

---

# Geometry of the Loss Landscape in Overparameterized Neural Networks: Symmetries and Invariances

---

Berfin Şimşek<sup>1,2</sup> François Ged<sup>1</sup> Arthur Jacot<sup>1</sup> Francesco Spadaro<sup>1</sup>  
Clément Hongler<sup>1\*</sup> Wulfram Gerstner<sup>2\*</sup> Johanni Brea<sup>2\*</sup>

## Abstract

We study how permutation symmetries in overparameterized multi-layer neural networks generate ‘symmetry-induced’ critical points. Assuming a network with  $L$  layers of minimal widths  $r_1^*, \dots, r_{L-1}^*$  reaches a zero-loss minimum at  $r_1^*! \cdots r_{L-1}^*!$  isolated points that are permutations of one another, we show that adding one extra neuron to each layer is sufficient to connect all these previously discrete minima into a single manifold. For a two-layer overparameterized network of width  $r^* + h =: m$  we explicitly describe the manifold of global minima: it consists of  $T(r^*, m)$  affine subspaces of dimension at least  $h$  that are connected to one another. For a network of width  $m$ , we identify the number  $G(r, m)$  of affine subspaces containing only symmetry-induced critical points that are related to the critical points of a smaller network of width  $r < r^*$ . Via a combinatorial analysis, we derive closed-form formulas for  $T$  and  $G$  and show that the number of symmetry-induced critical subspaces dominates the number of affine subspaces forming the global minima manifold in the mildly overparameterized regime (small  $h$ ) and vice versa in the vastly overparameterized regime ( $h \gg r^*$ ). Our results provide new insights into the minimization of the non-convex loss function of overparameterized neural networks.

## 1. Introduction

Neural network landscapes were traditionally thought of as highly non-convex landscapes, where non-global critical points may harm gradient-descent by slowing it down (due to saddles) or making it stop in local minima. Earlier works have argued in favor of a proliferation of saddles in high-dimensional neural network landscapes through an analogy with random error functions (Dauphin et al., 2014). On the other hand, practical neural network landscapes are found to exhibit surprising properties, such as the connectivity of global minima (Draxler et al., 2018; Garipov et al., 2018) and the convergence to a global minimum in the so-called overparameterized regime (Jacot et al., 2018), thereby ruling out proliferating saddles as a problem in this regime. Yet, in mildly overparameterized networks, gradient descent may find a global minimum only for a small fraction of random initializations (Sagun et al., 2014; Chizat & Bach, 2018; Frankle & Carbin, 2018).

In this work, we study the width-dependent scaling of the number of symmetry-induced critical points and the connectivity of global minima by exploiting the permutation symmetry and further invariances of the network parameterization. The permutation symmetry introduces an invariance to a permutation in parameterization that is characteristic for many machine learning models beyond neural networks, such as mixture models, multiple kernel learning, or matrix factorization.

Further invariances in a neural network of width  $m$  induce equal loss manifolds such that all points in the manifold are equivalent to a single point in a narrower network of width  $r < m$ . The mapping approach from a point in parameter space of the narrower network to a parameter manifold of the full network is particularly useful for the study of critical points as critical points of the narrow network turn into symmetry-induced critical subspaces of the full one. In particular, a global minimum of the narrow network turns into a collection of global minima subspaces that are connected to one another.

---

\*Shared senior authorship <sup>1</sup>Chair of Statistical Field Theory, École Polytechnique Fédérale de Lausanne, Switzerland <sup>2</sup>Laboratory of Computational Neuroscience, École Polytechnique Fédérale de Lausanne, Switzerland. Correspondence to: Berfin Şimşek <berfin.simsek@epfl.ch>.

