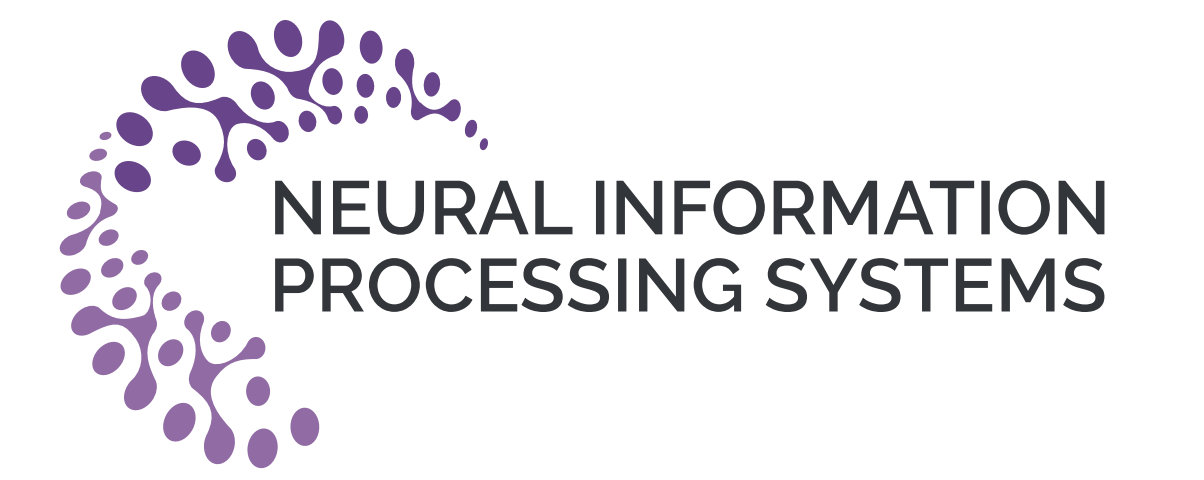


SGD vs GD : Rank Deficiency in Linear Networks

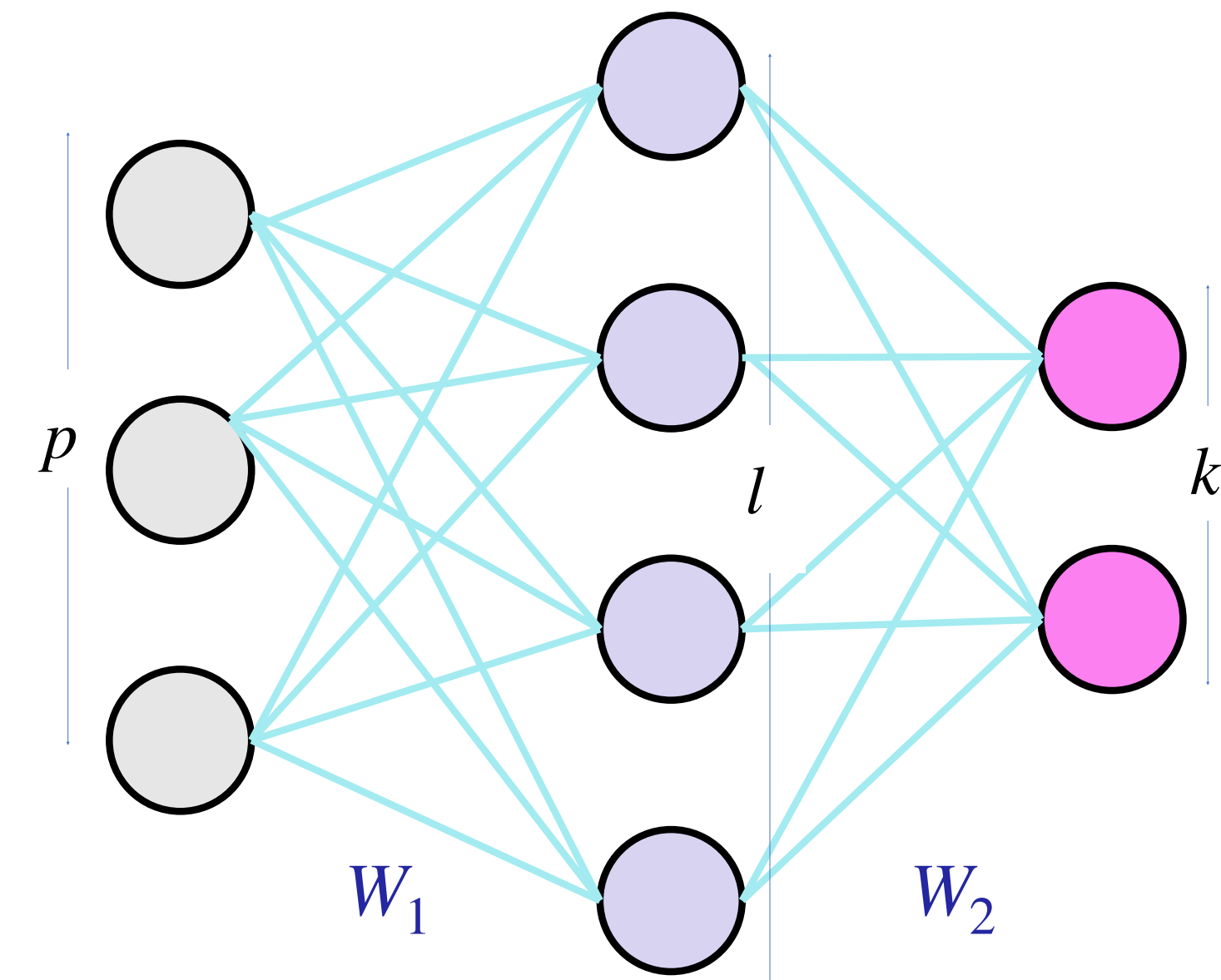
Aditya Varre, Margarita Sagitova, Nicolas Flammarion @ Theory of Machine Learning Lab (TML), Lausanne



