

Exploring the Loss Landscape of Regularized Neural Networks via Convex Duality

Sungyoon Kim¹, Aaron Mishkin², Mert Pilanci¹

¹Electrical Engineering Department

²Computer Science Department
Stanford University

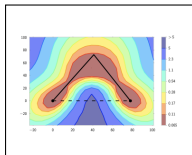
ICLR 2025

April 24th, 2025

Loss Landscape and Global Minima

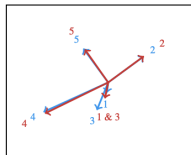
Loss Landscape and Global Minima

Mode Connectivity



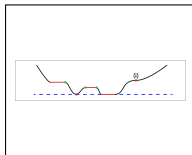
Garipov et al. (2018)

Permutation Symmetry



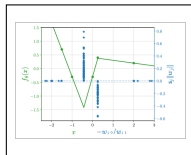
Brea et al. (2019)

Benign Landscape



Vidal et al. (2022)

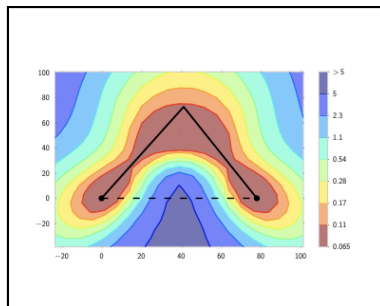
Unique Interpolator



Boursier and Flammarion (2023)

Loss Landscape and Global Minima

Mode Connectivity

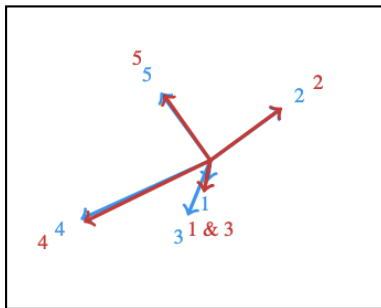


Garipov et al. (2018)

- ▶ Two different solutions are connected by a very simple curve.

Loss Landscape and Global Minima

Permutation Symmetry

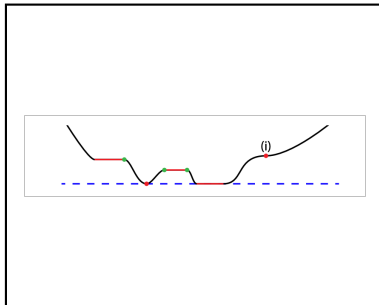


Brea et al. (2019)

- ▶ Permutations of an optimal neural network is still optimal.
- ▶ They are connected with a smooth path with low training loss.

Loss Landscape and Global Minima

Benign Landscape

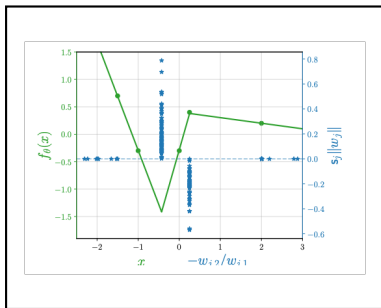


Vidal et al. (2022)

- For sufficiently wide neural networks, there is always a decreasing path to a global optimum.

Loss Landscape and Global Minima

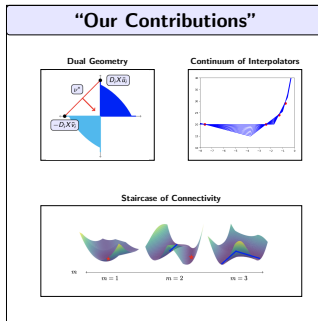
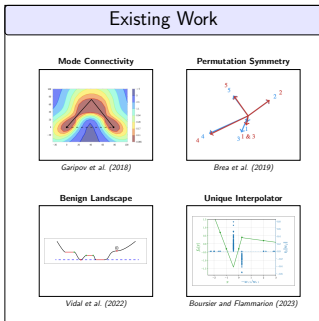
Unique Interpolator



Boursier and Flammarion (2023)

- Penalizing the bias and using free skip connections (e.g. an unregularized linear neuron).

