

BRPO: Batch Residual Policy Optimization

Sungryull Sohn^{*,1,2}, Yinlam Chow^{*,1}, Jayden Ooi^{*,1},
Ofir Nachum¹, Honglak Lee^{1,2}, Ed Chi¹ and Craig Boutilier¹

¹Google Research

²University of Michigan

srsohn@umich.edu, {yinlamchow, jayden, ofirnachum, honglak, edchi, cboutilier}@google.com

Abstract

In batch reinforcement learning (RL), one often constrains a learned policy to be close to the behavior (data-generating) policy, e.g., by constraining the learned action distribution to differ from the behavior policy by some maximum degree that is the *same at each state*. This can cause batch RL to be overly conservative, unable to exploit large policy changes at frequently-visited, high-confidence states without risking poor performance at sparsely-visited states. To remedy this, we propose *residual policies*, where the allowable deviation of the learned policy is *state-action-dependent*. We derive a new for RL method, *BRPO*, which learns both the policy and allowable deviation that jointly maximize a lower bound on policy performance. We show that BRPO achieves the state-of-the-art performance in a number of tasks.

1 Introduction

Deep reinforcement learning (RL) methods are increasingly successful in domains such as games [Mnih *et al.*, 2013], and robotic manipulation [Nachum *et al.*, 2019]. Much of this success relies on the ability to collect new data through *online* interactions with the environment during training, often relying on simulation. Unfortunately, this approach is impractical in many real-world applications where faithful simulators are rare, and in which active data collection through interactions with the environment is costly, time consuming, and risky.

Batch (or offline) RL [Lange *et al.*, 2012] is an emerging research direction that aims to circumvent the need for online data collection, instead learning a new policy using only offline trajectories generated by some behavior policy (e.g., the currently deployed policy in some application domain). In principle, any off-policy RL algorithm (e.g., DDPG [Lillicrap *et al.*, 2015], DDQN [Hasselt *et al.*, 2016]) may be used in this batch (or more accurately, “offline”) fashion; but in practice, such methods have been shown to fail to learn when presented with arbitrary, static, off-policy data. This can arise for several reasons: lack of exploration [Lange *et al.*, 2012], generalization error on out-of-distribution samples in value

estimation [Kumar *et al.*, 2019], or high-variance policy gradients induced by covariate shift [Mahmood *et al.*, 2014].

Various techniques have been proposed to address these issues, many of which can be interpreted as constraining or regularizing the learned policy to be *close* to the behavior policy [Fujimoto *et al.*, 2018; Kumar *et al.*, 2019] (see further discussion below). While these batch RL methods show promise, none provide improvement guarantees relative to the behavior policy. In domains for which batch RL is well-suited (e.g., due to the risks of active data collection), such guarantees can be critical to deployment of the resulting RL policies.

In this work, we use the well-established methodology of *conservative policy improvement (CPI)* [Kakade and Langford, 2002] to develop a theoretically principled use of behavior-regularized RL in the batch setting. Specifically, we parameterize the learned policy as a *residual policy*, in which a base (behavior) policy is combined linearly with a learned *candidate policy* using a mixing factor called the *confidence*. Such residual policies are motivated by several practical considerations. First, one often has access to offline data or logs generated by a deployed base policy which is known to perform reasonably well. The offline data can be used by an RL method to learn a candidate policy with better *predicted* performance, but if confidence in parts of that prediction is weak, relying on the base policy may be desirable. The base policy may also incorporate soft business constraints or some form of interpretability. Our residual policies blend the two in a learned, non-uniform fashion. When deploying a new policy, we use the CPI framework to derive updates that learn both the candidate policy and the confidence that jointly maximize a lower bound on performance improvement relative to the behavior policy. Crucially, while traditional applications of CPI, such as TRPO [Schulman *et al.*, 2015], use a constant or state-independent confidence, our performance bounds and learning rules are based on *state-action-dependent* confidences—this gives rise to bounds that are less conservative than their CPI counterparts.

In Sec. 2, we formalize residual policies and in Sec. 3 analyze a novel *difference-value function*. Sec. 4 holds our main result, a tighter lower bound on policy improvement for our residual approach (vs. CPI and TRPO). We derive the BRPO algorithm in Sec. 5 to jointly learn the candidate policy and confidence; experiments in Sec. 6 show its effectiveness.

*Equal contribution

2 Preliminaries

We consider a *Markov decision process (MDP)* $\mathcal{M} = \langle S; A; R; T; P_0 \rangle$, with state space S , action space A , reward function R , transition kernel T , and initial state distribution P_0 . A policy π interacts with the environment, starting at $s_0 \sim P_0$. At step t , the policy samples an action a_t from a distribution $\pi(\cdot|s_t)$ over A and applies. The environment emits a reward $r_t = R(s_t; a_t) \in [0; R_{\max}]$ and next state $s_{t+1} \sim T(\cdot|s_t; a_t)$. In this work, we consider discounted infinite-horizon problems with discount factor $\gamma \in [0; 1)$.

Let $\Pi = \{ \pi : S \times A \rightarrow [0; 1] \mid \sum_a \pi(a|s) = 1 \}$ be the set of Markovian stationary policies. The expected (discounted) cumulative return of policy π , is $J(\pi) := \mathbb{E}_{T; \pi} [\sum_{t=0}^{\infty} \gamma^t R(s_t; a_t) \mid s_0 \sim P_0]$. Our aim is to find an optimal policy $\pi^* \in \arg \max_{\pi} J(\pi)$. In reinforcement learning (RL), we must do so without knowledge of $R; T$, using only trajectory data generated from the environment (see below) or access to a simulator (see above).

We consider pure offline or *batch RL*, where the learner has access to a fixed data set (or *batch*) of state-actions-reward-next-state samples $B = \{(s; a; r; s')\}$, generated by a (known) *behavior policy* $\pi_b(\cdot|s)$. No additional data collection is permitted. We denote by d the γ -discounted occupation measure of the MDP w.r.t. π .

In this work, we study the problem of *residual policy optimization (RPO)* in the batch setting. Given the behavior policy $\pi_b(a|s)$, we would like to learn a *candidate policy* $\pi_c(a|s)$ and a *state-action confidence* $\beta(s; a)$, such that the final *residual policy* $\pi(a|s) = (1 - \beta(s; a)) \cdot \pi_c(a|s) + \beta(s; a) \cdot \pi_b(a|s)$ maximizes total return. As discussed above, this type of mixture allows one to exploit an existing, ‘‘well-performing’’ behavior policy. Intuitively, $\beta(s; a)$ should capture how much we can trust π_c at each $s; a$ pair, given the available data. To ensure that the residual policy is a probability distribution at every state $s \in S$, we constrain the confidence to lie in the set $\beta = \{ \beta : S \times A \rightarrow [0; 1] \mid \sum_a \beta(s; a) (\pi_c(a|s) - \pi_b(a|s)) = 0 \}$.

Related work. Similar to the above policy formulation, CPI [Kakade and Langford, 2002] also develops a policy mixing methodology that guarantees performance improvement when the confidence β is a constant. However, CPI is an online algorithm, and it learns the candidate policy independently of (not jointly with) the mixing factor; thus, extension of CPI to offline, batch setting is unclear. Other existing work also deals with online residual policy learning without jointly learning mixing factors [Johannink *et al.*, 2019; Silver *et al.*, 2018]. Common applications of CPI may treat β as a hyper-parameter, which specifies the maximum total-variation distance between the learned and behavior policy distributions (see standard proxies in [Schulman *et al.*, 2015; Pirota *et al.*, 2013] for details). *Batch-constrained Q-learning (BCQ)* [Fujimoto *et al.*, 2018; Fujimoto *et al.*, 2019] incorporates the behavior policy when defining the admissible action set in Q-learning for selecting the highest-valued actions that are similar to data samples in the batch. *BEAR* [Kumar *et al.*, 2019] is motivated as a means to control the accumulation of out-of-distribution value errors; but its main algorithmic contribution is realized by adding a regularizer to the loss that measures the kernel maximum mean dis-

crepancy (MMD) [Gretton *et al.*, 2007] between the learned and behavior policies similar to KL-control [Jaques *et al.*, 2019]. Algorithms such as SPI [Ghavamzadeh *et al.*, 2016] and SPIBB [Laroche and Trichelair, 2017] bootstraps the learned policy with the behavior policy when the uncertainty in the update for current state-action pair is high, where the uncertainty is measured by the visitation frequency of state-action pairs in the batch data. While these methods work well in some applications it is unclear if they have any performance guarantees.

