



The Hidden Convex Optimization Landscape of Two-Layer ReLU Networks

Victor Mercklé¹, Franck Iutzeler², Ievgen Redko³

¹ LJK, Univ. Grenoble Alpes ² IMT, Univ. Toulouse, CNRS ³ Paris Noah's Ark lab



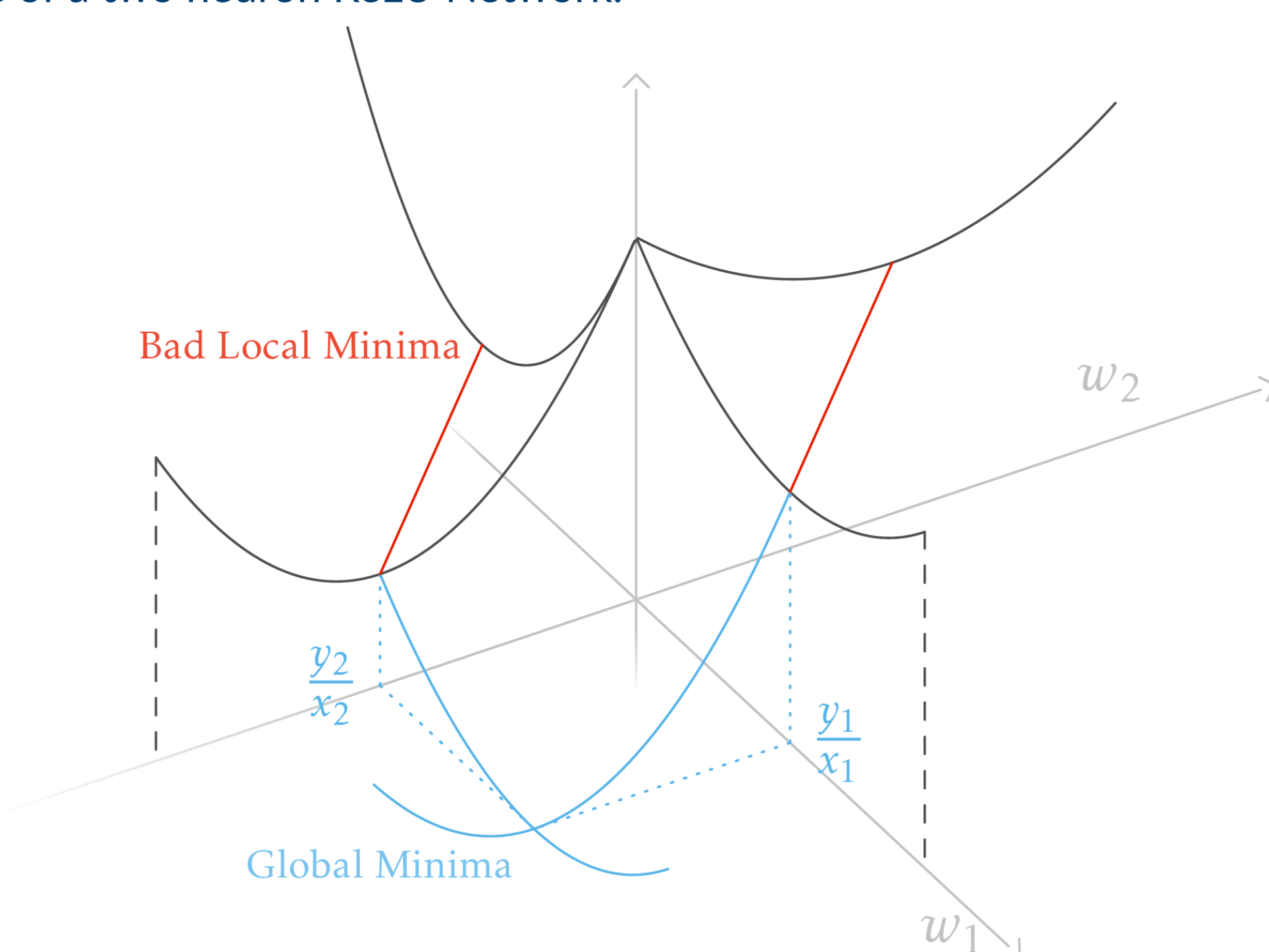
Online Blog Post

TL;DR

- Training a ReLU network is a **non-convex optimization** problem with many local minima which can be **hard** to escape in practice
- The **global minimum** is the unique solution of a derived **convex** optimization problem
- **Convex** problems have more predictable dynamics, and are **easier** to solve

Is this actually convex?

