
Towards Hyperparameter-free Policy Selection for Offline Reinforcement Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

How to select between policies and value functions produced by different training algorithms in offline reinforcement learning (RL)—which is crucial for hyperparameter tuning—is an important open question. Existing approaches based on off-policy evaluation (OPE) often require additional function approximation and hence hyperparameters, creating a chicken-and-egg situation. In this paper, we design hyperparameter-free algorithms for policy selection based on BVFT [XJ20], a recent theoretical advance in value-function selection, and demonstrate their effectiveness in discrete-action benchmarks such as Atari. To address performance degradation due to poor critics in continuous-action domains, we further combine BVFT with OPE to get the best of both worlds, and obtain a hyperparameter-tuning method for Q -function based OPE with theoretical guarantees as a side product.

1 Introduction and Related Works

Learning a good policy from historical data without interactive access to the actual environment, or offline (batch) reinforcement learning (RL), is a promising approach to applying RL to real-world scenarios when high-fidelity simulators are not available [LKTF20]. Despite the fast development in the training algorithms, a burning question that remains wide open is how to tune their hyperparameters, sometimes known as the offline policy selection problem [Pai+20; YDNTS20; Fu+21].

Standard approaches reduce the problem to off-policy evaluation (OPE), which estimates the expected return of the candidate policies and choose accordingly. Unfortunately, OPE itself is a difficult problem, and standard estimators such as importance sampling suffer exponential (in horizon) variance [LMS15; JL16]. While polynomial-variance estimators exist, either using TD (e.g., Fitted-Q Evaluation, or FQE [LVY19]) or marginalized importance sampling [LLTZ18; NCDL19; UHJ20], they require additional function approximation, inducing yet another set of hyperparameters (e.g., the neural-net architecture) which need to be carefully chosen. [Pai+20] recently conclude that FQE can be effective for offline policy selection, but “an important remaining challenge is how to choose hyperparameters for FQE”. (Incidentally, we are able to address this question as a side product of our approach in Section 5). In other words, to tune hyperparameters for training we need to tune the hyperparameters for OPE, creating a chicken-and-egg situation. To this end, we want to ask:

Can we design effective *hyperparameter-free* methods for offline policy selection?

The question has been investigated in the theoretical literature [FS11], mostly reformulated so that we select indirectly among value functions instead of policies to trade-off directness for tractability. More precisely, [FS11] imagines that training algorithms produce candidate Q -functions Q_1, Q_2, \dots, Q_m , which is a reasonable assumption as most offline algorithms produce value functions as a side product. The goal is to select $Q_i \approx Q^*$ —assuming one exists—so that the induced greedy policy, π_{Q_i} , is near-optimal. While $\|Q - Q^*\|$ is only a surrogate for the performance of π_Q , the hope is that whether $Q \approx Q^*$ can be more easily verified from holdout data without additional function approximation,

37 possibly by estimating the Bellman error (or residual) $\|Q - \mathcal{T}Q\|$. Unfortunately, $\|Q - \mathcal{T}Q\|$ is not
 38 amendable to statistical estimation in stochastic environments [SB18]. The naïve estimator which
 39 squares the TD error (see “1-sample BR” in Proposition 3) suffers the infamous *double-sampling*
 40 *bias* [Bai95], and debiasing approaches demand additional function approximation (and hence
 41 hyperparameters) [ASM08; FS11]. In prototypical real-world applications, hyperparameter-free
 42 heuristics such as picking the highest Q [GGMVS20] are often used despite the lack of theoretical
 43 guarantees, which we will compare to in our experiments.

44 In this paper, we attack the problem based on a recent theoretical breakthrough in value-function
 45 selection: [XJ20] propose a theoretical algorithm, BVFT, which provides a workaround to the double
 46 sampling issue without requiring additional function approximation; they estimate a form of projected
 47 Bellman error as a surrogate for $\|Q - Q^*\|$, where the function class for projection is created out of
 48 the candidate Q ’s themselves. See Section 3 for details. Our contributions are 2-fold:

- 49 1. We design a practical implementation of BVFT based on novel theoretical observations, removing
 50 its last hyperparameter which determines a discretization resolution. We empirically demonstrate
 51 that BVFT enjoys promising performance in discrete-action benchmarks such as Atari games,
 52 sometimes using 20x less data than required by FQE-based policy selection.
- 53 2. The vanilla BVFT suffers performance degradation in continuous-action benchmarks, where the
 54 training algorithms often have an actor-critic structure and output a (π, Q) pair where Q is far
 55 away from Q^* for various reasons. To address this challenge, we propose BVFT-PE, a variant
 56 of BVFT that allows us to select among (π, Q) pairs and pick one where $Q \approx Q^\pi$ and π yields
 57 a high return. To further handle the issue that Q from the critic is often a poor fit of Q^π , we
 58 propose to use *multiple* OPE algorithms to re-fit Q^π , and run BVFT-PE among the produced (π, Q)
 59 pairs. While the OPE algorithms often have many hyperparameters that need to be set, BVFT-PE
 60 automatically chooses between them, leaving very few to no hyperparameters untunable. This
 61 allows us to combine the strengths of OPE and BVFT and get the best of both worlds. We also
 62 show additional results that BVFT-PE can be used for hyperparameter tuning in Q -function-based
 63 OPE and provide theoretical guarantees, which is of independent interest.

64 2 Preliminaries

65 **Markov Decision Processes (MDPs)** In RL, we often model the environment as an MDP, specified
 66 by its state space \mathcal{S} , action space \mathcal{A} , reward function $R : \mathcal{S} \times \mathcal{A} \rightarrow [0, R_{\max}]$, transition function
 67 $P : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ ($\Delta(\cdot)$ is the probability simplex), discount factor $\gamma \in [0, 1]$, and a initial state
 68 distribution d_0 . We assume $\mathcal{S} \times \mathcal{A}$ is finite but can be arbitrarily large. A deterministic policy $\pi : \mathcal{S} \rightarrow$
 69 \mathcal{A} induces a random trajectory $s_0, a_0, r_0, s_1, a_1, r_1, \dots$ where $s_0 \sim d_0, a_t = \pi(s_t), r_t = R(s_t, a_t)$,
 70 and $s_{t+1} \sim P(\cdot | s_t, a_t), \forall t$. We measure the performance of π using $J(\pi) := \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r_t | \pi]$. In
 71 discounted MDPs, there always exists an optimal policy π^* that maximizes $J(\cdot)$ for all starting states.
 72 It is the greedy policy of the optimal Q -function, Q^* , i.e., $\pi^* = \pi_{Q^*} := (s \mapsto \arg \max_a Q^*(s, a))$.
 73 Q^* is the fixed point of Bellman optimality equation, $Q^* = \mathcal{T}Q^*$, where $\forall f \in \mathbb{R}^{\mathcal{S} \times \mathcal{A}}, (\mathcal{T}f)(s, a) :=$
 74 $R(s, a) + \gamma \mathbb{E}_{s' \sim P(\cdot | s, a)} [\max_{a'} f(s', a')]$. A related important concept is Q^π , which tells us the
 75 expected return of π when the trajectory starts from a specific state-action pair.

