# Empirical Phase Diagram for Three-layer Neural Networks with infinite Width

**Hanxu Zhou[1], Qixuan Zhou[1], Zhenyuan Jin[1], Tao Luo[1,2], Yaoyu Zhang[1,3], Zhi-Qin John Xu[1],***

(*corresponding author)

[1]School of Mathematical Sciences, Institute of Natural Sciences, MOE-LSC and Qing Yuan Research Institute, Shanghai Jiao Tong University
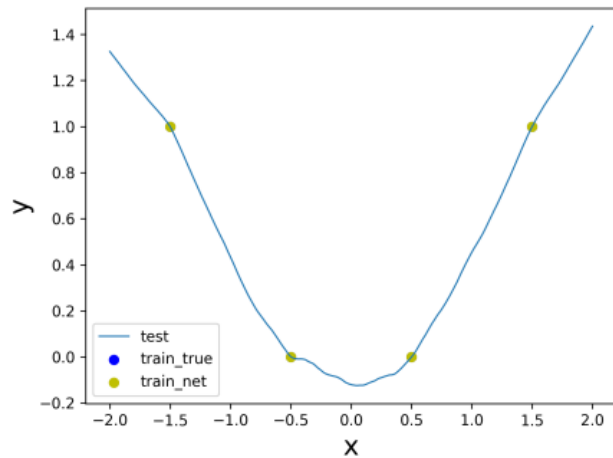
[2]CMA-Shanghai

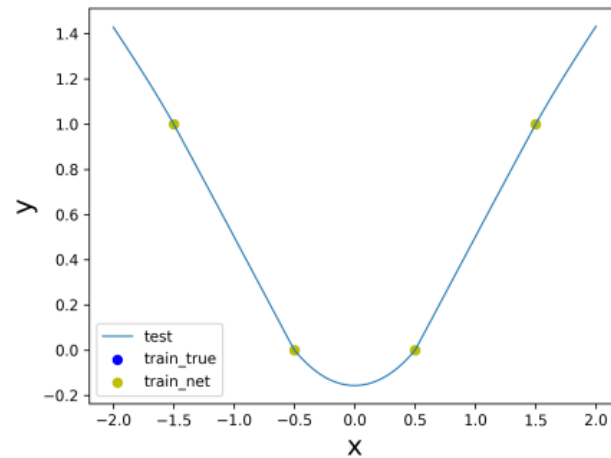[3]Shanghai Center for Brain Science and Brain-Inspired Technology

**2022.10.11**

**The output of different initialization methods has differentiated properties**



| NTK | Xavier | Condensed |

- Learning four data points by three-layer ReLU NNs with different initialization methods.
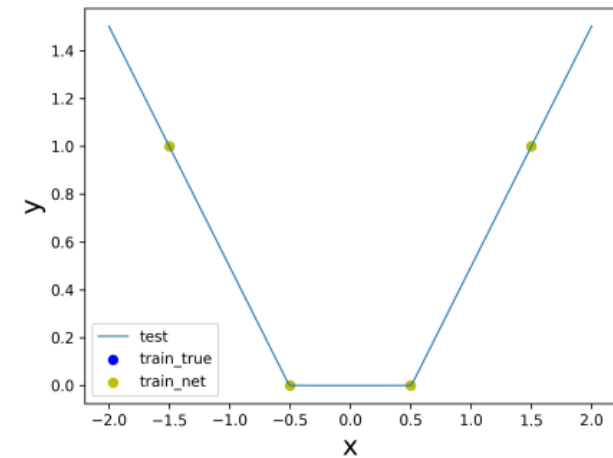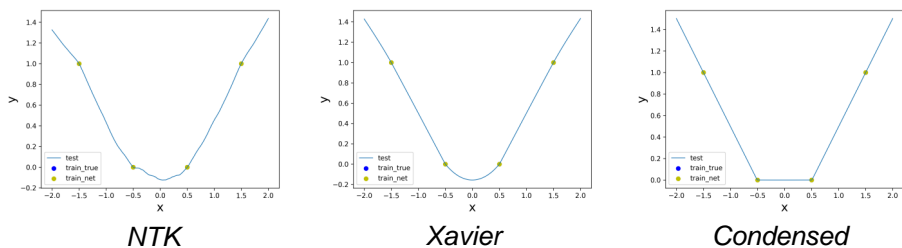
## The output of different initialization methods has differentiated properties