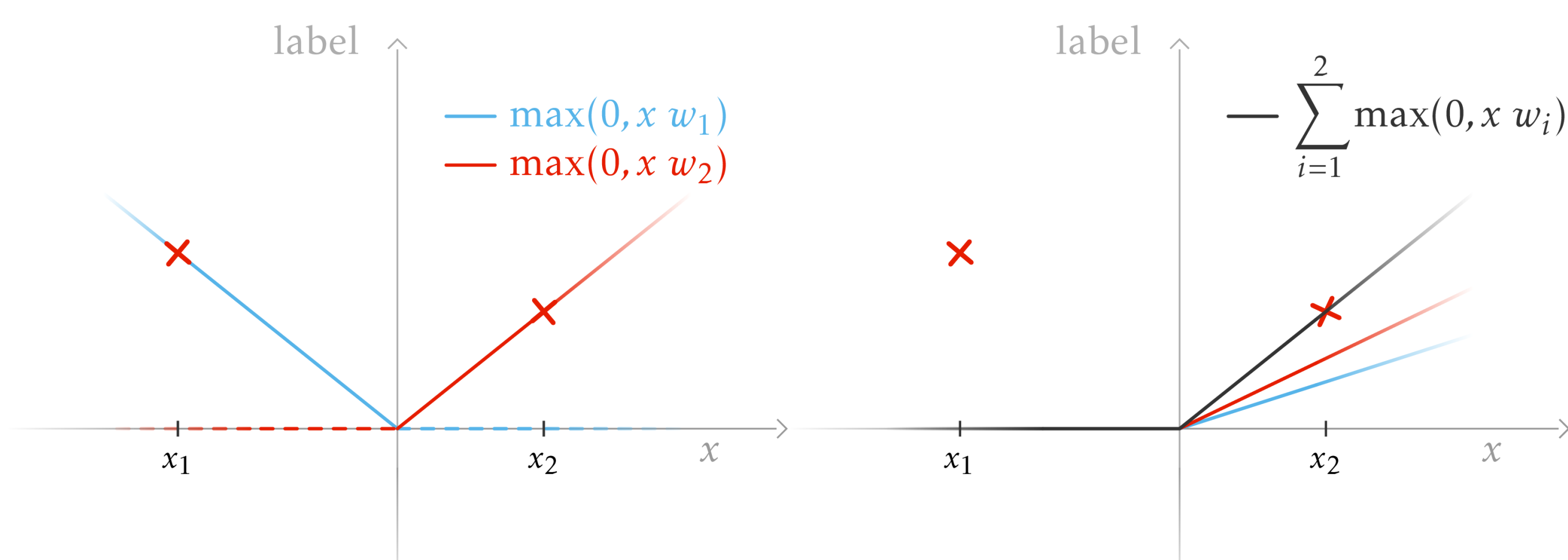
Loss landscape of a two neuron ReLU Network:



$$\mathcal{L}(w_1, w_2) = (\max(0, x_1 w_1) + \max(0, x_1 w_2) - y_1)^2 + (\max(0, x_2 w_1) + \max(0, x_2 w_2) - y_2)^2$$

- Two \times data points $x_1, x_2 \in \mathbb{R}$ with labels $y_1, y_2 \in \mathbb{R}$
- Two ReLU neurons $w_1, w_2 \in \mathbb{R}$, and second layer fixed to 1

Outputs of two networks found by gradient descent:

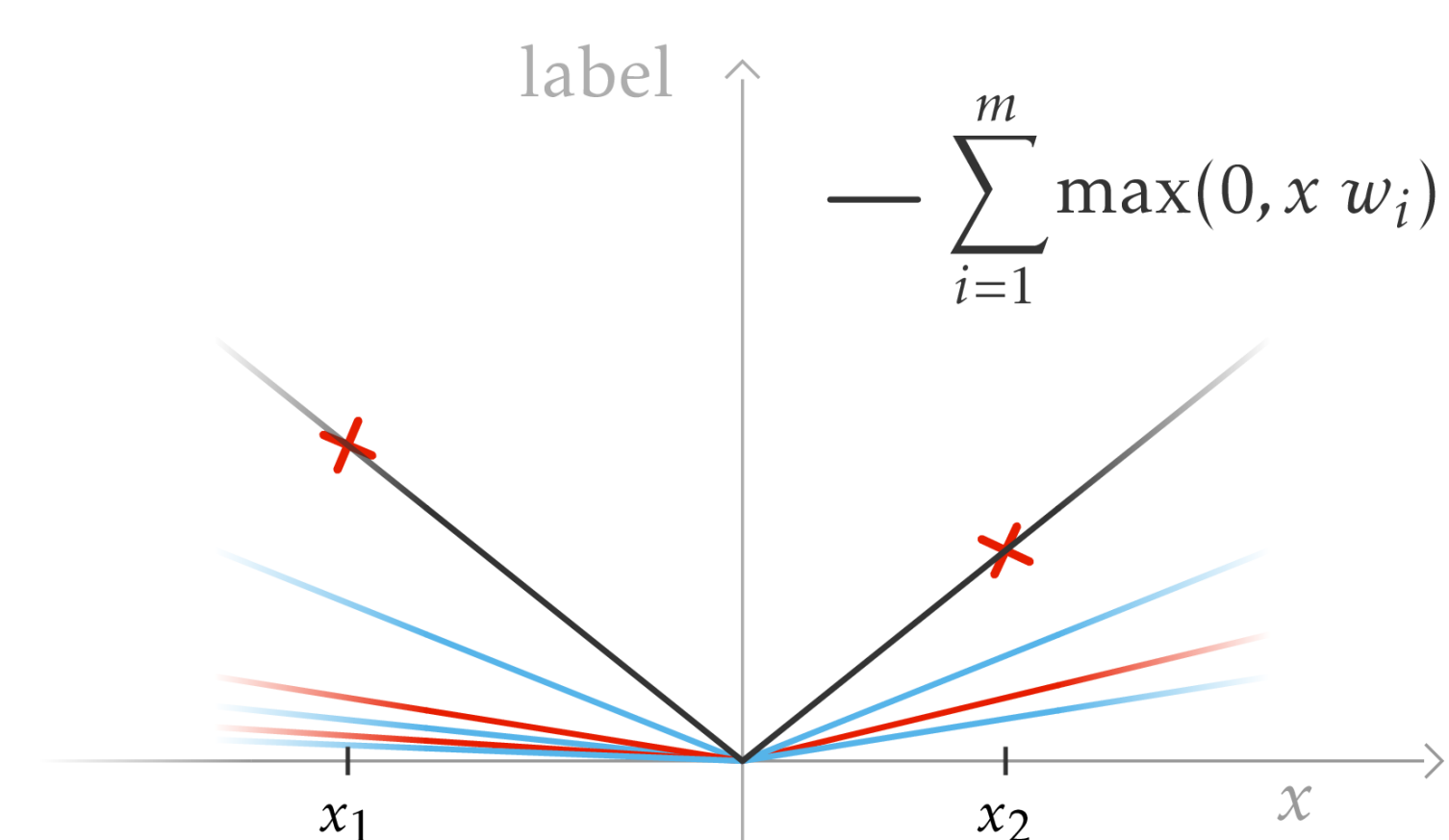


- $\begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}$ is the activation pattern of the red neuron that only activate x_2
- The blue neuron has a different activation pattern in the two local minima