76 **Offline Data** In offline RL, we are given a dataset of (s, a, r, s') tuples and cannot directly interact
 77 with the MDP. For the theoretical part of the paper, we assume the standard *offline sampling protocol*,
 78 that the tuples are generated i.i.d. as $(s, a) \sim \mu, r = R(s, a), s' \sim P(\cdot | s, a)$. With a slight abuse of
 79 notation, we also use $\mathbb{E}_\mu[\cdot]$ to denote the expectation over (s, a, r, s') sampled as above.

80 **Policy Selection/Ranking** We will use the following unified framework for policy selection throughout
 81 the paper: Suppose training algorithms (or the same algorithm with different hyperparameters)
 82 produce multiple (π, Q) pairs, $\{(\pi_i, Q_i)\}_{i=1}^m$, and our goal is to select a policy with good perfor-
 83 mance. The relationship between π and Q can differ in different contexts: for example, when training
 84 algorithms try to fit Q^* and induce a greedy policy, we have $\pi_i = \pi_{Q_i}$, and we only need to work
 85 with $\{Q_i\}_{i=1}^m$ as they contain all the relevant information. In the case where training algorithms have
 86 an actor-critic structure, π and Q are separate quantities and need to be reasoned about together. In
 87 real applications, the next step in the pipeline is to deploy the policy in the real system for online
 88 evaluation, and since we may have the resources to test more than 1 policies, we require all algorithms
 89 to produce a *ranking* over $\{\pi_i\}_{i=1}^m$, often by sorting the policies in ascending order w.r.t. a loss.

90 **2.1 Experiment Setup**

91 To avoid interrupting the flow of intertwined theoretical reasoning and empirical evaluation in the
 92 rest of the paper, we briefly describe our experiment setup here, with details deferred to Appendix C.

93 **Environments and datasets** We perform empirical evaluation on OpenAI Gym [Bro+16], Atari
 94 games [BNVB13], and Mujoco [TET12]. Taxi [Die00] is used for sanity check. We use standard
 95 offline datasets when available (RLUnplugged [Gul+21] for Atari, and D4RL [FKNTL21] for
 96 MuJoCo), and generate our own otherwise by mixing a trained expert policy with 30% chance of
 97 acting suboptimally. [Fu+21] have proposed a new benchmark for offline policy selection, which is
 98 also based on RLUnplugged/D4RL which we use and has only become available very recently. Also
 99 this benchmark (and [VLJY19]) focuses on OPE and policy selection without value functions, which
 100 does not exactly fit our purposes. We leave the evaluation on their benchmarks to future work.

101 **Training algorithms** We use several different training algorithms to generate the candidate models,
 102 including offline algorithms such as BCQ [FCGP19] and CQL [KZTL20]. To test the robustness of
 103 the policy-selection methods w.r.t. how the candidate policies are generated, we also use algorithms
 104 that learn from online interactions such as DQN [Mni+15] in some domains, though policy selection
 105 is always performed using separate offline datasets. This scenario is also of interest in its own right,
 106 as one can imagine training on (possibly imperfect) simulators via online algorithms and using
 107 limited realworld offline data for policy selection. For each algorithm, we consider different neural
 108 architectures, learning rates, and learning steps as hyperparameters to produce multiple candidate
 109 policies (and value functions) for selection; see Table 1 in Appendix C for details.

110 **Performance metrics** In each experiment, we run different policy-selection methods and evaluate
 111 the produced policy rankings using the following two metrics adapted from [YDNTS20]:

112 **top- k normalized regret** We take the best policy within the top- k recommended by the ranking,
 113 and calculate its gap compared to the best policy among all candidates. To make the value more
 114 interpretable, we normalize the regret by the gap between the best and the worst candidate policies.

115 **top- k precision** We take the k best policies in terms of their groundtruth values, and return the
 116 proportion of them appearing in the top- k policies in the ranking.

117 This process is repeated for 200 or 300 runs, with randomness coming from re-sampling a subset of
 118 the dataset for policy selection (usually of size 50,000; FQE needs much more data and we do not
 119 run it on random subsets) and sampling $m = 10$ or 15 policies from all candidates for comparison.
 120 All figures report the mean with error bars of twice the standard errors, i.e., 95% confidence intervals.

121 **3 Background: Batch Value-Function Tournament (BVFT)**

122 We briefly introduce the theoretical basis of our approach. As mentioned in Section 1, the Bellman
 123 error $\|Q - \mathcal{T}Q\|$ is an appealing quantity because $Q = Q^* \Leftrightarrow \|Q - \mathcal{T}Q\|_\infty = 0$, but it is not
 124 amendable to statistical estimation in stochastic environments due to the double-sampling bias
 125 [Bai95; ASM08; FS11]. Given this caveat, a closely related quantity has been extensively studied in
 126 the literature: given function class $\mathcal{G} \subset [0, \frac{R_{\max}}{1-\gamma}]^{\mathcal{S} \times \mathcal{A}}$, the (mean-squared) *projected Bellman error*
 127 of Q w.r.t. \mathcal{G} is defined as

$$\|Q - \mathcal{T}_{\mathcal{G}}Q\|_{2,\mu}^2, \text{ where } \mathcal{T}_{\mathcal{G}}Q := \arg \min_{g \in \mathcal{G}} \mathbb{E}_\mu[(g(s, a) - r - \gamma \max_{a'} Q(s', a'))^2]. \quad (1)$$

128 Here $\|\cdot\|_{2,\mu}^2 = \mathbb{E}_\mu[(\cdot)^2]$, and the dependence of $\mathcal{T}_{\mathcal{G}}$ on μ is suppressed for readability. This quantity
 129 can be straightforwardly estimated from data when \mathcal{G} has bounded statistical complexity; we just need
 130 to replace all $\mathbb{E}_\mu[\cdot]$ with their finite-sample approximation. The property of the projected Bellman
 131 error largely depends on the choice of \mathcal{G} , which needs to satisfy certain conditions for $\|Q - \mathcal{T}_{\mathcal{G}}^\mu Q\|$
 132 to be a good surrogate for $\|Q - Q^*\|$. The seminal work of [Gor95] has provided the following
 133 sufficient condition: (see [XJ20, Proposition 4] for a formal proof)

134 **Proposition 1.** *If (1) \mathcal{G} is a piecewise-constant class, and (2) $Q^* \in \mathcal{G}$, then $\mathcal{T}_{\mathcal{G}}$ is γ -contraction
 135 under $\|\cdot\|_\infty$, and $Q = Q^* \Leftrightarrow \|Q - \mathcal{T}_{\mathcal{G}}Q\|_{2,\mu} = 0$ if μ is fully supported on $\mathcal{S} \times \mathcal{A}$.*

136 Being piecewise constant means that there exists a partitioning of $\mathcal{S} \times \mathcal{A}$, and \mathcal{G} consists of all
 137 members of $[0, \frac{R_{\max}}{1-\gamma}]^{\mathcal{S} \times \mathcal{A}}$ that remain constant within each partition. A partitioning whose induced
 138 \mathcal{G} satisfies (1) and (2) is also closely related to Q^* -preserving state abstractions [LWL06].

