

Learning Dynamics of Deep Linear Networks Beyond the Edge of Stability

Avrajit Ghosh^{*1}, Soo Min Kwon^{*2}, Rongrong Wang¹, Saiprasad Ravishankar¹, Qing Qu²

¹ Computational Mathematics Science and Engineering, Michigan State University, ² Department of EECS, University of Michigan



Motivation

$$\mathbf{w}(t+1) = \mathbf{w}(t) - \eta \cdot \nabla_{\mathbf{w}} f(\mathbf{w})$$

- How does the step-size $\eta > 0$ affect learning dynamics?
- What kind of solutions are found when η is large?
- GD typically occurs at Edge of Stability (EoS) [1], where sharpness $\|\nabla^2 f(\mathbf{w})\|_2$ hovers about $2/\eta$.

Our Contributions

We focus on deep linear networks (DLNs) at EoS and show that:

- Oscillations only occur within a low-dimensional subspace.
- Subspace dimension depends on the step-size.
- DLN oscillates towards and periodically about the flattest (balanced) minima.
- Conservation law in DLNs (balancing) breaks at EoS.

Loss Optimization and Initialization

- Deep Matrix Factorization Loss (where $\text{rank}(\mathbf{M}_*) = r$):

$$\underset{\Theta}{\text{argmin}} f(\Theta) = \frac{1}{2} \|\mathbf{W}_L \cdot \dots \cdot \mathbf{W}_1 - \mathbf{M}_*\|_F^2.$$

- Update using GD with $\eta > 0$:

$$\mathbf{W}_\ell(t) = \mathbf{W}_\ell(t-1) - \eta \cdot \nabla_{\mathbf{W}_\ell} f(\Theta(t-1)), \quad \forall \ell \in [L].$$

- Initialization ($\alpha > 0$ is small):

$$\mathbf{W}_L(0) = \mathbf{0}, \quad \mathbf{W}_\ell(0) = \alpha \mathbf{I}_d, \quad \forall \ell \in [L-1],$$

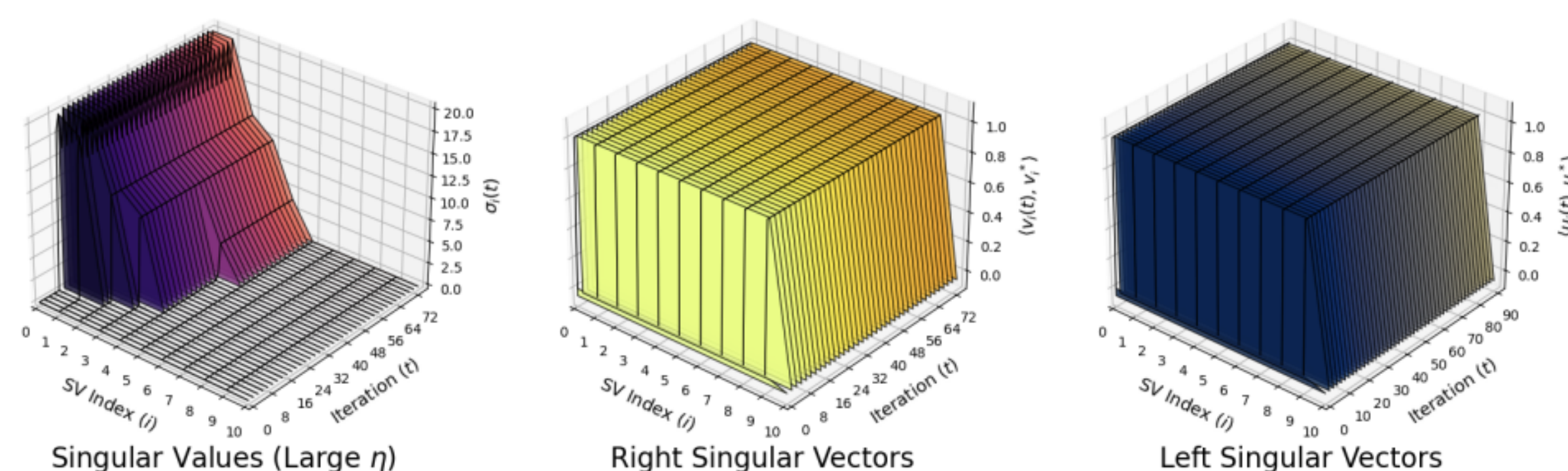
Implicit Bias in Singular Vectors of DLN:

Singular vectors remain static:

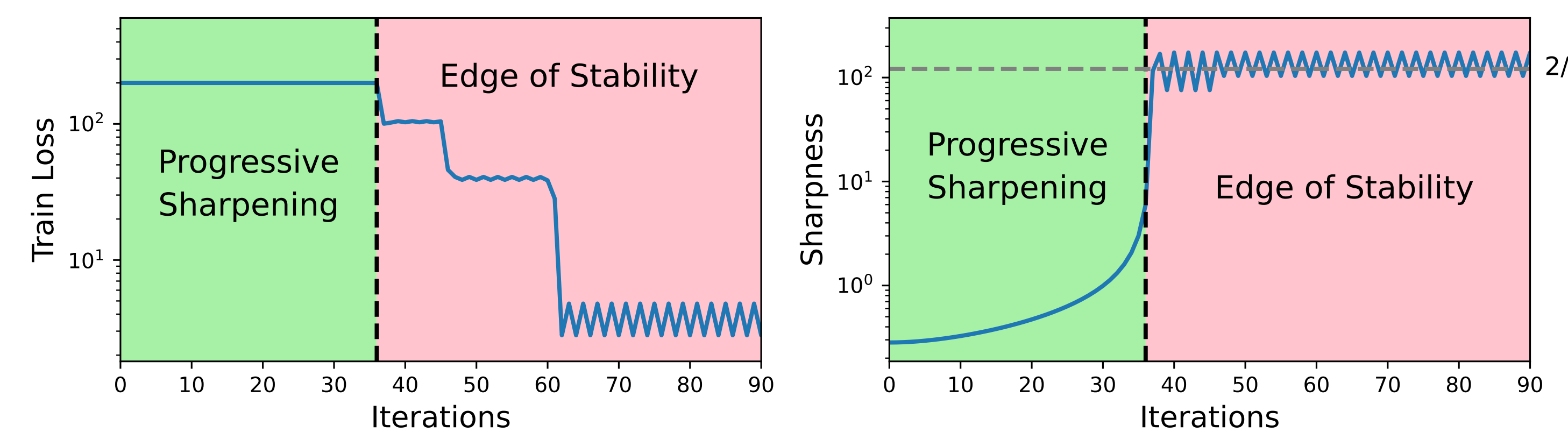
$$\mathbf{W}_L(t) = \mathbf{U}_* \Sigma_L(t) \mathbf{V}_*^\top, \quad \mathbf{W}_\ell(t) = \mathbf{V}_* \Sigma_\ell(t) \mathbf{V}_*^\top.$$

This simplifies the loss:

$$\frac{1}{2} \|\mathbf{W}_{L:1}(t) - \mathbf{M}^*\|_F^2 = \frac{1}{2} \sum_{i=1}^r (\sigma_i(\Sigma_{L:1}(t)) - \sigma_{*,i})^2, \quad (1)$$



Progressive Sharpening and EoS in DLNs



Broken Conservation Law: Balancing

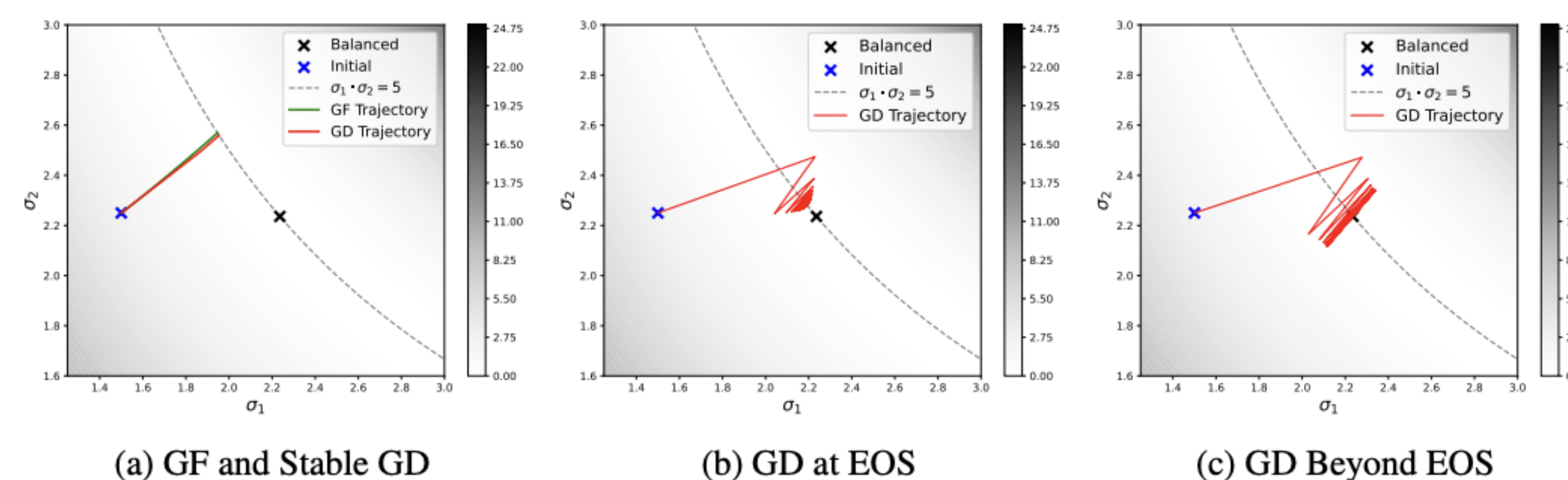
Theorem (Balancing of Singular Values). Consider the simplified loss in (1) and define $S_i := L\sigma_{*,i}^{2-\frac{2}{L}}$. If we run GD with $\eta > \frac{2}{S_i}$ and α satisfies