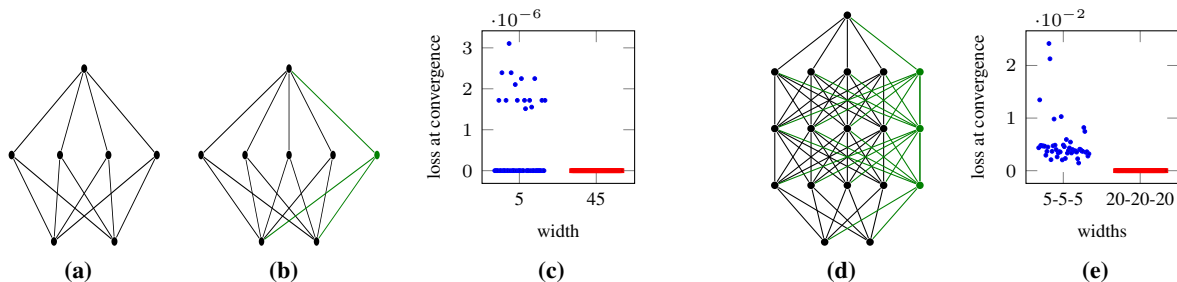


Figure 1. Graph of (a) a minimal network of width 4 (teacher) and (b) a mildly overparameterized student network of width 5. (c) With 50 random initializations, mildly overparameterized networks (blue) find a global minimum for only a fraction of initializations, whereas vastly overparameterized networks (red, width 45) consistently find a global minimum. (d) Graph of student network with three hidden layer learning from a teacher with widths 4 – 4 – 4. (e) Vastly overparameterized networks (red) consistently find a global minimum whereas mildly overparameterized networks (blue) typically do not.

### 1.1. Main Contributions

1. Suppose an  $L$ -layer Artificial Neural Network (ANN) with hidden layer widths  $r_1^*, \dots, r_{L-1}^*$  reaches a unique (up to permutation) zero-loss global minimum (we call such a network *minimal* if it cannot achieve zero loss if any neuron is removed). The permutation symmetries give rise to  $r_1^! \cdots r_{L-1}^!$  equivalent discrete global minima. We show that adding one neuron to each layer is sufficient to connect these global minima into a single zero-loss manifold.
2. For a two-layer overparameterized network of width  $m = r^* + h$ , we describe the geometry of the global minima manifold precisely: it consists of a union of a number  $T(r^*, m)$  of affine subspaces of dimension  $\geq h$  and it is connected. Furthermore, we show that the global minima manifold contains *all* the zero-loss points for smooth activation functions satisfying a technical condition and in the presence of infinitely many data points with full support of the input space.
3. The symmetries of the network generate symmetry-induced critical points, such as saddle points, which may prevent the convergence to a global minimum (see Figure 1). We find a surprising scaling relation between the number of subspaces formed by the symmetry-induced critical points and the number of subspaces making up the global minima:
  - When the number of additional neurons satisfies  $h \ll r^*$  (i.e. at the beginning of the overparameterized regime), the number of subspaces that make up the global minima manifold is much *smaller* than the number of subspaces that make up the symmetry-induced critical points. In this sense, there is a proliferation of saddles and the global minima manifold is ‘tiny’.
  - Conversely, when  $h \gg r^*$  (i.e. we are far into or within the overparameterized regime), the number of subspaces that make up the global minima

manifold is much *greater* than the number of subspaces that make up the symmetry-induced (non-global) critical points. In this sense the global minima manifold is ‘huge’.

4. One may worry that, by adding  $h$  neurons, a saddle of a network of width  $r$  could transform into a local minimum. However, we show that this is not the case and a saddle point in the smaller network transforms into symmetry-induced saddle points.

### 1.2. Related Work

A number of recent works have explored the typical path taken by a gradient-based optimizer. For very wide ANNs, the gradient flow converges to a global minimum in spite of the non-convexity of the loss (Jacot et al., 2018; Du et al., 2018; Chizat & Bach, 2018; Arora et al., 2019; Du et al., 2019; Lee et al., 2019a; 2020). First-order gradient algorithms provably escape strict saddles (Jin et al., 2017; Lee et al., 2019b), although they can face an exponential slowdown around these saddles (Du et al., 2017). For pruned ANNs, the training with typical (random) initialization does not reach any global minimum, in spite of their presence in the landscape (Frankle & Carbin, 2018).

