

1 We thank the reviewers for their comments and insightful reviews.

2 === R1 ===

3 **Best approximator in max-norm.** Prior literature (e.g., Remi Munos, "Error Bounds for Approximate Value  
4 Iteration", AAAI-05, Remi Munos, Csaba Szepesvari, *Finite-Time Bounds for Fitted Value Iteration*, JMLR 2008)  
5 studied errors measured in weighted p-norm. This approach is often preferable when using "standard" parametric  
6 regression algorithms, as the weighted p-norm can be directly minimized at learning time (e.g., in least-squares  
7 regression, the  $\ell_2$ -norm is minimized). Nonetheless, this gives rise to the so-called concentrability coefficients, which  
8 may be very large (even unbounded). In our case, the interpolation algorithm directly controls the max-norm (see e.g.,  
9 Lemma 1) and this allows us to derive the analysis directly in max-norm and avoid concentrability coefficients.

10 === R2 ===

11 **Trade-off amplification and anchor points.** In the worst case even the best trade-off between  $\bar{C}$  and  $K$  may indeed  
12 lead to an exponential complexity. This is not surprising as in the worst case the inherent Bellman error may  
13 be unbounded and standard AVI tends to diverge. Recent work (Jinglin Chen, Nan Jiang, *Information-Theoretic*  
14 *Considerations in Batch Reinforcement Learning*, ICML 2019, Conjecture 8) has even conjectured an exponential  
15 lower bound in case of unbounded Bellman error. As a consequence, there might exist a fundamental barrier to  
16 obtaining polynomial sample complexity in the worst case. Nonetheless, in many other cases a good trade-off between  
17 amplification and anchor points may correspond to a much smaller sample complexity (e.g., the condition in Prop.1  
18 could be achieved by a polynomial number of anchor points). Indeed in our experiments even in the case of unbounded  
19 Bellman error, the extrapolation can be controlled, and thus we can obtain satisfactory solutions by slightly increasing  
20 the number of anchor points without requiring an exponential number of them (see e.g., the experiments on the Tsitsiklis  
21 and Van Roy domain and the linear bandit).

22 **Construction of anchor points.** In the paper we propose a first heuristic algorithm to automatically construct a set  
23 of anchor points (see beginning of page 6). While we do not have any guarantee for the method, in our preliminary  
24 experiments it seems effective in building a compact set with small amplification factor.

25 **Incremental construction from  $H$  down to 1.** A good choice for the anchor points depends only on the linear  
26 architecture, i.e., it is computed exclusively on the basis of the feature map  $\phi_t(\cdot, \cdot)$ . Thus if one already knows a good  
27 set of anchor points at the final timestep  $t = H$  and the feature map  $\phi_t(\cdot, \cdot)$  does not vary a lot for  $t = 1, \dots, H$  then  
28 the set of anchor points for the timestep  $t = H$  is also a good choice for prior timesteps  $t = 1, \dots, H - 1$ .

29 **Infinite horizon.** The algorithm does not need any modification (other than the addition of the discount factor) to deal  
30 with the infinite horizon case, with the additional benefit that the identification of the support points can be done once at  
31 the beginning as opposed to every timestep (as it only depends on the feature map  $\phi(\cdot, \cdot)$ ). However, one would need to  
32 change the analysis. Note that Yang and Wang's analysis operates in the much easier setting of zero Bellman error.

33 **Continuous state space.** The main algorithm can be applied to the continuous state space case without any modification,  
34 and the choice of discrete MDPs in the experiments is for illustrative purposes.

35 === R3 ===

36 **Computational complexity of the heuristic.** In its most naive form and without further structure, one would need to  
37 loop through the state action pairs; for large or continuous state-actions spaces one should sample the state-action pairs  
38 to reduce the complexity. This should suffice as we only require approximate convex hulls. For computing the  $\theta$ 's, this  
39 is a linear program and the best known computational complexity can be found in the Arxiv paper "Solving Linear  
40 Programs in the Current Matrix Multiplication Time"; other methods, like interior point methods, may be used.

41 **Experiments.** We can report more thorough statistics for the anchor points (i.e., number and positions, and resulting  
42 extrapolation coefficients) as those are computed inside the program. Some of these are reported in figure 3 for that  
43 example, but we can add them for the other examples as well.

44 **Real Life Experiments** As the reviewers points out, we have chosen examples where realizability holds and the  
45 choice was deliberate to reduce the number of confounding factors: the approximation error may in general affect the  
46 comparison between methods, and might obscure the key underlying processes. We agree with the reviewer that a  
47 comparison with averagers like k-NN or kernel based would be practically interesting, but it raises questions of how  
48 to best define the function class (for example, the position of the points in a nearest neighbour procedure, the type of  
49 radial basis function, etc.) encoded by those settings as a fair comparison to our setting. Therefore in our work we focus  
50 on fixing the function approximation class (general linear value functions) and evaluate the impact of the algorithm  
51 used to fit the function class. A key benefit of our approach is that it does not modify the underlying (linear) feature  
52 representation, allowing the user to use the linear representation with, for example, approximate value iteration, and  
53 should this fail, the user can switch to our algorithm and progressively increase the number of support points while  
54 keeping the same feature representation. Finally, we can easily add more variations of the chain experiments with  
55 varying horizon (for example  $H = 50, 100, 200$ ) and feature representation as noted by the reviewer.