
A Connection between One-Step RL and Critic Regularization in Reinforcement Learning

Benjamin Eysenbach^{1,2} Matthieu Geist¹ Sergey Levine^{1,3} Ruslan Salakhutdinov²

Abstract

As with any machine learning problem with limited data, effective offline RL algorithms require careful regularization to avoid overfitting. One class of methods, known as one-step RL, perform just one step of policy improvement. These methods, which include advantage-weighted regression and conditional behavioral cloning, are thus simple and stable, but can have limited asymptotic performance. A second class of methods, known as critic regularization, perform many steps of policy improvement with a regularized objective. These methods typically require more compute but have appealing lower-bound guarantees. In this paper, we draw a connection between these methods: applying a multi-step critic regularization method with a regularization coefficient of 1 yields the same policy as one-step RL. While our theoretical results require assumptions (e.g., deterministic dynamics), our experiments nevertheless show that our analysis makes accurate, testable predictions about practical offline RL methods (CQL and one-step RL) with commonly-used hyperparameters.

1. Introduction

Reinforcement learning (RL) algorithms tend to perform better when regularized, especially when given access to only limited data, and especially in batch (i.e., offline) settings where the agent is unable to collect new experience. While RL algorithms can be regularized using the same tools as in supervised learning (e.g., weight decay, dropout), we will use “regularization” to refer to those techniques unique to the RL setting. Such regularization methods include policy regularization (penalizing the policy for sampling out-of-

¹Google Research ²Carnegie Mellon University ³UC Berkeley. Correspondence to: Benjamin Eysenbach <bey-senba@cs.cmu.edu>.

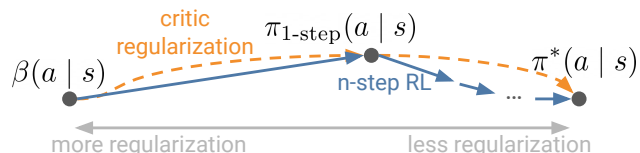


Figure 1: Both n -step RL and critic regularization can interpolate between behavioral cloning (left) and un-regularized RL (right) by varying the regularization parameter. Endpoints of these regularization paths are the same. We prove that these methods also obtain the same policy for an intermediate degree of regularization.

distribution action) and value regularization (penalizing the critic for making large predictions). Research on these sorts of regularization has grown significantly in recent years, yet theoretical work studying the tradeoffs between regularization methods remains limited (Vieillard et al., 2020).

Many RL methods perform regularization and can be classified by whether they perform one or many steps of policy improvement. *One-step RL* methods (Brandfonbrener et al., 2021; Peng et al., 2019; Peters & Schaal, 2007; Peters et al., 2010) perform one step of policy iteration, updating the policy to choose actions the are best according to the Q-function of the behavioral policy. The policy is often regularized to not deviate far from the behavioral policy. In theory, policy iteration can take a large number of iterations ($\tilde{O}(|\mathcal{S}||\mathcal{A}|/(1-\gamma))$) (Scherrer, 2013)) to converge, so one-step RL (one step of policy iteration) fails to find the optimal policy on most tasks. Empirically, policy iteration often converges in a smaller number of iterations (Sutton & Barto, 2018, Sec. 4.3), and the policy after just a single iteration can sometimes achieve performance comparable to multi-step RL methods (Brandfonbrener et al., 2021). *Critic regularization* methods modify the training of the value function such that it predicts smaller returns for unseen actions (Kumar et al., 2020; Chebotar et al., 2021; Yu et al., 2021; Hatch et al., 2022; Nachum et al., 2019; An et al., 2021; Bai et al.; Buckman et al., 2021). Intuitively, such critic regularization causes the policy to avoid sampling unseen actions. In this paper, we will use “critic regularization” to specifically refer to multi-step methods that use critic regularization; multi-step methods that do not utilize critic regularization (e.g., KL control (Ziebart, 2010),

TD3+BC (Fujimoto & Gu, 2021)) are outside the scope of our analysis.

These RL regularization methods appear distinct. Critic regularization typically involves solving a two-player game, whereby a policy predicts actions with high values while the critic decreases the values predicted for those actions. Prior work (Kumar et al., 2020) has argued that this complexity is worthwhile because the regularization effect is being propagated across time via Bellman backups: decreasing the value at one state will also decrease the value at states that lead there.

In this paper, we show that a certain type of one-step RL is equivalent to a certain type of critic regularization, under some assumptions (see Fig. 1). The key idea is that, when using a certain TD loss, the regularized critic updates converge not to the true Q-values, but rather the Q-values multiplied by an importance weight. For the critic, these importance weights mean that the Q-values end up estimating the expected returns of the behavioral policy (Q^β , as in many one-step methods (Peters et al., 2010; Peters & Schaal, 2007; Peng et al., 2019; Brandfonbrener et al., 2021)), rather than the expected returns of the optimal policy (Q^π). For the actor, these importance weights mean that the logarithm of the Q-values includes a term that looks like a KL divergence. This connection allows us to make precise how critic regularization methods implicitly regularize the policy.

The key contribution of this paper is a construction showing that one-step RL produces the same policy as a multi-step critic regularization method, for a certain regularization coefficient. We also discuss similar connections in settings with varying degrees of regularization, in goal-conditioned settings, and in RL settings with success classifiers. The main assumption behind these results is that the critic is updated using an update rule based on the cross entropy loss, rather than the MSE loss. Our result is potentially surprising because, algorithmically and mechanistically, one-step RL and critic regularization are very different. Nonetheless, our analysis may help explain why prior work has found that one-step RL and critic regularization methods can perform similarly on some (Brandfonbrener et al., 2021; Emmons et al., 2021) (but not all (Kostrikov et al., 2021)) problems. Our results hint that one-step RL methods may be a simpler approach to achieving the theoretical guarantees typically associated with critic regularization methods. Our analysis does not say that practical implementations (which violate our assumptions) will always behave the same in practice; they do not (Brandfonbrener et al., 2021, Sec. 7). Nonetheless, our experiments show that our analysis makes accurate testable predictions about practical methods. While our results do not say whether users should regularize the actor or critic in practice, they hint that one-step RL methods may be a simpler way of achieving the theoretical and empirical

properties of critic regularization on RL tasks that require strong regularization.

2. Related Work

Regularization has been applied to RL in many different ways (Neu et al., 2017; Geist et al., 2019), and features prominently in offline RL methods (Lange et al., 2012; Levine et al., 2020). While RL algorithms can be regularized using the same techniques as in supervised learning (e.g., weight decay, dropout), our focus will be on regularization methods unique to the RL setting. Such RL-specific regularization methods can be roughly categorized into one-step RL methods and critic regularization methods.

One-step RL methods (Brandfonbrener et al., 2021; Gülçehre et al., 2020; Peters & Schaal, 2007; Peng et al., 2019; Peters et al., 2010; Wang et al., 2018) apply a single step of policy improvement to the behavioral policy. These methods first estimate the Q-values of the behavioral policy, either via regression or iterative Bellman updates. Then, these methods optimize the policy to maximize these Q-values minus an actor regularizer. Many goal-conditioned or task-conditioned imitation learning methods (Savinov et al., 2018; Ding et al., 2019; Sun et al., 2019; Ghosh et al., 2020; Paster et al., 2021; Yang et al., 2021; Srivastava et al., 2019; Kumar et al., 2019b; Chen et al., 2021; Lynch & Sermanet, 2021; Li et al., 2020; Eysenbach et al., 2020a) also fits into this mold (Eysenbach et al., 2022), yielding policies that maximize the Q-values of the behavioral policy while avoiding unseen actions. Note that non-conditional imitation learning methods do not perform policy improvement, and do not fit into this mold. One-step methods are typically simple to implement and computationally efficient.

Critic regularization methods instead modify the objective for the Q-function so that it predicts lower returns for unseen actions (Kumar et al., 2020; Chebotar et al., 2021; Yu et al., 2021; Hatch et al., 2022; Nachum et al., 2019; An et al., 2021; Bai et al.; Buckman et al., 2021). Critic regularization methods are typically more challenging to implement correctly and more computationally demanding (Kumar et al., 2020; Nachum et al., 2019; Bai et al.; An et al., 2021), but can lead to better results on some challenging problems (Kostrikov et al., 2021) Our analysis will show that one-step RL is equivalent to a certain type of critic regularization.

Some regularization methods do not fit exactly into these two categories. Methods like KL control regularize both the actor and the reward function (Geist et al., 2019; Ziebart, 2010; Haarnoja et al., 2018; Abdolmaleki et al., 2018; Wu et al., 2019; Jaques et al., 2019; Rezaeifar et al., 2022). Other methods only regularize the policy used in the critic updates (Fujimoto et al., 2019; Kumar et al., 2019a).

Our results are conceptually similar to prior work on regularization in the supervised learning setting. There, regularization methods like weight decay, spectral regularization, and early stopping all appear quite distinct, but are actually mathematical equivalent under some assumptions (Bakushinskii, 1967; Wahba, 1987; Fleming, 1990; Santos, 1996; Bauer et al., 2007). This result is surprising because the methods are so different: weight decay modifies the loss function, early stopping modifies the update procedure, and spectral normalization is a post-hoc correction.

3. Preliminaries

We start by defining the single-task RL problem, and then introduce prototypical examples of one-step RL and critic regularization. We then define an actor critic algorithm for use in our analysis.

3.1. Notation

We assume an MDP with states s , actions a , initial state distribution $p_0(s_0)$, dynamics $p(s' | s, a)$, and reward function $r(s, a)$. We assume episodes always have infinite length (i.e., there are no terminal states). Without loss of generality, we assume rewards are positive; adding a positive constant to all rewards can make them all positive without changing the optimal policy. We will learn a Markovian policy $\pi(a | s)$ to maximize the expected discounted sum of rewards:

$$\max_{\pi} \mathbb{E}_{\pi(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 \sim p_0(s_0) \right],$$

where $\pi(\tau) = p(s_0) \prod_{t=0}^{\infty} \pi(a_t | s_t) p(s_{t+1} | s_t, a_t)$ is the probability of policy π sampling an infinite-length trajectory $\tau = (s_0, a_0, \dots)$. We define Q-values for policy $\pi(a | s)$ as

$$Q^{\pi}(s, a) = \mathbb{E}_{\pi(\tau)} \left[\sum_{t=0}^{\infty} \gamma^t r(s_t, a_t) \mid s_0 = s, a_0 = a \right].$$

Because the rewards are positive, these Q-values are also positive, $Q^{\pi}(s, a) > 0$. Since we focus on the offline setting, we will consider two policies: $\beta(a | s)$ is the *behavioral* policy that collected the dataset, and $\pi(a | s)$ is the *online* policy output by the algorithm that attempts to maximize the rewards. We will use $p(s, a, s')$ to denote the empirical distribution of transitions in an offline dataset, and $p(s, a)$ and $p(s)$ denote the corresponding marginal distributions. The behavioral policy is defined as $\beta(a | s) = p(a | s)$.

3.2. Examples of Regularization in RL

While actor and critic regularization methods can be implemented in many ways, we introduce two prototypical examples below to make our discussion more concrete.

Example of one-step RL: Brandfonbrener et al. (2021).

One-step RL first estimates the Q-values of the behavioral policy ($Q^{\beta}(s, a)$), and then optimizes the policy to maximize the Q-values minus a actor regularizer. While the actor regularizer can take different forms and the Q-values can be learned via regression, we will use a reverse KL regularizer and TD-style critic update so that the objective is similar to critic regularization:

$$\max_{\pi} \mathbb{E}_{p(s)\pi(a|s)} \left[Q^{\beta}(s, a) + \lambda \log \frac{\beta(a | s)}{\pi(a | s)} \right] \quad (1)$$

where $Q^{\beta}(s, a) = \lim_{t \rightarrow \infty} Q_t(s, a)$ and

$$Q_{t+1} \leftarrow \arg \min_Q \mathbb{E}_{p(s,a)} \left[(Q(s, a) - y^{\beta, Q_t}(s, a))^2 \right]$$

$$y^{\beta, Q_t}(s, a) \triangleq r(s, a) + \gamma \mathbb{E}_{p(s'|s,a)} [Q_t(s', a')].$$

The scalar λ is the regularization coefficient and $\beta(a | s)$ is an estimate of the behavioral policy, typically learned via behavioral cloning. Like most TD methods (Haarnoja et al., 2018; Mnih et al., 2013; Fujimoto et al., 2018), the TD targets y are not considered learnable. In practice, most methods do not solve optimize the critic to convergence at each step, instead taking just a few gradient steps before updating the TD targets. This one-step critic loss is different from the multi-step critic losses used in other RL methods (e.g., TD3, SVG(0)) because it uses the TD target $y^{\beta, Q}(s, a)$ (corresponds to a fixed policy) rather than $y^{\pi, Q}(s, a)$ (corresponding to a sequence of learned policies). One-step RL amounts to performing one step of policy iteration, rather than full policy optimization. While truncating the iterations of policy iteration can be suboptimal, it can also be interpreted as a form of early stopping regularization.

Example of critic regularization: Kumar et al. (2020).

