1 We would like to thank the reviewers for your thoughtful feedback and comments which would undoubtedly make the
2 paper better. We will update our paper to reflect your comments, fix typos and include missing references. Here, we
3 aim to address the concerns of each reviewer.

4 **R4 & R5 – main contribution / insights:** Our goal was to devise an effective algorithm that works well on a diverse
5 set of offline RL problems. We consider this to be the main contribution of our paper: a simple algorithm which achieves
6 good results on a large variety of environments/dataset, including challenging high-dimensional, partially observed
7 problems from pixels. Although some components of our algorithm have been considered in related recent or concurrent
8 work, this has usually been the case in the context of more complicated formulations. A secondary contribution of our
9 work is thus a careful analysis and ablation that sheds light on the important attributes, and to identify a particularly
10 simple but effective combination. We will update the paper to make this more overt.

11 **R2 & R3 – motivation behind Eq. 3 and 4**: One intuition for Eq. 4 (as explained in the subsection titled regularized
12 policy iteration) works as follows. Through exponential weighting, we approximately regularize the policy by how
13 much it diverges from the behavioral policy as measured by KL. The regularization would constrain the policy to not
14 deviate too much from the support of the dataset thus ensuring effective offline learning. Eq. 4 is therefore chosen
15 because of its connection to KL divergence. Both Eq. 3 and 4 are motivated by the policy improvement theorem.
16 Whereas Eq. 3 seeks to improve the policy by choosing a better action to copy, Eq. 4 does this in a soft manner. We
17 chose Eq. 3 because of its direct link to the policy improvement theorem. There indeed are many other choices of
18 filtering functions and finding the best is an interesting direction of future research.

19 **R2 – reproducibility:** We have open-sourced the code for CRR on Github and the link will be made available.

20 **R2 – network architecture:** Due to the complexity of the locomotion datasets, we use relatively large networks. For
21 consistency reasons, we adopt the same network architecture for the control suite dataset though smaller networks
22 would suffice. We contend that other than being large in size, the networks used are not very unconventional in the
23 RL community. The convolutional nets are actually borrowed from Impala and ACME where we only modified the
24 activation functions. The residual MLPs are chosen for their expressive power.

25 **R3 – bootstrapping:** In the offline setting, the problem with bootstrapping arises when the current policy selects
26 actions not in the support of the dataset. When this happens, the bootstrapping target can be inaccurate due to the lack
27 of data. Since CRR's policies are learned by copying the actions of the dataset, the actions selected by CRR would
28 be close to that of the dataset thereby making the problem much less severe. The fact that CRR does better than BC
29 supports the claim the Q function is learned well. D4PG's poor performance (despite the same value learning rule)
30 suggests naive strategies in learning the policy is not sufficient since it is not constrained to the support of the dataset.

31 **R3 – weighting functions:** It is indeed the case that the two weighting functions produce different results. This finding
32 makes it clear that not all weighting functions are created equal and thus opens new doors to future research. As offline
33 RL algorithms cannot escape the existence of hyper-parameters, the choice of weighting functions is treated as one
34 additional hyper-parameter.

35 **R4 – advantage estimation:** We agree that adding constants which depend only on $s_t$ to $f$ in Eq. 2 (i.e. $f + g(s_t)$)
36 does not make a difference. However, $f$ is defined as a nonlinear transformation of the advantage estimates. $A_{max}$
37 underestimates the advantage where as $A_{mean}$ does not. Thus, using $A_{mean}$ vs. $A_{max}$ leads to different outcomes.

38 **R4 – convergence of Q:** Our algorithm performs standard policy evaluation to learn $Q$ for the current policy, effectively
39 as an inner loop (Eq. 1). Thus, if our algorithm converges to a pair $(\pi, Q)$, $Q$ will estimate the value of the policy, i.e.
40 $Q^\pi$ (restricted to the data in the batch and subject to the usual caveats of RL with nonlinear function approximation). In
41 the tabular setting, under some mild regularity conditions, CRR constitutes policy improvements in a restricted MDP
42 similar to that defined in BCQ [10]. Therefore each iteration of CRR will result in a better policy. Via similar arguments
43 for the convergence of policy iteration, $Q$ of CRR would converge to one corresponding to an optimal policy for the
44 restricted MDP. The theoretical results would be included in an updated version of the paper.

45 **R4 – additional baselines:** Thanks for raising the question. We would like to point out that our experiments already
46 include several baselines that are considered to be the state of the art at the moment of writing. In our experiments, we
47 also examine elements of the baselines and compare it to the equivalent bits in our method. For example, in Appendix
48 A.2 we evaluate the effect of K-step returns used by AWR and MARWIL.

49 **R4 – novelty:** Our main contribution is a simple and effective algorithm that identifies important ingredients of offline
50 learning algorithms (see response "Main contribution" above). We do not claim that the use of distributional Q learning
51 is a contribution per se. We will make this clearer in the text. CRR uses a different advantages estimate compared to
52 [35]. We include a toy example (in subsection titled "CRR vs. return-based methods.") that shows how the approach
53 of [35] can fail. We also show empirically (in appendix A.2) that our advantage estimation methods significantly
54 outperforms that adopted by [35].