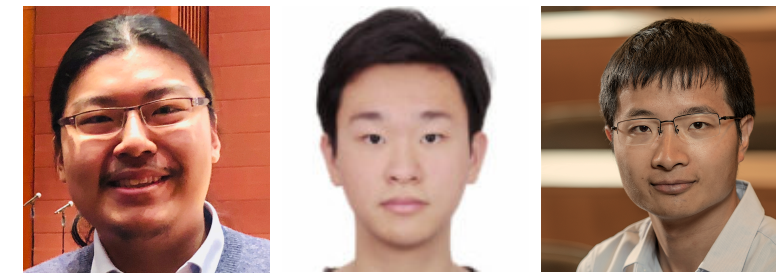
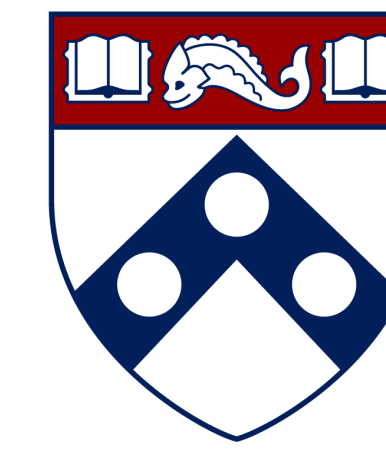


Imitating Deep Learning Dynamics via Locally Elastic Stochastic Differential Equations



Jiayao Zhang 张家耀
Hua Wang 王华
Weijie J. Su 苏炜杰
University of Pennsylvania
{zjiayao, wanghua, swj}@wharton.upenn.edu



Introduction

Local Elasticity: training on a sample has larger effects on samples similar to it than those dissimilar to it. For example, training on an image of cat has a greater effect on other images of cats than images of, say, dogs.

Understanding the training dynamics of deep learning models is perhaps a necessary step toward demystifying their effectiveness. In particular,

how do data from different classes gradually become *separable in their feature spaces* when training neural networks using stochastic gradient descent?

We take a *phenomenological approach* to model feature evolutions of neural net training using a set of stochastic differential equations (SDEs) that *each corresponds to a training sample*. Concretely, for binary classification, with m being time, superscripts class indices, subscripts sample indices, and (α, β) parameters measuring the strengths of *local elasticity*¹⁻³, we model

$$\begin{cases} \dot{x}_1^1(m) = x_1^1(m-1) + h \cdot \alpha x_1^1(m-1) + \text{noise}, \\ \dot{x}_2^2(m) = x_2^2(m-1) + h \cdot \beta x_1^1(m-1) + \text{noise}. \end{cases}$$

training sample from the j^{th} class at the $(m-1)$ -th step

Our main finding uncovers a **sharp phase transition** regarding the *intra-class* impact: if and only if the SDEs are *locally elastic* in the sense that the impact is more significant on samples from the same class as the input, the features of the training data are asymptotically linearly separable.

The LE-SDE/LE-ODE Model

We define the *generic* feature vector $\tilde{\mathbf{x}}(t) = (\tilde{x}^k(t))_{k=1}^K \in \mathbb{R}^{Kp}$ as the concatenation of p -dimensional feature vectors from K classes, and model its dynamics as $d\tilde{\mathbf{x}}(t) = \mathbf{M}(t)\tilde{\mathbf{x}}(t)dt + \Sigma(t)d\mathbf{B}_t$, where the drift $\mathbf{M}(t) = (\mathbf{E}(t) \otimes \mathbf{P}) \circ \mathbf{H}$ consists of the LE matrix $\mathbf{E}(t)$ that encodes the strengths of local elasticity (analogous to the magnitudes of α and β), the similarity matrix \mathbf{H} that encodes the direction in which features interacts (analogous to the phase of α and β), and a sampling matrix \mathbf{P} modeling randomnesses from such as mini-batch sampling and label corruption, and imbalanced datasets. Here, \mathbf{B}_t denotes standard Brownian motion.