Regression with Linear Network

Setup: The $(x_i)_{i=1}^n$ from \mathbb{R}^p and w.l.o.g assume the labels are generated by $y_i = U_*^T x_i \in \mathbb{R}^k$ and $U_* \in \mathbb{R}^{d \times k}$.

A linear network representing function f

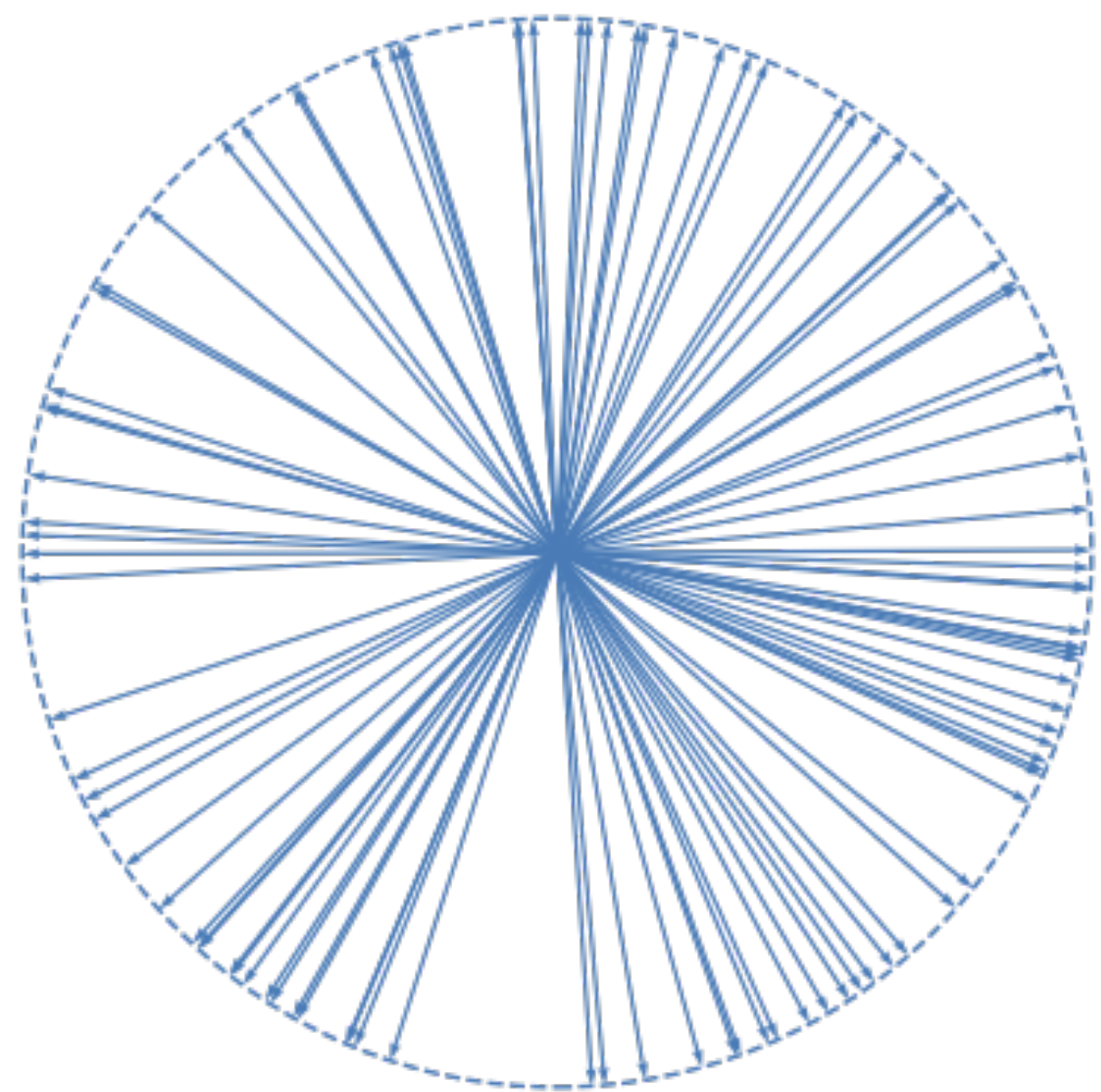


Linear Network: Let hidden layer be $W_1 \in \mathbb{R}^{d \times l}$ and weight layer $W_2 \in \mathbb{R}^{l \times k}$. The network represents the function