Our Contributions



We extend our knowledge on the loss landscape and global minima of neural networks via **the convex optimization perspective**.

We use an equivalent convex optimization problem to...

- ▶ discuss **novel geometric insights** of the global minima
- ▶ **phase transition in the connectivity** as the width changes
- ▶ construct a **continuum of optimal interpolators** with regularized bias.

Extensions to vector-valued networks, parallel three-layer networks, etc.

In this talk...

We use an equivalent convex optimization problem to...

- ▶ discuss **novel geometric insights** of the global minima
- ▶ **phase transition in the connectivity** as the width changes
- ▶ construct a continuum of optimal interpolators with regularized bias.

Extensions to vector-valued networks, parallel three-layer networks, et cetera.

Background: Convex Neural Networks

- ▶ Let $X \in \mathbb{R}^{n \times d}$ be the data matrix, $y \in \mathbb{R}^n$ be the labels, $u_j \in \mathbb{R}^d$, $\alpha_j \in \mathbb{R}$ for $j = 1, 2, \dots, m$, and $\beta > 0$.
- ▶ We use $(\cdot)_+$ to denote the ReLU activation.
- ▶ Also, $L : \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}_{\geq 0}$ is a convex loss function.

Our primary interest: “two-layer”, “scalar-output”, “ReLU - activation”, “weight decay”

Background: Convex Neural Networks

The objective of interest can be written as

$$\min_{u_j \in \mathbb{R}^d, \alpha_j \in \mathbb{R}} L\left(\sum_{j=1}^m (Xu_j)_+ \alpha_j, y\right) + \frac{\beta}{2} \sum_{j=1}^m (\|u_j\|_2^2 + \|\alpha_j\|_2^2).$$

Notation

We will use the term “nonconvex objective” to refer to the above problem. Also, the optimal objective will be noted as p_{noncvx}^ and the set of optimal parameters will be noted as Θ_m^* .*

The Convex Reformulation

Pilanci and Ergen (2020) showed that if $m \geq m^*$ for some critical width m^* , there exists an **equivalent convex optimization problem**.

- ▶ When we denote p_{cvx}^* as the optimal objective of the convex problem,

$$p_{\text{cvx}}^* = p_{\text{noncvx}}^*$$

- ▶ There exists a mapping between the optimal solution of the nonconvex objective and the equivalent convex problem (Pilanci and Ergen (2020), Wang et al. (2021)).

The Convex Reformulation

The equivalent convex optimization problem is written as

$$\min_{u_i, v_i} L\left(\sum_{i=1}^P D_i X(u_i - v_i), y\right) + \beta \sum_{i=1}^P (\|u_i\|_2 + \|v_i\|_2),$$

subject to constraints $(2D_i - I)Xu_i \geq 0$, $(2D_i - I)Xv_i \geq 0$.

- ▶ Here, $D_i = \text{Diag}(1(Xh \geq 0))$ denote all possible "hyperplane arrangement patterns"
- ▶ Intuition: in the constraint set, ReLU becomes linear and

$$(Xu_i)_+ = D_i Xu_i, \quad (Xv_i)_+ = D_i Xv_i.$$

The Convex Reformulation

The equivalent convex optimization problem is written as

$$\min_{u_i, v_i} L\left(\sum_{i=1}^P D_i X(u_i - v_i), y\right) + \beta \sum_{i=1}^P (\|u_i\|_2 + \|v_i\|_2),$$

subject to constraints $(2D_i - I)Xu_i \geq 0$, $(2D_i - I)Xv_i \geq 0$.

Notation

We will use the term “convex reformulation” to refer to the above problem. Also, the optimal objective will be noted as p_{cvx}^ and the set of optimal parameters will be noted as \mathcal{P}^* .*

The Dual Problem

The dual of the convex reformulation can be written as

$$\max_{\nu \in \mathbb{R}^n} -L^*(\nu) \quad \text{subject to} \quad |\nu^T(Xu)_+| \leq \beta \quad \forall \|u\|_2 \leq 1.$$

Here, L^* is the Fenchel conjugate of $L(\cdot, y)$.