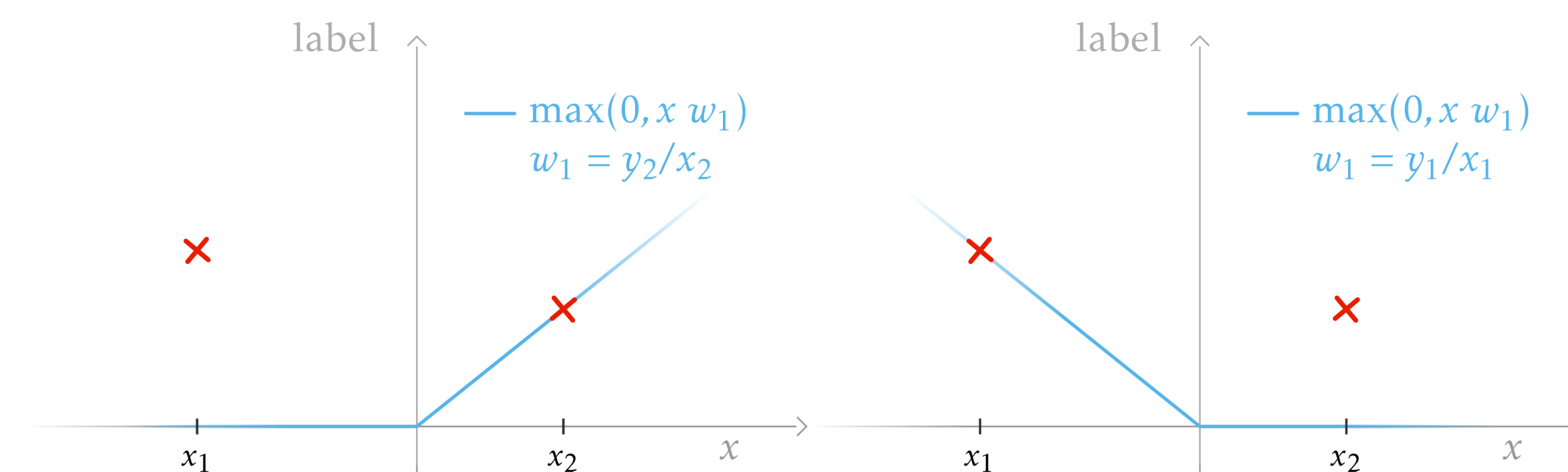
Equivalent convex problem:

$$\mathcal{L}(u_1, u_2) = \left\| \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} u_1 + \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} u_2 - \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} \right\|_2^2$$

- This convex loss has the same optimal as the non convex loss for any $m \geq 2$
- We can get the two optimal neurons from the convex solution



Minimum found by gradient descent with 10 neurons



The two local minima of a single neuron training

General Case and Result

Setting

$$\min_{W \in \mathbb{R}^{d \times m}, \alpha \in \mathbb{R}^m} \left\| \sum_{i=1}^m \max(0, X w_i) \alpha_i - y \right\|_2^2 + \lambda \sum_{i=1}^m \|w_i\|_2^2 + \alpha_i^2 \quad (1)$$

- n samples: $x_j \in \mathbb{R}^d$ (rows of X with labels $y_j \in \mathbb{R}, j = 1, \dots, n$)
- m neurons: First layer $w_i \in \mathbb{R}^d$, second layer $\alpha_i \in \mathbb{R}, i = 1, \dots, m$
- $\lambda \geq 0$ regularization, $\gamma > 0$, step size

Theorem (simplified) (Wang et al., 2022b)

- $D_i \in \{0, 1\}^{n \times n}$, one activation pattern
- $\mathcal{K}_i = \{v \in \mathbb{R}^d, \mathbb{1}_{X w_i \geq 0} = D_i\}$, convex cones

$$\min_{u_i, v_i \in \mathcal{K}_i} \left\| \sum_{i=1}^m D_i X (u_i - v_i) - y \right\|_2^2 + \lambda \sum_{i=1}^m \|u_i\|_2 + \|v_i\|_2 \quad (2)$$