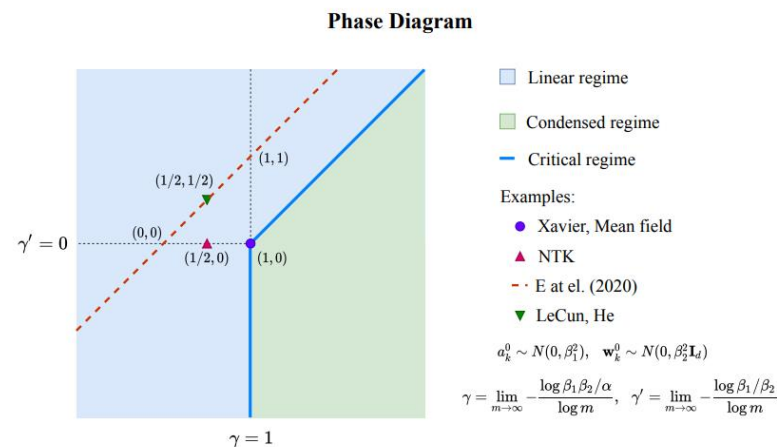

NTK            Xavier            Condensed

- Learning four data points by three-layer ReLU NNs with different initialization methods.

Phase Diagram for Two-layer ReLU Neural Networks at Infinite-width Limit

Tao Luo[#], Zhi-Qin John Xu[#], Zheng Ma, Yaoyu Zhang[*]



**Phase Diagram**

- Linear regime
- Condensed regime
- — Critical regime

Examples:
- Xavier, Mean field
- ▲ NTK
- — · E at el. (2020)
- ▼ LeCun, He

$a_k^0 \sim N(0, \beta_1^2), \quad \mathbf{w}_k^0 \sim N(0, \beta_2^2 \mathbf{I}_d)$

$\gamma = \lim_{m \to \infty} -\frac{\log \beta_1 \beta_2 / \alpha}{\log m}, \quad \gamma' = \lim_{m \to \infty} -\frac{\log \beta_1 / \beta_2}{\log m}$
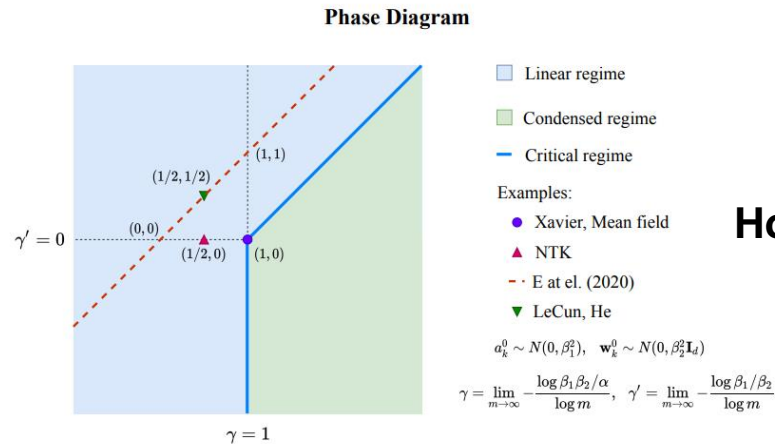
## The output of different initialization methods has differentiated properties

Phase Diagram for Two-layer ReLU Neural Networks at Infinite-width Limit

Tao Luo[#], Zhi-Qin John Xu[#], Zheng Ma, Yaoyu Zhang[*]



**Phase Diagram**

- Linear regime
- Condensed regime
- Critical regime

Examples:
- Xavier, Mean field
- NTK
- E at el. (2020)
- LeCun, He

$$a_k^0 \sim N(0, \beta_1^2), \quad \mathbf{w}_k^0 \sim N(0, \beta_2^2 \mathbf{I}_d)$$

$$\gamma = \lim_{m \to \infty} -\frac{\log \beta_1 \beta_2 / \alpha}{\log m}, \quad \gamma' = \lim_{m \to \infty} -\frac{\log \beta_1 / \beta_2}{\log m}$$

