

1 We thank all the reviewers for their comments and suggestions.

2 **Reviewer #1:** Thank you for the positive review. We will make the corrections you have pointed out.

3 **Reviewer #2:** We will add discussion of these points to clarify the theoretical contributions of the paper:

- 4 1. Yes, Theorem 9 is non-adaptive (i.e., requires knowing  $\sigma_g$  to tune the wavelet threshold and generator network
- 5 size). In practice, GANs need extensive tuning to perform well, so constructing an “adaptive GAN” could be
- 6 useful, but it is not clear to us how to do so. For now, we leave this as future work.
- 7 2. Yes, Theorem 9 relies heavily on Theorem 5 and prior work on the approximation ability of fully-connected
- 8 ReLU networks. Hence, our main technical contributions are in Theorems 4, 5, and 7. While proving Theorem
- 9 9 is straightforward given these results, we nevertheless feel that it is not obvious (and is worth explicitly
- 10 sharing with the NeurIPS community) that these results have implications for GANs.
- 11 3. Yes, this is an important issue. We provide an example of an acknowledged challenge (spatial adaptivity) that
- 12 neural networks can overcome, distinguishing them from some established estimators.

13 **Reviewer #4: Paper is poorly written/hard to follow... notation is unclear/terminology is not defined.**

14 We carefully reread the sections the reviewer mentions (lines 45-48, Section 2.2, Section 5), and we were unable to find

15 the undefined notation. We would appreciate if the reviewer could specify which aspects of the organization can be

16 improved, and which notation/terminology can be clarified – we would be happy to make these improvements.

17 **formal problem statement is vague:** The formal problem is to lower and upper bound the general minimax rate

18  $M(B_{p_d, q_d}^{\sigma_d}, B_{p_g, q_g}^{\sigma_g})$  (defined in (2)), as well as on its linear counterpart (defined in (4)). We did not understand why the

19 reviewer found this vague and would appreciate if the reviewer could clarify this – we’ll be happy to clarify these.

20 **lines 125-156, authors never explain how exactly their results generalize all these known results:**

21 Each loss listed in lines 125-156 corresponds to a particular value (or set of values) of  $\sigma_d$ ,  $p_d$ , and  $q_d$ . Our results

22 (specifically, Theorems 4 and 5) “generalize” these known rates in that they hold simultaneously for many different

23 values of  $\sigma_d$ ,  $p_d$ , and  $q_d$ , including but not limited to those listed on lines 125-156. We’ll add this explanation to paper.

24 **...as if parts of the paper were taken directly from a textbook without a proper introduction of terminology.**

25 To our knowledge, convergence rates under all these losses have not been studied in a unified setting. Many of the

26 results cited here are from within the last year, and we hope this section can help consolidate this very recent work.

27 Due to space constraints, we omit some common alternative definitions of these losses (e.g., optimal-transport def. of

28 Wasserstein metric), but they are equivalently defined in terms of Besov spaces in the way we give (e.g., Wasserstein

29 metric is equivalent to  $d_{B_{\infty, \infty}^1}$ ). We include common alternative notation (e.g.,  $C^1(1)$  on line 132) without definition to

30 help relate to these defs with which the reader may be familiar, but emphasize that these alternative notations are *not*

31 needed to understand the paper. If reviewer finds this confusing, we can remove these notations or define in appendix.

32 **...need to clearly demonstrate relevance of new results for GANs and nonparametric density estimation.**

33 While the paper’s focus is theoretical, its results are relevant to both density estimation and GAN literatures:

34 **Relevance for density estimation:** Usually, density estimation is not performed in isolation, but rather as a sub-routine.

35 Hence, importance of our theory for density estimation is best seen in downstream applications. Two examples:

- 36 1. In distributionally robust optimization, a bound on the convergence rate of density estimation is directly used
- 37 to tune the optimizer (see, e.g., Esfahani & Kuhn (2015) or Staib & Jegelka (2019)).
- 38 2. Risk bounds for density estimation can be used to derive risk bounds for estimating simpler properties of
- 39 densities (such as smooth functionals; see, e.g., Kandasamy et. al (NIPS, 2015)).

40 Thus, our results might enable the creation of both new tools and new theoretical analyses, based on Besov IPMs.

41 **Relevance for GANs:** Our results are among the first finite-sample guarantees for GANs for a large family of

42 distributions, for which we show that well-optimized GANs are minimax optimal – results that were previously

43 unknown to the community. Even many simpler theoretical properties of GANs are unknown; e.g., weak consistency

44 was only recently studied (in [31]), without investigating convergence rates.

45 Besides establishing basic theoretical properties of GANs, some ways our work might contribute to GANs include:

- 46 1. Suggesting why GANs can perform density estimation well in high dimensions (by implicitly using a weak
- 47 loss, under which minimax rates might not be exponentially bad with dimension).
- 48 2. Suggesting why GANs can strictly outperform many (i.e., linear) classical density estimators.