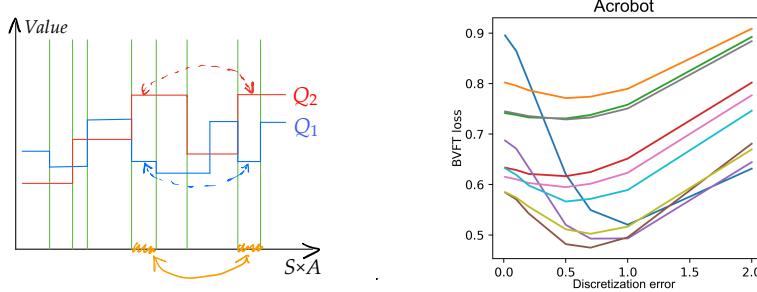


Figure 1: **Left:** Visualization of the partitioning that induces $\mathcal{G}_{1,2}$. The 2 subsets of $\mathcal{S} \times \mathcal{A}$ marked orange belong to the same partition despite being separated from each other, because Q_1 and Q_2 are constant across them. Since $\mathcal{S} \times \mathcal{A}$ is partitioned according to Q 's output, the number of partitions is independent of $|\mathcal{S} \times \mathcal{A}|$, allowing BVFT to scale to arbitrarily complex state-action spaces. **Right:** BVFT-loss vs. discretization error ϵ_{dct} of 10 candidate Q 's of a typical run in Acrobot, all having the U-shape predicted by our theoretical reasoning.

139 The problem is that it is very difficult to find \mathcal{G} that satisfies the above 2 criteria,¹ and it will just
 140 be another set of hyperparameters to tune if we leave the design of \mathcal{G} to the user. [XJ20] offers a
 141 resolution in the context of selecting Q^* from $\{Q_i\}_{i=1}^m$: consider the base case of $m = 2$ where
 142 we only need to select between Q_1 and Q_2 . If one of them is Q^* (which can be relaxed) but we do
 143 not know which, we can still create $\mathcal{G}_{1,2}$ that satisfies both criteria of Proposition 1 as the minimal
 144 piecewise-constant class such that $\{Q_1, Q_2\} \subset \mathcal{G}_{1,2}$. If the output of each Q_i takes at most N
 145 possible values, $\mathcal{G}_{1,2}$ will be induced by partitioning $\mathcal{S} \times \mathcal{A}$ into at most N^2 regions; see Figure 1L.
 146 [XJ20] further shows that this idea extends to arbitrary m via pairwise comparison (“tournament”):

147 **Proposition 2** (Simplification of [XJ20]; see Appendix B.1 for a proof sketch). *If $Q^* \in \{Q_i\}_{i=1}^m$ and μ is fully supported on $\mathcal{S} \times \mathcal{A}$, then $Q_i = Q^* \Leftrightarrow Q_i$ having 0 BVFT-loss, where*

$$\text{BVFT-loss}(Q_i ; \{Q_j\}_{j=1}^m) := \max_j \|Q_i - \mathcal{T}_{\mathcal{G}_{i,j}} Q_i\|_{2,\mu}. \quad (2)$$

149 For each Q_i , BVFT calculates its projected Bellman error w.r.t. $\mathcal{G}_{i,j}$ —where $\mathcal{G}_{i,j}$ is created just like
 150 $\mathcal{G}_{1,2}$ but using Q_i itself and every other Q_j —and scores Q_i using the worst-case projected error.
 151 See [XJ20] for the complete theory that accounts for $Q^* \notin \{Q_i\}_{i=1}^m$ and finite-sample effects. For
 152 readability we will stick to the above simplified reasoning.

153 **Computation** The empirical version of BVFT-loss can be computed exactly in closed form. The
 154 central step is the calculation of $\mathcal{T}_{\mathcal{G}} Q$. Recall that \mathcal{G} is piecewise constant and induced by a $\mathcal{S} \times \mathcal{A}$
 155 partitioning. For any (\tilde{s}, \tilde{a}) , $(\mathcal{T}_{\mathcal{G}} Q)(\tilde{s}, \tilde{a})$ is simply the average of $r + \gamma \max_{a'} Q(s', a')$ over all
 156 data points (s, a, r, s') where (s, a) falls in the same partition as (\tilde{s}, \tilde{a}) . Implemented with running
 157 averages, the computational complexity is $O(m^2 n)$ for m candidate functions and n data points in
 158 addition to $(|\mathcal{A}| + 1)mn$ candidate-function evaluations (i.e., caching $Q_i(s, a)$ and $\max_{a'} Q_i(s', a')$).

159 4 BVFT with Automatic Resolution Selection

160 Despite the appealing theoretical properties, it is unclear if BVFT can be converted into a practical
 161 algorithm. The sample complexity of estimating BVFT-loss depends on the complexity of $\mathcal{G}_{i,j}$,
 162 which is controlled by N^2 where N is the number of possible values in $[0, \frac{R_{\max}}{1-\gamma}]$ each Q_i can take.
 163 To handle the issue that N can be very large or even infinite, [XJ20] discretizes the output of $\{Q_i\}_{i=1}^m$
 164 up to ϵ_{dct} error before producing $\{\mathcal{G}_{i,j}\}$, where ϵ_{dct} needs to be carefully chosen to trade-off between
 165 discretization errors and the complexity of \mathcal{G} (which is now $O(1/\epsilon_{\text{dct}}^2)$). The theoretical value of ϵ_{dct}
 166 (as in [XJ20, Theorem 2]) not only relies on unknown properties of the MDP and the data, but is also
 167 very small which makes BVFT-loss expensive to estimate. Indeed, the previous theory predicts that
 168 BVFT-loss becomes an unstable statistic when $\epsilon_{\text{dct}} \rightarrow 0$ due to the unbounded complexity of $\mathcal{G}_{i,j}$.

¹The ideal \mathcal{G} can be created similarly to Figure 1 according to knowledge of Q^* . In Appendix D we show that our approach closely tracks the skyline of using the ideal \mathcal{G} in a tabular domain where we can compute Q^* .

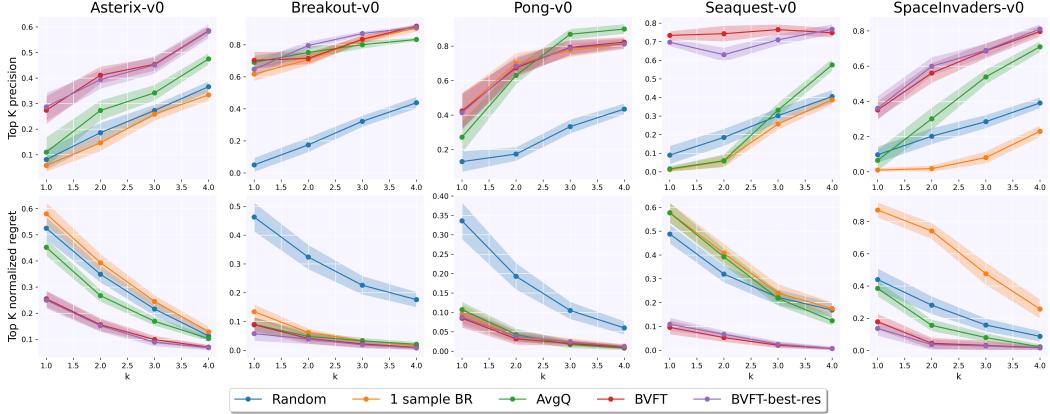


Figure 2: Top- k metrics of policy rankings vs. k in Atari. Row 1 shows top- k precision (the higher the better), and Row 2 shows top- k regret (the lower the better). Training algorithms are BCQ with different hyperparameters. The dataset for policy selection has 50,000 transition, which is an order of magnitude less than needed by FQE in Atari (see FQE in Enduro [VLJY19], as well as Figure 4).