**How about the more general case?**

**Difficulty:**

- Multi-layer structure
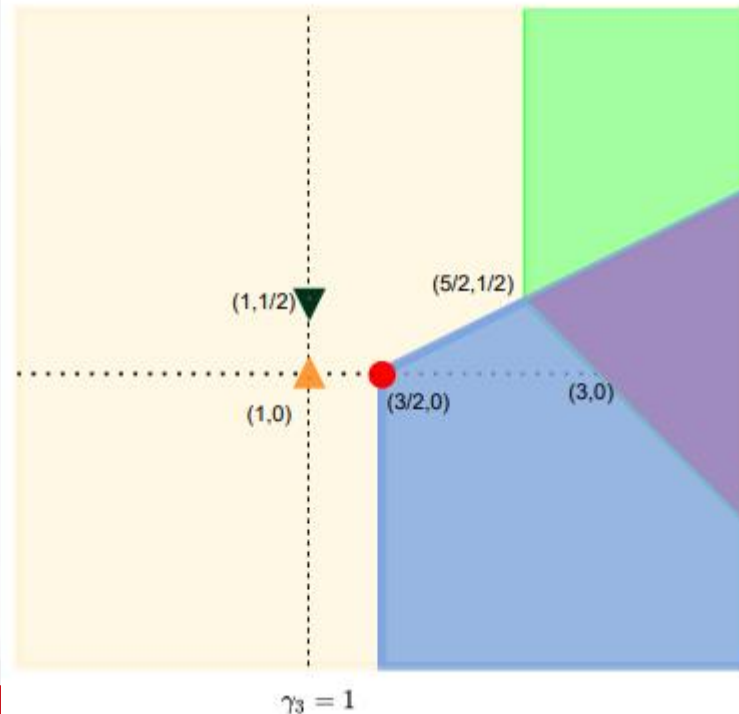- Non-linearity
- Distinct characteristics

**Curiosity:**

- Different frow two-layer
- Distinct dynamics in one NN

**This study:** make a step towards drawing a phase diagram for three-layer ReLU NNs with infinite width

- Figure out **key quantities** and **divide** the dynamics into:
  - a linear regime
  - a condensed regime
  - a critical regime.
- Identify the **condensation** as the strong non-linear signature behavior
- Suggest a complicated **dynamical regimes** consisting of three possible regimes, together with their mixture.

**A three-layer NN with $m$ hidden neurons for each layer is,**

$$f_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{\alpha} \boldsymbol{a}^T \sigma\big(\boldsymbol{W}^{[2]}\sigma(\boldsymbol{W}^{[1]}\boldsymbol{x})\big)$$

where, $\boldsymbol{x} = [\boldsymbol{x}^T, 1]^T$, $\boldsymbol{W}^{[1]} = [\boldsymbol{W}^{[1]}, b_k^{[1]}]^T$, $\bar{\boldsymbol{x}} = \sigma(\boldsymbol{W}^{[1]}\boldsymbol{x})$, $\bar{\boldsymbol{x}} = [\bar{\boldsymbol{x}}^T, 1]^T$, $\boldsymbol{W}^{[2]} = [\boldsymbol{W}^{[2]}, b_k^{[2]}]^T$, and $\boldsymbol{a}_k^0 \sim \mathcal{N}(0, \beta_3^2)$, $\boldsymbol{W}_{kk'}^{[2],0} \sim \mathcal{N}(0, \beta_2^2)$, $\boldsymbol{W}_{kk'}^{[1],0} \sim \mathcal{N}(0, \beta_1^2)$,

**The empirical risk is,**

$$R_S(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^{n} (f_{\boldsymbol{\theta}}(\boldsymbol{x}_i) - y_i)^2$$

**A three-layer NN with $m$ hidden neurons for each layer is,**

$$f_{\boldsymbol{\theta}}(\boldsymbol{x}) = \frac{1}{\alpha}\boldsymbol{a}^T\sigma\big(\boldsymbol{W}^{[2]}\sigma(\boldsymbol{W}^{[1]}\boldsymbol{x})\big), \qquad a_k^0 \sim \mathcal{N}(0,\beta_3^2), \ W_{kk'}^{[2],0} \sim \mathcal{N}(0,\beta_2^2), \ W_{kk'}^{[1],0} \sim \mathcal{N}(0,\beta_1^2),$$

where, $\boldsymbol{x} = [\boldsymbol{x}^T, 1]^T$, $\boldsymbol{W}^{[1]} = [\boldsymbol{W}^{[1]}, b_k^{[1]}]^T$, $\bar{\boldsymbol{x}} = \sigma(\boldsymbol{W}^{[1]}\boldsymbol{x})$, $\bar{\boldsymbol{x}} = [\bar{\boldsymbol{x}}^T, 1]^T$, $\boldsymbol{W}^{[2]} = [\boldsymbol{W}^{[2]}, b_k^{[2]}]^T$.