CQL (Kumar et al., 2020) modifies the standard Bellman loss to include an additional term that decreases the values predicted for unseen actions. The actor objective is to maximize Q values; some CQL implementations also regularize the actor loss (Hoffman et al., 2020; Kumar et al., 2020)). The objectives can then be written as

$$\max_{\pi} \mathbb{E}_{p(s)\pi(a|s)} [Q^{\pi}(s, a)] \quad (2)$$

where $Q^{\pi}(s, a) = \lim_{t \rightarrow \infty} Q_t(s, a)$ and

$$Q_{t+1} = \arg \min_Q \mathbb{E}_{p(s,a)} \left[(Q(s, a) - y^{\pi, Q_t}(s, a))^2 \right]$$

$$+ \lambda \left(\mathbb{E}_{p(s)\pi(a|s)} [Q(s, a)] - \mathbb{E}_{p(s)\beta(a|s)} [Q(s, a)] \right).$$

The second term decreases the Q-values for unseen actions (those sampled from $\pi(a | s)$) while the third term increases the values predicted for seen actions (those sampled from

the behavioral policy $\beta(a | s)$). Unlike standard temporal difference methods, the CQL updates resemble a competitive game between the actor and the critic. In practice, this cyclic dependency can create unstable learning (Kumar et al., 2020; Hoffman et al., 2020).

3.3. How are these methods connected?

Prior work has observed that one-step methods and critic regularization methods perform similarly on many (Fujimoto & Gu, 2021; Emmons et al., 2021) (but not all (Kostrikov et al., 2021)) tasks. Despite the differences in objectives and implementations of these two methods (and, more broadly, the actor/critic regularization methods for which they are prototypes), are there deeper, unifying connections between the methods?

In the next section, we introduce a different actor-critic method that will allow us to draw a connection between one-step RL and critic regularization. We experimentally validate this equivalence in Sec. 5.1. Despite its difference from practically-used methods, such as one-step RL and CQL, we will show that it makes accurate predictions about the behavior of these practical methods (Sec. 5.2 and 5.3).

3.4. Classifier Actor Critic

To support our analysis, we will introduce a new actor-critic algorithm. This algorithm is similar to prior work, but trains the critic using a cross entropy loss instead of an MSE loss. We introduce this algorithm not because we expect it to perform better than existing actor-critic methods, but rather because it allows us to make precise a connection between actor and critic regularization. This method treats the value function like a classifier, so we will call it *classifier actor critic*. We will then introduce actor-regularized and critic-regularized versions of this method. The subsequent section (Sec. 4) will show that these two regularized methods learn the same policy.

The key to our analysis will be to treat Q-values like probabilities, so we define the critic loss in terms of a cross-entropy loss, similar to prior work (Kalashnikov et al., 2018; Eysenbach et al., 2021). Recalling that Q-values are positive (Sec. 3.1), we transform the Q-values to have the correct range by using $\frac{Q}{Q+1} \in [0, 1]$. We will minimize the cross-entropy loss applied to the transformed Q-values:

$$\begin{aligned} & \mathbb{E}_{p(s,a)} \left[\mathcal{CE} \left(\frac{Q(s,a)}{Q(s,a)+1}; \frac{y^{\pi, Q_t}(s,a)}{y^{\pi, Q_t}(s,a)+1} \right) \right] \\ &= -\mathbb{E}_{p(s,a)} \left[\frac{y^{\pi, Q_t}(s,a)}{y^{\pi, Q_t}(s,a)+1} \log \frac{Q(s,a)}{Q(s,a)+1} \right. \\ & \quad \left. + \frac{1}{y^{\pi, Q_t}(s,a)+1} \log \frac{1}{Q(s,a)+1} \right] \end{aligned}$$

$$\stackrel{\text{const.}}{=} -\mathbb{E}_{p(s,a)} \left[\underbrace{y^{\pi, Q_t}(s,a) \log \frac{Q(s,a)}{Q(s,a)+1} + \log \frac{1}{Q(s,a)+1}}_{\triangleq \mathcal{L}_{\text{critic}}(Q, y^{\pi, Q_t})} \right]. \quad (3)$$

In the last line we scale both the positive and negative term by $y^{\pi, Q_t}(s,a) + 1$, a choice that does not change the optimal classifier but reduces notational clutter. When the TD target can be computed exactly, solving this optimization problem results in performing one SARSA update: $Q(s,a) \leftarrow r(s,a) + \gamma Q(s',a')$ (see Lemma 4.1). Thus, by solving this optimization problem many times, each time using the previous Q-value to compute the TD targets, we will converge to the correct Q-values (see Lemma 4.1). The actor objective is to maximize the expected *log* of the Q-values:

$$\max_{\pi} \mathcal{L}_{\text{actor}}(\pi) \triangleq \mathbb{E}_{p(s)\pi(a|s)} [\log(Q^{\pi}(s,a))] \quad (4)$$

where $Q^{\pi}(s,a) = \lim_{t \rightarrow \infty} Q_t(s,a)$ and

$$Q_{t+1} = \arg \min_Q \mathcal{L}_{\text{critic}}(Q, y^{\pi, Q_t}).$$

While most actor-critic methods do not use the logarithm transformation, prior work on conditional behavioral cloning (e.g., (Savinov et al., 2018; Ding et al., 2019; Sun et al., 2019; Ghosh et al., 2020; Srivastava et al., 2019)) implicitly includes this transformation (Eysenbach et al., 2022). In the absence of additional regularization, the optimal policy $\pi(a | s) = \mathbb{1}(a = \arg \max_{a'} Q(s, a'))$ is the same as the optimal policy for the standard actor objective (without the logarithm). We next introduce a one-step version of this method, as well as a critic regularization variant that resembles CQL. While we will implicitly use a regularization coefficient of 1 below, Appendix B.1 discusses versions of classifier actor critic with varying degrees of regularization.

One-step RL. To make classifier actor critic resemble one-step RL (Brandfonbrener et al., 2021), we make two changes: estimating the value of the behavioral policy and adding a regularization term to the actor objective. To estimate the value of the behavioral policy, we modify the critic loss to sample the next action a' from the behavioral policy (i.e., we use $y^{\beta, Q_t}(s,a)$ rather than $y^{\pi, Q_t}(s,a)$). We also regularize the policy by adding a relative entropy term to the actor loss, analogous to the reverse KL penalty used in one-step RL:

$$\max_{\pi} \mathbb{E}_{p(s)\pi(a|s)} [\log Q^{\beta}(s,a) + \log \beta(a | s) - \log \pi(a | s)] \quad (5)$$

where $Q^{\beta}(s,a) = \lim_{t \rightarrow \infty} Q_t(s,a)$ and

$$Q_{t+1} = \arg \min_Q \mathcal{L}_{\text{critic}}(Q, y^{\beta, Q_t}).$$

In tabular settings, this critic objective estimates the Q-values for $\beta(a | s)$ (Appendix Lemma 4.1).

Critic regularization. To emulate CQL, we modify the critic loss (Eq. 3) by adding a penalty term that decreases the values for unseen actions. Whereas CQL applies this penalty to the Q-values directly, we will apply it to the logarithm of the Q-values:¹

$$\max_{\pi} \mathbb{E}_{p(s)\pi(a|s)} [\log Q_r^{\pi}(s, a)] \quad (6)$$

where $Q_r^{\pi}(s, a) = \lim_{t \rightarrow \infty} Q_t(s, a)$ and $Q_{t+1}(s, a) = \arg \min_Q \mathcal{L}_{\text{critic}}^r(Q, y^{\pi, Q_t})$:

$$\begin{aligned} \mathcal{L}_{\text{critic}}^r(Q, y^{\pi, Q_t}) &\triangleq \mathcal{L}_{\text{critic}}(Q, y^{\pi, Q_t}) \\ &+ \lambda \left(\mathbb{E}_{p(s)} [\log(Q(s, a) + 1)] - \mathbb{E}_{\beta(a|s)} [\log(Q(s, a) + 1)] \right). \end{aligned}$$

4. A Connection between One-Step RL and Critic Regularization

This section provides our main result, which is that actor and critic regularization yield the same policy under some settings. The key to proving this connection will be to analyze the Q-values learned by critic regularization. While we mainly focus on the single-task setting, Sec. 4.1 describes how similar results also apply to other settings, including goal-conditioned RL, example-based control, and settings with smaller degrees of regularization. All proofs are in Appendix A.

To relate one-step RL to critic regularization, we start by analyzing the Q-values learned by both methods. We first show that the classifier critic converges to the correct Q-values:

Lemma 4.1. *Assume that states and actions are tabular (discrete and finite), that rewards are positive, and that TD targets can be computed exactly (without sampling). Incrementally update the critic by solving a sequence of optimization problems:*

$$Q_{t+1} \leftarrow \arg \min_Q \mathcal{L}_{\text{critic}}(Q, y^{\pi, Q_t}).$$

This sequence of Q-functions will converge to Q^{π} :

$$\lim_{t \rightarrow \infty} Q_t(s, a) = Q^{\pi}(s, a) \text{ for all states } s \text{ and actions } a.$$

Because one-step RL trains the critic using $\mathcal{L}_{\text{critic}}(Q, y^{\beta, Q})$, it learns Q-values corresponding to $Q^{\beta}(s, a)$. When regularization is added to the critic updates, it learns different Q-values. Perhaps surprisingly, this regularization means that our estimates for the value of policy $\pi(a | s)$ look like the value of the original behavioral policy:

¹From a dimensional analysis perspective (Huntley, 1967), this choice makes sense because it allows the penalty term to have the same “units” as the critic loss: \log Q-values. A second motivation for regularizing the logarithm is that the actor loss uses a logarithm.

Lemma 4.2. *Assume that states and actions are tabular (discrete and finite), that rewards are positive, and that TD targets can be computed exactly (without sampling). Incrementally update the critic by minimizing a sequence of regularized critic losses using policy π and hyperparameter $\lambda = 1$:*

$$Q_{t+1} \leftarrow \arg \min_Q \mathcal{L}_{\text{critic}}^r(Q, y^{\pi, Q_t}).$$

In the limit, this sequence of Q-functions will converge to the Q-values for the behavioral policy ($\beta(a | s)$), weighted by the ratio of the behavioral and online policies:

$$\lim_{t \rightarrow \infty} Q_t(s, a) = \frac{Q^{\beta}(s, a)\beta(a | s)}{\pi(a | s)}$$

for all states s and actions a .

Proof sketch. The ratio $\frac{\beta(a|s)}{\pi(a|s)}$ above is an importance weight. Ordinarily, a TD backup for policy $\pi(a | s)$ would entail sampling an action $a \sim \pi(a | s)$. However, this importance weight means that TD backup is effectively performed by sampling an action $a \sim \beta(a | s)$. Such a TD backup resembles the TD backup for $\beta(a | s)$. The full proof is in Appendix A. \square

Intuitively, this result says that critic regularization reweights the Q-values to assign higher values to in-distribution actions, where $\beta(a | s)$ is large. An unexpected part of this result is that the Q-values correspond to the behavioral policy. In other words, critic regularization added to a multi-step RL method (one using $y^{\pi, Q_t}(s, a)$) yields the same critic as a one-step RL method (one using $y^{\beta, Q_t}(s, a)$). Our main result is a direct corollary of this Lemma:

Theorem 4.3. *Let a behavioral policy $\beta(a | s)$ be given and let $Q^{\beta}(s, a)$ be the corresponding value function. Let $\pi(a | s)$ be an arbitrary policy (typically learned) with support constrained to $\beta(a | s)$ (i.e., $\pi(a | s) > 0 \implies \beta(a | s) > 0$). Let $Q_r^{\pi}(s, a)$ be the critic obtained by the regularized critic update (Eq. 6) to this policy with $\lambda = 1$. Then critic regularization results in the same policy as one-step RL:*

$$\mathbb{E}_{\pi(a|s)} [\log Q_r^{\pi}(s, a)] = \mathbb{E}_{\pi(a|s)} \left[\log Q^{\beta}(s, a) + \log \frac{\beta(a | s)}{\log \pi(a | s)} \right]$$

for all states s .

Since both forms of regularization result in the same objective for the actor, they must produce the same policy in the end. While prior work has mentioned that critic regularization implicitly regularizes the policy (Yu et al., 2021), this result shows that under the assumptions stated above, the implicit regularization of critic regularization results in the exact same policy learning objective as one-step RL. This equivalence holds when $\lambda = 1$, and not necessarily for other regularization coefficients. Appendix B.1 shows how a variant of this result that includes an additional regularization mechanism does apply to different regularization

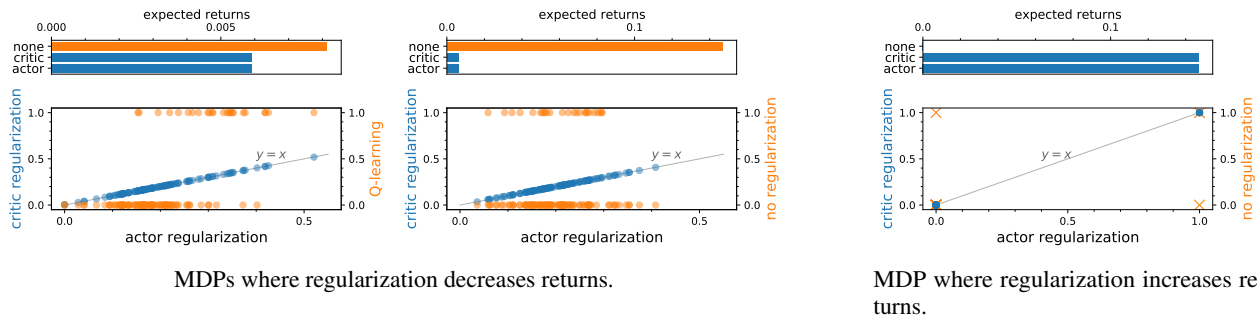


Figure 2: **Actor and critic regularization produce identical policies.** Across three tabular settings, we plot the action probabilities $\pi(a | s)$ for the policies produced by one-step RL and critic-regularized classifier actor-critic ($R^2 \geq 0.999$). We also plot the action probabilities for a policy learned by an unregularized policy to confirm that the equivalence between one-step RL and critic regularization is not a coincidence.

coefficients. This connection between one step RL and critic regularization concerns their *objective functions*, not the *procedures* used to optimize those objective functions. Indeed, because practical offline RL algorithms sometimes use different optimization procedures (e.g., TD vs. MC estimates of $Q^\beta(s, a)$), they will incur errors in estimating $Q^\beta(s, a)$, violating Theorem 4.3’s assumption that these Q-values are estimated exactly.

