

1 We thank the reviewers for their helpful comments and enthusiasm.

2 **Reviewer #2:** Thank you for your comments. Regarding your points:

3 **1. Motivation:** We would like to clarify our aims: our main goal is to study the connections and properties of MMD
4 DRO, to reveal its potential benefits, and aid a better general understanding of the DRO landscape with different
5 divergence measures.

6 We do not claim that Wasserstein or phi-divergence uncertainty are bad or yield poor out-of-sample bounds; we merely
7 highlight differences between them and MMD DRO. For example, in discussing phi-divergences, we state in line 88
8 that one cannot obtain a generalization bound via Principle 2.1 alone (to be fair, we will add a reminder to the reader
9 that e.g. Namkoong and Duchi (2017) achieve generalization bounds via other means). For Wasserstein, we focus
10 in lines 34-36 and lines 89-99 on other complications, such as difficulty of optimization, and that many of the upper
11 bounds are only asymptotic. Indeed, the convergence results of Blanchet, Kang, and Murthy (2016) (in Section 3 of
12 their paper) are also asymptotic in nature.

13 **2. Upper bounds instead of exact reformulation, and assuming ℓ_f is in an RKHS:** We again emphasize that there
14 are tradeoffs between all three DRO approaches. While discarding the non-negativity constraint in MMD DRO may
15 weaken the bounds, in return we obtain a simple non-asymptotic upper bound (unlike Wasserstein).

16 Moreover, we contend that the assumption that ℓ_f is in an RKHS \mathcal{H} is not so restrictive after all. If the kernel k is
17 universal, as is the case for many kernels used in practice such as Gaussian and Laplace kernels, we can readily extend
18 our results to all bounded continuous functions as described below. We will add this clarification to the paper:

19 Suppose the loss ℓ_f of our predictor f is any bounded continuous function on a compact metric space \mathcal{X} . By definition
20 [Muandet et al., Definition 3.3] if k is a universal kernel on \mathcal{X} (associated with the RKHS \mathcal{H}), then for any $\epsilon > 0$, there
21 is some $\ell' \in \mathcal{H}$ with $\sup_{x \in \mathcal{X}} |\ell_f(x) - \ell'(x)| < \epsilon$. It follows that for any measure \mathbb{P} , we can bound the expectation of
22 $\ell_f(x)$ by that of ℓ' : $\mathbb{E}_{x \sim \mathbb{P}}[\ell_f(x)] < \mathbb{E}_{x \sim \mathbb{P}}[\ell'(x)] + \epsilon$. Then, we can apply our results to $\ell' \in \mathcal{H}$.

23 **3. MMD DRO is a more conservative upper bound (Theorem 5.1):** Separate from the task of producing a valid (and
24 hopefully tight) upper bound is the task of designing a regularizer that is practically useful. And stronger regularizers
25 are often better. One drawback of the previously considered chi-squared DRO/variance regularization is that the
26 regularization ceases to have any effect when the training data can be fit perfectly e.g. in deep learning (since then the
27 loss for each datapoint is zero, and so the variance of the loss on the dataset is also zero). In such a regime, stronger
28 penalties such as the RKHS norm continue to be meaningful.

29 **Reviewer #3:** Thank you for your support.

30 **Re: discrete approximation of MMD DRO uncertain set may not contain the population:** Yes, any such discrete
31 approximation can have similar issues. We present it mainly to link variance regularization and MMD DRO, e.g. as in
32 Theorem 5.1.

33 **Reviewer #4:** Thank you for your feedback and support.

34 Before addressing your main comments in detail, we emphasize that at a high level we hope to present MMD DRO
35 as an alternative worth studying, with complementary properties to existing techniques and rich connections. In that
36 context,

37 **(a) Wasserstein convergence with fewer assumptions:** Thank you for pointing us to the references on non-Euclidean
38 Wasserstein convergence; we are happy to mention them in the camera ready. Regarding your comment about assuming
39 ℓ_f is in an RKHS, please see point 2 of our response to Reviewer #2.

40 **(b) Faster convergence:** The point you make about norms cancelling with rates is fair. We mention Wasserstein's
41 $O(n^{-1/d})$ rate in the paper because it is relevant to the application of Principle 2.1. However, as discussed in point 1 of
42 our response to Reviewer #2, we don't mean the remark about $O(n^{-1/2})$ vs $O(n^{-1/d})$ to claim the MMD results were
43 always better. Instead, the different convergence properties of different distances motivates studying different DRO
44 formulations. We will edit the paper to make this clearer.