Equivalence: The non convex (1) and the convex (2) problems have the same optimal loss when there are enough neurons ($m \geq n + 1$)

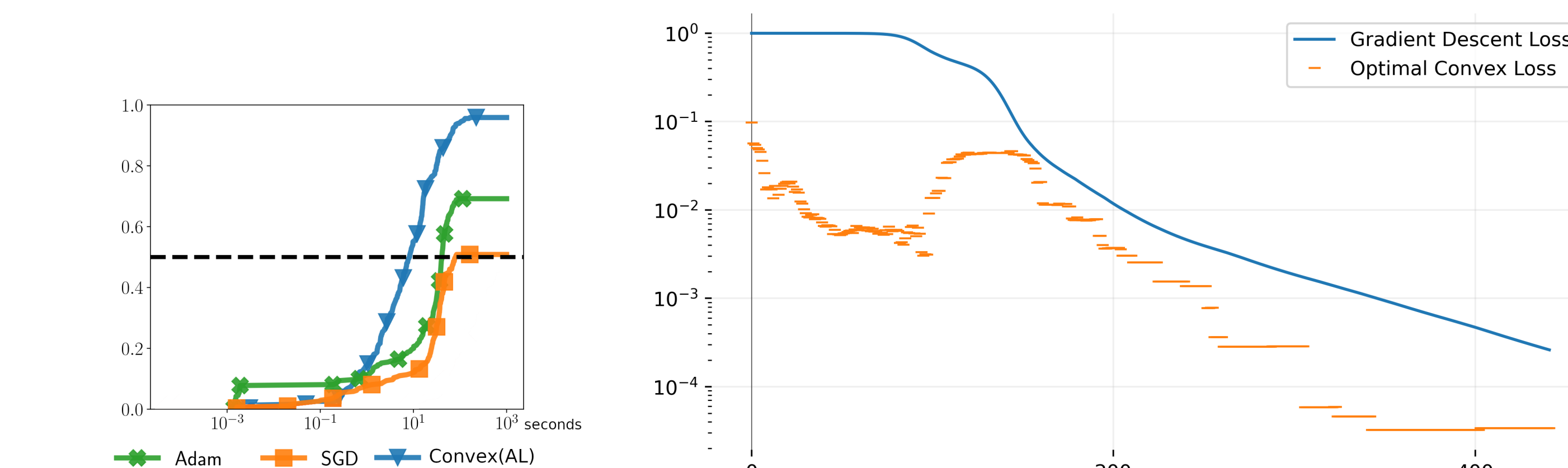
Characterization: We can find the optimal neurons for the non convex problem using the solution of the convex problem using a simple mapping:

$$\begin{aligned} (w_i, \alpha_i) &= \left(\frac{u_i}{\sqrt{\|u_i\|_2}}, \sqrt{\|u_i\|_2} \right) && \text{if } u_i \text{ is non-zero} \\ (w_i, \alpha_i) &= \left(\frac{v_i}{\sqrt{\|v_i\|_2}}, -\sqrt{\|v_i\|_2} \right) && \text{if } v_i \text{ is non-zero} \end{aligned}$$