3 The Difference-value Function

We begin by defining and characterizing the *difference-value function*, a concept we exploit in the derivation of our batch RPO method in Secs. 4 and 5. For any $s \in S$, let $V^{\pi_b}(s)$ and $V^{\pi_c}(s)$ be the value functions induced by policies π_b and π_c , respectively. Using the structure of the residual policy, we establish two characterizations of the *difference-value function* $V^{\pi}(s) := V^{\pi_c}(s) - V^{\pi_b}(s)$.

Lemma 1. *Let $A(s; a) := Q(s; a) - V^{\pi_b}(s)$ be the advantage function w.r.t. residual policy π_b , where Q is the state-action value. The difference-value is $V^{\pi}(s) = \mathbb{E}_{T; \pi} [\sum_{t=0}^{\infty} \gamma^t \hat{A}(s_t; a_t) \mid s_0 = s]$; where $\hat{A}(s; a) = \sum_{a_2 A} \pi_b(a_2|s) \cdot \pi_c(a|s) \cdot \frac{\pi_b(a_2|s)}{\pi_b(a|s)} \cdot A(s; a)$ is the residual reward that depends on π_c and difference of candidate policy and behavior policy π_b .*

This result establishes that the difference value is essentially a value function w.r.t. the residual reward. Moreover, it is proportional to the advantage of the target policy, the confidence, and the difference of policies. While the difference value can be estimated from behavior data batch B , this formulation requires knowledge of the advantage function A w.r.t. the target policy, which must be re-learned at every β -update in an off-policy fashion. Fortunately, we can show that the difference value can also be expressed as a function of the advantage w.r.t. the *behavior policy* π_b :

Theorem 2. *Let $A(s; a) := Q(s; a) - V^{\pi_b}(s)$ be the advantage function induced by π_b , in which Q is the state-action value. The difference-value is given by $V^{\pi}(s) = \mathbb{E}_{T; \pi} [\sum_{t=0}^{\infty} \gamma^t A(s_t; a_t) \mid s_0 = s]$, where*

$$A^{\pi}(s) = \sum_{a \in A} \pi_b(a|s) \pi_c(a|s) \frac{\pi_b(a|s)}{\pi_b(a|s)} A(s; a)$$

is the residual reward that depends on π_c and difference of candidate policy π_c and behavior policy π_b .

In our RPO approach, we exploit the nature of the difference-value function to solve the maximization w.r.t. the confidence and candidate policy: $(\beta(s; \cdot); \pi_c(\cdot|s)) \in \arg \max_{\beta, \pi_c} \sum_{s \in S} d(s) V^{\pi}(s), \forall s \in S$. Since $\beta(s; \cdot) = 0$ implies $V^{\pi}(s) = 0$, the optimal difference-value function $V^{\pi}(s) := \max_{\beta, \pi_c} \sum_{s \in S} d(s) V^{\pi}(s)$ is always lower-bounded by 0. We motivate computing $(\beta; \pi_c)$ with the above difference-value formulation rather than as a standard RL problem as follows. In the tabular case, optimizing $(\beta; \pi_c)$ with either formulation gives an identical result. However,

both the difference-value function in Theorem 2 and the standard RL objective require sampling data generated by the updated policy π . In the batch setting, when fresh samples are unavailable, learning $(\pi; \pi)$ with off-policy data may incur instability due to high generalization error [Kumar *et al.*, 2019]. While this can be alleviated by adopting the CPI methodology, applying CPI directly to RL can be overly conservative [Schulman *et al.*, 2015]. By contrast, we leverage the special structure of the difference-value function (e.g., non-negativity) below, using this new formulation together with CPI to derive a less conservative RPO algorithm.

4 Batch Residual Policy Optimization

We now develop an RPO algorithm that has stable learning performance in the batch setting and *performance improvement guarantees*. For the sake of brevity, in the following we only present the main results on performance guarantees of RPO. Proofs of these results can be found in the appendix of the extended paper. We begin with the following baseline result, directly applying Corollary 1 of the TRPO result to RPO to ensure the residual policy π performs no worse than π .

Lemma 3. *For any value function $U : S \rightarrow \mathbb{R}$, the difference-return satisfies $J - J \geq \frac{1}{1-\gamma} \mathbb{E}_{U; \pi; \pi} - \frac{1}{(1-\gamma)^2} \cdot U; \pi; \pi \cdot \mathbb{E}_S \frac{1}{d} \frac{1}{2} D_{\text{KL}}(\pi(s) \| \pi(s))$; where the surrogate objective and the penalty weight are*

$$\mathbb{E}_{U; \pi; \pi} := \mathbb{E}_{(s; a; s') \sim d} (s; a) \frac{(a|s)}{(a|s)} U(s; a; s'); \\ U; \pi; \pi := \max_s \mathbb{E}_{\pi} [U(s; a; s')];$$

where $U(s; a; s') := R(s; a) + U(s') - U(s)$.

When $U = V$, one has $\mathbb{E}_{\pi} [U(s; a; s')] = 0, \forall s \in S$, which implies that the inequality is tight—this lemma then coincides Lemma 1. While this CPI result forms the basis of many RL algorithms (e.g., TRPO, PPO), in many cases it is very loose since $U; \pi; \pi$ is a maximum over all states. Thus, using this bound for policy optimization may be *overly conservative*, i.e., algorithms which rely on this bound must take very small policy improvement steps, especially when the penalty weight $U; \pi; \pi$ is large, i.e., $|U; \pi; \pi| = (1-\gamma) \gg |\mathbb{E}_{U; \pi; \pi}|$. While this approach may be reasonable in online settings—when collection of new data (with an updated behavior policy $\pi \leftarrow \pi$) is allowed—in the batch setting it is challenging to overcome such conservatism.

To address this issue, we develop a CPI method that is specifically tied to the difference-value formulation, and uses a state-action-dependent confidence $d(s; a)$. We first derive the following theorem, which bounds the difference returns that are generated by π and π .

Theorem 4. *The difference return of $(\pi; \pi)$ satisfies*

$$J - J \geq \frac{1}{1-\gamma} L^0; \pi; \pi - \frac{1}{1-\gamma} \cdot L^{00}; \pi; \pi \cdot \max_{s_0 \in S} L^{000}; \pi; \pi (s_0);$$

where the surrogate objective function, regularization, and

penalty weight are given by

$$L^0; \pi; \pi := \mathbb{E}_{(s; a) \sim d} (s; a) \frac{(a|s)}{(a|s)} A(s; a) \\ L^{00}; \pi; \pi := \mathbb{E}_{(s; a) \sim d} (s; a) \frac{j(a|s)}{(a|s)} \frac{(a|s)j}{(a|s)} \\ L^{000}; \pi; \pi (s_0) := \mathbb{E}_{(s; a) \sim d(s_0)} (s; a) \frac{j(a|s)}{(a|s)} \frac{(a|s)j}{(a|s)} jA(s; a)j$$

respectively, in which $d(s_0)$ is the discounted occupancy measure w.r.t. π given initial state s_0 .

Unlike the difference-value formulations in Lemma 1 and Theorem 2, which require the knowledge of advantage function A or the trajectory samples generated by π , the lower bound in Theorem 4 is comprised only of terms that can be estimated directly using the data batch B (i.e., data generated by π). This makes it a natural objective function for batch RL. Notice also that the surrogate objective, the regularization, and the penalty weight in the lower bound are each proportional to the confidence and to the relative difference of the candidate and behavior policies. However, the \max operator requires state enumeration to compute this lower bound, which is intractable when S is large or uncountable.