**The gradient flow of $\boldsymbol{\theta} = \mathrm{vec}\{\boldsymbol{a}, \boldsymbol{W}^{[2]}, \boldsymbol{W}^{[1]}\}$,**

$$\frac{\mathrm{d}\boldsymbol{a}}{\mathrm{d}t} = -\frac{1}{n}\sum_{i=1}^{n}\frac{1}{\alpha}\sigma\big(\boldsymbol{W}^{[2]}\sigma(\boldsymbol{W}^{[1]}\boldsymbol{x}_i)\big)e_i,$$

$$\frac{\mathrm{d}\boldsymbol{W}^{[2]}}{\mathrm{d}t} = -\frac{1}{n}\sum_{i=1}^{n}\frac{1}{\alpha}\boldsymbol{a}\odot\sigma'\big(\boldsymbol{W}^{[2]}\sigma(\boldsymbol{W}^{[1]}\boldsymbol{x}_i)\big)\sigma(\boldsymbol{W}^{[1]}\boldsymbol{x}_i)^{\mathbf{T}}e_i,$$

$$\frac{\mathrm{d}\boldsymbol{W}^{[1]}}{\mathrm{d}t} = -\frac{1}{n}\sum_{i=1}^{n}\frac{1}{\alpha}\boldsymbol{W}^{[2]\mathbf{T}}\big(\boldsymbol{a}\odot\sigma'(\boldsymbol{W}^{[2]}\sigma(\boldsymbol{W}^{[1]}\boldsymbol{x}_i))\big)\odot\sigma'(\boldsymbol{W}^{[1]}\boldsymbol{x}_i)\boldsymbol{x}_i^{\mathbf{T}}e_i,$$

where $e_i = \left(\frac{1}{\alpha}\boldsymbol{a}^T\sigma\big(\boldsymbol{W}^{[2]}\sigma(\boldsymbol{W}^{[1]}\boldsymbol{x})\big) - y_i\right)$ , the operation $\odot$ is the Hadamard product.

**The gradient flow of $\theta = \mathrm{vec}\{a, W^{[2]}, W^{[1]}\}$,**

$$\frac{\mathrm{d}a}{\mathrm{d}t} = -\frac{1}{n}\sum_{i=1}^{n}\frac{1}{\alpha}\sigma(W^{[2]}\sigma(W^{[1]}x_i))e_i,$$

$$\frac{\mathrm{d}W^{[2]}}{\mathrm{d}t} = -\frac{1}{n}\sum_{i=1}^{n}\frac{1}{\alpha}a \odot \sigma'(W^{[2]}\sigma(W^{[1]}x_i))\sigma(W^{[1]}x_i)^{\mathbf{T}}e_i,$$

$$\frac{\mathrm{d}W^{[1]}}{\mathrm{d}t} = -\frac{1}{n}\sum_{i=1}^{n}\frac{1}{\alpha}W^{[2]\mathbf{T}}(a \odot \sigma'(W^{[2]}\sigma(W^{[1]}x_i))) \odot \sigma'(W^{[1]}x_i)x_i^{\mathbf{T}}e_i,$$

$$\sigma(au) = a\sigma(u), \sigma'(au) = \sigma'(u)$$

**The normalized gradient flow of $\theta$,**

$$\frac{\mathrm{d}\overline{a}}{\mathrm{d}\overline{t}} = -\left(\frac{1}{n}\sum_{i=1}^{n}\kappa_3\sigma(\overline{W}^{[2]}\sigma(\overline{W}^{[1]}x_i))\right)e_i,$$

$$\frac{\mathrm{d}\overline{W}^{[2]}}{\mathrm{d}\overline{t}} = -\kappa_1^2\left(\frac{1}{n}\sum_{i=1}^{n}\kappa_3\overline{a}\odot\sigma'(\overline{W}^{[2]}\sigma(\overline{W}^{[1]}x_i))\sigma(\overline{W}^{[1]}x_i)^{\mathbf{T}}\right)e_i,$$

$$\frac{\mathrm{d}\overline{W}^{[1]}}{\mathrm{d}\overline{t}} = -\kappa_2^2\left(\frac{1}{n}\sum_{i=1}^{n}\kappa_3\overline{W}^{[2]\mathbf{T}}(\overline{a}\odot\sigma'(\overline{W}^{[2]}\sigma(\overline{W}^{[1]}x_i)))\odot\sigma'(\overline{W}^{[1]}x_i)x_i^{\mathbf{T}}\right)e_i,$$

where $\overline{a} = \frac{1}{\beta_3}a$, $\overline{W}^{[2]} = \frac{1}{\beta_2}W^{[2]}$, $\overline{W}^{[1]} = \frac{1}{\beta_1}W^{[1]}$, $\kappa_1 = \frac{\beta_3}{\beta_2}$, $\kappa_2 = \frac{\beta_3}{\beta_1}$, $\kappa_3 = \frac{\beta_1\beta_2\beta_3}{\alpha}$, $t = \left(\alpha\prod_{i=1}^{3}\kappa_i\right)^{-\frac{2}{3}t}$.