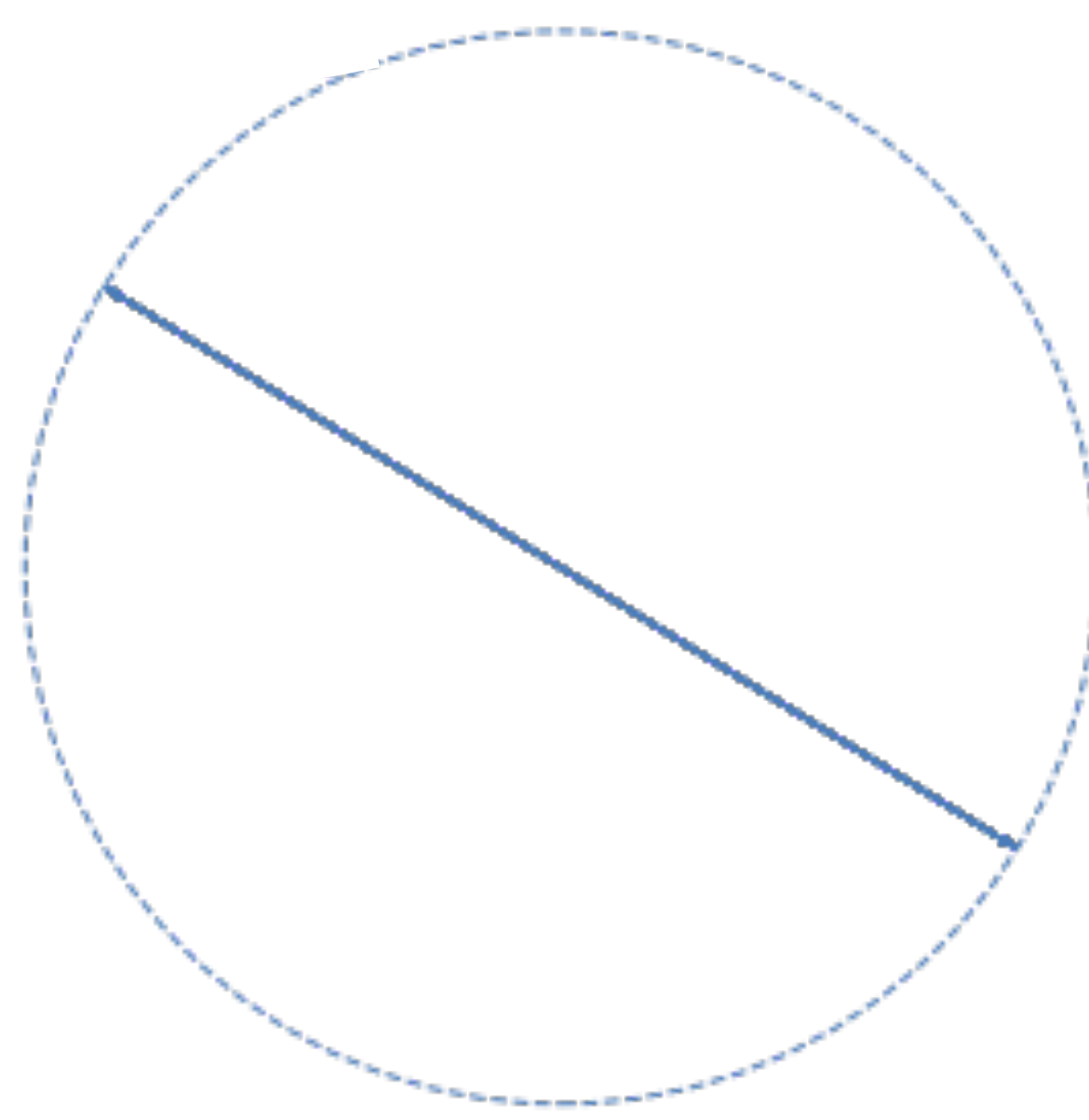
$$f(x) = W_2^T W_1^T x.$$

Square Loss: $L(W_1, W_2) = \frac{1}{2} \|Y - XW_1W_2\|^2$ where X, Y denotes the data in a matrix form.

A simple illustration: Let us consider the case with $k = 1$, we train GD and SGD with same initialization



(a) Gradient Descent



(b) Stochastic Gradient Descent

The directions of columns of W_1 at the end of training linear network.

SGD and simple structures:

- a) Limiting dynamics of SGD - Blanc et. al. 2020, Damian et. al. 2020, Li et. al. 2021.
- b) SGD on diagonal networks and bias to sparse predictor- Pesme et. al. 2021, Pillaud-Vivien et.al. 2022
- c) Empirical works in deep networks - Andruischenko et. al. 2023, Chen et. al. 2023.

How does **stochasticity** enable the emergence of **simple structures**?

Gradient Flow: Gradient flow on linear networks can be written as

$$\begin{aligned} dW_1 &= -\nabla_{W_1} L(W_1, W_2) dt = X^T (Y - XW_1W_2) W_2^T dt, \\ dW_2 &= -\nabla_{W_2} L(W_1, W_2) dt = W_1^T X^T (Y - XW_1W_2) dt. \end{aligned}$$