- ▶ Denote the optimal objective of the dual problem as d_{cvx}^* .
- ▶ If $m \geq m^*$,

$$p_{\text{cvx}}^* = d_{\text{cvx}}^* = p_{\text{noncvx}}^*.$$

Optimal Set Characterization

Mishkin and Pilanci (2023) show that for strictly convex L ,

- ▶ \mathcal{P}^* is a polyhedral set
- ▶ For $(u_i, v_i)_{i=1}^P, (u'_i, v'_i)_{i=1}^P \in \mathcal{P}^*$, if both u_i, u'_i are nonzero, they are positive scalings of each other.
- ▶ Θ_m^* can be characterized up to permutation/splitting symmetries.
- ▶ Our work largely builds upon this characterization, and many concepts needed for proof were adapted.

Optimal Polytope and the Dual Optimum

Theorem (The Optimal Polytope, informal)

Suppose L is a strictly convex loss function. The directions of optimal parameters of the convex problem, noted as \bar{u}_i, \bar{v}_i , **are uniquely determined from the dual optimum ν^*** . Moreover, the solution set is the **polytope**,

$$\mathcal{P}^* = \left\{ (c_i \bar{u}_i, d_i \bar{v}_i)_{i=1}^P \mid c_i, d_i \geq 0 \quad \forall i \in [P], \quad \sum_{i=1}^P D_i X \bar{u}_i c_i - D_i X \bar{v}_i d_i = y^* \right\}$$

for the unique optimal model fit y^* .

- ▶ Either $u_i^* = 0$ ($v_i^* = 0$), or
- ▶ They are positive scalings of the solution to the optimization problem,

$$\max_{\|u\|_2 \leq 1} |(\nu^*)^T (Xu)_+|,$$

Geometric Intuition: The Rectified Ellipsoid

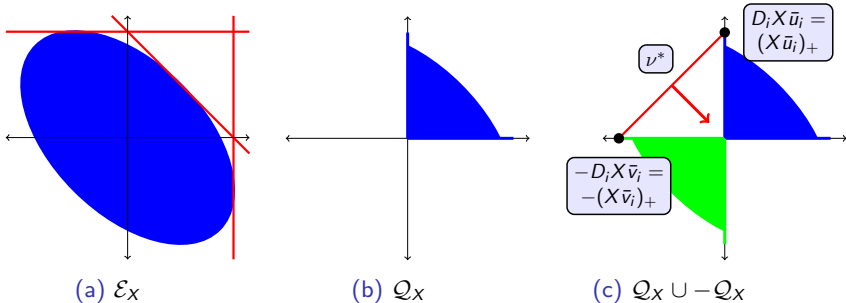
Definition (Pilanci and Ergen (2020))

The rectified ellipsoid is defined as the set

$$\mathcal{Q}_X = \left\{ (Xu)_+ \mid \|u\|_2 \leq 1 \right\}.$$

- It is the image of the ReLU mapping of the ellipsoid $\mathcal{E}_X = \{Xu \mid \|u\|_2 \leq 1\}$.

Geometric Intuition: The Rectified Ellipsoid



- ▶ The dual variable ν^* **decides the “face”** where the points $\{D_i X \bar{u}_i\}_{i=1}^P \cup \{-D_i X \bar{v}_i\}_{i=1}^P$ lies on.
- ▶ Blue set corresponds to \bar{u}_i , green set corresponds to \bar{v}_i .