**The normalized gradient flow of $\theta$,**

$$\frac{\mathrm{d}\overline{a}}{\mathrm{d}\overline{t}} = -\left(\frac{1}{n}\sum_{i=1}^{n}\kappa_3\sigma(\overline{W}^{[2]}\sigma(\overline{W}^{[1]}x_i))\right)e_i,$$

$$\frac{\mathrm{d}\overline{W}^{[2]}}{\mathrm{d}\overline{t}} = -\kappa_1^2\left(\frac{1}{n}\sum_{i=1}^{n}\kappa_3\overline{a}\odot\sigma'(\overline{W}^{[2]}\sigma(\overline{W}^{[1]}x_i))\sigma(\overline{W}^{[1]}x_i)^{\mathbf{T}}\right)e_i,$$

$$\frac{\mathrm{d}\overline{W}^{[1]}}{\mathrm{d}\overline{t}} = -\kappa_2^2\left(\frac{1}{n}\sum_{i=1}^{n}\kappa_3\overline{W}^{[2]\mathbf{T}}(\overline{a}\odot\sigma'(\overline{W}^{[2]}\sigma(\overline{W}^{[1]}x_i)))\odot\sigma'(\overline{W}^{[1]}x_i)x_i^{\mathbf{T}}\right)e_i,$$

**The scaling parameters and infinite-width limit,**

$$\kappa_1 = \frac{\beta_3}{\beta_2},\ \kappa_2 = \frac{\beta_3}{\beta_1},\ \kappa_3 = \frac{\beta_1\beta_2\beta_3}{\alpha},\ \overline{t} = \left(\alpha\prod_{i=1}^{3}\kappa_i\right)^{-\frac{2}{3}}t,$$

**Assumption 3.1**: $m_1 = m_2 = m$
**Assumption 3.2**: $\beta_2 = B\beta_3$

$$\gamma_1 = \lim_{m\to\infty}-\frac{\log\kappa_1}{\log m} = 0,\ \gamma_2 = \lim_{m\to\infty}-\frac{\log\kappa_2}{\log m},\ \gamma_3 = \lim_{m\to\infty}-\frac{\log\kappa_3}{\log m}$$

**The scaling parameters and infinite-width limit,**

$$\kappa_1 = \frac{\beta_3}{\beta_2}, \, \kappa_2 = \frac{\beta_3}{\beta_1}, \, \kappa_3 = \frac{\beta_1\beta_2\beta_3}{\alpha}, \, \bar{t} = \left(\alpha \prod_{i=1}^3 \kappa_i\right)^{-\frac{2}{3}} t, \, \gamma_2 = \lim_{m\to\infty} -\frac{\log\kappa_2}{\log m}, \gamma_3 = \lim_{m\to\infty} -\frac{\log\kappa_3}{\log m}$$

**Some common initialization methods**
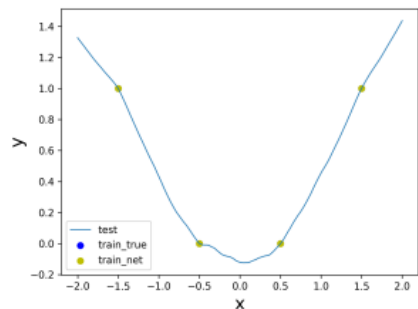
Table 1: Common initialization methods with their scaling parameters

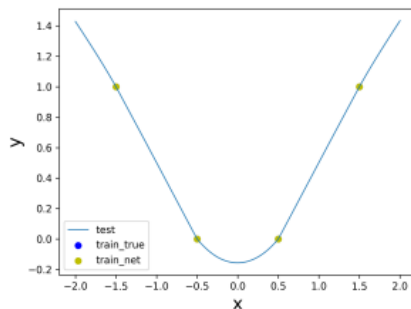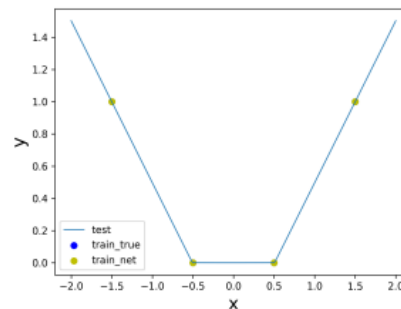| Name | $\alpha$ | $a$ | $W^{[2]}$ | $W^{[1]}$ | $\kappa_2$ | $\kappa_3$ | $\gamma_2$ | $\gamma_3$ |
|---|---|---|---|---|---|---|---|---|
| NTK<br>Jacot et al. (2018) | $\sqrt{m_1 m_2}$ | $1$ | $1$ | $1$ | $1$ | $\sqrt{\frac{1}{m_1 m_2}}$ | $0$ | $1$ |
| Lecun<br>LeCun et al. (2012) | $1$ | $\sqrt{\frac{1}{m_2}}$ | $\sqrt{\frac{1}{m_1}}$ | $\sqrt{\frac{1}{d}}$ | $\sqrt{\frac{d}{m_2}}$ | $\sqrt{\frac{1}{m_1 m_2 d}}$ | $\frac{1}{2}$ | $1$ |
| He<br>He et al. (2015) | $1$ | $\sqrt{\frac{2}{m_2}}$ | $\sqrt{\frac{2}{m_1}}$ | $\sqrt{\frac{2}{d}}$ | $\sqrt{\frac{d}{m_2}}$ | $\sqrt{\frac{8}{m_1 m_2 d}}$ | $\frac{1}{2}$ | $1$ |
| Xavier<br>Glorot and Bengio (2010) | $1$ | $\sqrt{\frac{2}{m_2+1}}$ | $\sqrt{\frac{2}{m_1+m_2}}$ | $\sqrt{\frac{2}{d+m_1}}$ | $\sqrt{\frac{d+m_1}{m_2+1}}$ | $\sqrt{\frac{8/(m_1+m_2)}{(m_2+1)(d+m_1)}}$ | $0$ | $\frac{3}{2}$ |