Let the equivalent linear predictor $\beta = W_1W_2$, $\nabla L(\beta) = X^T(X\beta - Y)$.

Block Matrix: Using a block matrix $\theta = [W_1^T | W_2] \in \mathbb{R}^{(p+k) \times (p+k)}$.

$$d\theta = \theta J dt \quad \text{where} \quad J := \begin{bmatrix} 0_{p \times p} & -\nabla L(\beta) \\ -\nabla L(\beta)^T & 0_{k \times k} \end{bmatrix},$$

reveals the inherent **multiplicative nature** of the gradient in linear networks.

Stochastic Gradient Flow: We choose a continuous version of Label noise gradient descent. With $B_t \in \mathbb{R}^{n \times k}$ and δ being the standard Brownian motion and noise level,

$$d\theta = \theta [Jdt + \sqrt{\eta\delta} d\xi] \quad d\xi := \begin{bmatrix} 0_{p \times p} & X^T dB_t \\ (X^T dB_t)^T & 0_{k \times k} \end{bmatrix},$$

Dichotomy via Determinant

Theorem : For the SGF, the following property holds for the evolution of determinant

$$d(\det(\theta^T \theta)) = -2\eta\delta k \text{Tr}\{(X^T X)\} \det(\theta^T \theta) dt.$$

Hence, $\det(\theta^T \theta) = \det(\theta_0^T \theta_0) \exp\{-2\delta k \text{Tr}\{(X^T X)\} t\}$ where θ_0 is the initialization.