The simplest LE matrix consists of two values $\alpha(t)$ and $\beta(t)$. In this case, in the large sample limit, we obtain the following LE-ODE for the mean features, where \circ denotes the Hadamard product (with a slight abuse of notation):

$$d\bar{\mathbf{x}}(t) = \mathbf{M}(t)\bar{\mathbf{x}}(t)dt = ((\mathbf{E}(t) \otimes \mathbf{P}) \circ \mathbf{H})\bar{\mathbf{x}}(t)dt.$$

When $\mathbf{P} = \mathbf{1}\mathbf{1}^T/K$, we write the LE-ODE as the follows where $\bar{\mathbf{x}} = \mathbb{E}_{\text{data}}\tilde{\mathbf{x}}$:

$$d \begin{bmatrix} \bar{x}_1^1 \\ \vdots \\ \bar{x}_p^1 \\ \vdots \\ \bar{x}_1^K \\ \vdots \\ \bar{x}_p^K \end{bmatrix} = \frac{1}{K} \begin{bmatrix} \alpha(t) & \beta(t) & \dots & \beta(t) \\ \beta(t) & \alpha(t) & \dots & \beta(t) \\ \vdots & \vdots & \ddots & \vdots \\ \beta(t) & \dots & \dots & \alpha(t) \end{bmatrix} \circ \begin{bmatrix} H_{11} & H_{12} & \dots & H_{1K} \\ H_{21} & H_{22} & \dots & H_{2K} \\ \vdots & \vdots & \ddots & \vdots \\ H_{K1} & H_{K2} & \dots & H_{KK} \end{bmatrix} \begin{bmatrix} \bar{x}_1^1 \\ \vdots \\ \bar{x}_p^1 \\ \vdots \\ \bar{x}_1^K \\ \vdots \\ \bar{x}_p^K \end{bmatrix} dt.$$

The Separation Theorem

Our main contributions include the following separation theorem.

Theorem (Separation of LE-SDE)

Suppose $\gamma(t) = \alpha(t) - \beta(t) > 0$, assume $\mathbf{H} = (H_{ij})$ is positive semi-definite (PSD) with positive diagonal entries. As $t \rightarrow \infty$, we have

- if $\gamma(t) = \omega(1/t)$, the features are *separable with probability tending to 1*;
- if $\gamma(t) = o(1/t)$, and the number of per-class-feature n tending to ∞ at an arbitrarily slow rate, the features are asymptotically *pairwise separable with probability 0*.

Here, $\gamma(t) = \omega(1/t)$ stands for $\gamma(t) \gg 1/t$ as $t \rightarrow \infty$. For example, $1/t^{0.5} = \omega(1/t)$ and $(t \ln t)^{-1} = o(1/t)$ as $t \rightarrow \infty$.

Visualizing Phase Transition

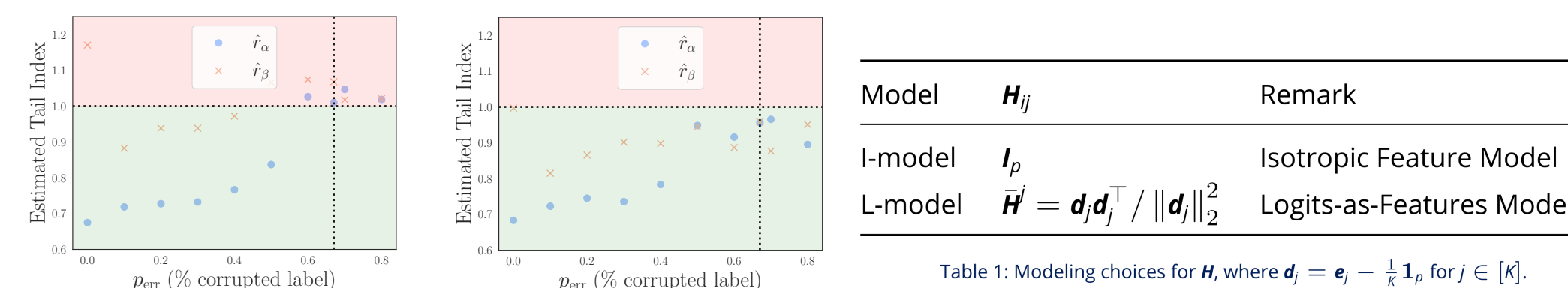


Figure 1(a): Estimated tail indices under I-model. Figure 1(b): Estimated tail indices under L-model.