The Staircase of Connectivity: Motivation

- ▶ Denote $\text{card}((u_i, v_i)_{i=1}^P)$ as the number of nonzero vectors in $\{u_i\}_{i=1}^P \cup \{v_i\}_{i=1}^P$.
- ▶ A solution map exists between Θ_m^* and the cardinality - constrained set

$$\mathcal{P}_m^* = \left\{ (u_i, v_i)_{i=1}^P \in \mathcal{P}^* \mid \text{card}((u_i, v_i)_{i=1}^P) \leq m \right\}.$$

The Staircase of Connectivity: Motivation

- ▶ Though \mathcal{P}^* is a connected set, \mathcal{P}_m^* **might not be connected** due to cardinality constraints - phase transitional behavior in connectivity!
- ▶ As Θ_m^* and \mathcal{P}_m^* are related by the solution map, connectivity properties of Θ_m^* can be deduced from that of \mathcal{P}_m^* .

The Staircase of Connectivity

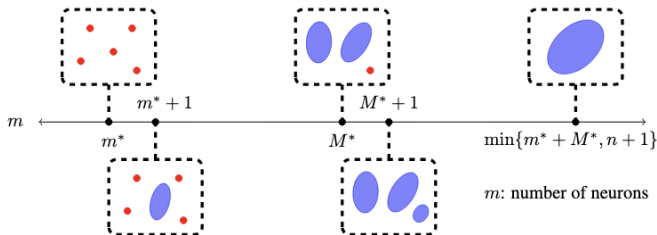


Figure: A schematic for the staircase of connectivity

- ▶ Red dots : isolated points, blue sets : connected components with more than one point.
- ▶ Critical widths m^* , M^* governs the phase transition.

The Staircase of Connectivity

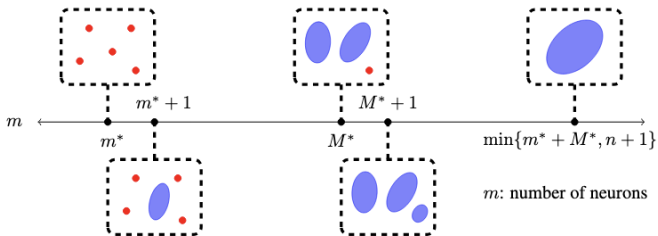


Figure: A schematic for staircase of connectivity

- ▶ When $m = m^*$, Θ_m^* is a set of finite isolated points.

The Staircase of Connectivity

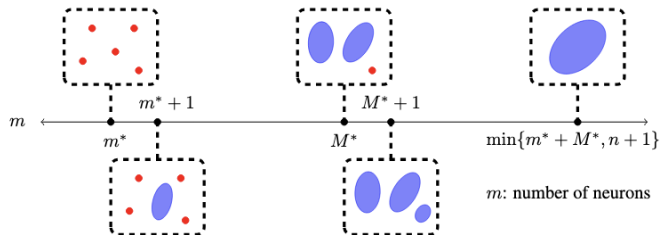


Figure: A schematic for staircase of connectivity

- ▶ When $m \geq m^* + 1$, there exists a path between two different optimal solutions in Θ_m^*

The Staircase of Connectivity

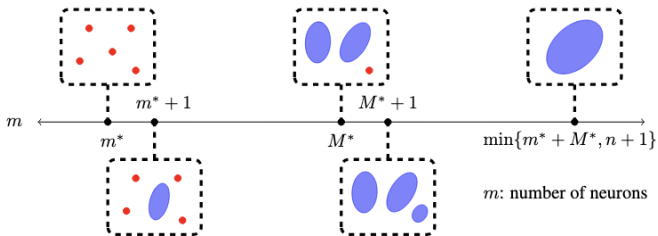


Figure: A schematic for staircase of connectivity

- ▶ When $m = M^*$, there exists an isolated point in Θ_m^* .

The Staircase of Connectivity

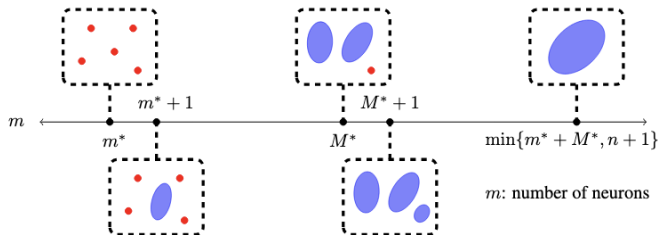


Figure: A schematic for staircase of connectivity

- ▶ When $m \geq M^* + 1$, there is no isolated point in Θ_m^* .
- ▶ Moreover, for an optimal solution $(w_i, \alpha_i)_{i=1}^m$, any permutation $(w_{\sigma(i)}, \alpha_{\sigma(i)})_{i=1}^m$ has a path inside Θ_m^* that connects the two solutions.

