
Chaotic Dynamics are Intrinsic to Neural Network Training with SGD

Anonymous Author(s)

Affiliation

Address

email

1 Proof for reviewer tUai

2 Consider a neural network trained by SGD, whose output at time step t is given by a function
3 $\mathbf{f}(\mathbf{w}^{(t)}; \mathbf{z}^{(t)})$ parametrized by weights $\mathbf{w}^{(t)}$ from the reals, subject to inputs $\mathbf{z}^{(t)}$ from the reals. For
4 more clarity, the inputs $\mathbf{z}^{(t)}$ will be omitted in the rest of the proof. Otherwise, we assume the same
5 formalism as in the main paper.

6 In our paper, we defined the finite-time Lyapunov matrix after a single time step as

$$\mathbf{\Lambda}^{(t+1,t)} = \frac{1}{2} \ln \mathbf{Y}^{(t+1,t)T} \mathbf{Y}^{(t+1,t)}, \quad (1)$$

7 where $\mathbf{Y}^{(t+1,t)}$ is the tangent map after a single time step given by

$$\mathbf{Y}^{(t+1,t)} = \mathbf{J}_F^{(t)} = \left\{ \frac{\partial w_i^{(t+1)}}{\partial w_j^{(t)}} \right\}_{ij}. \quad (2)$$

8 The Neural Tangent Kernel is defined as $K^{(t)} = \nabla \mathbf{f}(\mathbf{w}^{(t)}) \nabla \mathbf{f}(\mathbf{w}^{(t)})^T$, where

$$\nabla \mathbf{f}^{(t)} = \left\{ \frac{\partial f_i^{(t)}}{\partial w_j^{(t)}} \right\}_{ij} \quad (3)$$

9 Using the chain rule, the expression for the tangent map can be rewritten as

$$Y_{ij}^{(t+1,t)} = \frac{\partial w_i^{(t+1)}}{\partial w_j^{(t)}} = \sum_k \frac{\partial w_i^{(t+1)}}{\partial f_k^{(t)}} \frac{\partial f_k^{(t)}}{\partial w_j^{(t)}} \quad (4)$$

10 Using the update rule $w_i^{(t+1)} = w_i^{(t)} - \gamma g_i^{(t)}$ for SGD (no momentum), we find

$$Y_{ij}^{(t+1,t)} = \sum_k \frac{\partial}{\partial f_k^{(t)}} (w_i^{(t)} - \gamma g_i^{(t)}) \frac{\partial f_k^{(t)}}{\partial w_j^{(t)}} = \sum_k -\gamma \frac{\partial g_i^{(t)}}{\partial f_k^{(t)}} \frac{\partial f_k^{(t)}}{\partial w_j^{(t)}}. \quad (5)$$

11 Finally, using the definition of the gradient as $g_i^{(t)} = \frac{\partial \mathcal{L}(\mathbf{w}^{(t)})}{\partial w_i^{(t)}}$, we obtain

$$Y_{ij}^{(t+1,t)} = -\gamma \sum_k \frac{\partial^2 \mathcal{L}(\mathbf{w}^{(t)})}{\partial w_i^{(t)} \partial f_k^{(t)}} \frac{\partial f_k^{(t)}}{\partial w_j^{(t)}}, \quad (6)$$

12 or in matrix notation

$$\mathbf{Y}^{(t+1,t)} = -\gamma \mathbf{H}_{wy}^{(t)} \cdot \nabla \mathbf{f}(\mathbf{w}^{(t)}), \quad (7)$$

13 where $\mathbf{H}_{wy}^{(t)}$ denotes the Parameter-Output-Hessian at time step t .