**Empirical phase diagram**

**Intuitive experiments of synthetic data**
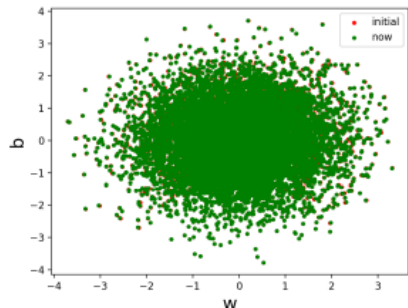


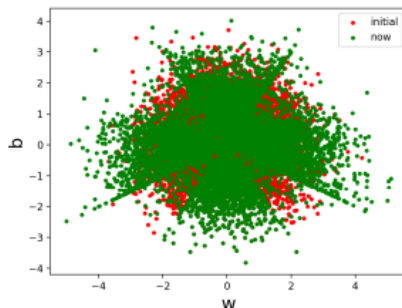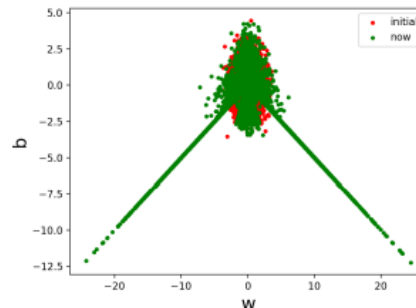(a) $\gamma_3 = 1.0$     (b) $\gamma_3 = 1.5$     (c) $\gamma_3 = 2.0$
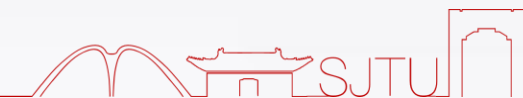
(d) $\gamma_3 = 1.0$     (e) $\gamma_3 = 1.5$     (f) $\gamma_3 = 2.0$

Learning four data points by three-layer ReLU NNs with $m = 10000$ and $\gamma_2 = 0$. The scatter plots in the second row are $\left\{ W_K^{[1]} \right\}_{k=1}^{m} = \left\{ (w_k^{[1]}, b_k^{[1]}) \right\}_{k=1}^{m}$, where red plots represent initial position and green plots represent final position.
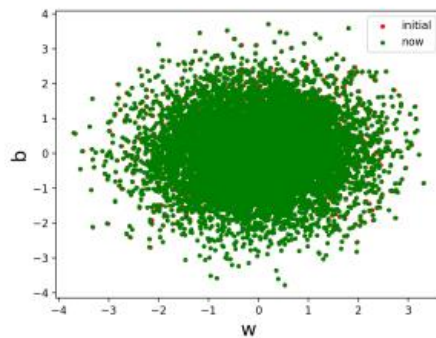
**Regime identification and separation**
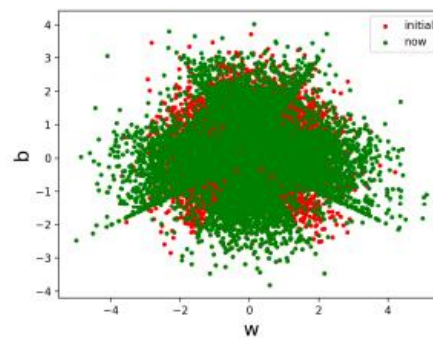
- Relative distance,