Limitations. Our theoretical analysis makes assumptions that may not always hold in practice. For example, our results use a critic loss based on the cross entropy loss, while most (but not all (Kalashnikov et al., 2018; Eysenbach et al., 2020b)) practical methods use the MSE. Our analysis assumes that critic regularization arrives at an equilibrium, and ignores errors introduced by function approximation and sampling. Nonetheless, our theoretical results will make accurate predictions about prior offline RL methods.

4.1. Extensions of the Analysis

We extend this analysis in three ways. *First*, we also show that a similar connection can be established for lesser degrees of regularization ($\lambda < 1$) (see Appendix B.1). *Second*, we show that a similar connection holds for RL problems defined via success examples (Pinto & Gupta, 2016; Tung et al., 2018; Kalashnikov et al., 2021; Singh et al., 2019; Zolna et al., 2020; Calandra et al., 2017; Eysenbach et al., 2021) These results use existing actor-critic method, rather than classifier actor critic (see Appendix C). *Third*, we extend our analysis to multi-task settings by looking at goal-conditioned RL problems. Taken together, these extensions show that the connection between actor and critic regularization extends to other commonly-studied problem settings.

5. Numerical Simulations

Our numerical simulations study whether the theoretical connection between actor and critic regularization holds

empirically. The first experiments (Sec. 5.1) will use classifier actor-critic, and we will expect the equivalence to hold exactly in this setting. We then study whether this connection still holds for practical prior methods (one-step RL and CQL), which violate our assumptions. We study these commonly-used methods in both tabular settings (Sec. 5.2) and on a benchmark offline RL task with continuous states and actions (Sec. 5.3). We do not expect these methods to always be the same (see, e.g., Kostrikov et al. (2021, Table 1)), and we will focus our experiments on critic regularization with moderate regularization coefficients. See Appendix F for details and hyperparameters for the experiments. Code for the tabular experiments is available online.²

5.1. Exact Equivalence for Classifier Actor Critic

Our first experiment aims to validate our theoretical result under the required assumptions: when using classifier actor-critic as the RL algorithm, and when using a tabular environment. We use a 5×5 deterministic gridworld with 5 actions (up/down/left/right/nothing). We describe the reward function and other details in Appendix F. To ensure that critic regularization converges to a fixed point and to avoid oscillatory learning dynamics, we update the policy using an exponential moving average. We also include (unregularized) classifier actor-critic to confirm that regularization is important in some settings.

We compare these three methods in three environments. The first setting (Fig. 2 (left)) checks our theory that one-step RL and critic regularization should obtain the same policy. The second setting (Fig. 2 (center)) shows that one-step RL and critic regularization learn the same (suboptimal) policy in settings where using the Q-values for the behavioral policy lead to a suboptimal policy. The final setting is designed so that regularization increases the expected returns. The dataset is a single trajectory from the initial state to

²Code: <https://github.com/ben-eynsbach/ac-connection>

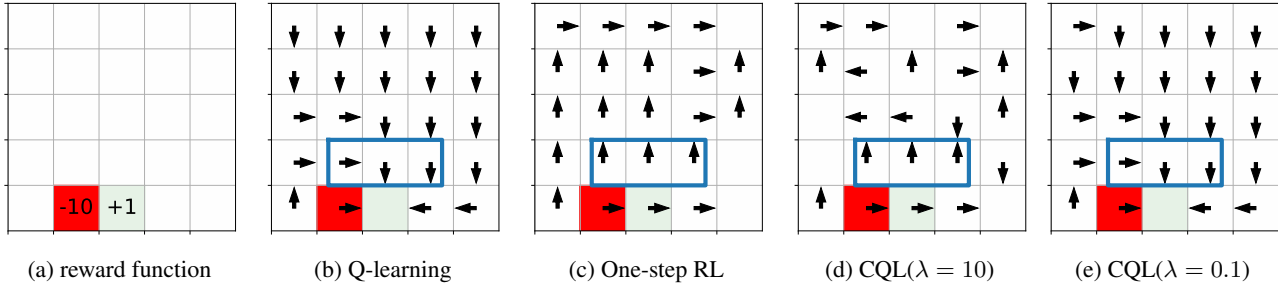


Figure 3: **CQL can behave like one-step RL.** We design a gridworld (a) so that one-step RL (c) learns a suboptimal policy. For the three cells highlighted in blue, the optimal policy (b) navigates towards the high-reward state (green) while the one-step RL policy (c) navigates away from the high-reward state. (d) CQL with a large regularization coefficient exhibits the same suboptimal behavior as one-step RL, taking actions that lead away from the high-reward states. (e) CQL with a small regularization coefficient behaves like Q-learning. For clarity, we only show the argmax action in each state; we omit the arrow when the argmax action is “do nothing”.

the goal. With such limited data, unregularized classifier actor critic overestimates the Q-values at unseen actions, learning a policy that mistakenly takes these actions. In contrast, the regularized approaches learn to imitate the expert trajectory. Fig. 2 (right) shows that both forms of regularization produce the optimal policy. In summary, these tabular experiments validate our theoretical results, including in settings where regularization is useful and harmful. These experiments also demonstrate that the actor-critic method introduced in Sec. 3.4 does converge (Lemma 4.1).

5.2. Predictions about Prior Methods: Tabular Setting

Based on our theoretical analysis, we predict that practical implementations of one-step RL and critic regularization will exhibit similar behavior, for a certain critic regularization coefficient. This section studies the tabular setting, and the following section will use a continuous control benchmark. For critic regularization, we used CQL (Kumar et al., 2020) together with soft value iteration; following (Brandfonbrener et al., 2021), we implement one-step RL (reverse KL) using Q-learning.

We designed a deterministic gridworld so one-step RL would fail to learn the optimal policy (see Fig. 3 (left)). If CQL interpolates between the behavioral policy (random) and the optimal policy, then the argmax action would always be the same as the action for π^* . Based on our analysis, we make a different prediction: that CQL will learn a policy similar to the one-step RL policy. We show results in Fig. 3, just showing the argmax action for visual clarity. The CQL policy takes actions away from both the high-reward state and the low reward state, similar to the behavioral policy but different from both the behavioral policy and the optimal policy. This experiment suggests that CQL can exhibit behavior similar to one-step RL. Of course, this effect is mediated by the regularization strength: a larger regularization coefficient would cause CQL to learn a random policy, and

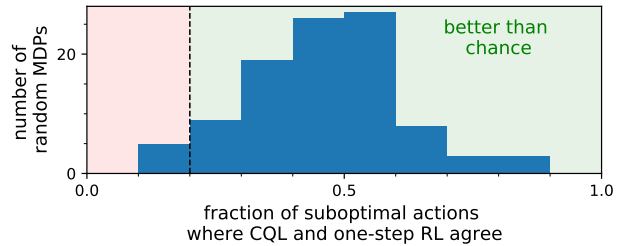


Figure 4: CQL and one-step RL take similar actions on most MDPs that resemble Fig. 3

a coefficient of 0 would make CQL identical to Q-learning. We extend Fig. 3 to include classifier actor critic and regularized variants in Appendix Fig. 9.

How often does one-step RL approximate CQL? To show that the results in Fig. 3 are not cherry-picked, we repeated this experiment using 100 MDPs that are structurally similar to that in Fig. 3, but where the locations of the high-reward and low reward state are randomized. In each randomly generated MDP, we determine whether CQL exhibits behavior similar to one-step RL by looking at the states where CQL takes actions that differ from the reward-maximizing actions (as determined by running Q-learning with unlimited data). Since there are five total actions, a random policy would have a similarity score of 20%. As shown in Fig. 4, the similarity score is significantly higher than chance for the vast majority of MDPs, showing that one-step RL and CQL($\lambda = 10$) produce similar policies on most such gridworlds.

When does one-step RL approximate CQL? Because one-step RL is highly regularized (policy iteration is truncated after just one step), one might imagine that it would be most similar to CQL with a very large regularization coefficient. To study this, we use the same environment

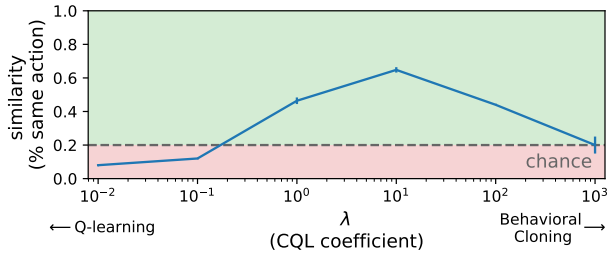


Figure 5: One-step RL is most similar to CQL with a moderate regularization coefficient.

(Fig. 3) and measure the fraction of states where one-step RL and CQL choose the same argmax action. As shown in Fig. 5, one-step RL is most similar to CQL with *moderate* regularization ($\lambda = 10$), and is less similar to CQL with a very strong regularization.

5.3. Predictions about Prior Methods: Continuous Control Setting

Our final set of experiments studies whether our theoretical results can make accurate testable predictions about practically-used regularization methods in a setting where they are commonly used: offline RL benchmarks with continuous states and actions. For these experiments, we will use well-tuned implementations of CQL and one-step RL from Hoffman et al. (2020), using the default hyperparameters without modification. While our theoretical results do not apply directly to these practical methods, which violate the assumptions in our analysis, they nonetheless raise the questions of whether these methods perform similarly in practice. We made one change to the one-step RL implementation to make the comparison more fair: because CQL learns two Q functions and takes the minimum (a trick introduced in Fujimoto et al. (2018)), we applied this same parametrization to the one-step RL implementation. Since offline RL methods can perform differently on datasets of varying quality (Wang et al., 2020; Fujimoto & Gu, 2021; Paine et al., 2020; Wang et al., 2021; Fujimoto et al., 2019), we will repeat our experiments on four datasets from the D4RL benchmark (Fu et al., 2020).

Lower bounds on Q-values. One oft-cited benefit of critic regularization is that it has guarantees about value-estimation (Kumar et al., 2020): under appropriate assumptions, the learned value function will underestimate the discounted expected returns of the policy. Because our analysis shows a connection between one-step RL and critic regularization, it raises the question of whether one-step RL methods have similar value-estimation properties. Taken at face value, this hypothesis seems obvious: the behavioral critic estimates the value of the behavioral policy, so it should underestimate the value of any policy that is better

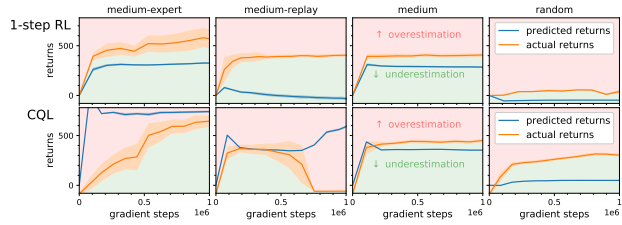


Figure 6: **Q-value under/over-estimation.** (Top) Experiments on benchmark datasets of varying quality show that one-step RL underestimates the Q-values. (Bottom) Despite the theoretical guarantees about critic regularization (CQL) yielding underestimates, in practice we observe that the values learned via critic regularization can sometimes overestimate the actual returns. We plot the mean and standard deviation across five random seeds. Note that the Q-values are equivalent to the value function $V^\pi(s)$.

than the behavioral policy. Despite this, the lower bound property of methods like one-step RL are rarely discussed, suggesting that it has yet to be widely appreciated.

Fig. 6 shows both these predicted and actual (discounted) returns throughout the course of training. The results for one-step RL confirm our theoretical prediction on $4/4$ datasets: the Q-values from one-step RL underestimate the actual returns. In contrast, we observe that critic regularization overestimates the true returns on $2/4$ environments, perhaps because the regularization coefficients used to achieve good returns in practice are too weak to guarantee the lower bound property, and perhaps because the theoretical guarantees are only guaranteed to hold at convergence. In total, these experiments confirm our theoretical predictions that one-step RL will result in Q-values that are underestimations, while also questioning the claim that critic regularization methods are always preferable for ensuring underestimation.

Critic regularization causes actor regularization. Our analysis in Sec. 4 not only suggests that one-step RL methods might inherit properties of critic regularization (as studied in the previous section), but also suggests that critic regularization methods may behave like one-step methods. In particular, while critic regularization methods such as CQL do not explicitly regularize their actor, we hypothesize that they *implicitly* regularize the actor (Lemma 4.2), similar to how one-step RL methods explicitly regularize the actor.

We measure the MSE between the action in the dataset and the action predicted by the learned policy. Fig. 7 shows the results. Some CQL implementations, including ours, “warm-start” the actor by applying a behavioral cloning loss for 50,000 iterations; we omit these initial gradient steps from our plots so that any effect is caused solely by the critic regularization. On $4/4$ datasets, we observe that the

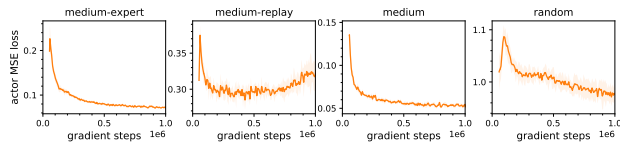


Figure 7: **Critic regularization causes actor regularization.** Performing critic regularization via CQL implicitly results in actor regularization, similar to one-step RL: the MSE between the predicted actions and the dataset actions decreases. We plot the mean and standard deviation across five random seeds.