169 To resolve this issue, we draw an interesting connection between BVFT-loss and the naïve “1-sample”
170 estimator for $\|Q - \mathcal{T}Q\|$, which reveals the unexpected behavior of BVFT-loss in the regime of
171 $\epsilon_{\text{dct}} \rightarrow 0$ that differs from the previous theoretical predictions:

172 **Proposition 3.** Consider $\mathbb{E}_\mu[(Q(s, a) - r - \gamma \max_{a'} Q(s', a'))^2]$, the naïve and biased estimator
173 for $\|Q - \mathcal{T}Q\|_{2,\mu}^2$, which we call “1-sample BR (Bellman residual)”. If each candidate Q_i never
174 predicts the same value for any two (s, a) pairs seen in the dataset, then with $\epsilon_{\text{dct}} = 0$, the empirical
175 version of BVFT-loss($Q_i; \{Q_j\}_{j=1}^m$) and 1-sample BR of Q_i coincide (up to squaring).

176 *Proof.* Given a dataset D consisting of (s, a, r, s') tuples, the empirical version of 1-sample
177 BR is $\frac{1}{|D|} \sum_{(s, a, r, s') \in D} (Q(s, a) - r - \gamma \max_{a'} Q(s', a'))^2$, and that of BVFT-loss squared is
178 $\max_j \frac{1}{|D|} \sum_{(s, a, r, s') \in D} (Q(s, a) - (\widehat{\mathcal{T}}_{G_{i,j}} Q)(s, a))^2$, where $\widehat{\mathcal{T}}_{G_{i,j}}$ is the empirical version of $\mathcal{T}_{G_{i,j}}$
179 based on the same dataset. It suffices to show that for any data point (s, a, r, s') , $(\widehat{\mathcal{T}}_{G_{i,j}} Q)(s, a) =$
180 $r + \gamma \max_{a'} Q_i(s', a')$, which follows immediately from $\epsilon_{\text{dct}} = 0$ and Q_i never predicting the exact
181 same value twice, since $G_{i,j}$ does not provide any aggregation over data points in this case and $\widehat{\mathcal{T}}_{G_{i,j}}$
182 coincides with the 1-sample Bellman update (c.f. the paragraph on computation in Section 3). \square

183 This result implies that, 1-sample BR, which is a reasonable objective and coincides with $\|Q - \mathcal{T}Q\|$
184 in deterministic environments, provides a safeguard to BVFT-loss when ϵ_{dct} is too small. Since
185 the double-sampling bias of 1-sample BR is positive [Bai95], we expect BVFT-loss to gradually
186 decrease as ϵ_{dct} increases, hit a minimum, and increase again due to discretization errors when ϵ_{dct}
187 becomes too large.² Indeed, this is precisely what we observe empirically; see Figure 1R.

188 Based on this novel observation, we propose to search for a grid of discretization errors in BVFT
189 and pick the resolution that minimizes the loss (Eq.(2)); see pseudocode in Appendix A. In the
190 experiments we will always use this rule to automatically select the discretization resolution for BVFT
191 and its variants. The remaining questions can only be answered empirically: Does BVFT ever exhibit
192 more interesting behavior than 1-sample BR (especially given their intimate relationship), is our
193 resolution selection rule a good one, and how does BVFT compare to other baselines?

194 **Empirical Evaluation** To answer these questions, we empirically compare BVFT, 1-sample BR,
195 and another simple hyperparameter-free heuristic, AvgQ [GGMVS20], which simply ranks $\{Q_i\}_{i=1}^m$
196 based on their average value on the data. “Random” ranks the policies in a completely random
197 manner. BVFT and 1-sample BR rank $\{Q_i\}_{i=1}^m$ in ascending order of their loss functions, respectively.
198 The experiments are done on 5 Atari games (Figure 2) and 4 Gym control problems (Figure 3; action
199 space is made discrete in Pendulum).

²This can be violated in some extreme cases; see Appendix A for details.

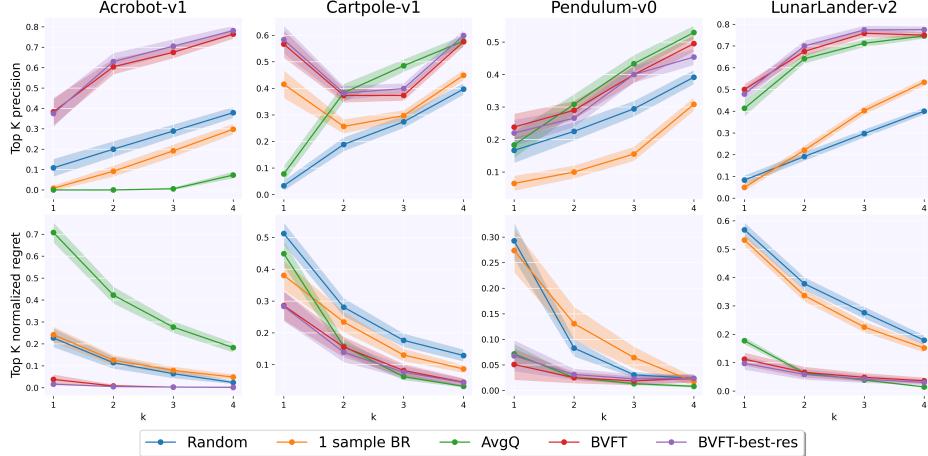


Figure 3: Results in Gym control problems. The dataset for policy selection has 50,000 transitions.

200 • *Does BVFT ever exhibit more interesting behavior than 1-sample BR?* Perhaps surprisingly, BVFT
 201 deviates significantly from the behavior of 1-sample BR, and almost always outperforms the latter by
 202 a wide margin. This is particularly interesting in CartPole and Pendulum, where the deterministic
 203 dynamics make 1-sample BR an **unbiased** estimate of $\|Q - \mathcal{T}Q\|$, and we expected it to perform
 204 well. Contrary to our expectation, 1-sample BR performs poorly even in these domains, where BVFT
 205 often performs much better. In fact, we observe similar phenomenon in a version of Taxi where we
 206 can compute $\|Q - \mathcal{T}Q\|$ as a skyline; see Appendix D.

207 • *Is our resolution selection rule effective?* To examine the effectiveness of our resolution selection
 208 rule, we compare to a skyline of BVFT itself (“BVFT-best-res”), where the best fixed resolution is
 209 chosen based on average statistics in the hindsight. As we can see, BVFT with automatic resolu-
 210 tion selection closely tracks the performance of BVFT-best-res—to the extent that they are mostly
 211 indistinguishable—indicating that our resolution selection rule is near-optimal.

212 • *Comparison to AvgQ.* The simple heuristic of picking Q with the highest predicted value—which is
 213 sometimes used in prototypical applications [GGMVS20]—can be surprisingly effective sometimes,
 214 such as in Pendulum, LunarLander, and Pong. BVFT performs equally well in these domains.
 215 Moreover, as results from other domains reveal, AvgQ is very unstable and can fail catastrophically,
 216 such as in Acrobot and Seaquest, and is much less robust compared to BVFT.

217 One may wonder if AvgQ works better with *pessimistic* training algorithms: if every Q_i is an under-
 218 estimation of the true return, then maximizing Q will be a well-justified heuristic. However, such
 219 pessimism is not always guaranteed especially when the function approximation is misspecified. In
 220 Appendix D we show additional experiments in Atari with a pessimistic algorithm CQL, and the
 221 results are qualitatively similar to Figure 2.

