

1 We thank the reviewers for the positive feedback and valuable comments. Our response follows.

2 **1) Computing smooth Wasserstein distance (SWD) [R3]:** This is a great and important point. We will add a discus-
3 sion about computation to the revision. As R3 mentions, the placeholder MC algorithm we used, in theory, has an
4 exponential in d runtime (though no major issues arose in practice). Our submission studies statistical aspects of SWD,
5 but we are actively working on efficient algorithms that compute it. To avoid sampling the kernel, we developed a
6 method to compute $W_1^{(\sigma)}(P_n, Q_m) = \sup_{f \in \text{Lip}_1} P_n(f * \varphi_\sigma) - Q_m(f * \varphi_\sigma)$ under NN parameterization of f , with the high-
7 dimensional convolution $f * \varphi_\sigma$ implemented in closed form. Accuracy and complexity of the alg. are being explored.

8 **2) Choice of SD [R3,R4]:** We will revise the presentation to thoroughly motivate SWD for inference. Briefly, SWD
9 inherits much of its structure from classic W_1 [27], while enjoying faster empirical convergence. [27] shows: (i) SWD
10 metrizes the same topology as W_1 ; (ii) it Γ -converges to W_1 as $\sigma \rightarrow 0$ (with convergence of optimizers); (iii) it is stable
11 in σ , $|W_1^{(\sigma)}(P, Q) - W_1^{(\tau)}(P, Q)| \leq 2\sqrt{d}|\sigma^2 - \tau^2|$; (iv) $W_1^{(\sigma)}$ is continuous and nonincreasing in σ . These suggest that
12 SWD preserves the compatibility of W_1 for inference, while having superior statistical properties. Re generalization,
13 Supp. A.1 states concentration inequalities for SWD, which imply generalization bounds for SMDE in the sense of [14]
14 (lines 248-255). SWD achieves fast $n^{-1/2}$ generalization, as opposed to the $n^{-1/d}$ generalization gap of W_1 .

15 **3) Effect/choice of σ [R2,R3,R4]:** Item 2 lists properties of $W_1^{(\sigma)}$ as a function of σ . From a statistical standpoint, the
16 expected empirical SWD scales as $\sigma^{-d/2}n^{-1/2}$ (see [26] and Thm. 1 herein). We are currently working on relaxing the
17 prefactor under the manifold hypothesis, i.e., showing that it depends only on the manifold’s (rather than the ambient)
18 dimension. R4 is correct that characterizing the dependence on σ in limit distribution results is nontrivial, and we
19 will comment on that in the revision, as suggested. In practice, having an excessively large σ may result in generative
20 models that are slow to learn small-scale structures in the target distribution (though in theory they should be recovered
21 eventually). That said, SWD is remarkably stable over moderate changes in σ and is thus not hard to tune.

22 **4) Sliced W_1 (SW_1) vs. $W_1^{(\sigma)}$ [R4]:** We will gladly expand the SW_1 discussion as follows. SW_1 and SWD are different
23 approaches for alleviating the curse of dimensionality (CoD) of W_1 . SW_1 uses 1D projections, while SWD levels out
24 local irregularities in the distributions by convolving with \mathcal{N}_σ . SW_1 and SWD share many similarities, e.g., both are
25 metrics and (topologically) equivalent to W_1 . SW_1 is easily computed using the 1D formula, while computational
26 aspects of SWD are still in the works (Item 1). SWD might be preferable when comparing to regular W_1 , as it is within an
27 additive $2\sigma\sqrt{d}$ gap from W_1 [27] (see Item 2). Comparison results for SW_1 seem weaker, assuming compact support and
28 involving implicit dimension-dependent constants (cf., e.g., Lem. 5.1.4. in <http://cvgmt.sns.it/paper/2341/>).

29 **5) Entropic OT (EOT) vs. $W_1^{(\sigma)}$ [R3]:** Generally, we do not view SWD and EOT as competing techniques. It is
30 beneficial for the community to have several methods for dealing with the CoD, especially since tradeoffs often emerge.
31 We agree that a thorough literature review of EOT is appropriate, and will cover the following in the revision: (i) EOT
32 is not a metric, while SWD retains the metric structure of W_1 ; (ii) $n^{-1/2}$ rate for EOT is proven only for smooth costs
33 (excluding W_1) with compactly supported distributions (arXiv:1810.02733) or squared cost with subgaussian distributions
34 (arXiv:1905.11882), while for SWD the $n^{-1/2}$ rate holds under mild polynomial moment conditions; (iii) EOT CLT
35 (arXiv:1905.11882) is similar to (arXiv:1705.01299) but markedly different from ours. Notably, they derive the result
36 for two-sample populations with the unknown centering constants $\mathbb{E}[S_\epsilon(P_n, Q)]$ or $\mathbb{E}[S_\epsilon(P_n, Q_m)]$, which differ from
37 $S_\epsilon(P, Q)$. (iv) A major virtue of EOT lies in fast algs., and matching those for SWD is a central goal going forward.

38 **6) $W_1^{(\sigma)}$ vs. MMD and $\sigma \rightarrow \infty$ limit [R2,R3]:** MMDs do not suffer from the CoD, due to low complexity of RKHS
39 function class. A main motivator for studying SWD is to alleviate the Wasserstein CoD, while preserving the metric
40 structure. In light of Item 2, we are thus interested in $\sigma < \infty$. Also, [Cor. 2.4, arXiv:2005.00738] shows that
41 $\lim_{\sigma \rightarrow \infty} W_1^{(\sigma)}(P, Q) = |\mathbb{E}_P[X] - \mathbb{E}_Q[Y]|$, which is not informative as a discrepancy measure between distributions.

42 **7) Weed-Berthet [R3]:** Their focus is different from ours: they study density estimation under W_p and do not deal
43 with limit distributions or MDE. Re computation, please see Item 1. Still, we are glad to cite this paper in the revision.

44 **8) Applying analysis to EOT [R3]:** Our proof technique to derive the asymptotic distribution relies on expressing
45 SWD *exactly* as the supremum of an empirical process indexed by smoothed Lip_1 functions. As EOT cannot be written
46 in this manner (for some popper function class), our proof technique does not directly extend to EOT.

47 **9) Fig. 1 [R2]:** The purpose of the point clouds is to show that the point clouds closely overlap. This implies that indeed
48 a limiting distribution is emerging. We will add additional description to make this much more clear in the revision.

49 **10) GAN method [R3], parametric experiments [R4], and appeal to ML audience [R6]:** When referring to GAN-
50 type methods, does R3 mean the generator’s parameterization or the adversarial training? The former is covered by the
51 MDE theory by viewing the weights/biases of the generator NN as θ (adversarial training aspects are indeed beyond the
52 scope of this work). Our last experiment in fact implements a SWD GAN: while the data comes from a parametric
53 distribution (GMM with $\exp(d)$ modes) we match it with a (generator) NN transformation of a latent variable. Per R6’s
54 suggestion, we will revise the presentation of this experiment to better communicate it to ML audiences. In general we
55 will put more emphasis on NN-based ML applications, in addition to our focus on parametric setups.