The Staircase of Connectivity

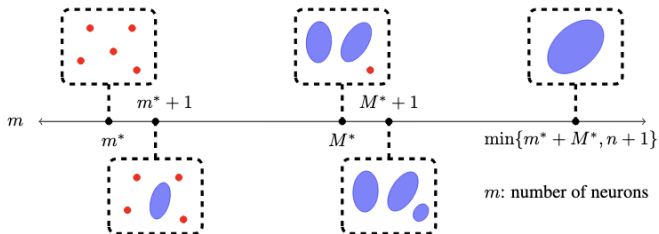


Figure: A schematic for staircase of connectivity

- ▶ When $m \geq \min\{M^* + m^*, n + 1\}$, Θ_m^* is connected.

Relations with Existing Landscape Results

- ▶ Haefele and Vidal (2017) shows that when $m \geq n + 1$, there is no spurious local minima. We further characterize that **all sublevel sets are connected**.
- ▶ Nguyen (2021) shows that when there is no regularization, Θ_m^* is connected when $m \geq n + 1$. **We extend the result to the regularized case.**
- ▶ Our analysis is also tightly connected to Simsek et al. (2021), who add a neuron to connect permutations of optimal solutions.

Conclusion

- ▶ We derived novel characterizations of the loss landscape and global minima of neural networks by leveraging tools from convex optimization.
- ▶ An extension of these results to different nonconvex problems that has convex reformulations could be an interesting future direction.

Poster session: Hall 3 + Hall 2B #350, today 3pm – 5:30pm

Bibliography

- Boursier, E. and Flammarion, N. (2023). Penalising the biases in norm regularisation enforces sparsity. *Advances in Neural Information Processing Systems*, 36:57795–57824.
- Brea, J., Simsek, B., Illing, B., and Gerstner, W. (2019). Weight-space symmetry in deep networks gives rise to permutation saddles, connected by equal-loss valleys across the loss landscape. *arXiv preprint arXiv:1907.02911*.
- Garipov, T., Izmailov, P., Podoprikin, D., Vetrov, D. P., and Wilson, A. G. (2018). Loss surfaces, mode connectivity, and fast ensembling of dnns. *Advances in neural information processing systems*, 31.
- Haeffele, B. D. and Vidal, R. (2017). Global optimality in neural network training. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7331–7339.
- Mishkin, A. and Pilanci, M. (2023). Optimal sets and solution paths of relu networks. In *International Conference on Machine Learning*, pages 24888–24924. PMLR.
- Nguyen, Q. (2021). A note on connectivity of sublevel sets in deep learning. *arXiv preprint arXiv:2101.08576*.
- Pilanci, M. and Ergen, T. (2020). Neural networks are convex regularizers: Exact polynomial-time convex optimization formulations for two-layer networks. In *International Conference on Machine Learning*, pages 7695–7705. PMLR.
- Simsek, B., Ged, F., Jacot, A., Spadaro, F., Hongler, C., Gerstner, W., and Brea, J. (2021). Geometry of the loss landscape in overparameterized neural networks: Symmetries and invariances. In *International Conference on Machine Learning*, pages 9722–9732. PMLR.
- Vidal, R., Zhu, Z., and Haeffele, B. D. (2022). Optimization landscape of neural networks. *Mathematical Aspects of Deep Learning*, page 200.
- Wang, Y., Lacotte, J., and Pilanci, M. (2021). The hidden convex optimization landscape of regularized two-layer relu networks: an exact characterization of optimal solutions. In *International Conference on Learning Representations*.