222 **Comparison to OPE** Despite the lack of no known method for tuning OPE’s hyperparameters
 223 except for “cheating” in the simulator using online roll-outs (in Section 5 we will combine BVFT with
 224 OPE to address this issue), which makes it very difficult to have fair comparisons with OPE-based
 225 policy-selection methods, it is still instructive to have a rough sense of how our method compares to
 226 OPE. In this and the next sections, we choose Fitted Q-Evaluation (FQE) as a representative OPE
 227 algorithm.³ Figure 4L shows the ranking metrics vs. sample size in 2 Atari games. We tune the
 228 neural architecture for FQE by choosing the training architecture that produces the *best* policy in
 229 Asterix.⁴ While FQE has equally good performance compared to BVFT in the large-sample regime
 230 (10^6 transitions, which is typically required for FQE in Atari), the performance degradation is severe
 231 when sample size decreases. In comparison, BVFT is much more sample-efficient and provides almost
 232 the same level of performance with about 20x less samples, which is also what we used in Figure 2.

³Despite recent exciting developments in marginalized importance sampling (MIS), FQE often shows strong empirical performance thanks to its simplicity [VLJY19; Pai+20; Fu+21], whereas MIS is not as off-the-shelf due to its difficult optimization, so we defer the comparison as well as combining BVFT with MIS to future work.

⁴The recent OPE benchmarks [Pai+20; Fu+21] do not include Atari, so we choose our own hyperparameters.

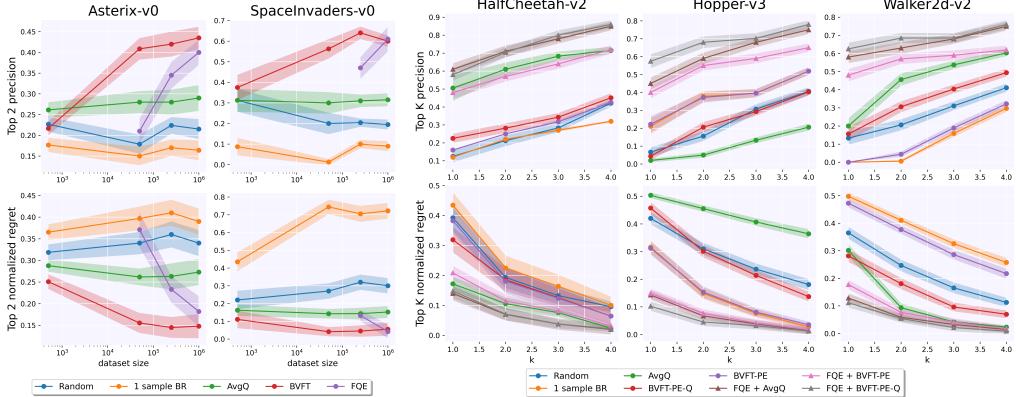


Figure 4: **Left:** Top-2 metrics vs. sample size in 2 Atari games, with FQE as an additional baseline. We did not experiment with FQE on small sample sizes because the trend of performance degradation is clear. **Right:** Top- k metrics vs. k across 3 Mujoco domains with continuous actions. All methods use a dataset of size 50,000, except that the FQE component of any method uses 10^6 .

233 **Additional Results** In Appendix D we examine the sensitivity of different methods to the dataset
 234 size and exploratoriness. BVFT remains advantageous compared to the baselines, is insensitive to data
 235 exploratoriness, and enjoys improved performance when more data becomes available unlike other
 236 hyperparameter-free baselines.

237 5 Fighting Poor Value-Function Estimates with BVFT + OPE

238 The experiments so far are on discrete-action domains, and the candidate policies are always greedy
 239 w.r.t. the Q -function. However, in continuous-action domains, actor-critic-type algorithms are often
 240 used, where a critic Q is trained to evaluate a parameterized policy π , and the actor π is in turn
 241 improved based on Q . The trained (π, Q) pair in general does not satisfy $\pi = \pi_Q$, so it seems unwise
 242 to ignore π and only focus on Q in policy selection. On a related note, $\arg \max_{a'} \cdot$ is often expensive if
 243 not infeasible to calculate in continuous action spaces, so we need to make changes to BVFT anyway.

244 More importantly, the success of BVFT relies on the existence of a reasonable approximation of
 245 Q^* among the candidate functions, which does not always hold especially for actor-critic algo-
 246 rithms. Indeed, in the Mujoco experiments shown in Figure 4R, we observe poor performance of all
 247 hyperparameter-free methods, including variants of BVFT for continuous-action domains we develop
 248 later, and no method can consistently outperform the trivial baseline of random ranking. We address
 249 this issue of poor critics in the rest of this section. The issue is quite complicated and has many
 250 contributing factors—as we will explain subsequently—requiring us to take a multi-step approach to
 251 address one factor at a time.

252 **Step 1: BVFT-PE and BVFT-PE-Q for Joint Selection of (π, Q)** We first address the issue that
 253 $\arg \max_{a'} \cdot$ may be infeasible to calculate in continuous-action domains, and that both π and Q need
 254 to be taken into consideration for actor-critic algorithms. To address this issue, we propose a variant
 255 of BVFT, called BVFT-PE, with the following loss function:

$$\text{BVFT-PE-loss}((\pi_i, Q_i); \{(\pi_j, Q_j)\}_{j=1}^m) := \max_j \|Q_i - \mathcal{T}_{\mathcal{G}_{i,j}}^{\pi_i} Q_i\|_{2,\mu}, \quad (3)$$

256 where $\mathcal{G}_{i,j}$ is exactly the same as in BVFT, and $\mathcal{T}_{\mathcal{G}_{i,j}}^{\pi_i}$ is the same as Eq.(1), except that $\max_{a'} Q(s', a')$
 257 should be replaced by $Q(s', \pi(s'))$. (In actor-critic algorithms, π is often a stochastic policy, in which
 258 case the term means $\mathbb{E}_{a' \sim \pi(\cdot | s')} [Q(s', a')]$.) This loss function naturally generalizes the BVFT-loss:
 259 when each π_i is the greedy policy of Q_i , we immediately have $\mathcal{T}_{\mathcal{G}_{i,j}}^{\pi_i} Q_i = \mathcal{T}_{\mathcal{G}_{i,j}} Q_i$, thus recovering
 260 BVFT-loss as a special case.