MSE between the CQL policy’s actions and the actions in the datasets decreases throughout training. Perhaps the one exception is on the `medium-replay` dataset, where the MSE eventually starts to increase after $5e5$ gradient steps. While directly regularizing the actor leads to MSE errors that are $\sim 3\times$ smaller, this plot nevertheless provides evidence that critic regularization indirectly regularizes the actor.

6. Conclusion

In this paper, we drew a connection between two seemingly-distinct RL regularization methods: one-step RL and critic regularization. While our analysis made assumptions that are typically violated in practice, it nonetheless made accurate, testable predictions about practical methods with commonly-used hyperparameters: critic regularization methods can behave like one-step methods, and vice versa.

Acknowledgements. We thank Aviral Kumar, George Tucker, David Brandfonbrener, and Scott Fujimoto for discussions and feedback on drafts of the paper. We thank Young Geng and Ilya Kostrikov for sharing code. We thank Chongyi Zheng for spotting a bug that prompted this project. This material is supported by the Fannie and John Hertz Foundation, NSF GRFP (DGE1745016), ONR N000142312368, DARPA/AFRL FA87502321015, and the Office of Naval Research

References

- Abdolmaleki, A., Springenberg, J. T., Tassa, Y., Munos, R., Heess, N., and Riedmiller, M. Maximum a posteriori policy optimisation. In *International Conference on Learning Representations*, 2018.
- Abramovich, S. and Persson, L.-E. Some new estimates of the ‘Jensen gap’. *Journal of Inequalities and Applications*, 2016(1):1–9, 2016.
- Agarwal, A., Jiang, N., Kakade, S. M., and Sun, W. Reinforcement learning: Theory and algorithms. *CS Dept., UW Seattle, Seattle, WA, USA, Tech. Rep.*, 2019.
- An, G., Moon, S., Kim, J.-H., and Song, H. O. Uncertainty-based offline reinforcement learning with diversified Q-ensemble. *Advances in neural information processing systems*, 34:7436–7447, 2021.
- Bai, C., Wang, L., Yang, Z., Deng, Z.-H., Garg, A., Liu, P., and Wang, Z. Pessimistic bootstrapping for uncertainty-driven offline reinforcement learning. In *International Conference on Learning Representations*.
- Bakushinskii, A. B. A general method of constructing regularizing algorithms for a linear incorrect equation in hilbert space. *Zhurnal Vychislitel’noi Matematiki i Matematicheskoi Fiziki*, 7(3):672–677, 1967.
- Bauer, F., Pereverzev, S., and Rosasco, L. On regularization algorithms in learning theory. *Journal of complexity*, 23(1):52–72, 2007.
- Brandfonbrener, D., Whitney, W., Ranganath, R., and Bruna, J. Offline RL without off-policy evaluation. *Advances in Neural Information Processing Systems*, 34:4933–4946, 2021.
- Buckman, J., Gelada, C., and Bellemare, M. G. The importance of pessimism in fixed-dataset policy optimization. In *International Conference on Learning Representations*, 2021.
- Calandra, R., Owens, A., Upadhyaya, M., Yuan, W., Lin, J., Adelson, E. H., and Levine, S. The feeling of success: Does touch sensing help predict grasp outcomes? In *Conference on Robot Learning*, pp. 314–323. PMLR, 2017.
- Chebotar, Y., Hausman, K., Lu, Y., Xiao, T., Kalashnikov, D., Varley, J., Irpan, A., Eysenbach, B., Julian, R. C., Finn, C., et al. Actionable models: Unsupervised offline reinforcement learning of robotic skills. In *International Conference on Machine Learning*, pp. 1518–1528. PMLR, 2021.
- Chen, L., Lu, K., Rajeswaran, A., Lee, K., Grover, A., Laskin, M., Abbeel, P., Srinivas, A., and Mordatch, I. Decision transformer: Reinforcement learning via sequence modeling. *Advances in neural information processing systems*, 34:15084–15097, 2021.
- Ding, Y., Florensa, C., Abbeel, P., and Phielipp, M. Goal-conditioned imitation learning. *Advances in neural information processing systems*, 32, 2019.
- Emmons, S., Eysenbach, B., Kostrikov, I., and Levine, S. Rvs: What is essential for offline rl via supervised learning? In *International Conference on Learning Representations*, 2021.

- Eysenbach, B., Geng, X., Levine, S., and Salakhutdinov, R. Rewriting history with inverse rl: Hindsight inference for policy improvement. In *Advances in Neural Information Processing Systems*, 2020a.
- Eysenbach, B., Salakhutdinov, R., and Levine, S. C-learning: Learning to achieve goals via recursive classification. In *International Conference on Learning Representations*, 2020b.
- Eysenbach, B., Levine, S., and Salakhutdinov, R. Replacing rewards with examples: Example-based policy search via recursive classification. *Advances in Neural Information Processing Systems*, 34, 2021.
- Eysenbach, B., Udatha, S., Salakhutdinov, R. R., and Levine, S. Imitating past successes can be very suboptimal. *Advances in Neural Information Processing Systems*, 35: 6047–6059, 2022.
- Finn, C., Levine, S., and Abbeel, P. Guided cost learning: Deep inverse optimal control via policy optimization. In *International conference on machine learning*, pp. 49–58. PMLR, 2016.
- Fleming, H. E. Equivalence of regularization and truncated iteration in the solution of ill-posed image reconstruction problems. *Linear Algebra and its applications*, 130:133–150, 1990.
- Fu, J., Singh, A., Ghosh, D., Yang, L., and Levine, S. Variational inverse control with events: A general framework for data-driven reward definition. *Advances in neural information processing systems*, 31, 2018.
- Fu, J., Kumar, A., Nachum, O., Tucker, G., and Levine, S. D4rl: Datasets for deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- Fujimoto, S. and Gu, S. S. A minimalist approach to offline reinforcement learning. *Advances in neural information processing systems*, 34:20132–20145, 2021.
- Fujimoto, S., Hoof, H., and Meger, D. Addressing function approximation error in actor-critic methods. In *International conference on machine learning*, pp. 1587–1596. PMLR, 2018.
- Fujimoto, S., Meger, D., and Precup, D. Off-policy deep reinforcement learning without exploration. In *International conference on machine learning*, pp. 2052–2062. PMLR, 2019.
- Gao, X., Sitharam, M., and Roitberg, A. E. Bounds on the Jensen gap, and implications for mean-concentrated distributions. *arXiv preprint arXiv:1712.05267*, 2017.
- Geist, M., Scherrer, B., and Pietquin, O. A theory of regularized Markov decision processes. In *International Conference on Machine Learning*, pp. 2160–2169. PMLR, 2019.
- Ghosh, D., Gupta, A., Reddy, A., Fu, J., Devin, C. M., Eysenbach, B., and Levine, S. Learning to reach goals via iterated supervised learning. In *International Conference on Learning Representations*, 2020.
- Gülçehre, Ç., Wang, Z., Novikov, A., Le Paine, T., Colmenarejo, S. G., Zolna, K., Agarwal, R., Merel, J., Mankowitz, D. J., Paduraru, C., et al. RL unplugged: Benchmarks for offline reinforcement learning. 2020.
- Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pp. 1861–1870. PMLR, 2018.
- Hatch, K., Yu, T., Rafailov, R., and Finn, C. Example-based offline reinforcement learning without rewards. *Proceedings of Machine Learning Research vol*, 144:1–17, 2022.
- Hazan, E., Kakade, S., Singh, K., and Van Soest, A. Provably efficient maximum entropy exploration. In *International Conference on Machine Learning*, pp. 2681–2691. PMLR, 2019.
- Hoffman, M., Shahriari, B., Aslanides, J., Barth-Maron, G., Behbahani, F., Norman, T., Abdolmaleki, A., Cassirer, A., Yang, F., Baumli, K., Henderson, S., Novikov, A., Colmenarejo, S. G., Cabi, S., Gulcehre, C., Paine, T. L., Cowie, A., Wang, Z., Piot, B., and de Freitas, N. Acme: A research framework for distributed reinforcement learning. *arXiv preprint arXiv:2006.00979*, 2020.
- Huntley, H. E. *Dimensional analysis*. Dover publications, 1967.
- Jaques, N., Ghandeharioun, A., Shen, J. H., Ferguson, C., Lapedriza, A., Jones, N., Gu, S., and Picard, R. Way off-policy batch deep reinforcement learning of implicit human preferences in dialog. *arXiv preprint arXiv:1907.00456*, 2019.
- Kalashnikov, D., Irpan, A., Pastor, P., Ibarz, J., Herzog, A., Jang, E., Quillen, D., Holly, E., Kalakrishnan, M., Vanhoucke, V., et al. Scalable deep reinforcement learning for vision-based robotic manipulation. In *Conference on Robot Learning*, pp. 651–673. PMLR, 2018.
- Kalashnikov, D., Varley, J., Chebotar, Y., Swanson, B., Jonschkowski, R., Finn, C., Levine, S., and Hausman, K. Scaling up multi-task robotic reinforcement learning. In *5th Annual Conference on Robot Learning*, 2021.

- Kostrikov, I., Nair, A., and Levine, S. Offline reinforcement learning with implicit Q-learning. In *International Conference on Learning Representations*, 2021.
- Kumar, A., Fu, J., Soh, M., Tucker, G., and Levine, S. Stabilizing off-policy Q-learning via bootstrapping error reduction. *Advances in Neural Information Processing Systems*, 32, 2019a.
- Kumar, A., Peng, X. B., and Levine, S. Reward-conditioned policies. *arXiv preprint arXiv:1912.13465*, 2019b.
- Kumar, A., Zhou, A., Tucker, G., and Levine, S. Conservative Q-learning for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 33: 1179–1191, 2020.
- Lange, S., Gabel, T., and Riedmiller, M. Batch reinforcement learning. In *Reinforcement learning*, pp. 45–73. Springer, 2012.
- Levine, S., Kumar, A., Tucker, G., and Fu, J. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- Li, A., Pinto, L., and Abbeel, P. Generalized hindsight for reinforcement learning. *Advances in neural information processing systems*, 33:7754–7767, 2020.
- Lu, C., Ball, P., Parker-Holder, J., Osborne, M., and Roberts, S. J. Revisiting design choices in offline model based reinforcement learning. In *International Conference on Learning Representations*, 2021.
- Lynch, C. and Sermanet, P. Language conditioned imitation learning over unstructured data. *Proceedings of Robotics: Science and Systems*, 2021.
- Lyu, J., Ma, X., Yan, J., and Li, X. Efficient continuous control with double actors and regularized critics. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 7655–7663, 2022.
- Mnih, V., Kavukcuoglu, K., Silver, D., Graves, A., Antonoglou, I., Wierstra, D., and Riedmiller, M. A. Playing atari with deep reinforcement learning. *ArXiv*, abs/1312.5602, 2013.
- Nachum, O., Dai, B., Kostrikov, I., Chow, Y., Li, L., and Schuurmans, D. Algaedice: Policy gradient from arbitrary experience. *arXiv preprint arXiv:1912.02074*, 2019.
- Nash, J. Non-cooperative games. *Annals of mathematics*, pp. 286–295, 1951.
- Neu, G., Jonsson, A., and Gómez, V. A unified view of entropy-regularized Markov decision processes. *arXiv preprint arXiv:1705.07798*, 2017.
- Paine, T. L., Paduraru, C., Michi, A., Gulcehre, C., Zolna, K., Novikov, A., Wang, Z., and de Freitas, N. Hyperparameter selection for offline reinforcement learning. *arXiv preprint arXiv:2007.09055*, 2020.
- Paster, K., McIlraith, S. A., and Ba, J. Planning from pixels using inverse dynamics models. In *International Conference on Learning Representations*, 2021.
- Peng, X. B., Kumar, A., Zhang, G., and Levine, S. Advantage-weighted regression: Simple and scalable off-policy reinforcement learning. *arXiv preprint arXiv:1910.00177*, 2019.
- Peters, J. and Schaal, S. Reinforcement learning by reward-weighted regression for operational space control. In *Proceedings of the 24th international conference on Machine learning*, pp. 745–750, 2007.
- Peters, J., Mulling, K., and Altun, Y. Relative entropy policy search. In *Twenty-Fourth AAAI Conference on Artificial Intelligence*, 2010.
- Pinto, L. and Gupta, A. Supersizing self-supervision: Learning to grasp from 50k tries and 700 robot hours. In *2016 IEEE international conference on robotics and automation (ICRA)*, pp. 3406–3413. IEEE, 2016.
- Rezaeifar, S., Dadashi, R., Vieillard, N., Hussenot, L., Bachem, O., Pietquin, O., and Geist, M. Offline reinforcement learning as anti-exploration. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 8106–8114, 2022.
- Santos, R. J. Equivalence of regularization and truncated iteration for general ill-posed problems. *Linear algebra and its applications*, 236:25–33, 1996.
- Savinov, N., Dosovitskiy, A., and Koltun, V. Semi-parametric topological memory for navigation. In *International Conference on Learning Representations*, 2018.
- Scherrer, B. Improved and generalized upper bounds on the complexity of policy iteration. *Advances in Neural Information Processing Systems*, 26, 2013.
- Siegel, N., Springenberg, J. T., Berkenkamp, F., Abdolmaleki, A., Neunert, M., Lampe, T., Hafner, R., Heess, N., and Riedmiller, M. Keep doing what worked: Behavior modelling priors for offline reinforcement learning. In *International Conference on Learning Representations*.
- Singh, A., Yang, L., Hartikainen, K., Finn, C., and Levine, S. End-to-end robotic reinforcement learning without reward engineering. *arXiv preprint arXiv:1904.07854*, 2019.

