

1 We thank the reviewers for thoughtful reviews and encouraging comments. We respond only to questions and concerns.

2 (R1) “In the introduction ...”: Good suggestion. Will refer specifically to CDRD, CDRK in Table 1a. The point is  
3 extending the cited importance sampling methods to RL would lead to convergence rate that *deteriorates* in horizon.

4 (R1) “Something is said in Remark 2 ...”: Our analysis handles known and unknown behavior policy and simply take  
5 as a condition the nuisance estimation rate. Knowing the behavior policy can help estimate nuisances. In experiments  
6 we consider known behavior policy as is common in offline RL. We still need to estimate  $w_t^K$  even if behavior is known.

7 (R1) “Out of the 9 combinations, ...”: We focused on these as they represent the two extremes. We can easily provide in  
8 the supplement a general analysis of all combinations under the intersection of the assumptions need for each extreme.

9 (R1) “In the beginning of section 4, ...”: Unfortunately, no, as smoothness conditions are necessary in continuous  
10 action space, else the finite data may be unrepresentative of the infinite possible unseen actions. This is the same as in  
11 density estimation. We will comment on this.

12 (R1) “Although half of the paper ...”: With limited space, we thought actual *learning* would be of greatest interest. We  
13 will follow your suggestion and add to the supplement experiments for policy evaluation and varying  $H$ . Note this is  
14 easy with the submitted code – we will just run it and report the results.

15 (R1) “The paper is well written ... In general, it would be interesting ...”: The derivation is different and more  
16 complicated than [3] as in [3] the density ratio exists and is used directly. On the other hand we need to analyze the  
17 errors due to the kernelization, which complicates the analysis as it introduces slower leading terms with rate that  
18 depends on horizon dimensionality.

19 (R1) “Typos: ... isn’t the inverse?”: Thanks for catching. Yes; the latter is a typo; it is the reverse.

20 (R2) “The experimental setting ...”: The qualitative results are the same as we vary these. We will run the (submitted)  
21 code with a range of parameters and include additional plots in supplement.

22 (R2) “The overall structure are well planned and the paper is well written ...”: We will add reminders of notation when  
23 used for first time much after problem setup and add short descriptions of steps in equations in appendix.

24 (R2) “Line 72–73 ...”: Yes; it’s a typo; “latter” and “former” should be exchanged.

25 (R2) “Line 81–82 ...”: By Radon-Nikodym thm, exists  $f$  such that  $\mathbb{E}_{\pi_e}[g(a) \mid s] = \mathbb{E}_{\pi_b}[f(a)g(a) \mid s]$  for all  $g$   
26 measurable if and only if  $\pi_e(\cdot \mid s)$  is absolutely continuous wrt  $\pi_b(\cdot \mid s)$  (this is for each  $s$ ). However, if  $\pi_e(\cdot \mid s)$   
27 is discrete, it is *not* enough for behavior to have positive density at its atoms. E.g., Dirac at 0.5 is *not* abs cts wrt the  
28 uniform distribution on  $[0, 1]$  even though latter has density 1 at 0.5. Recall  $\mu \ll \nu$  means ( $\mu(A) > 0 \Rightarrow \nu(A) >$   
29  $0, \forall A$  measurable), so if  $\mu \ll \nu$  and  $\mu(\{0.5\}) = 1$  then  $\nu(\{0.5\}) > 0$ , i.e., has an atom at 0.5. Will add this example.

30 (R3) “The novelty of ...”: We respectfully disagree. Not only do we make it doubly robust *and* extend it to RL, we also  
31 analyze it and give rates under lax conditions and show the naïve extension yields very bad rates as horizon grows.

32 (R3) “The evaluation is ...”: Indeed while we avoid curse of dimension in state space and horizon, we may suffer  
33 from it in action space, since we kernelize actions. In many practical offline RL settings, however, states are complex  
34 (e.g., many+rich health indicators) and actions simple and often one-dimensional (e.g., insulin dosing/timing). We will  
35 comment on this and run the (submitted) code on growing action dimension and add to supplement to visualize this.

36 (R4) “The expected return ...”: In OPE, MSE is actually the metric of interest. In offline learning, policy value is  
37 indeed of interest and our policy gradient experiments (Sec 5) showcase how low-error gradients lead to high-value  
38 learning. Moreover, as Remark 8 mentions, error bound on gradients can be combined with standard gradient ascent  
39 analysis to get value regret guarantees: we simply replace the stochastic-policy gradient error bounds of Kallus &  
40 Uehara ’20 with our new ones in Thms 11–13 therein; while this is straightforward use of existing work, we’ll flesh this  
41 out more explicitly in supplement for completeness.

42 (R4) “The only concern ...”: The primary contribution is a theoretical study of rates and the simple experiments are  
43 intended only to illustrate the new theory. Extensive experimentation is beyond the scope of such a short paper with so  
44 many new results already. We can nonetheless easily run our (submitted) code on the Warfarin dosing experiment of  
45 Kallus & Zhou ’18 (the code is also public) and add to the supplement.

46 (R4) “The paper is well-organized and well-written ... more intuition explanations ... explain more on the comparison  
47 results in Table 1”: Thank you; we will use this feedback to further improve the clarity. We will add *more* in-words  
48 explanations of each result and we will move Table 1.