$$\mathrm{RD}(\boldsymbol{W}^{[1]}) = \frac{\|\boldsymbol{\theta}^*_{\boldsymbol{W}_1} - \boldsymbol{\theta}_{\boldsymbol{W}_1}(0)\|_2}{\|\boldsymbol{\theta}_{\boldsymbol{W}_1}(0)\|_2}, \quad \mathrm{RD}(\boldsymbol{W}^{[2]}) = \frac{\|\boldsymbol{\theta}^*_{\boldsymbol{W}_2} - \boldsymbol{\theta}_{\boldsymbol{W}_2}(0)\|_2}{\|\boldsymbol{\theta}_{\boldsymbol{W}_2}(0)\|_2},$$
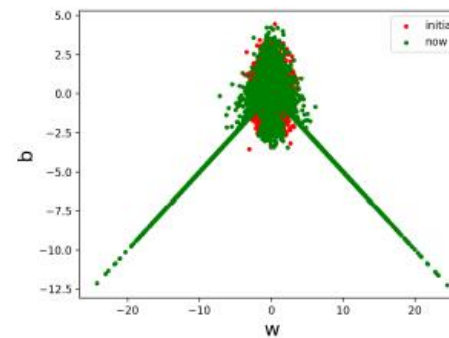
- We empirically consider that as $m \to \infty$,

  - Linear regime: $\displaystyle\sup_{t \in [0,+\infty)} \mathrm{RD}\left(\boldsymbol{W}^{[i]}(t)\right) \to 0, i = 1,2$

  - Condensed regime: $\displaystyle\sup_{t \in [0,+\infty)} \mathrm{RD}\left(\boldsymbol{W}^{[i]}(t)\right) \to +\infty \; i = 1,2$

  - Critical regime: $\displaystyle\sup_{t \in [0,+\infty)} \mathrm{RD}\left(\boldsymbol{W}^{[i]}(t)\right) \to O(1), i = 1,2$



(d) $\gamma_3 = 1.0$      (e) $\gamma_3 = 1.5$      (f) $\gamma_3 = 2.0$

**Regime identification and separation**

- Relative distance,
  $$\mathrm{RD}(\boldsymbol{W}^{[1]}) = \frac{\|\boldsymbol{\theta}^{*}_{\boldsymbol{W}_1} - \boldsymbol{\theta}_{\boldsymbol{W}_1}(0)\|_2}{\|\boldsymbol{\theta}_{\boldsymbol{W}_1}(0)\|_2}, \ \mathrm{RD}(\boldsymbol{W}^{[2]}) = \frac{\|\boldsymbol{\theta}^{*}_{\boldsymbol{W}_2} - \boldsymbol{\theta}_{\boldsymbol{W}_2}(0)\|_2}{\|\boldsymbol{\theta}_{\boldsymbol{W}_2}(0)\|_2},$$

- We empirically found that as $m \to \infty$,
  - Linear regime: $\underset{t \in [0,+\infty)}{\mathrm{Sup}} \ \mathrm{RD}\left(\boldsymbol{W}^{[i]}(t)\right) \to 0, i = 1,2$
  - Condensed regime: $\underset{t \in [0,+\infty)}{\mathrm{Sup}} \ \mathrm{RD}\left(\boldsymbol{W}^{[i]}(t)\right) \to +\infty \ i = 1,2$
  - Critical regime: $\underset{t \in [0,+\infty)}{\mathrm{Sup}} \ \mathrm{RD}\left(\boldsymbol{W}^{[i]}(t)\right) \to O(1), i = 1,2$
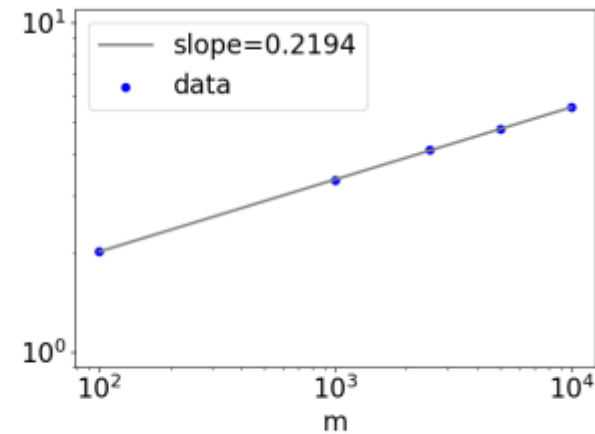
$\mathrm{RD}(\boldsymbol{W}^{[1]})$ v.s. m. Still learn four data points by three-layer ReLU NNs with different $\gamma_3$'s and $\gamma_2 = 0$.