- Srivastava, R. K., Shyam, P., Mutz, F., Jaśkowski, W., and Schmidhuber, J. Training agents using upside-down reinforcement learning. *arXiv preprint arXiv:1912.02877*, 2019.
- Sun, H., Li, Z., Liu, X., Zhou, B., and Lin, D. Policy continuation with hindsight inverse dynamics. *Advances in Neural Information Processing Systems*, 32, 2019.
- Sutton, R. S. and Barto, A. G. *Reinforcement learning: An introduction*. MIT press, 2018.
- Tung, H.-Y., Harley, A. W., Huang, L.-K., and Fragkiadaki, K. Reward learning from narrated demonstrations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7004–7013, 2018.
- Vieillard, N., Kozuno, T., Scherrer, B., Pietquin, O., Munos, R., and Geist, M. Leverage the average: an analysis of kl regularization in reinforcement learning. *Advances in Neural Information Processing Systems*, 33:12163–12174, 2020.
- Villafior, A., Dolan, J., and Schneider, J. Fine-tuning offline reinforcement learning with model-based policy optimization. *Offline Reinforcement Learning Workshop at Neural Information Processing Systems*, 2020.
- Wahba, G. Three topics in ill-posed problems. In *Inverse and ill-posed problems*, pp. 37–51. Elsevier, 1987.
- Wang, Q., Xiong, J., Han, L., Liu, H., Zhang, T., et al. Exponentially weighted imitation learning for batched historical data. *Advances in Neural Information Processing Systems*, 31, 2018.
- Wang, R., Wu, Y., Salakhutdinov, R., and Kakade, S. Instabilities of offline rl with pre-trained neural representation. In *International Conference on Machine Learning*, pp. 10948–10960. PMLR, 2021.
- Wang, Z., Novikov, A., Zolna, K., Merel, J. S., Springenberg, J. T., Reed, S. E., Shahriari, B., Siegel, N., Gulcehre, C., Heess, N., et al. Critic regularized regression. *Advances in Neural Information Processing Systems*, 33: 7768–7778, 2020.
- Wen, J., Kumar, S., Gummadi, R., and Schuurmans, D. Characterizing the gap between actor-critic and policy gradient. In *International Conference on Machine Learning*, pp. 11101–11111. PMLR, 2021.
- Wu, Y., Tucker, G., and Nachum, O. Behavior regularized offline reinforcement learning. *arXiv preprint arXiv:1911.11361*, 2019.
- Yang, R., Fang, M., Han, L., Du, Y., Luo, F., and Li, X. MHER: Model-based hindsight experience replay. *ArXiv*, abs/2107.00306, 2021.
- Yu, T., Kumar, A., Rafailov, R., Rajeswaran, A., Levine, S., and Finn, C. Combo: Conservative offline model-based policy optimization. *Advances in neural information processing systems*, 34:28954–28967, 2021.
- Zhou, W., Bajracharya, S., and Held, D. Plas: Latent action space for offline reinforcement learning. In *Conference on Robot Learning*, pp. 1719–1735. PMLR, 2021.
- Ziebart, B. D. *Modeling purposeful adaptive behavior with the principle of maximum causal entropy*. Carnegie Mellon University, 2010.
- Zolna, K., Novikov, A., Konyushkova, K., Gulcehre, C., Wang, Z., Aytar, Y., Denil, M., de Freitas, N., and Reed, S. Offline learning from demonstrations and unlabeled experience. *arXiv preprint arXiv:2011.13885*, 2020.

Appendices

In the Appendices, we provide proofs of the theoretical results (Appendix C), extend the analysis to other RL settings (Appendices C–D), and then provide details of the experiments (Appendix F).

A. Proofs

A.1. Proof of Lemma 4.1

Proof sketch. Lemma 4.1 shows that classifier actor critic converges. The key idea of the proof will be to show that the incremental updates for classifier actor critic are exactly the same as the incremental updates for Q-learning. Q-learning converges, so an algorithm that performs the same incremental updates as Q-learning must also converge.

Proof. As the cross entropy loss is minimized when the predictions equal the labels, updates for $\mathcal{L}_{\text{critic}}(Q, \pi)$ can be written as $\frac{Q(s,a)}{Q(s,a)+1} \leftarrow \frac{y^{\pi, Q_t}(s,a)}{y^{\pi, Q_t}(s,a)+1}$. If the updates are performed by averaging over all possible next states (e.g., in the tabular setting), these updates are equivalent to directly updating $Q(s, a) \leftarrow y^{\pi, Q_t}(s, a) = r(s, a) + \gamma \mathbb{E}_{p(s'|s, a)\pi(a'|s')} [Q_t(s', a')]$, which is the standard policy evaluation update for policy $\pi(a | s)$. Thus, we can invoke the standard result that policy evaluation converges to Q^π (Agarwal et al., 2019, Theorem 1.14.) to argue that updates for $\mathcal{L}_{\text{critic}}$ likewise converge to Q^π . \square

In this proof, the TD targets were the expectation over the next state and next action. If Eq. 3 were optimized using a single-sample estimate of this expectation, $\mathbf{y} = r(s, a) + \gamma Q_t(s', a')$, then the updates would be biased:

$$\frac{Q(s, a)}{Q(s, a) + 1} \leftarrow \mathbb{E} \left[\frac{\mathbf{y}}{\mathbf{y} + 1} \right] \leq \frac{\mathbb{E}[\mathbf{y}]}{\mathbb{E}[\mathbf{y}] + 1} = \frac{y^{\pi, Q_t}(s, a)}{y^{\pi, Q_t}(s, a) + 1}.$$

In settings with stochastic transitions or policies, these updates would result in estimating a lower bound on $Q^\pi(s, a)$.

A.2. Proof of Lemma 4.2 and Theorem 4.3

Proof. Our proof proceeds in three steps. First, we derive the update equations for the regularized critic update. That is, if we maintained a table of Q-values, what would the new value for $Q(s, a)$ be? Second, we show that these updates are equivalent to performing policy evaluation on a *re-parametrized* critic $\tilde{Q}(s, a) = Q(s, a) \frac{\pi(a|s)}{\beta(a|s)}$. We invoke the standard results for policy evaluation to prove convergence that $\tilde{Q}(s, a)$ converges. Finally, we undo the reparametrization to obtain convergence results for $Q(s, a)$.

Step 0. We start by rearranging the regularized critic objective:

$$\begin{aligned} \mathcal{L}_{\text{critic}}^r(Q, y^{\pi, Q_t}) &\triangleq \mathcal{L}_{\text{critic}}(Q, y^{\pi, Q_t}) + \left(\mathbb{E}_{p(s)\pi(a|s)} [\log(Q(s, a) + 1)] - \mathbb{E}_{p(s)\beta(a|s)} [\log(Q(s, a) + 1)] \right) \\ &= -\mathbb{E}_{p(s,a)} \left[y^{\pi, Q_t}(s, a) \log \frac{Q(s, a)}{Q(s, a) + 1} + \log \frac{1}{Q(s, a) + 1} \right] \\ &\quad + \left(\mathbb{E}_{p(s)\pi(a|s)} [\log(Q(s, a) + 1)] - \mathbb{E}_{p(s)\beta(a|s)} [\log(Q(s, a) + 1)] \right) \\ &= -\mathbb{E}_{p(s,a)} \left[y^{\pi, Q_t}(s, a) \log \frac{Q(s, a)}{Q(s, a) + 1} + \log \frac{1}{Q(s, a) + 1} \right] \\ &\quad - \left(\mathbb{E}_{p(s)\pi(a|s)} \left[\log \frac{1}{Q(s, a) + 1} \right] + \mathbb{E}_{p(s)\beta(a|s)} \left[\log \frac{1}{Q(s, a) + 1} \right] \right) \\ &= -\mathbb{E}_{p(s,a)} \left[y^{\pi, Q_t}(s, a) \log \frac{Q(s, a)}{Q(s, a) + 1} \right] + \mathbb{E}_{p(s)\beta(a|s)} \left[\log \frac{1}{Q(s, a) + 1} \right]. \end{aligned}$$

For the cancellation on the third line, we used the fact that $p(s, a) = p(s)\beta(a | s)$.

Step 1. To start, note that the regularized critic update is equivalent to a weighted classification loss: positive examples are sampled $(s, a) \sim p(s)\beta(a | s)$ and receive weight $\frac{y^{\pi, Q_t}(s, a)}{y^{\pi, Q_t}(s, a)+1}$, and negative examples are sampled $(s, a) \sim p(s)\pi(a | s)$

and receive weight $\frac{1}{y^{\pi, Q_t(s,a)+1}}$. The Bayes' optimal classifier is given by

$$\frac{Q(s, a)}{Q(s, a) + 1} = \frac{\frac{y^{\pi, Q_t(s,a)}}{y^{\pi, Q_t(s,a)+1}} p(s) \beta(a | s)}{\frac{y^{\pi, Q_t(s,a)}}{y^{\pi, Q_t(s,a)+1}} p(s) \beta(a | s) + \frac{1}{y^{\pi, Q_t(s,a)+1}} p(s) \pi(a | s)} = \frac{y^{\pi, Q_t(s,a)} \beta(a | s)}{y^{\pi, Q_t(s,a)} \beta(a | s) + \pi(a | s)}.$$

Solving for $Q(s, a)$ on the left hand side, the optimal value for $Q(s, a)$ is given by

$$Q(s, a) = y^{\pi, Q_t(s,a)} \frac{\beta(a | s)}{\pi(a | s)} = (r(s, a) + \mathbb{E}_{p(s'|s,a)\pi(a'|s')} [Q_t(s', a')]) \frac{\beta(a | s)}{\pi(a | s)}. \quad (7)$$

This equation tells us what each update for the regularized critic loss does.

Step 2. To analyze these updates, we define $\tilde{Q}(s, a) \triangleq Q_t(s, a) \frac{\pi(a|s)}{\beta(a|s)}$. Then these updates can be written using $\tilde{Q}(s, a)$ as

$$\tilde{Q}(s, a) \frac{\beta(a | s)}{\pi(a | s)} = \left(r(s, a) + \mathbb{E}_{p(s'|s,a)\pi(a'|s')} \left[\tilde{Q}(s', a') \frac{\beta(a' | s')}{\pi(a' | s')} \right] \right) \frac{\beta(a | s)}{\pi(a | s)}, \quad (8)$$

which can be simplified to

$$\tilde{Q}(s, a) = r(s, a) + \mathbb{E}_{p(s'|s,a)\beta(a'|s')} \left[\tilde{Q}(s', a') \right]. \quad (9)$$

Note that the ratio $\frac{\beta(a'|s')}{\pi(a'|s')}$ inside the expectation acts like an importance weight, so that the expectation over $\pi(a' | s')$ becomes an expectation over $\beta(a' | s')$. Thus, the regularized critic updates are equivalent to perform policy evaluation on $\tilde{Q}(s, a)$. An immediate consequence is that the regularized critic updates converge, and they converge to $\tilde{Q}^*(s, a) = Q^\beta(s, a)$.

Step 3. Finally, we translate these convergence results for $\tilde{Q}(s, a)$ into convergence results for $Q(s, a)$. Written in terms of the original Q-values, we see that the optimal critic for the regularized critic update is

$$Q^*(s, a) = \tilde{Q}^*(s, a) \frac{\beta(a | s)}{\pi(a | s)} = Q^\beta(s, a) \frac{\beta(a | s)}{\pi(a | s)}. \quad (10)$$

This completes the proof of Lemma 4.2. □

We now prove Theorem 4.3 by applying a logarithm:

Proof.

$$\log Q^*(s, a) = \log \left(Q^\beta(s, a) \frac{\beta(a | s)}{\pi(a | s)} \right) = \log Q^\beta(s, a) + \log \beta(a | s) - \log \pi(a | s).$$

□

We note that our proof does not account for stochastic and function approximation errors. However, if we assume that the TD updates are deterministic (e.g., as they are in deterministic MDPs), then the updates for classifier actor-critic are identical to those of Q-learning (Lemma 4.1). Thus, it immediately inherits any theoretical results regarding the propagation of errors for Q-learning.

While this Theorem 4.3 shows that one-step RL and critic regularization have the same fixed point, it does not say how many transitions or gradient updates are required to reach those fixed points.

A.3. Why use the cross-entropy loss?

Our proof of Theorem 4.3 helps explain why classifier actor-critic use the cross entropy loss for the critic loss, rather than the MSE loss. Precisely, our analysis requires that the optimal Q function be a *ratio*, $\tilde{Q}(s, a) = \frac{Q(s,a)\pi(a|s)}{\beta(a|s)}$. The cross entropy loss can readily estimate ratios. For example, the optimal classifier for data drawn from $p(x)$ and $q(x)$ is $C(x) = \frac{p(x)}{p(x)+q(x)}$,

so the ratio can be expressed as $\frac{C(x)}{1-C(x)} = \frac{p(x)}{q(x)}$. However, fitting a function $C(x)$ to data drawn from (say) a 1:1 mixture of $p(x)$ and $q(x)$ would result in $C(x) = \frac{1}{2}p(x) + \frac{1}{2}q(x)$, which we cannot transform to express the ratio $\frac{p(x)}{q(x)}$ as a function of $C(x)$.