We address this by introducing a slack variable $\lambda \geq 0$ to replace the \max -operator with suitable constraints. This allows the bound on the difference return to be rewritten as: $J - J \geq \frac{1}{1-\gamma} L^0; \pi; \pi - \min_{\lambda \geq 0} L^{000}; \pi; \pi (s_0); \lambda \frac{1}{(1-\gamma)^2} L^{00}; \pi; \pi$.

Consider the Lagrangian of the lower bound:

$$\frac{L^0; \pi; \pi}{1-\gamma} - \min_{\lambda \geq 0} \max_{(s) \geq 0, \forall s} \frac{L^{00}; \pi; \pi}{(1-\gamma)^2} \times \lambda (s) (L^{000}; \pi; \pi (s));$$

To simplify this saddle-point problem, we restrict the Lagrange multiplier to be $\lambda(s) = \lambda \cdot P_0(s) \geq 0$, where $\lambda \geq 0$ is a scalar multiplier. Using this approximation and the strong duality of linear programming [Boyd and Vandenberghe, 2004] over primal-dual variables $(\lambda; \lambda)$, the saddle-point problem on $(\lambda; \lambda)$ can be re-written as

$$L; \pi; \pi := \max_{\lambda \geq 0} \min_{\lambda \geq 0} \frac{L^0; \pi; \pi - L^{00}; \pi; \pi \lambda + L^{000}; \pi; \pi}{1-\gamma} \\ = \frac{1}{1-\gamma} L^0; \pi; \pi - \frac{1}{1-\gamma} L^{00}; \pi; \pi \lambda + L^{000}; \pi; \pi \lambda \quad (1)$$

where $L^{000}; \pi; \pi = \mathbb{E}_S [L^{000}; \pi; \pi (s)]$. The equality is based on the KKT condition on $(\lambda; \lambda)$. Notice that the only difference between the CPI lower bound in Theorem 4 and the objective function $\mathcal{L}; \pi; \pi$ is that the \max operator is replaced by expectation w.r.t the initial distribution.

With certain assumptions on the approximation error of the Lagrange multiplier parametrization $\lambda(s) \approx P_0(s)$, we can characterize the gap between the original CPI objective function in Theorem 4 and $\mathcal{L}; \pi; \pi$. One approach is to look into the KKT condition of the original saddle-point problem and bound the sub-optimality gap introduced by this Lagrange parameterization. Similar derivations can be found in the analysis of approximate linear programming (ALP) algorithms [Abbasi-Yadkori *et al.*, 2019; Farias and Roy, 2003].

Compared with the vanilla CPI result from Lemma 3, there are two characteristics in problem (1) that make the optimization w.r.t. $\mathcal{L}; \pi; \pi$ less conservative. First, the penalty weight

L_{π}^{000} here is smaller than L_{π} in Lemma 3, which means that the corresponding objective has less incentive to force π to be close to π^* . Second, compared with entropy regularization in vanilla CPI, here the regularization and penalty weights are both linear in $2 \in [0, 1]^A$; thus, unlike vanilla CPI, whose objective is linear in π , our objective is quadratic in π —this modification ensures the optimal value is not a generate extreme point.

5 The BRPO Algorithm

We now develop the BRPO algorithm, for which the general pseudo-code is given in Algorithm 1. Recall that if the candidate policy π and con dence π^* are jointly optimized

$$(\pi; \pi^*) \underset{\pi; \pi^*}{2} \arg \max L_{\pi; \pi^*}; \quad (2)$$

then the residual policy $(a|s) = (1 - \beta(s; a)) \pi(a|s) + \beta(s; a) \pi^*(a|s)$ performs no worse than behavior policy π . Generally, solutions for problem (2) use a form of minorization-maximization (MM) Hunter and Lange, 2004, a class of methods that also includes expectation maximization. In the terminology of MM algorithms, $L_{\pi; \pi^*}$ is a surrogate function satisfying the following MM properties

$$J_{\pi} \leq J_{\pi^*} \leq L_{\pi; \pi^*}; \quad J_{\pi} = L_{\pi; \pi^*} = L_{\pi^*}; \quad \pi = \pi^* = 0; \quad (3)$$

which guarantees that it minorizes the difference-return J_{π} with equality at $\pi = \pi^*$ (with arbitrary π^*) or at $\pi^* = \pi$ (with arbitrary π). This algorithm is also reminiscent of proximal gradient methods. We optimize π and π^* in RPO with a simple two-step coordinate-ascent. Specifically, at iteration $k \geq 0; 1; \dots; K$, given con dence π_{k-1} , we first compute an updated candidate policy π_k , and with π_k fixed, we update π_{k+1} , i.e., $L_{\pi_k; \pi_{k+1}} \leq L_{\pi_{k-1}; \pi_k} \leq L_{\pi_k; \pi_{k+1}}$. When π_k and π_{k+1} are represented tabularly or with linear function approximators, under certain regularity assumptions (the Kurdyka-Lojasiewicz property Xu and Yin, 2013) coordinate ascent guarantees global convergence (to the limit point) for BRPO.

However, when more complex representations (e.g., neural networks) are used to parameterize these decision variables, this property no longer holds. While one may still compute $(\pi; \pi^*)$ with first-order methods (e.g., SGD), convergence to local optima is not guaranteed. To address this, we next further restrict the MM procedure to develop closed-form solutions for both the candidate policy and the con dence.

The Closed-form Candidate Policy. To effectively update the candidate policy when given the con dence π^* , we develop a closed-form solution for π . Our approach is based on maximizing the following objective, itself a more conservative version of the CPI lower bound:

$$\max_{\pi} \hat{L}_{\pi; \pi^*} := E_{s \sim d} E_a \left((s; a) \frac{(a|s)}{(a|s)} A \right) \frac{\max_{\pi} \int_{j \in A} (s; j) g(s) D_{KL}(\pi(j|s) \| \pi^*(j|s))}{2(1 - \beta(s))} \frac{1}{1 - \beta(s)}; \quad (4)$$

¹For example, when β is state-dependent (which automatically satisfies the equality constraints in), the linear objective in vanilla CPI makes the optimal value $(\int_{j \in A} a 0-1$ vector. Especially when $\frac{2}{1 - \beta(s)} E_{s \sim d} \frac{1}{2} D_{KL}(\pi(s) \| \pi^*(s))$ is large, then most entries of π become zero, i.e., will be very close to π^* .

where $g(s) = (1 + \log E_{\pi^*} [\exp(g(a|s)^2)]) > 0$ for any arbitrary non-negative function g . To show that $\hat{L}_{\pi; \pi^*}$ in (4) is an eligible lower bound (so that the corresponding solution is an MM), we need to show that it satisfies the properties in (3). When $\pi = \pi^*$, by the definition of $\hat{L}_{\pi; \pi^*}$, the second property holds. To show the first property, we first consider the following problem:

$$\max_{\pi} \frac{1}{2} L_{\pi; \pi^*}^0 - \frac{1}{1 - \beta(s)} E_{s \sim d} E_a \left(\frac{(a|s)}{(a|s)} A \right); \quad (5)$$

where $L_{\pi; \pi^*}^0(\pi)$ is given in Theorem 4, and

$$E_{s \sim d} E_a \left(\frac{(a|s)}{(a|s)} A \right) = \int_{j \in A} \int_{s \in \mathcal{S}} \frac{h_p(s, j) D_{KL}(\pi(j|s) \| \pi^*(j|s))}{h_q(s, j) D_{KL}(\pi(j|s) \| \pi^*(j|s))} ds; \quad (6)$$

The concavity of $\frac{h_p}{h_q}$ (i.e., $E_{s \sim d} \left[\frac{h_p}{h_q} \right] \geq \frac{h_p}{h_q}$) and monotonicity of expectation imply that the objective in (4) is a lower bound of that in (7) below. Furthermore, by the weighted Pinsker's inequality Bolley and Villani, 2005, we have: $0 \leq L_{\pi; \pi^*}^0 - L_{\pi; \pi^*} \leq 2 E_{s \sim d} \int_{j \in A} \int_{s \in \mathcal{S}} \frac{h_p(s, j) D_{KL}(\pi(j|s) \| \pi^*(j|s))}{h_q(s, j) D_{KL}(\pi(j|s) \| \pi^*(j|s))} ds$, which implies the objective in (5) is a lower-bound of that in (2) and validates the first MM property.