261 A closer inspection of Eq.(3) reveals an interesting fact: $\text{BVFT-PE-loss}((\pi_i, Q_i)) = 0$ as long
 262 as $Q_i = Q^{\pi_i}$, so BVFT-PE-loss is really a loss function for *policy evaluation*, hence the name
 263 BVFT-PE. (We will actually provide the theoretical guarantees of BVFT-PE for policy evaluation
 264 at the end of this section; see Theorem 4.) While the loss recovers BVFT-loss when $\pi = \pi_Q$,

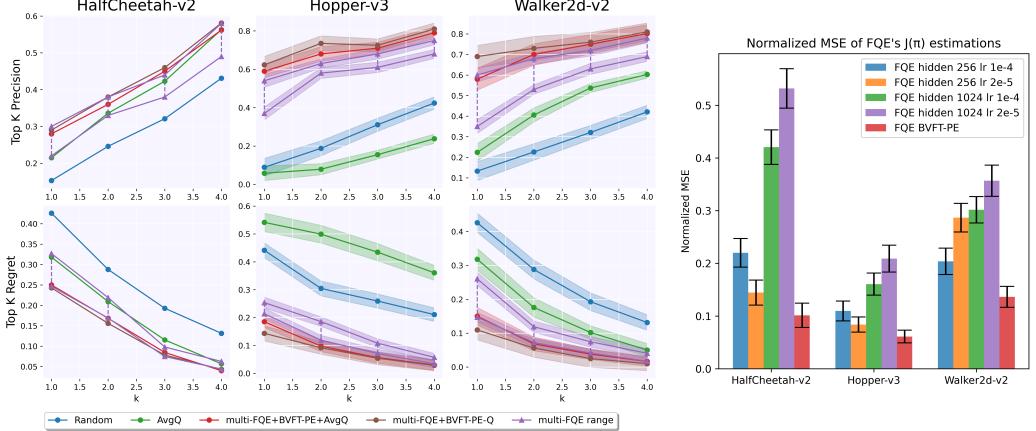


Figure 5: **Left:** Policy selection with multiple FQE instances in Mujoco. To avoid cluttering, we show the range of FQE performance across different hyperparameters (with upper and lower bounds connected by vertical dashed lines), and have removed error bars from HalfCheetah for better visibility (their widths are similar to the other two figures). The curves are cluttered in HalfCheetah and both our strategies perform similarly to the upper bound. **Right:** OPE accuracy for using BVFT-PE for hyperparameter tuning in FQE, where $J(\pi)$ is approximated by $\mathbb{E}_{s \sim d_0}[Q(s, \pi(s))]$.

more generally there can exist the degenerate cases of $Q = Q^\pi$ but π itself being a poor policy, where $\text{BVFT-PE-loss}((\pi, Q))$ is still 0. We address this issue by subtracting a Q term from the BVFT-PE-loss : $\text{BVFT-PE-Q} := \text{BVFT-PE-loss} - \lambda \mathbb{E}_\mu[Q]$, and the structure of this loss resembles a telescoping identity for $J(\pi)$ commonly used in the OPE literature; see Appendix A.1 for details.

Step 2: Re-fitting Q with OPE The above derivations address some of the basic theoretical issues, but still implicitly assume that at least one good π_i is paired with a $Q_i \approx Q^{\pi_i}$. In actor-critic algorithms, however, we often observe that the actor π converges way before the critic Q does, leaving us with a poorly estimated Q . Indeed, we have already seen in Figure 4 that BVFT-PE and BVFT-PE-Q are still not effective when applied to the candidate (π, Q) pairs.

A natural solution to the problem is to use OPE algorithms to refit Q^π to replace the critic Q , in order to provide a more accurate value function. We implement this using FQE as the OPE algorithm, though in principle we can use any OPE algorithm that provides an estimate of Q^π , including kernel loss [FLL19], MQL [UHJ20], or even a model-based method through planning. Figure 4 shows that BVFT-PE and BVFT-PE-Q enjoy strong empirical performance. However, the baseline of simply ranking the policies according to FQE itself (“FQE+AvgQ”) is equally effective, which suggests that most of the performance gains should be attributed to FQE. That said, even if one’s conclusion is that “OPE-based policy selection is more superior in continuous-action domains”, we are just all the way back to where we started in Section 1:

How can we tune the hyperparameters for OPE itself?

Step 3: Re-fitting Q with Multiple OPE Algorithms Our last idea is retrospectively simple, yet allows us to combine the strengths of BVFT and OPE and get the best of both worlds. Suppose we have L OPE algorithms which we wish to select from. We will simply run each OPE algorithm on each policy, producing $m \times L$ candidate (π, Q) pairs in total, in the form of $\{(\pi_i, Q_i^l)\}$, where Q_i^l is the estimation of Q^{π_i} by the l -th OPE algorithm. There are two strategies we can proceed with:

Strategy 1: Run BVFT-PE-Q on all $m \times L$ pairs of (π, Q) . Rank the m policies according to their highest position in the ranking produced by BVFT-PE-Q.

Strategy 2: Within each π_i , use BVFT-PE to select $Q_i^{l_i^*} \approx Q^\pi$, then rank the policies based on the predictions of $Q_i^{l_i^*}$ (e.g., using AvgQ).

Empirical Evaluation We use these strategies in the same experiment setting as Figure 4R, with multiple instances of FQE using different hyperparameters (see legend of Figure 5R). The results are shown in Figure 5L, where both strategies perform similarly to or better than the best FQE instance. We emphasize that this is highly nontrivial, because BVFT-PE and BVFT-PE-Q do not observe the *identities* of the FQE algorithms across different policies and runs (in fact, each FQE algorithm is represented as merely $2mn$ numbers), so matching the best FQE performance is strong evidence for

299 BVFT algorithms' model selection capabilities. Between the two BVFT strategies, Strategy 1 (using
300 BVFT-PE-Q) slightly outperforms Strategy 2, but comes with an additional hyperparameter λ ; we
301 tuned it on Hopper and use the same constant in all experiments. In comparison, Strategy 2 is only
302 slightly worse and does not have its own hyperparameters, which is an advantage.

303 **Hyperparameter Tuning in OPE** In Strategy 2, we are basically using OPE as the policy-selection
304 algorithm, and BVFT-PE as a subroutine for hyperparameter tuning for OPE itself. Since OPE
305 is an important component of the offline RL pipeline in its own right, our procedure for tuning
306 OPE algorithms using BVFT-PE is also of independent interest. To this end, we conduct additional
307 experiments in Mujoco to test the OPE accuracy of FQE with different hyperparameters. As Figure 5R
308 shows, FQE with hyperparameters tuned by BVFT-PE consistently outperforms the *best fixed* set
309 of hyperparameters across all 3 Mujoco domains. This implies that BVFT-PE can select the best
310 hyperparameters for each policy individually, which is an appealing property. We also supplement the
311 empirical results with a theoretical guarantee for selecting Q^π out of candidate functions for a fixed
312 π , which follows from similar proof techniques as [XJ20, Theorem 2]; see Appendix B for the proof.

313 **Theorem 4.** *Let C be the same as in [XJ20, Theorem 2], which characterizes the exploratoriness
314 of the data distribution. Consider any policy π and candidate Q -functions, $\{Q^l\}_{l=1}^L$, with $Q^l \in$
315 $[0, \frac{R_{\max}}{1-\gamma}]$. Let (π, \hat{Q}) be the pair that minimizes the empirical version of BVFT-PE-loss applied to
316 $\{(\pi, Q^l)\}_{l=1}^L$ with $\epsilon_{dct} = \frac{\epsilon R_{\max}}{8\sqrt{C}}$. Then, using a dataset of size $\tilde{O}\left(\frac{C^2 \ln(L/\delta)}{\epsilon^4 (1-\gamma)^4}\right)$ where \tilde{O} suppresses
317 logarithmic terms, w.p. $\geq 1 - \delta$, $\sup_{\nu: \|\nu/\mu\|_\infty \leq C} \|\hat{Q} - Q^\pi\|_{2,\nu} \leq \epsilon \cdot \frac{R_{\max}}{1-\gamma} + \frac{(2+4\sqrt{C}) \min_l \|Q^l - Q^\pi\|_\infty}{1-\gamma}$.*

318 As the theorem implies, when one of $\{Q^l\}_{l=1}^L$ is a good approximation of Q^π (i.e., $\min_l \|Q^l - Q^\pi\|_\infty$
319 is small), BVFT-PE is able to identify $\hat{Q} \approx Q^\pi$ with a polynomial sample complexity.