Results for other settings

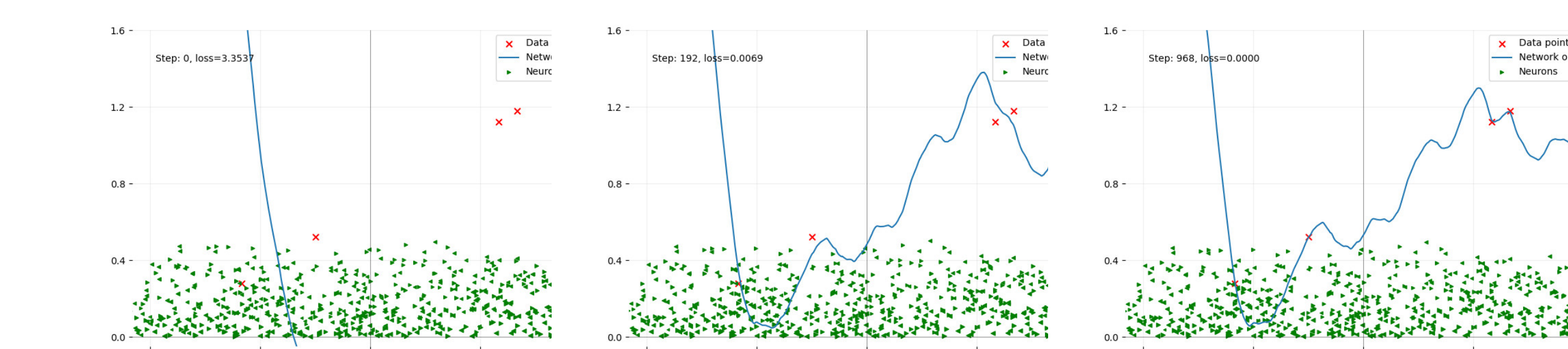
- Additional layers, by using all possible combinations of activation matrices (Ergen and Pilanci, 2021)
- Vector output, the regularization changes to require a nuclear norm
- Batch Normalization by replacing $D_i X$ with the first matrix in its Singular Value Decomposition
- Wasserstein Generative Adversarial Network, as a convex-concave game
- Parallel networks (Wang et al., 2022a)

Applications

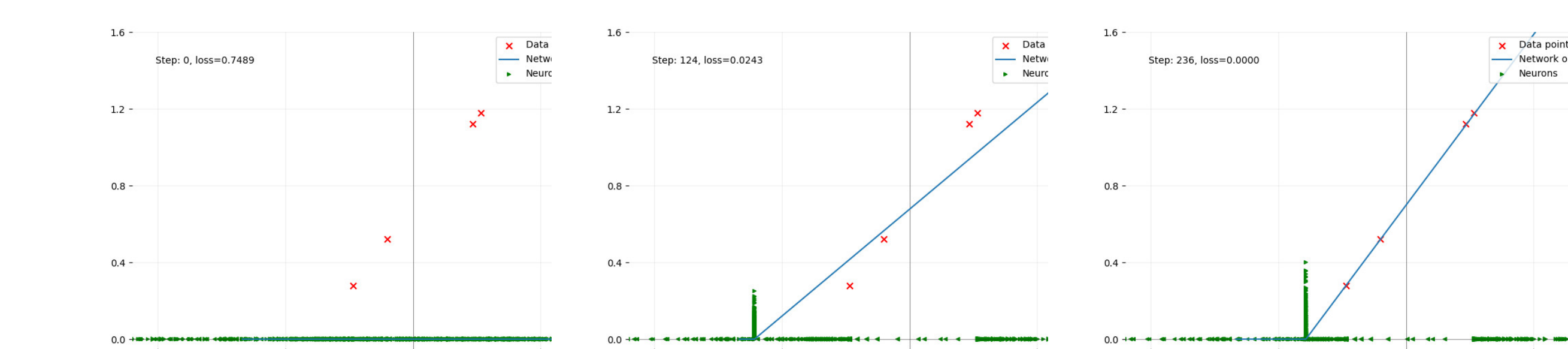


Solving the convex problem can be faster

Large scale initialization dynamic



Small scale initialization dynamic



References

- T. Ergen and M. Pilanci. Global optimality beyond two layers: Training deep relu networks via convex programs. In *International Conference on Machine Learning*, pages 2993–3003. PMLR, 2021.
- Y. Wang, T. Ergen, and M. Pilanci. Parallel deep neural networks have zero duality gap. In *The Eleventh International Conference on Learning Representations*, Sept. 2022a.
- Y. Wang, J. Lacotte, and M. Pilanci. The hidden convex optimization landscape of regularized two-layer relu networks: An exact characterization of optimal solutions. In *International Conference on Learning Representations, ICLR, 2022b*.