

1 We would like to thank the reviewers for their insightful comments. As has been noted by all four reviewers, our paper
2 presents a novel algorithm for batch reinforcement learning that is not only simpler and faster than the state-of-the-art
3 algorithms, but also performs significantly better. Because the BCQ and BEAR code are publicly available, we were
4 able to make a careful and comprehensive comparison of the performance of BAIL, BCQ, BEAR, MARWIL and naive
5 behavioral cloning using the MuJoCo benchmark. For our experiments, we created non-expert training batches in a
6 manner identical to what was done in the BCQ paper and included additional partially-trained training batches for the
7 environments Ant and Humanoid using SAC. Our experiments showed that BAIL wins for 20 of the 22 batches, with
8 overall performance 42% or more higher than the other algorithms. Moreover, BAIL is computationally 30-50 times
9 faster than BCQ and BEAR. We provide robust anonymized code for reproducibility. We will also make our datasets
10 publicly available for future benchmarking.

11 BAIL learns a value function by training a neural network to obtain the “upper envelope of the data”. To the best of our
12 knowledge, the notion of the upper envelope of a dataset is novel, and can possibly be applied to other RL problems in
13 the future.

14 **Response to Reviewer 2:** Thank you for your positive feedback with an accept recommendation. In terms of testing
15 on alternative domains, we are currently focused on MuJoCo, where Ant and Humanoid are the most challenging
16 environments. But we will consider testing on alternative benchmarks in future work.

17 **Response to Reviewer 4:** Thank you for your positive feedback with an accept recommendation. We will add to the
18 paper more intuition on why BAIL performs better than the other methods. Intuitively, BAIL performs better than BCQ
19 and BEAR because BCQ and BEAR rely on carefully tuned policy constraints to prevent the use of out-of-distribution
20 actions. A loose constraint can cause extrapolation error to accumulate quickly, and a tight constraint will prevent the
21 policy from choosing some of the good actions. BAIL, however, identifies and imitates the highest-performing actions
22 in the dataset, thus avoiding the need to carefully tune such a constraint. In terms of the computation, for BCQ and
23 BEAR, we used the authors’ implementations. These Q-learning based algorithms typically have more, larger networks,
24 and sample multiple candidate actions for each update. Thus it is not surprising that BAIL can be 30-50 times faster.

25 **Response to Reviewer 3:** Thank you for your mostly positive review, and for pointing out the papers (Oh et al,
26 2018) and (Vinyals et al, 2019), which we will add to the related work section. In our view, all DRL algorithms
27 are heuristics, and performance guarantees for schemes using neural-network function-approximators are rare. We
28 remark the experimental results presented in the main body indeed use “sub-optimal training data,” since those batches
29 were obtained from training data well before optimal performance is achieved (and in many cases using sub-optimal
30 algorithms for training). We will make this more clear in the revision. The results for the execution batches are
31 summarized in the main body and presented in detail in the appendix.

32 We decided to use L_2 regularization in the definition of the upper-envelope since it leads to a clean definition and
33 theory. In our computational experiments, we found L_2 regularization and early-stopping regularization to give similar
34 performance, with early stopping being faster. So we chose to use early stopping in our experiments. We note that there
35 is a well-known theory showing, under some conditions, the equivalence of early stopping to L_2 regularization (see the
36 Deep Learning book by Goodfellow et al., Section 7.8).

37 **Response to Reviewer 1:** Thank you for the feedback. It seems you are very concerned about our table results. We
38 would like to point out that they are definitely not “misleading”. In fact, the table is presented in this manner to make
39 fair comparisons, as clearly explained in the paper. On page 8, line 277, we explain: “For each batch, all the algorithms
40 that are within 10% of the highest average value are considered winners and are indicated in bold.” So it is possible for
41 multiple algorithms to be in bold in a table row. Indeed, when two algorithms have very similar scores, this should
42 be considered a tie, since their rankings can change easily when using a different random seed. To complement the
43 tables, on page 8, line 283, we also explain how we computed the average improvement ratio, which shows that BAIL
44 performs 42% better than BCQ, and 101% better than BC.

45 You ask “how close to expert does the training data have to be for BAIL to perform competitively?” This question
46 is answered starting at page 2, line 39, where we explained that we generated a large variety of batches, including
47 far-from-expert datasets with different random noise levels and with different learning agents. Our batch generation
48 procedure is also consistent with the procedures in prior work. Thus Table 1 answers your question by showing that
49 BAIL performs well on all the training datasets with various noise levels.

50 For the execution batches, we used the same hyper-parameters for the training set, and the results show that BAIL is
51 better than BEAR, similar to BCQ, and slightly weaker than naive behavioral cloning. If we change the hyper-parameters
52 we will get at least the same performance as behavioral cloning (which is the strongest in the execution batches). Thus,
53 in summary, BAIL is the clear winner for the training batches, and BAIL, BCQ, and naive behavioral cloning are
54 equally good for the execution batches. For the other minor points you suggested, thank you for pointing them out, and
55 we will modify them in our revision.