A.4. Validating the Theory

Our theoretical results suggest that one-step RL and critic regularization should be most similar with critic regularization is applied with a regularization coefficient of $\lambda = 1$. To test this hypothesis, we took the task from Fig. 2 (Left) and measured the similarity between one-step RL and critic-regularized classifier actor critic, for varying values of the critic regularization parameter. We measured the similarity of the policies obtained by the two methods by counting the fraction of states where the two methods choose the same (argmax) action. The results, shown in Fig. 8, validate our theoretical prediction that these methods should be most similar with $\lambda = 1$.

A.5. What about using the policy gradient?

Our analysis fundamentally requires using TD learning: the key step is that doing TD backups using one policy is equivalent to doing (modified) TD backups with a different policy. However, the actor updates for both methods could be implemented using a policy gradient or natural gradient, rather than a straight-through gradient estimator. Indeed, much of the work on one-step RL methods (Peng et al., 2019; Siegel et al.) uses an actor update that resembles a policy gradient or natural policy gradient (e.g., 1-step RL with a reverse KL penalty (Brandfonbrener et al., 2021)).

B. Varying the regularization coefficient

While our main analysis (Theorem 4.3) showed that regularization and critic regularization yield the same policy when these regularizers are applied with a certain strength, in practice the strength of regularization is controlled by a hyperparameter. This hyperparameter raises a question: *does the connection between one-step RL and critic regularization hold for different values of this hyperparameter?*

In this section, we show that there remains a precise connection between actor and critic regularization, even for different values of this hyperparameter. This result not only suggests that the connection is stronger than initially suggested by the main result. Proving this connection also helps highlight how many regularization methods can be cast from a similar mold.

B.1. A Regularization Coefficient.

We start by modifying the actor regularizer and critic regularizer introduced in Sec. 3.4 to include an additional hyperparameter.

Mixture policy. Both the actor and critic losses will make use of a mixture policy, $(1 - \lambda)\pi(a | s) + \lambda\beta(a | s)$, where $\lambda \in [0, 1]$ will be a hyperparameter. Larger values of λ yield a mixture policy that is closer to the behavioral policy; this will correspond to higher degrees of regularization. Mixtures of policies are commonly used in practice (Kumar et al., 2020, Appendix F), (Villaflor et al., 2020, Eq. 11), (Finn et al., 2016, Sec. 4.3) (Lyu et al., 2022) (Hazan et al., 2019, Eq. 2.5), even though it rarely appears in theoretical offline RL literature. Indeed, because critic regularization resembles a two-player zero-sum game, mixture policies might even be *required* to find a (Nash) equilibrium of the critic regularizer (Nash, 1951).

λ -weighted critic loss. With this concept of a mixture policy, we define the λ -weighted actor and critic regularizers. For the λ -weighted critic loss, we will change how the TD targets are computed. Instead of sampling the next action from π or β , we will sample the next action from a λ_{TD} -weighted combination of these two policies, reminiscent of how prior work

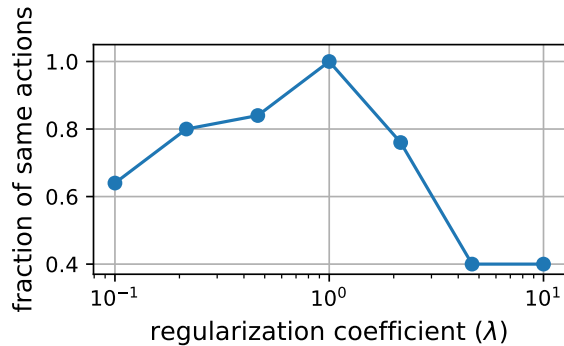


Figure 8: Under the assumptions of Theorem 4.3, one-step RL is most similar to critic regularization with a coefficient of $\lambda = 1$.

has regularized the actions sampled for the TD backup (Fujimoto et al., 2019; Zhou et al., 2021):

$$y^{\lambda_{\text{TD}}} \triangleq y^{(1-\lambda)\pi+\lambda\beta}(s, a) = r(s, a) + \gamma \mathbb{E}_{\substack{p(s'|s, a) \\ (1-\lambda_{\text{TD}})\pi(a|s)+\lambda_{\text{TD}}\beta(a|s)}} [Q(s', a')].$$

When introducing one-step RL in Sec. 3.4, we used $\lambda_{\text{TD}} = 1$.

Using this TD target, the λ -weighted critic loss can now be written as a combination of the un-regularized objective (Eq. 3) plus the regularized objective (Eq. 6):

$$\begin{aligned} \mathcal{L}_{\text{critic}}^r(Q, \lambda_{\text{critic}}) &\triangleq (1 - \lambda_{\text{critic}}) \left(-\mathbb{E}_{p(s, a)} \left[\frac{y^{\lambda_{\text{TD}}}(s, a)}{y^{\lambda_{\text{TD}}}(s, a) + 1} \log \frac{Q(s, a)}{Q(s, a) + 1} + \frac{1}{y^{\lambda_{\text{TD}}}(s, a) + 1} \log \frac{1}{Q(s, a) + 1} \right] \right) \\ &\quad + \lambda \left(-\mathbb{E}_{\substack{p(s, a) \\ a^- \sim \pi(\cdot|s)}} \left[\frac{y^{\lambda_{\text{TD}}}(s, a)}{y^{\lambda_{\text{TD}}}(s, a) + 1} \log \frac{Q(s, a)}{Q(s, a) + 1} + \frac{1}{y^{\lambda_{\text{TD}}}(s, a) + 1} \log \frac{1}{Q(s, a) + 1} \right] \right) \\ &= -\mathbb{E}_{\substack{p(s, a) \\ a^- \sim (1-\lambda_{\text{critic}})\pi(\cdot|s)+\lambda_{\text{critic}}\beta(\cdot|s)}} \left[\frac{y^{\lambda_{\text{TD}}}(s, a)}{y^{\lambda_{\text{TD}}}(s, a) + 1} \log \frac{Q(s, a)}{Q(s, a) + 1} + \frac{1}{y^{\lambda_{\text{TD}}}(s, a) + 1} \log \frac{1}{Q(s, a^-) + 1} \right]. \end{aligned} \quad (11)$$

The second line rewrites this objective: the first term looks the same as the original “positive” term in the critic objective, while the “negative” term uses actions sampled from a mixture of the current policy and the behavioral policy. When $\lambda_{\text{critic}} = 1$, we recover the regularized critic loss introduced in Sec. 3.4.

λ -weighted actor loss. Finally, the strength of the actor regularizer can be controlled by changing the reverse KL penalty. While it may seem like changing the reward scale would vary the strength of the actor loss, this is not the case for classifier actor critic because of the $\log(\cdot)$ in the actor loss. Instead, we will relax the reverse KL penalty between the learned policy $\pi(a | s)$ and the behavioral policy $\beta(a | s)$ so that only the mixture policy only needs to be close to behavioral policy:

$$\mathcal{L}_{\text{actor}}^r(\pi, \lambda_{\text{KL}}) \triangleq \mathbb{E}_{p(s)\pi(a|s)} [\log Q(s, a) + \log \beta(a | s) - \log ((1 - \lambda_{\text{KL}})\pi(a | s) + \lambda_{\text{KL}}\beta(a | s))]. \quad (12)$$

As indicated on the second line, replacing $\beta(a | s)$ with the mixture policy has an effect similar to that of decreasing the weight applied to the KL penalty. The approximation on the second line is determined by the Jensen Gap (Abramovich & Persson, 2016; Gao et al., 2017). When introducing one-step RL in Sec. 3.4, we used $\lambda_{\text{KL}} = 1$, together with $\lambda_{\text{TD}} = 1$.

In summary, the strength of the actor and critic regularizers can be controlled through additional hyperparameters ($\lambda_{\text{critic}}, \lambda_{\text{TD}}, \lambda_{\text{KL}}$). Indeed, it is typical for offline RL methods to require many hyperparameters (Brandfonbrener et al., 2021; Lu et al., 2021; Paine et al., 2020; Wu et al., 2019), and performance is sensitive to their settings. However, the close connection that we have shown between actor and critic regularizers allows us to decrease the number of hyperparameters.

B.2. Analysis

In our main result (Thm. 4.3), we showed that one-step RL and critic regularization are equivalent when $\lambda_{\text{critic}} = \lambda_{\text{TD}} = \lambda_{\text{KL}} = 1$. This is a large value for the regularization strength, and we now consider what happens for smaller degrees of regularization: is there still a connection between one-step RL and critic regularization?

The following theorem will prove that this is the case. In particular, applying critic regularization with coefficient λ_{critic} yields the same policy as applying one-step RL with $\lambda_{\text{TD}} = \lambda_{\text{KL}} = \lambda_{\text{critic}}$. That is, there is a very simple recipe for converting the hyperparameters for critic regularization into the hyperparameters for one-step RL.

Theorem B.1. *Let policy $\pi(a | s)$ be given, let $Q^\beta(s, a)$ be the Q -function of the behavioral policy, and let $Q_r^{\lambda_{\text{TD}}}(s, a, \lambda_{\text{critic}})$ be the critic obtained by the λ_{critic} -weighted regularized critic update (Eq. 11) using TD targets $y^{\lambda_{\text{TD}}}(s, a)$. If $\lambda_{\text{critic}} = \lambda_{\text{TD}} = \lambda_{\text{KL}}$, then the λ_{KL} -weighted actor loss (Eq. 12) is equivalent to the un-regularized policy objective using the regularized critic:*

$$\begin{aligned} &\mathbb{E}_{p(s)\pi(a|s)} [\log Q(s, a) + \log \beta(a | s) - \log ((1 - \lambda_{\text{KL}})\pi(a | s) + \lambda_{\text{KL}}\beta(a | s))] \\ &= \mathbb{E}_{\pi(a|s)} [\log Q_r^{\lambda_{\text{TD}}}(s, a, \lambda_{\text{critic}})] \quad \text{for all states } s. \end{aligned}$$

While we used the cross entropy loss for this result, it turns out that the result also holds for the more standard MSE loss (we omit the proof for brevity).

Limitations. Before presenting the proof in Sec. B.3, we discuss a few limitations of this result. Like the rest of the analysis in this paper, the form of the critic regularizer is different from that often used in practice. Additionally, our analysis assumes ignores many sources of errors (e.g., sampling, function approximation), and assumes that each objective is optimized exactly.

B.3. Proof of Theorem B.1

Proof. We start by defining the fixed point of the λ -weighted regularized critic loss. Like in the single-task setting, this loss resembles a weighted classification problem, so we can write down the Bayes' optimal classifier as

$$\begin{aligned} \frac{Q(s, a)}{Q(s, a) + 1} &= \frac{\frac{y^{\lambda_{\text{TD}}(s, a)}}{y^{\lambda_{\text{TD}}(s, a)} + 1} p(s) \beta(a | s)}{\frac{y^{\lambda_{\text{TD}}(s, a)}}{y^{\lambda_{\text{TD}}(s, a)} + 1} p(s) \beta(a | s) + \frac{1}{y^{\lambda_{\text{TD}}(s, a)} + 1} p(s) ((1 - \lambda_{\text{critic}}) \pi(a | s) + \lambda_{\text{critic}} \beta(a | s))} \\ &= \frac{y^{\lambda_{\text{TD}}(s, a)} \beta(a | s)}{y^{\lambda_{\text{TD}}(s, a)} \beta(a | s) + (1 - \lambda_{\text{critic}}) \pi(a | s) + \lambda_{\text{critic}} \beta(a | s)}. \end{aligned}$$

Solving for $Q(s, a)$ on the left hand side, the optimal value for $Q(s, a)$ is given by

$$\begin{aligned} Q(s, a) &= y^{\lambda_{\text{TD}}(s, a)} \frac{\beta(a | s)}{(1 - \lambda_{\text{critic}}) \pi(a | s) + \lambda_{\text{critic}} \beta(a | s)} \\ &= (r(s, a) + \mathbb{E}_{p(s' | s, a), a' \sim (1 - \lambda_{\text{TD}}), \pi(\cdot | s') + \lambda_{\text{TD}} \beta(\cdot | s)} [Q(s', a')]) \frac{\beta(a | s)}{(1 - \lambda_{\text{critic}}) \pi(a | s) + \lambda_{\text{critic}} \beta(a | s)}. \end{aligned} \quad (13)$$

Note that the next action a' is sampled from a mixture policy defined by λ_{TD} . This equation tells us what each update for the λ -weighted regularized critic loss does.