Now recall the optimization problem $\max_{\pi} \frac{1}{2} \hat{L}_{\pi; \pi^*}$. Since this optimization is over the state-action mapping the Interchangeability Lemma Shapiro et al., 2009 allows swapping the order of $E_{s \sim d}$ and \max_{π} . This implies that at each $s \in \mathcal{S}$ the candidate policy can be solved using:

$$\begin{aligned} & 2 \arg \max_{\pi} E_a \left(\frac{(a|s)}{(a|s)} A \right) (s) \log \frac{(a|s)}{(a|s)} \\ & = \arg \max_{\pi} E_a \left((s; a) A \right) (s) \log \frac{(a|s)}{(a|s)}; \quad (7) \end{aligned}$$

where $(s) = \max_{\pi} \int_{j \in A} (s; j) g(s) D_{KL}(\pi(j|s) \| \pi^*(j|s))$ is the state-dependent penalty weight of the relative entropy regularization. By the KKT condition of (7), the optimal candidate policy π^* has the form

$$(a|s) = \frac{(a|s) \exp \left(\frac{(s; a) A}{(s)} \right)}{E_{a^0} \left[\exp \left(\frac{(s; a^0) A}{(s)} \right) \right]}; \quad (8)$$

Notice that the optimal candidate policy is a relative softmax policy, which is a common solution policy for many entropy-regularized RL algorithms Haarnoja et al., 2018. Intuitively, when the mixing factor vanishes (i.e., $(s; a) = 0$), the candidate policy equals to the behavior policy, and with con dence we obtain the candidate policy by modifying the behavior policy via exponential twisting.

The Closed-form Con dence. Given candidate policy π , we derive efficient scheme for computing the con dence π^* that solves the MM problem $\max_{\pi^*} \frac{1}{2} L_{\pi; \pi^*}$. Recall that this optimization can be reformulated as a concave quadratic program (QP) with linear equality constraints, which has a unique optimal solution Faybusovich and Moore, 1997. However, since the decision variable (i.e., the con dence mapping) is in finite-dimensional, solving this QP is intractable without some assumptions about this mapping. To

resolve this issue, instead of using the surrogate objective L_{π}^0 in MM, we turn to its sample-based estimate. Specifically, given a batch of data $\mathcal{B} = \{f(s_i; a_i; r_i; s_i^0)_{i=1}^{jB}\}$ generated by the behavior policy, denote by

$$L_{\pi}^0 := \frac{1}{jB} \sum_{i=1}^{jB} \left(\sum_{a \in \mathcal{A}} \pi(a|s_i) A(s_i, a) \right)$$

$$L_{\pi}^{00} := \frac{1}{jB} \sum_{i=1}^{jB} \sum_{a \in \mathcal{A}} \pi(a|s_i) A(s_i, a)$$

$$L_{\pi}^{000} := \frac{1}{jB} \sum_{i=1}^{jB} \sum_{a \in \mathcal{A}} \pi(a|s_i) A(s_i, a)$$

the sample-average approximation (SAA) of functions L_{π}^0 , L_{π}^{00} , and L_{π}^{000} respectively, where

$$\bar{L}_{\pi}^0 = \frac{1}{j} \sum_{i=1}^j f(s_i; a_i) A(s_i, a_i); \quad (9)$$

$$\bar{L}_{\pi}^{00} = \frac{1}{j} \sum_{i=1}^j \sum_{a \in \mathcal{A}} \pi(a|s_i) A(s_i, a); \quad (10)$$

and $\bar{L}_{\pi}^{000} = \frac{1}{j} \sum_{i=1}^j \sum_{a \in \mathcal{A}} \pi(a|s_i) A(s_i, a)$ are jA_j jB_j -dimensional vectors, where each element is generated by a state sample from \mathcal{B} , and $\bar{L}_{\pi}^0 = \frac{1}{j} \sum_{i=1}^j f(s_i) g_{s_i, 2B}$ is a jA_j jB_j -dimensional decision vector, where each jA_j -dimensional element vector corresponds to the con dence w.r.t. state samples in \mathcal{B} . Since the expectation L_{π}^0 , L_{π}^{00} , and L_{π}^{000} is over the stationary distribution induced by the behavior policy, all the SAA functions are unbiased Monte-Carlo estimates of their population-based counterparts. We now denote

\bar{L}_{π}^0 , \bar{L}_{π}^{00} , \bar{L}_{π}^{000} as the SAA-MM objective and use this to solve for the con dence vector over the batch samples.

Now consider the following maximization problem:

$$\max_{\pi} \frac{1}{2} \left(\bar{L}_{\pi}^0 - \bar{L}_{\pi}^{00} \right) A; \quad \pi \in \mathcal{P}_{a2A} \quad (11)$$

where the feasible set $\mathcal{P}_{a2A} = \{ \pi \in \mathcal{P} \mid \pi(a|s_i) \in [0, 1]; \sum_{a \in \mathcal{A}} \pi(a|s_i) = 1; \forall s_i \in \mathcal{S} \}$ only imposes constraints on the states that appear in the batch

This finite-dimensional QP problem can be expressed in the following quadratic form:

$$\max_{\pi} \frac{1}{2} \pi^T M \pi + \pi^T c$$

where the symmetric matrix is given by

$$M := \frac{(D_{jA} - j_j + D_{jA}^>)}{jB_j(1 - \pi_j)}$$

$\pi_j = \sum_{i=1}^j \pi(a_i|s_i)$, and $D_{jA} = \text{diag}(f_j A_j g_{a2A; s2B})$ is a jB_j jA_j -diagonal matrix whose elements are the absolute advantage function. By definition, it is positive-semi-definite, hence the QP above is concave. Using its KKT condition, the unique optimal con dence vector over batch is given as

$$\bar{\pi} = \min_{\pi} \frac{1}{2} \pi^T M \pi + \pi^T c; \quad \pi \in \mathcal{P}_{a2A}; \quad (12)$$

where $M^{-1} = \text{blkdiag}(f_j A_j g_{a2A; s2B})$ is a jB_j jA_j -matrix, and the Lagrange multiplier, $\lambda \in \mathbb{R}^{jB_j}$ w.r.t. constraint $\pi \in \mathcal{P}_{a2A}$ is given by

$$\lambda = (M^{-1} c; \lambda) \quad (13)$$

We first construct the con dence function $\bar{\pi}(s; a)$ from the con dence vector $\bar{\pi}$ over \mathcal{B} , in the following tabular fashion:

$$\bar{\pi}(s; a) = \begin{cases} \bar{\pi}_{s;a} & \text{if } (s; a) \in \mathcal{B} \\ 0 & \text{otherwise} \end{cases}$$

While this construction preserves optimality w.r.t. the CPI objective (2), it may be overly conservative, because the policy equates to the behavior policy by setting $\bar{\pi} = 0$ at state-action pairs that are not in \mathcal{B} (i.e., no policy improvement). To alleviate this conservatism, we propose to learn a con dence function that generalizes to out-of-distribution samples.

Learning the Con dence. Given a con dence vector $\bar{\pi}$ corresponding to samples in batch \mathcal{B} , we learn the con dence function $\pi(s; a)$ in supervised fashion. To ensure that the con dence function satisfies the constraint: $\sum_{a \in \mathcal{A}} \pi(a|s) = 1$, i.e., we parameterize it as

$$\pi(s; a) := \frac{\psi(a|s)}{\sum_{a \in \mathcal{A}} \psi(a|s)}; \quad \psi(s; a) \in \mathcal{S} \times \mathcal{A}; \quad (14)$$

where ψ is a learnable policy mapping, such that $\min_{\psi} \sum_{(s; a) \in \mathcal{B}} \psi(a|s) \log \frac{\psi(a|s)}{\sum_{a \in \mathcal{A}} \psi(a|s)}$. We then learn ψ via the following KL distribution-matching objective [Rusuet al., 2015]:

$$\min_{\psi} \frac{1}{|\mathcal{B}|} \sum_{(s; a) \in \mathcal{B}} \psi(a|s) \log \frac{\psi(a|s)}{(1 - \bar{\pi}_{s;a}) \psi(a|s) + \bar{\pi}_{s;a}}$$

While this approach learns ψ by generalizing the con dence vector to out-of-distribution samples, when is a NN, one challenge is to enforce the constraint $\sum_{a \in \mathcal{A}} \psi(a|s) = 1$.