Comments :

- a) For GF, $\det(\theta^T \theta) = \det(\theta_0^T \theta_0) \rightarrow$ **rank preserving invariant.**
- b) For SGF, $\lim_{t \rightarrow \infty} \det(\theta(t)^T \theta(t)) \rightarrow 0 \rightarrow$ **SGF diminishes the rank** and simplifies the parameters.

Limitations :

- a) Diminishing only the **smallest singular value** is sufficient to reduce the determinant. Therefore the above result does not capture the complete simplification of the parameters
- b) The dichotomy is effective only when $\det(\theta_0^T \theta_0) \neq 0$, which occurs when $l \geq p + k$, thereby limiting the result to networks with large widths

Repulsive Force between Singularvalues

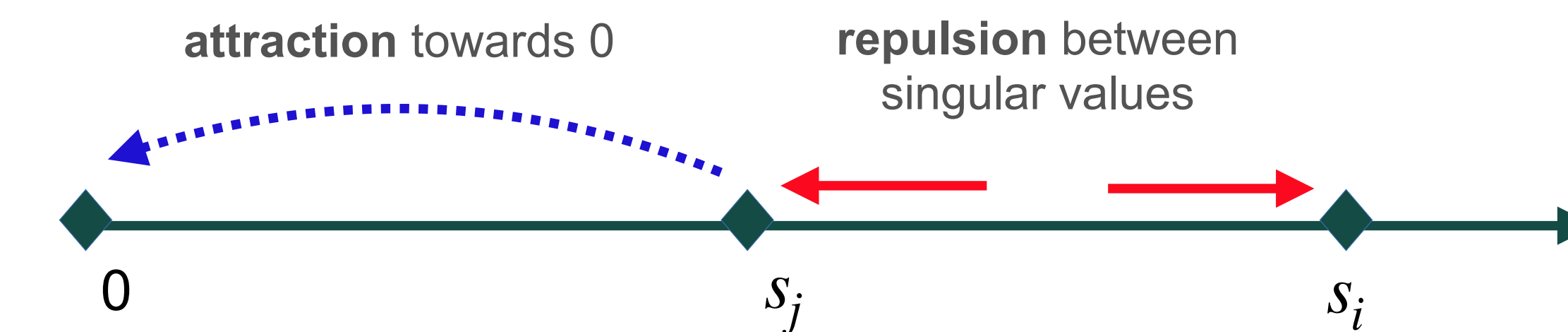
From the SDE on the matrix-valued process on θ , we can derive the SDE governing the singular values [Dyson, 1962, Bru, 1989].

Time rescaled SGF for scalar regression, $W_1 = W, W_2 = a, X = I$,
 $dW = dXa^T; \quad da = W^T dX$, where $dX = \frac{1}{\eta\delta} (Y - Wa) dt + dB_t$

SVD: Let $W = U\Sigma V^T$. Let s_1, s_2, \dots, s_l be the eigenvalues of $W^T W$ (squared singular values of W)

Theorem: Until the collision time, the evolution of eigenvalues is given by

$$d(s_i) = p c_i^2 dt + \sum_{j \neq i} \frac{s_j c_j^2 + s_i c_i^2}{s_i - s_j} dt + 2\sqrt{s_i c_i^2} (d\tilde{X})_i$$



Forces at play:

(a) A repulsion force between the eigenvalues due to the **skew-symmetric** term

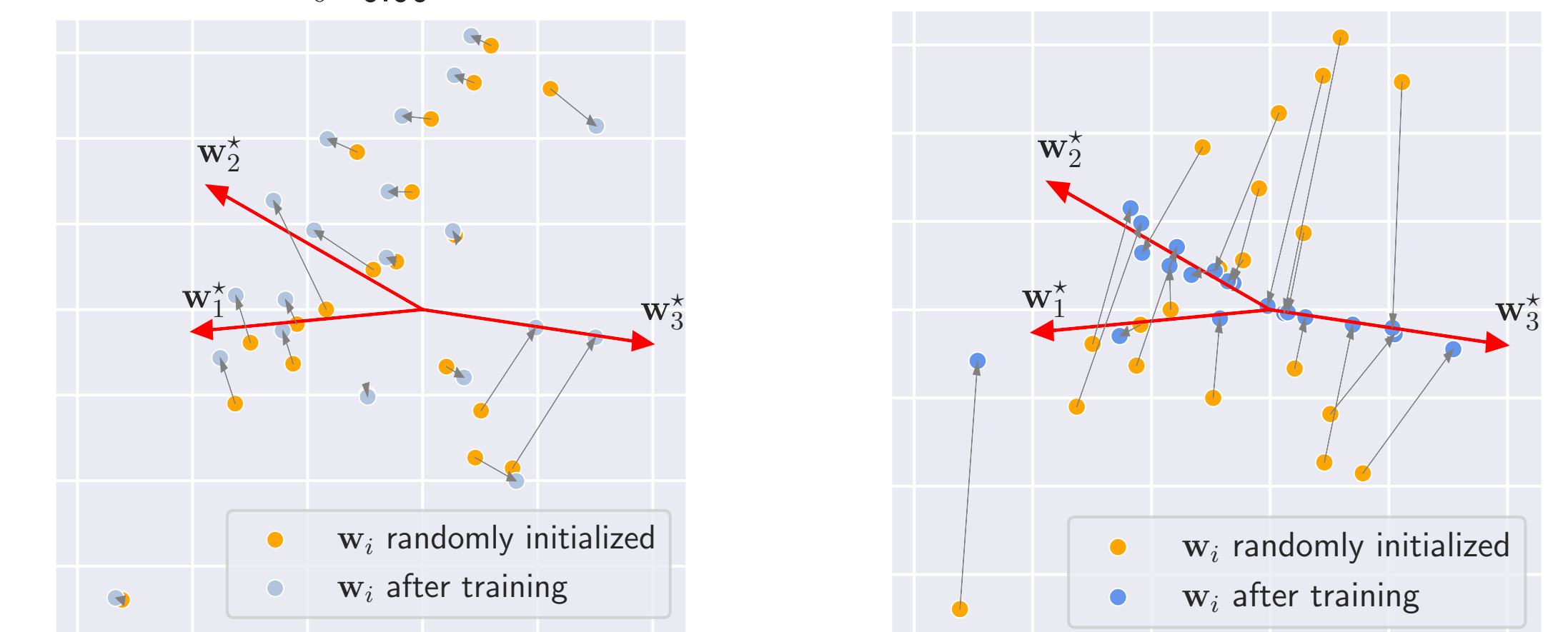
(b) A pull towards zero due to the Geometric Brownian motion

$$d \log(s_i) = p \frac{c_i^2}{s_i} dt - 2 \frac{c_i^2}{s_i} dt + \frac{1}{s_i} \sum_{j \neq i} \frac{s_j c_j^2 + s_i c_i^2}{s_i - s_j} dt + 2\sqrt{\frac{c_i^2}{s_i}} (d\tilde{X})_i$$

Empirics with ReLU: Let $k = 1$, we train one hidden layer ReLU network with GD and SGD from the same initialization.

Teacher : $y = \sum_{i=1}^3 \sigma(x^T w_i^*)$ **Student:** $y = \sum_{j=1}^m \sigma(w_j^T x)$

$\delta=0.00$ $\delta=1.00$



Neurons in ReLU networks aligns along the teacher directions when trained with a noise level $\delta = 1$.

References:

- a) Blanc et. al. Implicit regularization for deep neural networks driven by an ornstein-uhlenbeck like process, COLT 2020.
- b) Li et. al. What happens after sgd reaches zero loss?—a mathematical framework., ICLR 2021
- c) M.-F. Bru. Diffusions of perturbed principal component analysis 1989
- d) Chen et. al. Stochastic collapse: How gradient noise attracts SGD dynamics towards simpler subnetworks, NeurIPS 2023