To analyze these updates, we define

$$\tilde{Q}(s, a) \triangleq Q(s, a) \frac{(1 - \lambda_{\text{critic}}) \pi(a | s) + \lambda_{\text{critic}} \beta(a | s)}{\beta(a | s)}.$$

Like before, the ratio $\frac{\beta(a' | s')}{(1 - \lambda_{\text{TD}}) \pi(a' | s') + \lambda_{\text{TD}} \beta(a' | s')}$ can act like an importance weight. When $\lambda_{\text{TD}} = \lambda_{\text{critic}}$, then this importance weight cancels with the sampling distribution, providing the following identity:

$$\begin{aligned} &\mathbb{E}_{p(s' | s, a), a' \sim (1 - \lambda_{\text{TD}}), \pi(\cdot | s') + \lambda_{\text{TD}} \beta(\cdot | s)} [Q(s', a')] \\ &= \mathbb{E}_{p(s' | s, a), a' \sim (1 - \lambda_{\text{TD}}), \pi(\cdot | s') + \lambda_{\text{TD}} \beta(\cdot | s)} \left[\tilde{Q}(s, a) \frac{\beta(a | s)}{(1 - \lambda_{\text{critic}}) \pi(a | s) + \lambda_{\text{critic}} \beta(a | s)} \right] \\ &= \mathbb{E}_{p(s' | s, a), a' \sim \beta(\cdot | s')} [\tilde{Q}(s, a)]. \end{aligned}$$

Substituting this identity in Eq. 13, we can write the updates using $\tilde{Q}(s, a)$:

$$\begin{aligned} &\tilde{Q}(s, a) \frac{\beta(a | s)}{(1 - \lambda_{\text{critic}}) \pi(a | s) + \lambda_{\text{critic}} \beta(a | s)} \\ &= \left(r(s, a) + \mathbb{E}_{p(s' | s, a), a' \sim \beta(\cdot | s')} [\tilde{Q}(s, a)] \right) \frac{\beta(a | s)}{(1 - \lambda_{\text{critic}}) \pi(a | s) + \lambda_{\text{critic}} \beta(a | s)}, \end{aligned}$$

which can be simplified to

$$\tilde{Q}(s, a) = r(s, a) + \mathbb{E}_{p(s' | s, a), a' \sim \beta(\cdot | s')} [\tilde{Q}(s, a)].$$

We then translate these convergence results for $\tilde{Q}(s, a)$ into convergence results for $Q(s, a)$. Written in terms of the original Q-values, we see that the optimal critic for the regularized critic update is

$$Q^*(s, a) = Q^\beta(s, a) \frac{\beta(a | s)}{(1 - \lambda_{\text{critic}}) \pi(a | s) + \lambda_{\text{critic}} \beta(a | s)}. \quad (14)$$

Note that this holds for any value of $\lambda_{\text{critic}} = \lambda_{\text{TD}} \in [0, 1]$. This result suggests that two common forms of regularization, decreasing the values predicted at unseen actions and regularizing the actions used in the TD backup, can produce the same effect: a critic that estimates the Q-values of the behavioral policy (multiplied by some importance weight).

Finally, substitute this Q-function into the un-regularized actor loss, we see that the result is equivalent to the λ -weighted actor loss:

$$\mathbb{E}_{p(s)\pi(a|s)} [\log Q^*(s, a)] = \mathbb{E}_{p(s)\pi(a|s)} \left[\log Q^\beta(s, a) + \underbrace{\log \beta(a | s) - \log ((1 - \lambda_{\text{KL}})\pi(a | s) + \lambda_{\text{KL}}\beta(a | s))}_{\lambda\text{-weighted actor regularizer}} \right]$$

□

C. Regularization for Goal-Conditioned Problems

Like single-task RL problems, goal-conditioned RL problems have also been approached with both one-step methods (Ghosh et al., 2020; Ding et al., 2019; Sun et al., 2019) and critic regularization (Chebotar et al., 2021). In these problems, the aim is to learn a goal-conditioned policy $\pi(a | s, s_g)$ that maximizes the expected discounted sum of goal-conditioned rewards $r_g(s, a)$, where goals are sampled $s_g \sim p_g(s_g)$:

$$\max_{\pi} \mathbb{E}_{p_g(s_g)} \mathbb{E}_{\pi(\tau|s_g)} \left[\sum_{t=0}^{\infty} \gamma^t r_g(s_t, a_t) \right].$$

We will use the goal-conditioned reward function $r_g(s, a) = p(s' = s_g | s, a)$, which is defined in terms of the environment dynamics. In settings with discrete states, maximizing this reward function is equivalent to maximizing the sparse indicator reward function ($r_g(s, a) = \mathbb{1}(s_g = s)$).

In this section, we show that one-step RL and critic regularization are equivalent for a certain goal-conditioned actor-critic method. Unlike our analysis in the single-task setting, this analysis here uses an existing method, C-learning (Eysenbach et al., 2020b). C-learning is a TD method that already makes use of the cross entropy loss for training the critic:

$$\begin{aligned} \max_Q (1 - \gamma) \mathbb{E}_{p(s,a,s')} \left[\log \frac{Q(s, a, s_g = s')}{Q(s, a, s_g = s') + 1} \right] &+ \gamma \mathbb{E}_{p(s,a)p_g(s_g)} \left[y^{\pi, Q_t}(s, a, s) \log \frac{Q(s, a, s_g)}{Q(s, a, s_g) + 1} \right] \\ &+ \mathbb{E}_{p(s,a)p_g(s_g)} \left[\log \frac{1}{Q(s, a, s_g = s') + 1} \right], \end{aligned}$$

where $y^{\pi, Q_t}(s, a, s_g) = \mathbb{E}_{p(s'|s,a)\pi(a'|s',s_g)} [Q(s', a', s_g)]$ serves the role of the TD target.

The first two terms increase the Q-values while the last term decreases the Q-values. The actor is updated to maximize the Q-values. While this objective for the actor can be written in many ways, we will write it as maximizing a log ratio because it will allow us to draw a precise equivalence between actor and critic regularization:

$$\max_{\pi} \mathbb{E}_{p_g(s_g)p(s)\pi(a|s,s_g)} [\log Q(s, a, s_g)]$$

We will now consider variants of C-learning that incorporate actor and critic regularization.

One-step RL. We will consider a variant of C-learning that resembles one-step RL (Brandfonbrener et al., 2021). The critic update will be similar to before, but the next-actions sampled for the TD updates will be sampled from the *marginal* behavioral policy:

$$\begin{aligned} \max_Q (1 - \gamma) \mathbb{E}_{p(s,a,s')} \left[\log \frac{Q(s, a, s_g = s')}{Q(s, a, s_g = s') + 1} \right] &+ \gamma \mathbb{E}_{p(s,a)p_g(s_g)} \left[y^{\beta, Q_t}(s, a, s) \log \frac{Q(s, a, s_g)}{Q(s, a, s_g) + 1} \right] \\ &+ \mathbb{E}_{p(s,a)p_g(s_g)} \left[\log \frac{1}{Q(s, a, s_g = s') + 1} \right], \end{aligned}$$

where $y^{\beta, Q_t}(s, a, s_g) = \mathbb{E}_{p(s'|s,a)\beta(a'|s')} [Q_t(s', a', s_g)]$. The actor update will be modified to include a reverse KL divergence:

$$\max_{\pi} \mathbb{E}_{p(s)p_g(s_g)\pi(a|s,s_g)} [\log Q(s, a, s_g) + \log \beta(a | s) - \pi(a | s, s_g)]. \quad (15)$$

Note that we are regularizing the policy to be similar to the average behavioral policy, $\beta(a | s)$. Compared to regularization towards a goal-conditioned behavioral policy $\beta(a | s, s_g)$, this choice gives the policy additional flexibility: when trying to reach goal s_g , it is allowed to take actions that were not taken by $\beta(a | s, s_g)$, as long as they were taken by the behavioral policy when trying to reach some other goal s'_g .

Critic regularization. To regularize the critic, we will modify the “negative” term in the C-learning objective to use actions sampled from the policy:

$$\max_Q (1 - \gamma) \mathbb{E}_{p(s,a,s')} \left[\log \frac{Q(s, a, s_g = s')}{Q(s, a, s_g = s') + 1} \right] \quad (16)$$

$$+ \gamma \mathbb{E}_{p(s,a)p_g(s_g)} \left[y^{\pi, Q_t}(s, a, s_g) \log \frac{Q(s, a, s_g)}{Q(s, a, s_g) + 1} \right] \quad (17)$$

$$+ \mathbb{E}_{p(s)p_g(s_g), a \sim \pi(\cdot | s, s_g)} \left[\log \frac{1}{Q(s, a, s_g) + 1} \right]. \quad (18)$$

C.1. Analysis for Goal-Conditioned Problems

Like in the single-task setting, these two forms of regularization yield the same fixed points:

Theorem C.1. *Let policy $\pi(a | s, s_g)$ be given, let $Q^\beta(s, a, s_g)$ be the Q -values for the marginal behavioral policy $\beta(a | s)$ and let $Q_r^\pi(s, a, s_g)$ be the critic obtained by the regularized critic update (Eq. 18). Then performing regularized policy updates (Eq. 15) using the behavioral critic is equivalent to the un-regularized policy objective using the regularized critic:*

$$\mathbb{E}_{\pi(a|s,s_g)} [\log Q^\beta(s, a, s_g) + \log \beta(a | s) - \log \pi(a | s, s_g)] = \mathbb{E}_{\pi(a|s,s_g)} [\log Q_r^\pi(s, a, s_g)]$$

for all states s and goals s_g .

Proof. We start by determining the fixed point of critic-regularized C-learning. Like in the single-task setting, the C-learning objective resembles a weighted-classification problem, so we can write down the Bayes’ optimal classifier as

$$\frac{Q(s, a, s_g)}{Q(s, a, s_g) + 1} = \frac{((1 - \gamma)p(s' = s_g | s, a) + \gamma p(s = s_g)y(s', s_g))\beta(a | s)}{((1 - \gamma)p(s' = s_g | s, a) + \gamma p(s = s_g)y(s', s_g))\beta(a | s) + p(s_g)\pi(a | s, s_g)}.$$

Solving for $Q(s, a, s_g)$ on the left hand side, the optimal value for $Q(s, a, s_g)$ is given by

$$Q(s, a, s_g) = ((1 - \gamma)p(s' = s_g | s, a) + \gamma p(s = s_g)y(s', s_g)) \frac{\beta(a | s)}{\pi(a | s, s_g)}$$

This tells us what each critic-regularized C-learning update does.

To analyze these updates, we define $\tilde{Q}(s, a, s_g) \triangleq Q(s, a, s_g) \frac{\pi(a|s,s_g)}{\beta(a|s)}$. Then these updates can be written using $\tilde{Q}(s, a, s_g)$ as

$$\tilde{Q}(s, a, s_g) \frac{\beta(a | s)}{\pi(a | s, s_g)} = \left((1 - \gamma)p(s' = s_g | s, a) + \gamma \mathbb{E}_{p(s'|s,a)\pi(a'|s',s_g)} \left[\tilde{Q}(s', a', s_g) \frac{\beta(a' | s')}{\pi(a' | s', s_g)} \right] \right) \frac{\beta(a | s)}{\pi(a | s, s_g)}.$$

These updates can be simplified to

$$\tilde{Q}(s, a, s_g) = (1 - \gamma)p(s' = s_g | s, a) + \gamma \mathbb{E}_{p(s'|s,a)\beta(a'|s')} \left[\tilde{Q}(s', a', s_g) \right].$$

Like before, the ratio $\frac{\beta(a'|s')}{\pi(a'|s',s_g)}$ inside the expectation acts like an importance weight. Thus, the regularized critic updates are equivalent to perform policy evaluation on $\tilde{Q}(s, a, s_g)$. Note that this is estimating the probability that the *average* behavioral policy $\beta(a | s)$ reaches goal s_g ; this is *not* the probability that a goal-directed behavioral policy $\beta(a | s, s_g)$ reaches the goal.

Finally, we translate these convergence results for $\tilde{Q}(s, a, s_g)$ into convergence results for $Q(s, a, s_g)$. Written in terms of the original Q -values, we see that the optimal critic for the regularized critic update is

$$Q^*(s, a, s_g) = \tilde{Q}^*(s, a, s_g) \frac{\beta(a | s)}{\pi(a | s, s_g)} = Q^{\beta(\cdot | \cdot)}(s, a, s_g) \frac{\beta(a | s)}{\pi(a | s, s_g)}.$$

Thus, critic regularization implicitly regularizes the actor objective so that it is the same objective as one-step RL:

$$\begin{aligned} & \mathbb{E}_{p(s), s_g \sim p(s), \pi(a|s, s_g)} [\log Q^*(s, a, s_g)] \\ &= \mathbb{E}_{p(s), s_g \sim p(s), \pi(a|s, s_g)} \left[\log Q^{\beta(\cdot)}(s, a, s_g) + \log \beta(a | s) - \log \pi(a | s, s_g) \right]. \end{aligned}$$

□

D. Regularization for Example-based Control Problems

While specifying tasks in terms of reward functions is standard for MDPs, it can be difficult for real-world applications of RL. So, prior work has looked at specifying tasks by goal states (as in the previous section) or sets of states representing good outcomes (Pinto & Gupta, 2016; Tung et al., 2018; Fu et al., 2018). In addition to requiring more flexible and user-friendly forms of task specification, these algorithms targeted at real-world applications often demand regularization. In the same way that prior goal-conditioned RL algorithms have employed critic regularization, so too have prior example-based control algorithms (Singh et al., 2019; Hatch et al., 2022). In this section, we extend our analysis to regularization of an example-based control algorithm. Again, we will show that a certain form of critic regularization is equivalent to regularizing the actor.

We first define the problem of example-based control (Fu et al., 2018). In these problems, the agent is given a small collection of states $s \sim p_e(s)$, which are examples of successful outcomes. The aim is to learn a policy $\pi(a | s)$ that maximizes the probability of reaching a success state:

$$\max_{\pi} \mathbb{E}_{p(s_g)} \mathbb{E}_{\pi(\tau|s_g)} \left[\sum_{t=0}^{\infty} \gamma^t p_e(s_t) \right].$$

Note that this objective function is exactly equivalent to a reward-maximization problem, with a reward function $r(s, a) = p_e(s_t)$.