Instead, using an in-graph convex optimization [Namos and Kolter, 2017], we parameterize ψ with a NN with the following constraint-projection layer: $\mathcal{S} \rightarrow \mathcal{A}$ before the output:

$$\begin{aligned} \psi(s) & \in \arg \min_{\pi} \frac{1}{2} \sum_{a \in \mathcal{A}} \pi(a|s) \log \frac{\pi(a|s)}{\sum_{a \in \mathcal{A}} \pi(a|s)} \\ \text{s.t.} & \sum_{a \in \mathcal{A}} \pi(a|s) = 1; \end{aligned} \quad (15)$$

where, at any $s \in \mathcal{S}$, the jA_j -dimensional con dence vector label $\psi_{s;a} g_{a2A}$ is equal to $\bar{\pi}_{s;a} g_{a2A}$ chosen from the batch con dence vector $\bar{\pi}$ such that $\bar{\pi}_{s;a}$ is closest to $\psi_{s;a}$. Indeed, analogous to the closed-form solution in (12), this projection layer has a closed-form QP formulation with linear constraints: $\psi(s) = \min_{\pi} \frac{1}{2} \sum_{a \in \mathcal{A}} \pi(a|s) \log \frac{\pi(a|s)}{\sum_{a \in \mathcal{A}} \pi(a|s)}$; where Lagrange multiplier λ_s is given by $\lambda_s = (\sum_{a \in \mathcal{A}} \psi(a|s) \log \frac{\psi(a|s)}{\sum_{a \in \mathcal{A}} \psi(a|s)})^{-1} \sum_{a \in \mathcal{A}} \psi(a|s) \log \frac{\psi(a|s)}{\sum_{a \in \mathcal{A}} \psi(a|s)}$.

Although the ψ -update is theoretically justified, in practice, when the magnitude of λ_s becomes large (due to the conservatism of the weighted Pinsker inequality), the relative-softmax candidate policy (8) may be too close to the behavior policy π , impeding learning of the residual policy (i.e., $\pi - \pi$). To avoid this in practice, we can upper bound the temperature, i.e., $\psi(s) = \min_{\pi} \frac{1}{2} \sum_{a \in \mathcal{A}} \pi(a|s) \log \frac{\pi(a|s)}{\sum_{a \in \mathcal{A}} \pi(a|s)}$, or introduce a weak temperature-decay schedule, i.e. $\psi(s) = \frac{1}{k} \sum_{a \in \mathcal{A}} \pi(a|s) \log \frac{\pi(a|s)}{\sum_{a \in \mathcal{A}} \pi(a|s)}$, with a tunable $k \in [0, 1]$.

²If one restricts ψ to be only state-dependent, this constraint immediately holds.

Environment [#]	DQN		BRPO-C		BRPO (ours)		BCQ		KL-Q		SPIBB		BC		Behavior Policy
Acrobot-0.05	-91.2	9.1	-94.6	3.8	-91.9	9.0	-96.9	3.7	-93.0	2.6	-103.5	24.1	-102.3	5.0	-103.9
Acrobot-0.15	-83.1	5.2	-91.7	4.0	-86.1	10.1	-97.1	3.3	-92.1	3.2	-91.1	44.8	-113.1	5.6	-114.3
Acrobot-0.25	-83.4	3.9	-91.2	4.1	-85.3	4.8	-96.7	3.1	-90.0	2.9	-86.0	5.8	-124.1	7.0	-127.2
Acrobot-0.50	-84.3	22.6	-90.9	3.4	-83.7	16.6	-77.8	13.5	-84.5	3.8	-106.8	102.7	-173.7	8.1	-172.4
Acrobot-1.00	-208.9	174.8	-156.8	22.0	-121.7	10.2	-236.0	85.6	-227.5	148.1	-184.8	150.2	-498.3	1.7	-497.3
CartPole-0.05	82.7	0.5	220.8	117.0	336.3	122.6	255.4	11.1	323.0	13.5	28.8	1.2	205.6	19.6	219.1
CartPole-0.15	299.3	133.5	305.6	95.2	409.9	64.4	255.3	11.4	357.7	84.1	137.7	11.7	151.6	27.5	149.5
CartPole-0.25	368.5	129.3	405.1	74.4	316.8	64.1	247.4	128.7	441.4	79.8	305.2	119.7	103.0	20.4	101.9
CartPole-0.50	271.5	52.0	358.3	114.1	433.8	93.5	282.5	111.8	314.1	107.0	310.4	128.0	39.7	5.1	37.9
CartPole-1.00	118.3	0.3	458.6	51.5	369.0	42.3	194.0	25.1	209.7	48.4	147.1	0.1	22.6	1.5	21.9
LunarLander-0.05	236.4	177.6	35.6	61.7	88.2	32.0	81.5	14.9	84.4	26.3	-200.4	81.7	75.8	17.7	73.7
LunarLander-0.15	215.6	140.4	79.6	29.7	103.9	49.8	80.3	16.8	61.4	39.0	86.1	73.3	76.4	16.6	84.9
LunarLander-0.25	2.5	101.3	109.5	40.7	141.6	11.0	83.5	14.6	78.7	48.8	166.0	90.6	57.9	13.1	57.3
LunarLander-0.50	-104.6	68.3	42.5	71.4	101.0	39.6	-13.2	44.9	66.2	78.0	-134.6	17.1	-32.6	6.5	-36.0
LunarLander-1.00	-65.6	45.9	53.5	44.1	81.8	42.1	-69.1	44.0	-139.2	29.1	-107.1	94.4	-177.4	13.1	-182.6

Table 1: The mean and st. dev. of average return with the best hyperparameter configuration (with the top-2 results boldfaced). Full training curves are given in the appendix. For BRPO-C, the optimal confidence parameter is found by grid search.

Algorithm 1 BRPO algorithm

Require: B : batch data; Tunable parameter $\epsilon \in [0, 1]$

- 1: for $t = 1; \dots; N$ do
- 2: Sample mini-batch of transition $(s; a; r; d; s^0) \in B$
- 3: Compute \bar{r} from Eq. (12)
- 4: Update confidence by Eq. (15)
- 5: Update candidate policy by Eq. (8)
- 6: Construct target critic network $V_0(s^0) := (1 - \epsilon)E_{a^0} [Q_0(s^0, a^0)] + \max_{a^0} Q_0(s^0, a^0)$
- 7: Update $Q_0 \leftarrow \arg \max_a \frac{1}{2} (Q(s; a) + r + V_0(s^0))^2$
- 8: Update target network: $V_0 \leftarrow (1 - \epsilon)V_0 + \epsilon V_0(s^0)$

6 Experimental Results

To illustrate the effectiveness of BRPO, we compare against six baselines: DQN [Mnih et al., 2013], discrete BCQ [Fujimoto et al., 2018], KL-regularized Q-learning (KL-Q) [Jaques et al., 2018], SPIBB [Laroche and Trichelair, 2017], Behavior Cloning (BC) [Kober and Peters, 2010], and BRPO-C, which is a simplified version of BRPO that uses a constant (tunable) parameter as confidence weight. We do not consider ensemble models, thus do not include methods like BEAR [Kumar et al., 2019] among our baselines. CPI is also excluded since it is subsumed by BRPO-C with grid search on the confidence. It is also generally inferior to BRPO-C because candidate policy learning does not optimize the performance of the final mixture policy. We evaluated on three discrete-action OpenAI Gym tasks [Brockman et al., 2016]: Cartpole-v1, Lunarlander-v2, and Acrobot-v1.

The behavior policy in each environment is trained using standard DQN until it reaches 75% of optimal performance, similar to the process adopted in related work (Fujimoto et al., 2018). To assess how exploration and the quality of behavior policy affect learning, we generate five sets of data for each task by injecting different random exploration into the same behavior policy. Specifically, we add greedy exploration for $\epsilon = 1$ (fully random), 0.5, 0.25, 0.15, and 0.05, generating 100K transitions each for batch RL training.