320 6 Discussion and Conclusion

321 We present BVFT and its variants based on recent theoretical advances [XJ20], which are (nearly)
322 hyperparameter-free algorithms for policy selection in offline RL and empirically effective in discrete-
323 action benchmarks. When combined with OPE algorithms such as FQE, variants of BVFT are also
324 competitive in continuous-action benchmarks, and such a combination addresses the weaknesses
325 of BVFT (relying on the existence of good value functions among the candidates) and those of OPE
326 (having untunable hyperparameters) and gets the best of both worlds.

327 We conclude the paper with discussions of the limitations of our approach and open questions:

328 **Sample efficiency** Section 5 uses OPE to fit Q^π , which can be data intensive. A plausible solution is
329 to re-use the training data for OPE, as suggested by [Pai+20]. However, data reuse can cause serious
330 issues in realworld domains where the amount of training data taken for granted in deep RL is not
331 available. How to improve the sample efficiency of BVFT + OPE is an important open question.

332 **Computational complexity** The $O(m^2)$ complexity of BVFT can be prohibitive if we wish to
333 compare hundreds of models. A plausible solution is divide-and-conquer, i.e., eliminating bad models
334 by running BVFT in smaller groups. It will be interesting to see if this compromises performance.

335 **Applicability** It is important to understand what type of domains BVFT is particularly suited for.
336 Despite having provided partial answers (e.g., the discussion of discrete actions vs. continuous
337 actions), we still need a more thorough answer based on comprehensive evaluation across more
338 diverse domains and comparison to a wider range of baselines, which is beyond the scope of this
339 paper since we focus on algorithm development and some of our contributions are orthogonal to
340 existing approaches (e.g., BVFT + OPE). We look forward to investigating this question empirically in
341 the future with the help of the public benchmarks that have become available very recently [Fu+21].

342 **Data with insufficient coverage** While results in Appendix D show that BVFT is insensitive to
343 moderate changes in data exploratoriness, we will likely need to incorporate some form of pessimism—
344 which is shown to be important for training [KRNJ20; JYW20; RZMJR21]—when the data coverage
345 is seriously lacking. This can be, however, quite challenging, as BVFT uses a dynamic function space
346 for projection ($\mathcal{G}_{i,j}$) that varies from candidate to candidate, and states considered covered due to
347 generalization effects under one function space may be considered lacking data under a different one.
348 How to resolve this issue and design a pessimistic version of BVFT is an interesting question.

349 **References**

- 350 [ASM08] András Antos, Csaba Szepesvári, and Rémi Munos. “Learning near-optimal policies
 351 with Bellman-residual minimization based fitted policy iteration and a single sample
 352 path”. In: *Machine Learning* 71.1 (2008), pp. 89–129.
- 353 [Bai95] Leemon Baird. “Residual algorithms: Reinforcement learning with function approxi-
 354 mation”. In: *Machine Learning Proceedings 1995*. Elsevier, 1995, pp. 30–37.
- 355 [BNVB13] Marc G Bellemare, Yavar Naddaf, Joel Veness, and Michael Bowling. “The arcade
 356 learning environment: An evaluation platform for general agents”. In: *Journal of
 357 Artificial Intelligence Research* 47 (2013), pp. 253–279.
- 358 [Bro+16] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman,
 359 Jie Tang, and Wojciech Zaremba. *OpenAI Gym*. 2016. eprint: arXiv:1606.01540.
- 360 [CJ19] Jinglin Chen and Nan Jiang. “Information-Theoretic Considerations in Batch Re-
 361inforcement Learning”. In: *Proceedings of the 36th International Conference on
 362 Machine Learning*. 2019, pp. 1042–1051.
- 363 [Die00] Thomas G Dietterich. “Hierarchical reinforcement learning with the MAXQ value
 364 function decomposition”. In: *Journal of artificial intelligence research* 13 (2000),
 365 pp. 227–303.
- 366 [FCGP19] Scott Fujimoto, Edoardo Conti, Mohammad Ghavamzadeh, and Joelle Pineau. *Bench-
 367 marking Batch Deep Reinforcement Learning Algorithms*. 2019. arXiv: 1910.01708
 368 [`cs.LG`].
- 369 [FKNTL21] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. *D4RL:
 370 Datasets for Deep Data-Driven Reinforcement Learning*. 2021. arXiv: 2004.07219
 371 [`cs.LG`].
- 372 [FLL19] Yihao Feng, Lihong Li, and Qiang Liu. “A kernel loss for solving the bellman
 373 equation”. In: *Advances in Neural Information Processing Systems*. 2019, pp. 15430–
 374 15441.
- 375 [FMP19] Scott Fujimoto, David Meger, and Doina Precup. *Off-Policy Deep Reinforcement
 376 Learning without Exploration*. 2019. arXiv: 1812.02900 [`cs.LG`].
- 377 [FS11] Amir-massoud Farahmand and Csaba Szepesvári. “Model selection in reinforcement
 378 learning”. In: *Machine learning* 85.3 (2011), pp. 299–332.
- 379 [FSM10] Amir-massoud Farahmand, Csaba Szepesvári, and Rémi Munos. “Error Propagation
 380 for Approximate Policy and Value Iteration”. In: *Advances in Neural Information
 381 Processing Systems*. 2010, pp. 568–576.
- 382 [Fu+21] Justin Fu, Mohammad Norouzi, Ofir Nachum, George Tucker, Ziyu Wang, Alexander
 383 Novikov, Mengjiao Yang, Michael R Zhang, Yutian Chen, Aviral Kumar, et al.
 384 “Benchmarks for Deep Off-Policy Evaluation”. In: *arXiv preprint arXiv:2103.16596*
 385 (2021).
- 386 [GGMVS20] Diksha Garg, Priyanka Gupta, Pankaj Malhotra, Lovekesh Vig, and Gautam Shroff.
 387 “Batch-Constrained Distributional Reinforcement Learning for Session-based Recom-
 388 mendation”. In: *arXiv preprint arXiv:2012.08984* (2020).
- 389 [Gor95] Geoffrey J Gordon. “Stable function approximation in dynamic programming”. In:
 390 *Proceedings of the twelfth international conference on machine learning*. 1995,
 391 pp. 261–268.
- 392 [Gul+21] Caglar Gulcehre, Ziyu Wang, Alexander Novikov, Tom Le Paine, Sergio Gomez
 393 Colmenarejo, Konrad Zolna, Rishabh Agarwal, Josh Merel, Daniel Mankowitz, Cos-
 394 min Paduraru, Gabriel Dulac-Arnold, Jerry Li, Mohammad Norouzi, Matt Hoffman,
 395 Ofir Nachum, George Tucker, Nicolas Heess, and Nando de Freitas. *RL Unplugged: A
 396 Suite of Benchmarks for Offline Reinforcement Learning*. 2021. arXiv: 2006.13888
 397 [`cs.LG`].
- 398 [JL16] Nan Jiang and Lihong Li. “Doubly robust off-policy value evaluation for reinfor-
 399 cements learning”. In: *International Conference on Machine Learning*. PMLR. 2016,
 400 pp. 652–661.
- 401 [JYW20] Ying Jin, Zhuoran Yang, and Zhaoran Wang. “Is Pessimism Provably Efficient for
 402 Offline RL?” In: *arXiv preprint arXiv:2012.15085* (2020).