We estimate the LE matrix under *different modeling choices* of the similarity matrix H , as shown in Table 1, from simulations on GeoMNIST, a dataset consisting of three geometric shapes, using a variant of the AlexNet. We show in Figures 1 the estimated tail index r_α ($-\ln \alpha(t) \approx C - r_\alpha \ln t$) versus the corrupted label ratio p_{err} . As p_{err} increases, we expect LE effects diminish since the dataset is becoming more like random data.

- In the I-model (Figure 1(a)), a clear phase transition for the tail of $(\alpha - \beta)$ occurs around $p_{err} = 2/3$, when the dataset has completely random labels.
- Although in the L-model (Figure 1(b)), the phase transition for the tail is less obvious, note around $p_{err} = 2/3$ the index of β begins to dominate that of α .

Corollary: Connection with Neural Collapse

Neural collapse⁴⁻⁵ is a recent phenomenological finding on the geometry of logits of DNNs at convergence: they tend to form *equiangular tight frames* (ETFs). Let $B(t)$ be the definite integral of $\beta(\tau)$ from $\tau = 0$ to $\tau = t$, we have the following corollary.

Proposition (Neural Collapse of the LE-ODE)

Under L-model and the same setup as in Theorem 1, if $\gamma(t) > 0$ and there exists some $T > 0$ such that $B(t) < 0$ for $t \geq T$, then $\{\tilde{\mathbf{x}}^k(t) / \|\tilde{\mathbf{x}}^k(t)\|\}_{k=1}^K$ forms an ETF as $t \rightarrow \infty$.

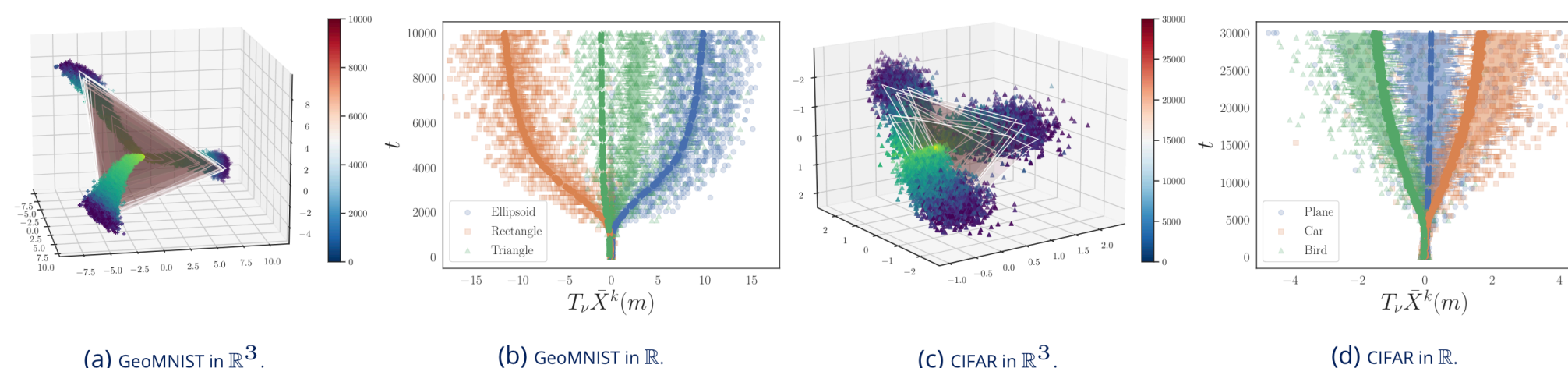


Figure 2: Logit evolution in training. More details are given in the caption of Figure 1 of our paper. Here note that the logits are tending to form an ETF as training progresses (c.f. (a) and (c)).

Imitating DNN Dynamics

Estimating the LE Matrix

Under the I-model and the L-model, the solutions to the LE-ODE can be solved exactly, whence we are able to estimate the LE matrix in terms of the parameters $\alpha(t)$ and $\beta(t)$ based on their cumulative functions $A(t)$ and $B(t)$.

