

1 We appreciate the valuable comments, which urged us to embody explicit connections to practices of learning. Apology
 2 that not all comments are replied here and our replies have to be short due to space, but they'll be fully addressed in a
 3 revision. We plead a reconsideration based on the improvement, as our contribution is truly innovative and nontrivial.

4 **Re: connection to learning, and when Cond.1&2 hold.** Here is an example (simplified and only briefly explained
 5 for length) in which the loss will be multiscale as considered in our paper: train a 2-layer neural network to fit data
 6 $\{x^k, y^k\}_k$, where the output $y^k = y_0^k + y_1^k + \xi^k$ admits a decomposition into large scale behavior $y_0^k = g_0(x^k)$,
 7 microscopic detail $y_1^k = \epsilon g_1(\epsilon x^k)$, and i.i.d. noise ξ^k . Assume g_0 and g_1 are regular enough so that universal
 8 approximation (UA) works and they can be approximated by wide enough neural networks with $\mathcal{O}(1)$ weights. Consider
 9 MSE loss $\sum_k \|y^k - \sum_i a_i \sigma(W_i x^k + b_i)\|^2$ with σ being the periodic activation in a recent progress [Implicit Neural
 10 Representations with Periodic Activation Functions, 2020]. Then there exists a minimizer and in its neighborhood the
 11 loss satisfies Cond.1 & 2: omit k WLOG, absorb bias into weight, and rewrite the loss as (denote by $\theta = [a_i, W_i]_i$)

$$f(\theta) = \left\| y_0 - \sum_{i \in I} a_i \sigma(W_i x) + \epsilon y_1 - \sum_{j \notin I} a_j \sigma(W_j x) \right\|^2 = \left\| g_0(x) - \sum_{i \in I} a_i \sigma(W_i x) \right\|^2 \\ + 2\epsilon \left\langle g_0(x) - \sum_{i \in I} a_i \sigma(W_i x), g_1(\epsilon x) - \sum_{j \notin I} a_j \sigma(W_j x) \right\rangle + \epsilon^2 \left\| g_1(\epsilon x) - \sum_{j \notin I} a_j \sigma(W_j x) \right\|^2$$

12 where I and I^c are sets of nodes, each large enough for UA to ensure vanishing loss. Renormalize by letting $\hat{x} = \epsilon x$ so
 13 that UA works for $g_1(\cdot)$, then the 2nd term rewrites as

$$2\epsilon \left\langle g_0(x) - \sum_{i \in I} a_i \sigma(W_i x), g_1(\hat{x}) - \sum_{j \notin I} a_j \sigma\left(\frac{W_j}{\epsilon} \hat{x}\right) \right\rangle.$$

14 This is in the form of $\epsilon \hat{f}_1(\theta/\epsilon, \theta)$ for some $\hat{f}_1(\phi, \varphi)$ that is quasiperiodic in ϕ (quasiperiodic because \hat{x} is multi-
 15 dim). The 3rd term rewrites similarly. Thus, we see $f(\theta) = f_0(\theta) + f_{1,\epsilon}(\theta)$ where f_0 is the 1st term and $f_{1,\epsilon}(\theta) =$
 16 $\epsilon \hat{f}_1(\theta/\epsilon, \theta) + \epsilon^2 \hat{f}_2(\theta/\epsilon, \theta)$ for some \hat{f}_1, \hat{f}_2 quasiperiodic in the 1st argument. Such $f_{1,\epsilon}$ satisfies Cond.1&2 due to its
 17 quasiperiodic small scale. \square

18 Like most theory papers, we also present numerical experiments in which our conclusions still hold although conditions
 19 for our theorems no longer apply. Thanks to the reviews the following will be added (and expanded):

20 **Neural network training.** We use fully connected 5-16-2 MLP to regress UCI Airfoil Self-Noise Data Set, with leaky
 21 ReLU, MSE as loss, and batch gradient. Fig.1 shows large LR again produces stochasticity as our paper studies.

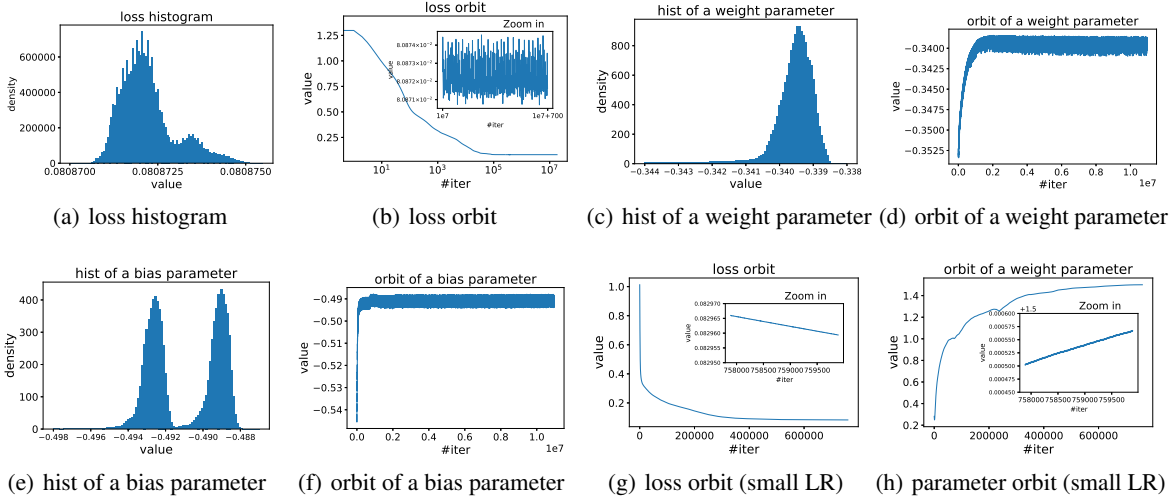


Figure 1: (a)-(f) use LR=0.0165 (large) and demonstrate stochasticity originated from chaos as GD converges to a statistical distribution rather than a local min. (g,h) use LR=0.001 (small) and GD converges to a local min.

22 **Re: $f_{1,\epsilon}$ satisfying Cond.1&2 is like a random variable; tautology?** $f_{1,\epsilon}$ does contribute like a r.v., but this needs to
 23 be proved, which is one of our main contributions – note both x and $f_{1,\epsilon}(x)$ are deterministic even under Cond.1&2!
 24 Cond.1&2 use auxiliary random variables to define the needed $f_{1,\epsilon}$, but $f_{1,\epsilon}$ is a deterministic function.

25 **Re: weaken isotropic noise assumption?** We don't require isotropic 'noise'. Kindly see e.g., Thm.2, which contains
 26 2 statements: (i) convergence to stochastic behavior for general covariance; (ii) explicit characterization of the limiting
 27 statistics when covariance is isotropic (note the same thing holds for SGD).

28 **Re: valid in multi-dim?** Apology that multi-dim. and nonconvex demonstrations were left in Appendix C.2, C.3.3, &
 29 C.5. This rebuttal also adds a neural network example, which is high-dim. & nonconvex, and our conclusion still holds.