All models use the same architecture for a given

environment—details (architectures, hyper-parameters, etc.) are described in the appendix of the extended paper. While training is entirely offline, policy performance is evaluated online using the simulator, at every 1000 training iterations. Each measurement is the average return over 100 evaluation episodes and 5 random seeds, and results are averaged over a sliding window of size 10.

Table 1 shows the average return of BRPO and the other baselines under the best hyper-parameter configurations in each task setting. Behavior policy performance decreases as ϵ increases, as expected, and BC matches that very closely. DQN performs poorly in the batch setting. Its performance improves as ϵ increases from 0.05 to 0.25, due to increased state-action coverage, but goes higher (0.5, 1.0), the state space coverage decreases again since the (near-) random policy is less likely to reach a state far away from the initial state.

BCQ, KL-Q and SPIBB follow the behavior policy in some ways, and showing different performance characteristics over the data sets. The underperformance relative to BRPO is more prominent for very low or very high, suggesting deficiency due to overly conservative updates or following the behavior policy too closely, when BRPO is able to learn.

Since BRPO exploits the statistics of each (s, a) pair in the batch data, it achieves good performance in almost all scenarios, outperforming the baselines. The stable performance and robustness across various scenarios make BRPO an appealing algorithm for batch/offline RL in real-world, where it is usually difficult to estimate the amount of exploration required prior to training, given access only to batch data.

7 Concluding Remarks

We have presented Batch Residual Policy Optimization (BRPO) for learning residual policies in batch RL settings. Inspired by CPI, we derived learning rules for jointly optimizing both the candidate policy and state-action dependent confidence mixture of a residual policy to maximize a conservative lower bound on policy performance. BRPO is thus more exploitative in areas of state space that are well-covered by the batch data and more conservative in others. While we have shown successful application of BRPO to various benchmarks, future work includes deriving finite-sample analysis of BRPO, and applying BRPO to more practical batch domains (e.g., robotic manipulation, recommendation systems).

³For algorithms designed for online settings, we modify data collection to sample only from online / batch data.

References

- [Abbasi-Yadkori et al., 2019] Y. Abbasi-Yadkori, P. Bartlett, X. Chen, and A. Malek. Large-scale markov decision problems via the linear programming dual. arXiv:1901.01992 2019.
- [Amos and Kolter, 2017] B. Amos and Z. Kolter. Optnet: Differentiable optimization as a layer in neural networks. In ICML, pages 136–145, 2017.
- [Bolley and Villani, 2005] F. Bolley and C. Villani. Weighted chi-square-kullback-pinsker inequalities and applications to transportation inequalities. *Annales de la Faculté des sciences de Toulouse: Mathématiques* volume 14, pages 331–352, 2005.
- [Boyd and Vandenberghe, 2004] S. Boyd and L. Vandenberghe. *Convex optimization* Cambridge university press, 2004.
- [Brockman et al., 2016] G. Brockman, V. Cheung, L. Petteersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. Openai gym, 2016.
- [Farias and Roy, 2003] P. De Farias and B. Van Roy. The linear programming approach to approximate dynamic programming. *Operations research* 51(6):850–865, 2003.
- [Faybusovich and Moore, 1997] L. Faybusovich and J. Moore. In finite-dimensional quadratic optimization: interior-point methods and control applications. *Applied Mathematics and Optimization* 36(1):43–66, 1997.
- [Fujimoto et al., 2018] S. Fujimoto, D. Meger, and D. Precup. Off-policy deep reinforcement learning without exploration. arXiv:1812.02900 2018.
- [Fujimoto et al., 2019] Scott Fujimoto, Edoardo Conti, Mohammad Ghavamzadeh, and Joelle Pineau. Benchmarking batch deep reinforcement learning algorithms. arXiv:1910.01708 2019.
- [Ghavamzadeh et al., 2016] M. Ghavamzadeh, M. Petrik, and Y. Chow. Safe policy improvement by minimizing robust baseline regret. *NIPS* 2016.
- [Gretton et al., 2007] A. Gretton, K. Borgwardt, M. Rasch, B. Schölkopf, and A. Smola. A kernel approach to comparing distributions. In *AAAI*, 2007.
- [Haarnoja et al., 2018] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. *ICML*, 2018.
- [Hasselt et al., 2016] H. Van Hasselt, A. Guez, and D. Silver. Deep reinforcement learning with double Q-learning. In *AAAI*, 2016.
- [Hunter and Lange, 2004] D. Hunter and K. Lange. A tutorial on EM algorithms. *The American Statistician* 58(1):30–37, 2004.
- [Jaques et al., 2019] Natasha Jaques, Asma Ghandeharioun, Judy Hanwen Shen, Craig Ferguson, Agata Lapedriza, Noah Jones, Shixiang Gu, and Rosalind Picard. Way off-policy batch deep reinforcement learning of implicit human preferences in dialog. arXiv:1907.00456 2019.
- [Johannink et al., 2019] T. Johannink, S. Bahl, A. Nair, J. Luo, A. Kumar, M. Loskyll, J. Ojea, E. Solowjow, and S. Levine. Residual reinforcement learning for robot control. In *ICRA*, pages 6023–6029. IEEE, 2019.
- [Kakade and Langford, 2002] S. Kakade and J. Langford. Approximately optimal approximate reinforcement learning. In *ICML*, pages 267–274, 2002.
- [Kober and Peters, 2010] W. Kober and J. Peters. Imitation and reinforcement learning. *IEEE Robotics & Automation Magazine* 17(2):55–62, 2010.
- [Kumaret al., 2019] A. Kumar, J. Fu, G. Tucker, and S. Levine. Stabilizing off-policy q-learning via bootstrapping error reduction. arXiv:1906.00949 2019.
- [Lange et al., 2012] S. Lange, T. Gabel, and M. Riedmiller. Batch reinforcement learning. *Reinforcement learning* pages 45–73. Springer, 2012.
- [Laroche and Trichelair, 2017] R. Laroche and P. Trichelair. Safe policy improvement with baseline bootstrapping. arXiv:1712.06924 2017.
- [Lillicrap et al., 2015] T. Lillicrap, J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra. Continuous control with deep reinforcement learning. arXiv:1509.02971 2015.
- [Mahmood et al., 2014] A. Mahmood, H. van Hasselt, and R. Sutton. Weighted importance sampling for off-policy learning with linear function approximation. *NIPS* pages 3014–3022, 2014.
- [Mnih et al., 2013] V. Mnih, K. Kavukcuoglu, D. Silver, A. Graves, I. Antonoglou, D. Wierstra, and M. Riedmiller. Playing atari with deep reinforcement learning. arXiv:1312.5602 2013.
- [Nachum et al., 2019] O. Nachum, M. Ahn, H. Ponte, S. Gu, and V. Kumar. Multi-agent manipulation via locomotion using hierarchical sim2real. *CoRL*, 2019.
- [Pirotta et al., 2013] Matteo Pirotta, Marcello Restelli, Alessio Pecorino, and Daniele Calandriello. Safe policy iteration. In *ICML*, pages 307–315, 2013.
- [Rusuet al., 2015] A. Rusu, S. Colmenarejo, C. Gulcehre, G. Desjardins, J. Kirkpatrick, R. Pascanu, V. Mnih, K. Kavukcuoglu, and R. Hadsell. Policy distillation. arXiv:1511.06295 2015.
- [Schulman et al., 2015] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz. Trust region policy optimization. In *ICML*, pages 1889–1897, 2015.
- [Shapiro et al., 2009] A. Shapiro, D. Dentcheva, and A. Ruszczyński. *Lectures on stochastic programming: modeling and theory* SIAM, 2009.
- [Silver et al., 2018] T. Silver, K. Allen, J. Tenenbaum, and L. Kaelbling. Residual policy learning. arXiv:1812.06298 2018.
- [Xu and Yin, 2013] Y. Xu and W. Yin. A block coordinate descent method for regularized multiconvex optimization with applications to nonnegative tensor factorization and completion. *SIAM Journal on imaging sciences* 6(3):1758–1789, 2013.