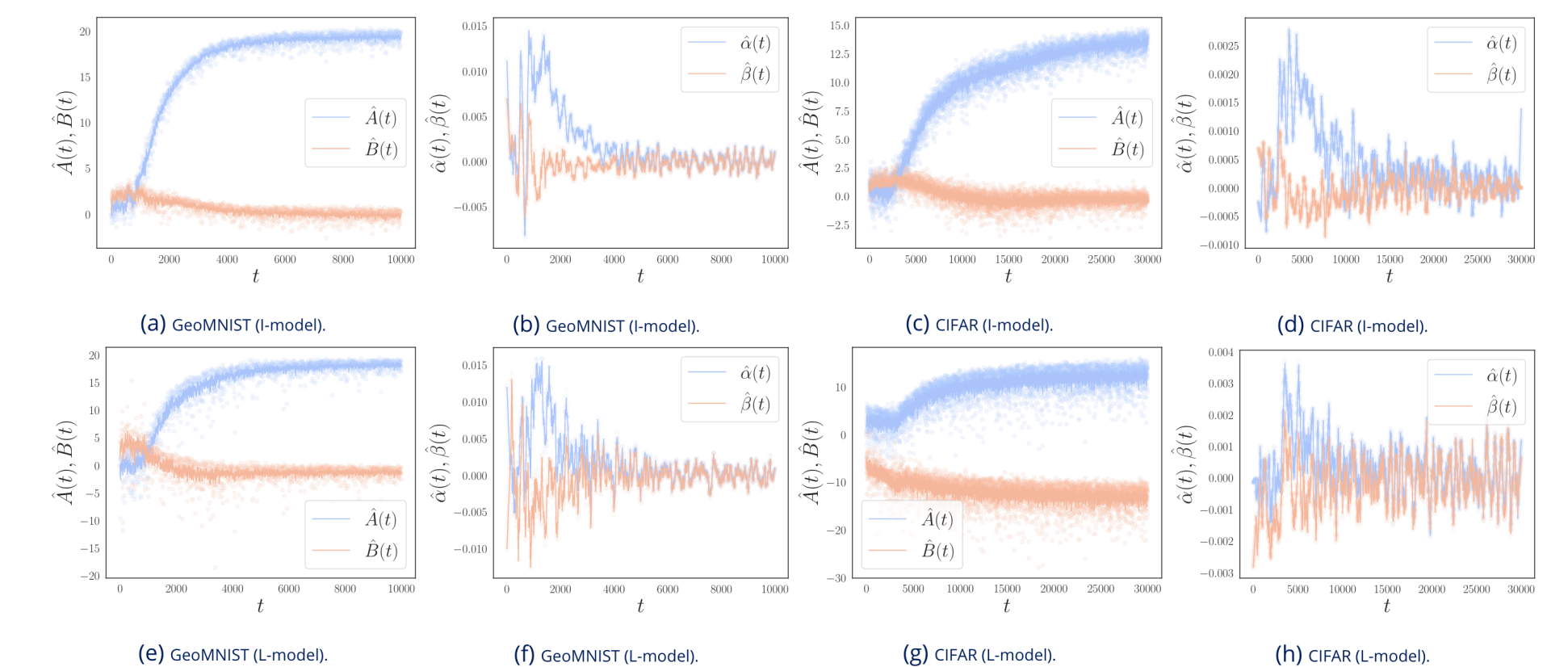


Figure 3: Estimated parameters for the LE matrix $E(t)$.

Imitating DNN Dynamics and Evaluations

With $E(t)$ estimated, we can simulate the LE-SDE using say forward Euler method to test if our model specification is reasonably correct. We assess the goodness-of-fit via the following relative difference (RD, the lower the better):

$$RD_k(t) := \frac{\|\tilde{\mathbf{x}}^k(t) - \tilde{\mathbf{y}}^k(t)\|_{\ell^k}}{(\|\tilde{\mathbf{x}}^k(t)\|_2 + \|\tilde{\mathbf{y}}^k(t)\|_2) / 2}.$$

