\bar{\omega}$ smoothing coefficient ($\nu_{\bar{\omega}}$)	0.995
$\bar{\phi}$ smoothing coefficient ($\nu_{\bar{\phi}}$)	0.9999
ν_A	0.01
$\nu_{A^{\text{MaxDiff}}}$	0.001
size of \mathcal{B}_K	1000
if-else update	$c = M$

Table 3: Hyperparameters of PPO in mujoco tasks

Parameter	Value
<i>PPO</i>	
optimizer	Adam
value function learning rate	10^{-3}
policy learning rate	3×10^{-4}
discount factor (γ)	0.99
number of hidden layers	2
number of hidden units per layer	256
number of samples per minibatch	100
activations	ReLU
Number of samples per update	80
Policy objective clipping coefficient	0.2

Table 4: Hyperparameters of CAC in pendulum task

Parameter	Value
<i>CAC</i>	
optimizer	Adam
learning rate	10^{-3}
discount factor (γ)	0.99
replay buffer size (\mathcal{B})	10^6
number of hidden layers	2
number of hidden units per layer	256
number of samples per minibatch	32
activations	ReLU
entropy coefficient (κ)	0.2
$\bar{\theta}$ smoothing coefficient	0.995
$\bar{\phi}$ smoothing coefficient	0.995
size of \mathcal{B}_K	1000
if-else update	$c = M$