- 403 [KRNJ20] Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten
 404 Joachims. “Morel: Model-based offline reinforcement learning”. In: *arXiv preprint*
 405 *arXiv:2005.05951* (2020).
- 406 [KZTL20] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. “Conservative q-
 407 learning for offline reinforcement learning”. In: *arXiv preprint arXiv:2006.04779*
 408 (2020).
- 409 [Lil+19] Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez,
 410 Yuval Tassa, David Silver, and Daan Wierstra. *Continuous control with deep rein-*
 411 *forcement learning*. 2019. arXiv: 1509.02971 [cs.LG].
- 412 [LKTF20] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. “Offline reinforcement
 413 learning: Tutorial, review, and perspectives on open problems”. In: *arXiv preprint*
 414 *arXiv:2005.01643* (2020).
- 415 [LLTZ18] Qiang Liu, Lihong Li, Ziyang Tang, and Dengyong Zhou. “Breaking the curse of
 416 horizon: Infinite-horizon off-policy estimation”. In: *Advances in Neural Information*
 417 *Processing Systems*. 2018, pp. 5361–5371.
- 418 [LMS15] Lihong Li, Rémi Munos, and Csaba Szepesvári. “Toward Minimax Off-policy Value
 419 Estimation”. In: *Proceedings of the 18th International Conference on Artificial Intel-*
 420 *ligence and Statistics*. 2015.
- 421 [LVY19] Hoang M. Le, Cameron Voloshin, and Yisong Yue. *Batch Policy Learning under*
 422 *Constraints*. 2019. arXiv: 1903.08738 [cs.LG].
- 423 [LWL06] Lihong Li, Thomas J Walsh, and Michael L Littman. “Towards a unified theory of
 424 state abstraction for MDPs”. In: *Proceedings of the 9th International Symposium on*
 425 *Artificial Intelligence and Mathematics*. 2006, pp. 531–539.
- 426 [Mni+15] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness,
 427 Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg
 428 Ostrovski, et al. “Human-level control through deep reinforcement learning”. In:
 429 *nature* 518.7540 (2015), pp. 529–533.
- 430 [NCDL19] Ofir Nachum, Yinlam Chow, Bo Dai, and Lihong Li. *DualDICE: Behavior-Agnostic*
 431 *Estimation of Discounted Stationary Distribution Corrections*. 2019. arXiv: 1906.
 432 04733 [cs.LG].
- 433 [Pai+20] Tom Le Paine, Cosmin Paduraru, Andrea Michi, Caglar Gulcehre, Konrad Zolna,
 434 Alexander Novikov, Ziyu Wang, and Nando de Freitas. “Hyperparameter selection
 435 for offline reinforcement learning”. In: *arXiv preprint arXiv:2007.09055* (2020).
- 436 [RZMJR21] Paria Rashidinejad, Banghua Zhu, Cong Ma, Jiantao Jiao, and Stuart Russell. “Bridg-
 437 ing Offline Reinforcement Learning and Imitation Learning: A Tale of Pessimism”.
 438 In: *arXiv preprint arXiv:2103.12021* (2021).
- 439 [SB18] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*.
 440 MIT press, 2018.
- 441 [TET12] Emanuel Todorov, Tom Erez, and Yuval Tassa. “MuJoCo: A physics engine for
 442 model-based control”. In: *2012 IEEE/RSJ International Conference on Intelligent*
 443 *Robots and Systems*. 2012, pp. 5026–5033. DOI: 10.1109/IROS.2012.6386109.
- 444 [UHJ20] Masatoshi Uehara, Jiawei Huang, and Nan Jiang. “Minimax weight and q-function
 445 learning for off-policy evaluation”. In: *International Conference on Machine Learning*.
 446 PMLR. 2020, pp. 9659–9668.
- 447 [UIJKSX21] Masatoshi Uehara, Masaaki Imaizumi, Nan Jiang, Nathan Kallus, Wen Sun,
 448 and Tengyang Xie. “Finite sample analysis of minimax offline reinforcement
 449 learning: Completeness, fast rates and first-order efficiency”. In: *arXiv preprint*
 450 *arXiv:2102.02981* (2021).
- 451 [VLJY19] Cameron Voloshin, Hoang M Le, Nan Jiang, and Yisong Yue. “Empirical study
 452 of off-policy policy evaluation for reinforcement learning”. In: *arXiv preprint*
 453 *arXiv:1911.06854* (2019).
- 454 [XJ20] Tengyang Xie and Nan Jiang. *Batch Value-function Approximation with Only Realiz-*
 455 *ability*. 2020. arXiv: 2008.04990 [cs.LG].
- 456 [YDNTS20] Mengjiao Yang, Bo Dai, Ofir Nachum, George Tucker, and Dale Schuurmans. “Offline
 457 Policy Selection under Uncertainty”. In: *arXiv preprint arXiv:2012.06919* (2020).

458 **Checklist**

- 459 1. For all authors...
- 460 (a) Do the main claims made in the abstract and introduction accurately reflect the paper's
461 contributions and scope? **[Yes]**
- 462 (b) Did you describe the limitations of your work? **[Yes]** See Section 6.
- 463 (c) Did you discuss any potential negative societal impacts of your work? **[No]**
- 464 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
465 them? **[Yes]**
- 466 2. If you are including theoretical results...
- 467 (a) Did you state the full set of assumptions of all theoretical results? **[Yes]**
- 468 (b) Did you include complete proofs of all theoretical results? **[No]** We provide complete
469 proofs of Proposition 3 and Theorem 4, which are new results. Earlier propositions are
470 simplifications of existing theoretical results to help readers understand the theoretical
471 intuitions, and we provide proof sketches and point the readers to the relevant literature
472 for the full versions of the theorems and proofs.
- 473 3. If you ran experiments...
- 474 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
475 mental results (either in the supplemental material or as a URL)? **[Yes]**
- 476 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
477 were chosen)? **[Yes]**
- 478 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
479 ments multiple times)? **[Yes]**
- 480 (d) Did you include the total amount of compute and the type of resources used (e.g., type
481 of GPUs, internal cluster, or cloud provider)? **[No]** We briefly mention computational
482 considerations of our methods in Section 6. Also, a large part of the computational
483 resources go into the training algorithms, which produces the candidate models our
484 algorithms will select from and is part of the experiment setup, which we did not keep
485 track of.
- 486 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 487 (a) If your work uses existing assets, did you cite the creators? **[Yes]**
- 488 (b) Did you mention the license of the assets? **[No]**
- 489 (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]**
- 490 (d) Did you discuss whether and how consent was obtained from people whose data you're
491 using/curating? **[N/A]**
- 492 (e) Did you discuss whether the data you are using/curating contains personally identifiable
493 information or offensive content? **[N/A]**
- 494 5. If you used crowdsourcing or conducted research with human subjects...
- 495 (a) Did you include the full text of instructions given to participants and screenshots, if
496 applicable? **[N/A]**
- 497 (b) Did you describe any potential participant risks, with links to Institutional Review
498 Board (IRB) approvals, if applicable? **[N/A]**
- 499 (c) Did you include the estimated hourly wage paid to participants and the total amount
500 spent on participant compensation? **[N/A]**