(a) $\gamma_3 = 0.9$  (b) $\gamma_3 = 1.5$  (c) $\gamma_3 = 2.1$

**Regime identification and separation**

$\mathrm{RD}(\boldsymbol{W}^{[1]})$ v.s. m. Still learn four data points by three-layer ReLU NNs with different $\gamma_3$'s and $\gamma_2 = 0$.
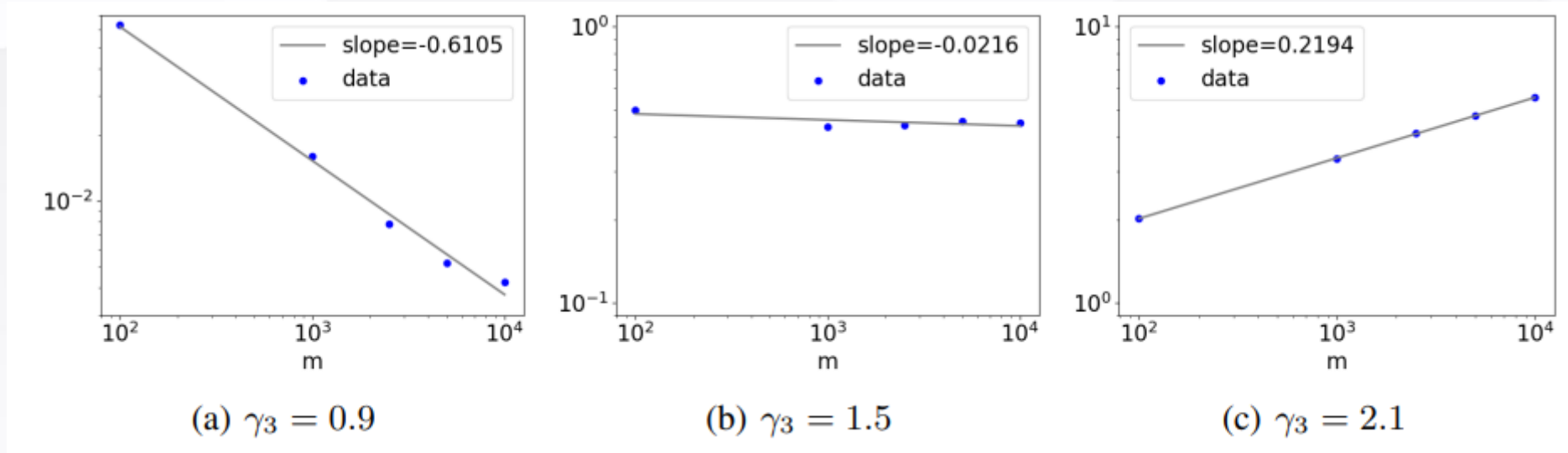


(a) $\gamma_3 = 0.9$       (b) $\gamma_3 = 1.5$       (c) $\gamma_3 = 2.1$

We quantify the growth of $\mathrm{RD}(\boldsymbol{W}^{[i]})$, $i = 1,2$, as $m \to \infty$, by defining,

$$S_{\boldsymbol{W}_{[i]}} = \lim_{m \to \infty} \frac{\log \mathrm{RD}(\boldsymbol{W}^{[i]})}{\log m}$$

- Linear regime:      $S_{\boldsymbol{W}_{[i]}} < 0$
- Condensed regime:      $S_{\boldsymbol{W}_{[i]}} > 0$
- Critical regime:      $S_{\boldsymbol{W}_{[i]}} = 0$

**Regime identification and separation**

We quantify the growth of $\mathrm{RD}(\boldsymbol{W}^{[i]})$, $i = 1,2$, as $m \to \infty$, by defining,

$$S_{\boldsymbol{W}^{[i]}} = \lim_{m \to \infty} \frac{\log \mathrm{RD}(\boldsymbol{W}^{[i]})}{\log m}$$

- Linear regime: $\qquad S_{\boldsymbol{W}^{[i]}} < 0$
- Condensed regime: $\qquad S_{\boldsymbol{W}^{[i]}} > 0$
- Critical regime: $\qquad S_{\boldsymbol{W}^{[i]}} = 0$

For synthetic data,



(a) $S_{\boldsymbol{W}_1}$      (b) $S_{\boldsymbol{W}_1}$      (c) $S_{\boldsymbol{W}_2}$

**Regime identification and separation**

For synthetic data,



(a) $S_{W_1}$       (b) $S_{W_1}$       (c) $S_{W_2}$

For mnist data,



(a) $S_{W_1}$ of mnist       (b) $S_{W_1}$ of mnist       (c) $S_{W_2}$ of mnist

- Characterize the linear, critical, and condensed regimes
- Identify the condensation as non-linear
- Figure out the relation between the training dynamics and initialization
- Draw the phase diagram
- Reveal different training dynamics within a neural network