$$0 < \alpha < \left(\ln \left(\frac{2\sqrt{2}}{\eta S_i} \right) \cdot \frac{\sigma_{*,i}^{4/L}}{L^2 \cdot 2^{\frac{2L-3}{L}}} \right)^{1/4},$$

then for all $\ell \in [L-1]$, there exists a $c \in (0, 1]$:

$$|\sigma_{L,i}^2(t+1) - \sigma_{L,i}^2(t)| < c \cdot |\sigma_{L,i}^2(t) - \sigma_{L,i}^2(t)|.$$

$$\text{Example: } f(\sigma_1, \sigma_2) = \frac{1}{2} (\sigma_2 \cdot \sigma_1 - \sigma_*)^2$$

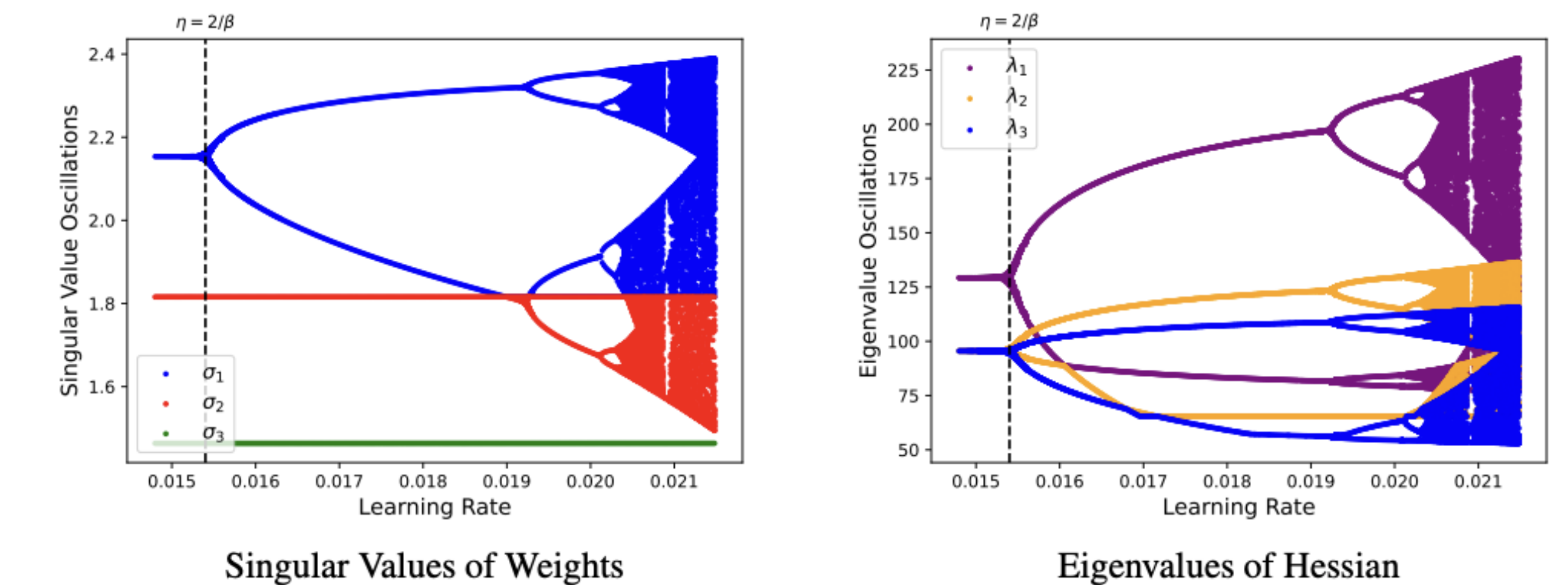


- Gradient flow ($\eta \rightarrow 0$): $|\sigma_1^2(t) - \sigma_2^2(t)|$ stays constant.

Define $S_i := L\sigma_{*,i}^{2-\frac{2}{L}}$ to be stability limit:

- Gradient descent at EoS ($\eta < 2/S_i$): $|\sigma_1^2(t) - \sigma_2^2(t)|$ monotonically reduces to $\mathcal{O}(\alpha)$, where α is initialization scale.
- Gradient descent beyond EoS ($\eta > 2/S_i$): $|\sigma_1^2(t) - \sigma_2^2(t)|$ monotonically reduces to zero due to oscillations!