A Proofs for Results in Section 3

A.1 Proof of Lemma 1

Before going into the derivation of this theorem, we first have the following technical result that studies the distance of the occupation measures that are induced by π and π' .

Lemma 5. The following expression holds for any state-next-state pair (s, s') :

$$(I - T) \pi' - (I - T) \pi = (I - T) \tau; \quad (I - T) \pi' - (I - T) \pi = \tau;$$

where $T(s, s')$ and $T'(s, s')$ represent the transition probabilities from state s to next-state s' following policy π and π' respectively, and for any state-next-state pair (s, s') , $\tau(s, s') = \sum_{a \in A} T(s, s'; a) (\pi'(a|s) - \pi(a|s)) \frac{V(s') - V(s)}{V(s') - V(s)}$.

Proof. Consider the following chain of equalities from matrix manipulations:

$$\begin{aligned} (I - T) \pi' - (I - T) \pi &= (I - T) \pi' - (I - T) \pi + \sum_{s, s' \in S} (T(s, s') - T'(s, s')) \pi(a|s) \frac{V(s') - V(s)}{V(s') - V(s)} \\ &= \sum_{s, s' \in S} (T(s, s') - T'(s, s')) \pi(a|s) \frac{V(s') - V(s)}{V(s') - V(s)} \end{aligned}$$

By multiplying the matrix $(I - T)$ on both sides of the above expression, it implies that

$$(I - T) \tau = \sum_{s, s' \in S} (T(s, s') - T'(s, s')) \pi(a|s) \frac{V(s') - V(s)}{V(s') - V(s)}$$

Using the definition of τ , completes the proof of this lemma. \square

Using the above result, for any initial state $s_0 \in S$, the value functions that are induced by π and π' have the following relationship:

$$\begin{aligned} V_\pi(s) - V_{\pi'}(s) &= \sum_{s_0=s} (I - T) \pi' - (I - T) \pi R + \sum_{s_0=s} (I - T) \pi' - (I - T) \pi (R_\# - R) \\ &= E_{T; \pi'} \sum_{t=0}^{\infty} \gamma^t (s_t; a_t) \frac{(a_t | s_t) - (a_t | s_t)'}{(a_t | s_t)} V(s_{t+1}) \mid s_0 = s \\ &\quad + E_{T; \pi} \sum_{t=0}^{\infty} \gamma^t (s_t; a_t) \frac{(a_t | s_t) - (a_t | s_t)'}{(a_t | s_t)} R(s_t; a_t) \mid s_0 = s \\ &= E_{T; \pi'} \sum_{t=0}^{\infty} \gamma^t (s_t; a_t) \frac{(a_t | s_t) - (a_t | s_t)'}{(a_t | s_t)} Q(s_t; a_t) \mid s_0 = s \\ &= E_{T; \pi'} \sum_{t=0}^{\infty} \gamma^t (s_t; a_t) \frac{(a_t | s_t) - (a_t | s_t)'}{(a_t | s_t)} A(s_t; a_t) \mid s_0 = s \end{aligned} \tag{16}$$

The second equality follows from the fact that $Q(s; a) = R(s; a) + \sum_{s' \in S} T(s, s'; a) V(s')$. The third equality follows from the result in Lemma 5 and the fact that for any state $s \in S$, $(I - T) \pi' - (I - T) \pi R(s) = V_\pi(s) - V_{\pi'}(s)$. The last equality is based on the fact of the confidence constraint that

$$E_{T; \pi'} \sum_{t=0}^{\infty} \gamma^t V(s_t) - E_{T; \pi} \sum_{t=0}^{\infty} \gamma^t V(s_t) = 0; \quad s_0 = s \in S$$

A.2 Proof of Theorem 2

Denote by V_π and $V_{\pi'}$ the vectors of value functions $V_\pi(s)$ and $V_{\pi'}(s)$ at every state $s \in S$ respectively. Re-writing the result in (16) in matrix form, it can be expressed as

$$V_\pi - V_{\pi'} = (I - T) \pi' - (I - T) \pi (R; \# + T; \pi' - T; \pi (V_\pi - V_{\pi'}));$$

where for any states $s \in S$, $R_{ij}(s) = \sum_{a \in A} P_{aj}(s|a) \frac{r_{aj}(s)}{P_{aj}(s|a)} R(s; a)$, and $T_{ij}(s^0|s) = \sum_{a \in A} P_{aj}(s^0|s) \frac{P_{aj}(s^0|s)}{P_{aj}(s|a)}$. This expression implies that

$$(I - (T + T_{ij}))^{-1} T_{ij} (V - V) = (I - T)^{-1} (R_{ij} + T_{ij} V);$$

which further implies that

$$V - V = (I - (T + T_{ij}))^{-1} T_{ij} (I - T)^{-1} (R_{ij} + T_{ij} V);$$

Here based on the definition of T_{ij} and the condense constraint, one can show that $(T + T_{ij})$ is a stochastic matrix (all the elements are non-negative, and $\sum_{j \in S} (T + T_{ij})(s^0|s) = 1, \forall s \in S$). Therefore the matrix $(I - (T + T_{ij}))^{-1} T_{ij}$ is invertible.

Using the matrix inversion lemma, one has the following equality:

$$(I - (T + T_{ij}))^{-1} T_{ij} = (I - (T + T_{ij}))^{-1} (I - T);$$

Therefore the difference of value function $V - V$ can further be expressed as

$$V - V = (I - T - T_{ij})^{-1} (R_{ij} + T_{ij} V);$$

In other words, at any state $s \in S$, the corresponding value function $V(s)$ is given by the following expression:

$$\begin{aligned} V(s) - V(s) &= E \sum_{t=0}^{\infty} (R_{ij} + T_{ij} V)(s_t) j T_{ij}^0; s_0 = s \\ &= E \sum_{t=0}^{\infty} Q_{ij}(s_t; a_t) j T_{ij}^0; s_0 = s \\ &= E \sum_{t=0}^{\infty} A_{ij}(s_t; a_t) j T_{ij}^0; s_0 = s; \end{aligned} \tag{17}$$

where the transition probability $T_{ij}^0(s^0|s) = (T + T_{ij})(s^0|s)$ is given by $\sum_{a \in A} P_{aj}(s^0|s) T(s^0|s; a) + (s; a) \frac{P_{aj}(s^0|s)}{P_{aj}(s|a)}$ at state-next-state pair (s, s^0) . By noticing that $T_{ij}^0(s^0|s)$ is indeed $T(s^0|s)$ (the transition probability that is induced by residual policy), the proof of Theorem 2 is completed.

B Proofs for Results in Section 4

B.1 Proof of Theorem 4

Define the state-action discounted stationary distribution w.r.t. an arbitrary policy $d(s; a) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P(s_t = s; a_t = a | s_0 = s_0, P_0)$ and its state-only counterpart $d(s) = \sum_{a \in A} d(s; a)$. Immediately one can write the difference of return (objective function of this problem) with the following chain of equalities/inequalities:

$$\begin{aligned} E \sum_{t=0}^{\infty} \gamma^t (A_{j;}(s_t) - T_{j;}(s_t; s_0)) &= \frac{1}{1 - \gamma} \sum_{s \in S} d(s) (A_{j;}(s) - T_{j;}(s)) \\ &= \frac{1}{1 - \gamma} \sum_{s \in S} d(s) (A_{j;}(s) - T_{j;}(s)) + (d(s) - d(s)) (A_{j;}(s) - T_{j;}(s)) \end{aligned}$$

Recall that $T_{j;}(s) = \sum_{a \in A} T(s; a) d(s; a) = \sum_{a \in A} T(s; a) \frac{(a|s)}{(a|s)}$. At any state $s \in S$, the difference of stationary distribution $d(s) - d(s)$ can be further expressed as

$$\begin{aligned} (d - d)(s) &= P_0^> (I - T)^{-1} (I - T)^{-1} (s) \\ &= P_0^> (I - T)^{-1} + (I - T)^{-1} T_{j;}; (I - (I - T)^{-1} T_{j;};)^{-1} \\ &= P_0^> (I - T)^{-1} T_{j;}; (I - (I - T)^{-1} T_{j;};)^{-1} (I - T)^{-1} (s) \end{aligned}$$