In this section, we show that one-step RL and critic regularization are equivalent for a certain example-based control algorithm. Unlike our analysis in the single-task setting, this analysis here uses an existing method, RCE (Eysenbach et al., 2021). RCE is a TD method that already makes use of the cross entropy loss for training the critic:

$$\max_Q (1 - \gamma) \mathbb{E}_{p_e(s)\beta(a|s)} \left[\log \frac{Q(s, a)}{Q(s, a) + 1} \right] + \mathbb{E}_{p(s, a)} \left[\gamma y^{\pi, Q_t}(s, a) \log \frac{Q(s, a)}{Q(s, a) + 1} + \log \frac{1}{Q(s, a) + 1} \right],$$

where $y^{\pi, Q_t}(s, a) = \mathbb{E}_{p(s'|s, a)\pi(a'|s')} [Q(s', a')]$ serves the role of the TD target. The first two terms increase the Q-values while the last term decreases the Q-values. The actor is updated to maximize the Q-values. While this objective for the actor can be written in many ways, we will write it as maximizing a log ratio because it will allow us to draw a precise equivalence between actor and critic regularization:

$$\max_{\pi} \mathbb{E}_{p(s)\pi(a|s)} [\log Q(s, a)]$$

We will now consider variants of RCE that incorporate actor and critic regularization.

One-step RL. We will consider a variant of RCE that resembles one-step RL (Brandfonbrener et al., 2021). The critic update will be similar to before, but the next-actions sampled for the TD updates will be sampled from the behavioral policy:

$$\max_Q (1 - \gamma) \mathbb{E}_{p_e(s)\beta(a|s)} \left[\log \frac{Q(s, a)}{Q(s, a) + 1} \right] + \mathbb{E}_{p(s, a)} \left[\gamma y^{\beta, Q_t}(s, a) \log \frac{Q(s, a)}{Q(s, a) + 1} + \log \frac{1}{Q(s, a) + 1} \right],$$

where $y^{\beta, Q_t}(s, a) = \mathbb{E}_{p(s'|s, a)\beta(a'|s')} [Q(s', a')]$. The actor update will be modified to include a reverse KL divergence:

$$\max_{\pi} \mathbb{E}_{p(s), \pi(a|s)} [\log Q(s, a) + \log \beta(a | s) - \pi(a | s)]. \quad (19)$$

Critic regularization. To regularize the critic, we will modify the “negative” term in the RCE objective to use actions sampled from the policy:

$$(1 - \gamma) \mathbb{E}_{p_e(s)\beta(a|s)} \left[\log \frac{Q(s, a)}{Q(s, a) + 1} \right] + \mathbb{E}_{p(s, a), a^- \sim \pi(\cdot|s)} \left[\gamma y^{\pi, Q_t}(s, a) \log \frac{Q(s, a)}{Q(s, a) + 1} + \log \frac{1}{Q(s, a^-) + 1} \right], \quad (20)$$

D.1. Analysis for Example-based Control Problems

Like in the single-task setting, these two forms of regularization yield the same fixed points:

Theorem D.1. *Let policy $\pi(a | s)$ be given, let $Q^\beta(s, a)$ be the Q -values for the behavioral policy $\beta(a | s)$ and let $Q_r^\pi(s, a)$ be the critic obtained by the regularized critic update (Eq. 20). Then performing regularized policy updates (Eq. 19) using the behavioral critic is equivalent to the un-regularized policy objective using the regularized critic:*

$$\mathbb{E}_{\pi(a|s)} [\log Q^\beta(s, a) + \log \beta(a | s) - \log \pi(a | s)] = \mathbb{E}_{\pi(a|s)} [\log Q_r^\pi(s, a)]$$

for all states s .

Proof. We start by determining the fixed point of critic-regularized RCE. Like in the single-task setting, The RCE objective resembles a weighted-classification problem, so we can write down the Bayes' optimal classifier as

$$\frac{Q(s, a)}{Q(s, a) + 1} = \frac{((1 - \gamma)p_e(s) + \gamma y^{\pi, Q_t}(s, a))\beta(a | s)}{((1 - \gamma)p_e(s) + \gamma y^{\pi, Q_t}(s, a))\beta(a | s) + \pi(a | s)}.$$

Solving for $Q(s, a)$ on the left hand side, the optimal value for $Q(s, a)$ is given by

$$Q(s, a) = ((1 - \gamma)p_e(s) + \gamma y^{\pi, Q_t}(s, a)) \frac{\beta(a | s)}{\pi(a | s)}$$

This tells us what each critic-regularized RCE update does.

To analyze these updates, we define $\tilde{Q}(s, a) \triangleq Q(s, a) \frac{\pi(a|s)}{\beta(a|s)}$. Then these updates can be written using $\tilde{Q}(s, a)$ as

$$\tilde{Q}(s, a) \frac{\beta(a | s)}{\pi(a | s)} = \left((1 - \gamma)p_e(s) + \gamma \mathbb{E}_{p(s'|s, a)\pi(a'|s')} \left[\tilde{Q}(s', a') \frac{\beta(a' | s')}{\pi(a' | s')} \right] \right) \frac{\beta(a | s)}{\pi(a | s)}.$$

These updates can be simplified to

$$\tilde{Q}(s, a) = (1 - \gamma)p_e(s) + \gamma \mathbb{E}_{p(s'|s, a)\beta(a'|s')} [\tilde{Q}(s', a')].$$

Like before, the ratio $\frac{\beta(a'|s')}{\pi(a'|s')}$ inside the expectation acts like an importance weight. Thus, the regularized critic updates are equivalent to perform policy evaluation on $\tilde{Q}(s, a)$.

Finally, we translate these convergence results for $\tilde{Q}(s, a)$ into convergence results for $Q(s, a)$. Written in terms of the original Q -values, we see that the optimal critic for the regularized critic update is

$$Q^*(s, a) = \tilde{Q}^*(s, a) \frac{\beta(a | s)}{\pi(a | s)} = Q^\beta(s, a) \frac{\beta(a | s)}{\pi(a | s)}.$$

Thus, critic regularization implicitly regularizes the actor objective so that it is the same objective as one-step RL:

$$\mathbb{E}_{p(s), \pi(a|s)} [\log Q^*(s, a)] = \mathbb{E}_{p(s), \pi(a|s)} [\log Q^\beta(s, a) + \log \beta(a | s) - \log \pi(a | s)].$$

□

E. Additional Experiments

E.1. Is classifier actor critic a good model of practically-used RL methods?

To study whether classifier actor critic is an accurate model for practically-used RL methods, we extended Fig. 3 by adding three additional methods: classifier actor critic, classifier actor critic with actor regularization (Equations 4 and 3). Comparing Q -learning with classifier actor critic, we see that both yield reward-maximizing policies (in line with Lemma 4.1), though these policies are different (i.e., they perform symmetry breaking in different ways). Comparing one-step RL with actor-regularized classifier actor critic, we observe that both methods take the same three actions at the states within the blue

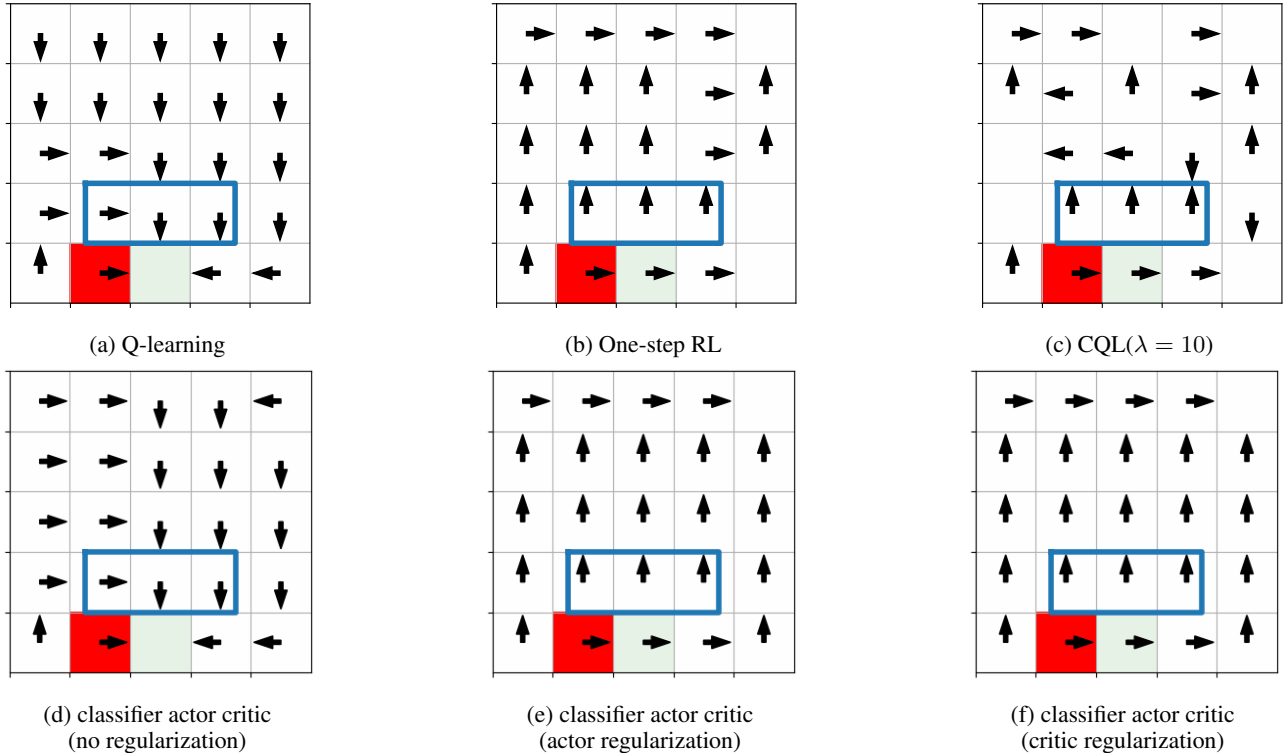


Figure 9: Extension of Fig. 3 with additional figures for classifier actor critic and regularized variants. As expected, classifier actor critic produces a reward-maximizing policy (Lemma 4.1), as does Q-learning. In line with Theorem 4.3, critic-regularized and actor-regularized classifier actor critic produce the same policies. Comparing (b) one-step RL (e) actor-regularized classifier actor critic, we see that the methods produce the same three actions within the blue box (states where we expect regularization to be important), but can produce different actions in states outside the blue box. The comparisons between (c) CQL and (d) critic-regularized classifier actor critic are similar, supporting the claim that classifier actor critic is only an approximate (no a perfect) model of practically-used offline RL methods.

box, states where we expect regularization to have a large effect. Outside this blue box, these two methods occasionally take different actions. Similarly, comparing CQL to critic-regularized classifier actor-critic, we observe that the methods take the same actions within the blue box, but occasionally take different actions outside the blue box. In line with Theorem 4.3, classifier actor-critic with critic regularization produces the exact same policy as classifier actor-critic with actor regularization. Taken together, these results provide empirical backing for our theoretical results, while also showing that classifier actor critic is only an approximate model of practical algorithms, not a perfect model.

F. Experimental Details

F.1. Tabular experiments

Implementing critic regularization for classifier actor critic. The objective for critic regularization in contrastive actor critic (Eq. 6) is nontrivial to optimize because of the cyclic dependency between the policy and the critic: simply alternating between optimizing the actor and the critic does not converge. In our experiments, we update the critic using an exponential moving average of the policy, as proposed in Wen et al. (2021). We found that this decision was sufficient for ensuring convergence. When applying CQL in the tabular setting (Figures 3 and 4), we did not do this because soft value iteration represents the policy implicitly in terms of the value function.

Fig. 2 (left) The initial state and goal state are located in opposite corners. The reward function is +1 for reaching the goal and 0 otherwise. We use a dataset of 20 trajectories, 50 steps each, collected by a random policy. We use $\gamma = 0.95$ and train for 20k full-batch updates, using a learning rate of $1e-2$. The Q table is randomly initialized using a standard normal

distribution.

Fig. 2 (center) The initial state and goal state are located in adjacent corners. The goal state has a reward of +3.5, the states between the initial state and goal state have a reward +1, and all other states (including the initial state) have a reward of +2. We use a dataset of 20 trajectories, 50 steps each, collected by a random policy. We use $\gamma = 0.95$ and train for 20k full-batch updates, using a learning rate of $1e-2$. The Q table is randomly initialized using a standard normal distribution.

Fig. 2 (right) The initial state and goal state are located in adjacent corners. The reward is +0.01 at the goal state and 0 otherwise. We use a dataset of 1 trajectories with 10 steps, which traces the following path:

$$[(0, 0), (1, 0), (1, 1), (1, 2), (1, 3), (1, 4), (0, 4), (0, 4), (0, 4), (0, 4)].$$

We use $\gamma = 0.95$ and train for 10k full-batch updates, using a learning rate of $1e-2$. The Q table is randomly initialized using a standard normal distribution.

Fig. 3 There is a bad state (reward of -10) next to the optimal state (reward of +1), so the behavioral policy navigates away from the optimal state. We generate 10 trajectories of length 100 from a uniform random policy. We use $\gamma = 0.95$ and train each method for 10k full-batch updates. The Q table is randomly initialized using a standard normal distribution. One-step RL performs SARSA updates while CQL performs soft value iteration (as suggested in the CQL paper).

Fig. 4 We generate 100 random variants of Fig. 3 by randomly sampling the high-reward state and low-reward state (without replacement). The datasets are generated in the same way.

Fig. 5 We use the same environment and dataset as in Fig. 3, but train the CQL agent with varying values of λ , each with 5 random seeds. We train the one-step RL agent for 5 random seeds. For each point on the X axis of Fig. 5, we compare compute 5×5 pairwise comparisons and report the mean and standard deviation.

E.2. Continuous control experiments

For the experiments in Figures 6 and 7, we used the implementation of one-step RL (reverse KL) and CQL provided by Hoffman et al. (2020). We choose this implementation because it is well tuned and uses similar hyperparameters for the two methods. As mentioned in the main text, the only change we made to the implementation was adding the twin-Q trick to one-step RL, such that it matched the critic architecture used by CQL. We did not change any of the other hyperparameters, including hyperparameters controlling the regularization strength.