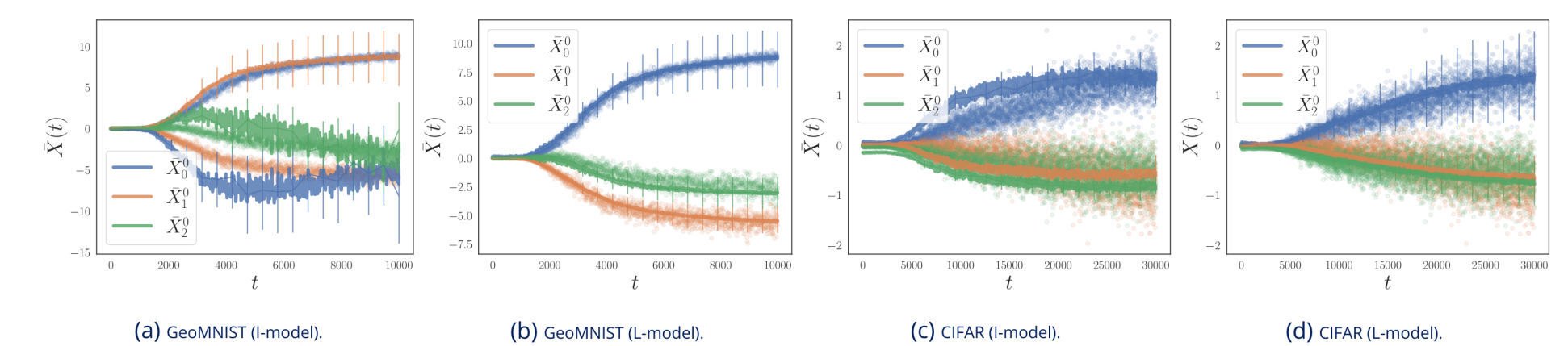


Figure 4: Simulated and genuine logistic trajectories.

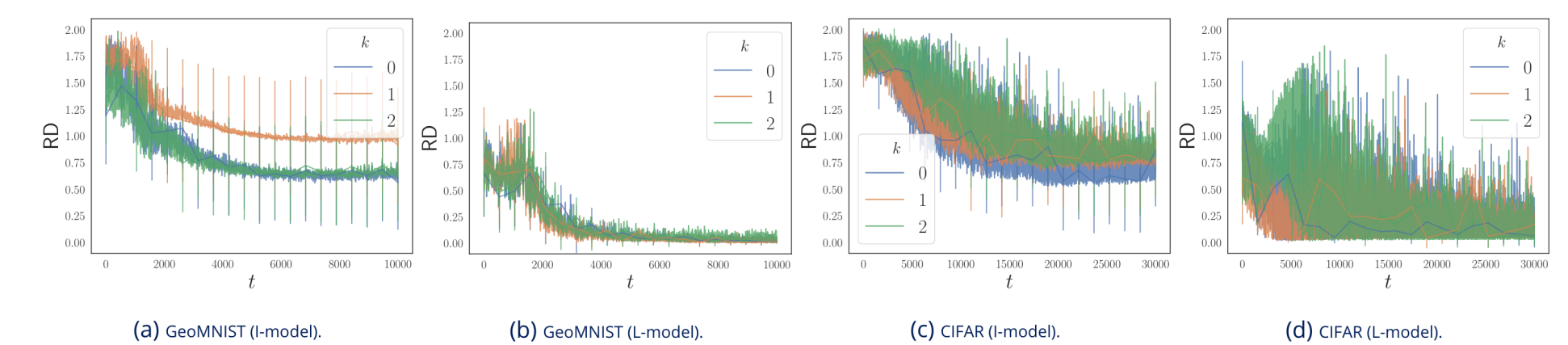


Figure 5: Relative difference (RD) between simulated and genuine trajectories.

References

- H. He and W. J. Su. *The local elasticity of neural networks*. ICLR, 2020.
- S. Chen, H. He, W. J. Su. *Label-aware neural tangent kernel: Toward better generalization and local elasticity*. NeurIPS, 2020.
- Z. Deng, H. He, W. J. Su. *Toward better generalization bounds with locally elastic stability*. ICML, 2021.
- V. Papayan, X. Y. Han, and D. L. Donoho. *Prevalence of neural collapse during the terminal phase of deep learning training*. PNAS, 2020.
- C. Fang, H. He, Q. Long, and W. J. Su. *Exploring deep neural networks via layer-peeled model: Minority collapse in imbalanced training*. PNAS, 2021.