Let $D = \sum_{s, s' \in S} P_0^> (I - T)^{-1} T_{j;}; (I - (I - T)^{-1} T_{j;};)^{-1} (I - T)^{-1} (s)$ be the occupation measure matrix induced by d . Combining the above arguments one has

$$\begin{aligned} & \sum_{i, j} (d - d); A_{j;}; ij \\ &= \sum_{i, j} P_0^> (I - T)^{-1} T_{j;}; (I - (I - T)^{-1} T_{j;};)^{-1} (I - T)^{-1}; A_{j;}; i \\ &= \frac{1}{1 - \gamma} \sum_{i, j} P_0^> D T_{j;}; ((I - D T_{j;};)^{-1}); D A_{j;}; i \\ &= \frac{1}{1 - \gamma} \sum_{i, j} P_0^> D T_{j;}; (I - (I + D T_{j;};))^{-1} D A_{j;}; i \end{aligned} \tag{18}$$

Now $I + D T_{j;};$ is a stochastic matrix, which is because for any state s ,

$$\begin{aligned} (D T_{j;}; e)(s) &= (1 - \gamma) E \sum_{t=0}^{\infty} \gamma^t T_{j;};(s^0; s) j T_{j;};; s_0 = s \\ &= (1 - \gamma) E \sum_{t=0}^{\infty} \gamma^t (s; a) (a|s) (a|s) j T_{j;};; s_0 = s \\ &= 0 \end{aligned}$$

Using this property one can upper bound the magnitude of each element of $(I - (I + D T_{j;};))^{-1}$ by $\frac{1}{1 - \gamma}$. Therefore, using Holder inequality one can further upper bound the expression in (18) as follows:

$$\begin{aligned} \sum_{i, j} (d - d); A_{j;}; ij &\leq \frac{1}{(1 - \gamma)^2} \sum_{s \in S} P_0^> D T_{j;}; \sum_{i, j} k D A_{j;}; k_j \\ &= \frac{1}{(1 - \gamma)^2} \sum_{s \in S; a \in A} d(s) (a|s) (s; a) \sum_{j} \frac{(a|s) (a|s) j}{(a|s)} \max_{s_0 \in S} \sum_{s_2 \in S} d(s; s_0) j A_{j;}; j \end{aligned}$$

Plugging in the definition of the discounted occupation measure d into the above expression, the proof of this theorem is completed.

C Practical Implementation Details of BRPO

In this section, we discuss several practical techniques to further boost training stability and effectiveness of BRPO.

Improving CPI with Optimal Advantage The derivation in Section 4 and Section 5 shows that optimizing finds a residual policy that performs no-worse than the behavior policy (modulo any Lagrangian approximation error). While we argue that this optimization is less conservative than existing methods (like TRPO) due to the state-action-dependent learned con dence, it might not be aggressive enough in leveraging the function approximation to generalize to unseen state and action. One major reason is that by design, only uses the long-term value of (in the form of A), in order to circumvent the issue of bad generalization to unseen state-action pair. This also makes policy improvement local. This is a fundamental challenge of batch RL, but can be relaxed depending on the domain. As a remedy to this issue, by a convex ensemble of the results from Lemma 1 and Theorem 2 (with any combination weight $\alpha \in [0, 1]$), notice that the difference-return also satisfies

$$J_{\pi} - J_{\pi_0} = \frac{1}{1} \mathbb{E}_{(s,a) \sim d} \left[\frac{(1-\alpha)}{1} L_{\pi_0}^{(0)} + \alpha \max_{s_0 \in \mathcal{S}} L_{\pi_0}^{(0,0)}(s_0) \right];$$

where

$$\mathbb{E}_{(s,a) \sim d} \left[\frac{(1-\alpha)}{1} L_{\pi_0}^{(0)} + \alpha \max_{s_0 \in \mathcal{S}} L_{\pi_0}^{(0,0)}(s_0) \right] := \mathbb{E}_{(s,a) \sim d} \left[\frac{(1-\alpha)}{1} W(s; a) + \alpha \max_{s_0 \in \mathcal{S}} W(s_0; a) \right]$$

with a weighted advantage function $W := (1-\alpha)A + \alpha A$. Therefore, without loss of generality one can replace $L_{\pi_0}^{(0)}$ with the weighted advantage function. Furthermore, to avoid estimating each policy update and assuming that CPI eventually finds π^* , one may directly estimate the optimal weighted advantage function

$$W(s; a) = Q_{\pi^*}(s; a) - V_{\pi^*}(s);$$

in which the value function Q_{π^*} is a Bellman fixed-point of $Q(s; a) = R(s; a) + \mathbb{P}_{s_0} T(s_0; s; a) V(s_0)$, with

$$V(s) = (1-\alpha) \mathbb{E}_a [Q(s; a)] + \alpha \max_a Q(s; a);$$

This approach of combining the optimal Bellman operator with the on-policy counterpart belongs to the general class of hybrid on/off-policy RL algorithms (Sutton & Schuermans, 2016) combining.

Therefore, we learn an advantage function that is a weighted combination of A and A . Using the batch data \mathcal{B} , the expected advantage can be learned with any critic-learning technique, such as SARSA (Sutton, 2018) reinforcement. We can learn A by DQN (Mnih, 2013) playing or other Q-learning algorithm. We provide pseudo-code of our BRPO algorithm in Algorithm 1.

D Experimental Details

This section describes more details about our experimental setup to evaluate the algorithms.

D.1 Behavior policy

We train the behavior policy using DQN, using architecture and hyper-parameters specified in Section D.2. The behavior policy was trained for each task until the performance reaches around 75% of the optimal performance since Fujimoto et al., 2019 and Kumar et al., 2019. Specifically, we trained the behavior policy for 100,000 steps for Lunarlander-v2, and 50,000 steps for Cartpole-v1 and Acrobot-v1. We used two-layers MLP with FC(32)-FC(16). The replay buffer size is 500,000 and batch size is 64. The performance of the behavior policies are given in Table 1.

D.2 Hyperparameters

For fair comparison, we generally used the same set of hyper-parameters and architecture across all methods and experiments, which are defined in Table 2 and Table 3. Similar to the behavior policy, we used two-layers MLP with FC(32)-FC(16) for all the critic agents and Behavioral cloning agent's policy. The final hyperparameters are found using grid search, with candidate set specified in Table 2 and Table 3.

Hyperparameters for BC, BCQ, SARSA, DQN, and KL-Q	Sweep range	Final value
Soft target update rate χ	-	0.5
Soft target update period	150, 500, 1500	500
Discount factor	-	0.99
Mini-batch size	-	64
Q-function learning rates	0.0003, 0.001, 0.002	0.001
Neural network optimizer	-	Adam
[BCQ] Behavior policy threshold (in [Fujimoto et al., 2019])	0.1, 0.2, 0.3, 0.4	0.3
[SPIBB] Bootstrapping set threshold	0.1, 0.2, 0.3, 0.4	0.2
[KL-Q] KL-regularization weight	0.01, 0.03, 0.1, 0.3	0.1

Table 2: The range of hyperparameters swept over and the final hyperparameters used for the baselines (BC, BCQ, SARSA, DQN, and KL-Q).

Hyperparameters for BRPO and BRPO-C	Sweep range	Final value
Soft target update rate χ	-	0.5
Soft target update period	150, 500, 1500	500
Discount factor	-	0.99
Mini-batch size	-	64
Q-function learning rates	0.0003, 0.001, 0.002	0.001
Neural network optimizer	-	Adam
[BRPO, BRPO-C] Mixing (for critic training)	0.1, 0.5, 0.9	0.9
[BRPO] Confidence (s; a) learning rates	0.0002, 0.0001	0.0001
[BRPO-C] Constant (for constant residual policy)	0.25, 0.33, 0.5, 0.66, 0.75	0.5

Table 3: The range of hyperparameters swept over and the final hyperparameters used for the proposed methods (BRPO and BRPO-C).

E Additional Results

Here are the learning curves for each environment and behavior policy over the course of batch training.