Period Doubling Route to Chaos in DLNs



We focus on the period-2 case! Here $\beta = L\sigma_*^{2-\frac{2}{L}}$ (i.e., the stability limit).

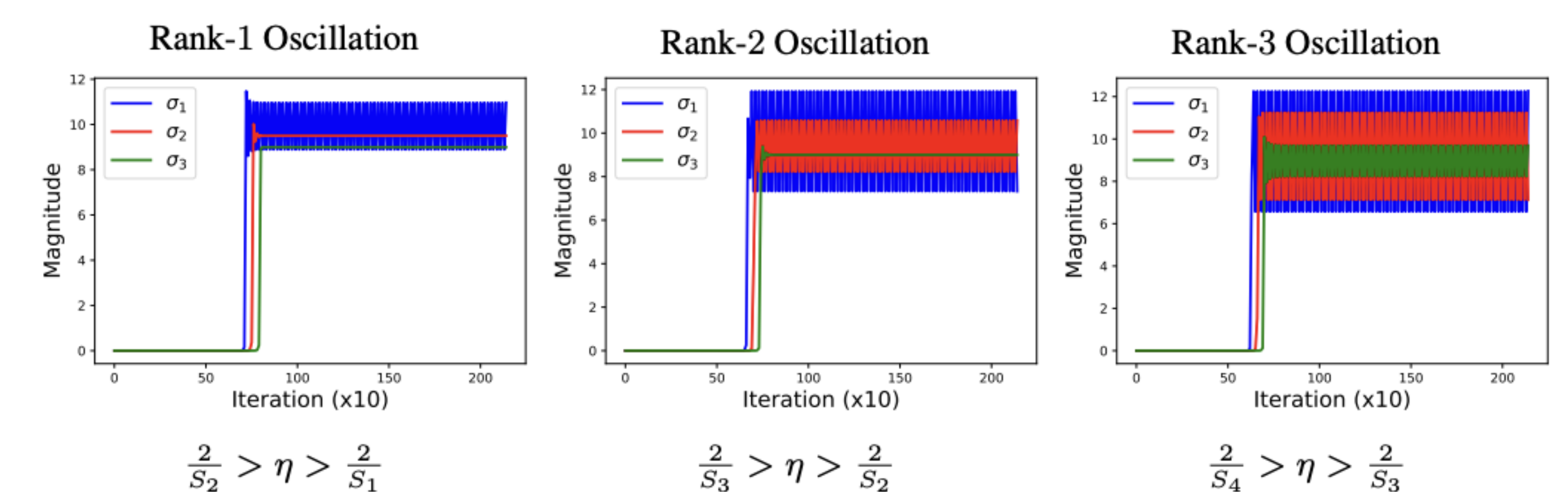
Periodic Subspace Oscillations Beyond EoS

Theorem (Rank- p Periodic Oscillations). Let $\mathbf{M}_* = \mathbf{U}_* \Sigma_* \mathbf{V}_*^\top$. Define $S_p := L\sigma_{*,p}^{2-\frac{2}{L}}$. If we run GD with $\eta > \frac{2}{S_p}$, then the top- p singular values of the end-to-end DLN oscillate in a 2-period orbit ($j \in \{1, 2\}$) around the balanced minimum and admit the following decomposition:

$$\mathbf{W}_{L:1} = \underbrace{\sum_{i=1}^p \rho_{i,j} \cdot \mathbf{u}_{*,i} \mathbf{v}_{*,i}^\top}_{\text{oscillation subspace}} + \underbrace{\sum_{k=p+1}^d \sigma_{*,k} \cdot \mathbf{u}_{*,k} \mathbf{v}_{*,k}^\top}_{\text{stationary subspace}}, \quad j \in \{1, 2\} \quad (2)$$

where $\rho_{i,1} \in (0, \sigma_{*,i}^{1/L})$ and $\rho_{i,2} \in (\sigma_{*,i}^{1/L}, (2\sigma_{*,i})^{1/L})$ are the two real roots of a $(2L-2)(2L-1)$ order polynomial.

$$\text{Example: } \text{rank}(\mathbf{M}_* \in \mathbb{R}^{10 \times 10}) = 3$$



- i -th singular value oscillates in a two-period orbit if $\eta > \frac{2}{S_i}$
- Oscillation range depends on how large the learning rate is!
- Does all Symmetry induced Conservation laws break at Edge of Stability?

References

[1] Cohen et al., "Gradient Descent on Neural Networks Typically Occurs at the Edge of Stability". ICLR 2021.