Another body of work focuses on the geometric investigation of neural network landscapes. Dauphin et al. (2014) suggested a proliferation of saddles in ANN landscapes through an analogy with high-dimensional Gaussian Processes. Other models have been proposed to understand the general structure of ANN landscapes inspired by statistical physics (Geiger et al., 2019), and via high-dimensional wedges (Fort & Jastrzebski, 2019). These model-based empirical works focus mainly on the Hessian spectrum at the critical points.

Another line of work suggests that global minima found by stochastic gradient descent are connected (i.e. there is a path linking arbitrary two minima along which the loss increases only negligibly) via simply parameterized low-

loss curves (Draxler et al., 2018; Garipov et al., 2018) or line segments (Sagun et al., 2017; Frankle et al., 2020; Fort et al., 2020). Theoretical work limited to ReLU-type activation functions, showed that in overparameterized networks, all global minima lie in a connected manifold (Freeman & Bruna, 2016; Nguyen, 2019), however without giving a geometrical description of this manifold. Cooper (2020) studied the geometry of a subset of the manifolds of critical points. Kuditiipudi et al. (2019) showed that the global minima for ReLU networks, for which *half* of the neurons can be dropped without incurring a significant increase in loss, are connected via piecewise linear paths of minimal cost.

In this paper, we show that adding or removing a *single* neuron radically changes the connectedness without any change in loss. We are the first to prove the connectivity of the global minima manifold for continuously differentiable activation functions. The focus on symmetries in our work is similar to that of (Fukumizu & Amari, 2000; Brea et al., 2019; Fukumizu et al., 2019) regarding the critical points coming from neuron replications. In an orthogonal direction, Kunin et al. (2020); Głuch & Urbanke (2021) present a catalog of symmetries appearing in deep networks, which however does not include the permutation symmetry. To the best of our knowledge, this work is the first to study the scaling of the number of critical points in ANN landscapes as a function of the overparameterization amount. A key challenge to overcome is the numerous equivalent arrangements of neurons inside the network.

**Notation.** For  $m \geq 1$ , set  $[m] = \{1, \dots, m\}$  and let  $S_m$  denote the symmetric group on  $m$  symbols, i.e. the set of permutations of  $[m]$ . For a permutation  $\pi \in S_m$  and  $D \geq 1$ , the map  $\mathcal{P}_\pi : \mathbb{R}^{Dm} \rightarrow \mathbb{R}^{Dm}$  permutes the units  $\vartheta_i \in \mathbb{R}^D$  of a point  $\boldsymbol{\theta} = (\vartheta_1, \dots, \vartheta_m)$  according to  $\pi$ , i.e.  $\mathcal{P}_\pi \boldsymbol{\theta} = (\vartheta_{\pi(1)}, \dots, \vartheta_{\pi(m)})$ ; we sometimes use  $\boldsymbol{\theta}_\pi := \mathcal{P}_\pi \boldsymbol{\theta}$ .

## 2. Symmetric Losses

Numerous machine learning models involve permutation-symmetric parameterizations: mixture models, matrix factorization, and neural networks. In this section, we abstract away the particular parameterization of these models and focus on the implications of permutation symmetry on the gradient flow. In particular, the discussion here is general and applies to ANNs which is the main focus of this paper.

**Definition 2.1.** A loss function  $L^m : \mathbb{R}^{Dm} \rightarrow \mathbb{R}$  is a **symmetric loss**<sup>1</sup> on  $m$  units if it is a  $C^1$  function and if for any  $\pi \in S_m$  and any  $\boldsymbol{\theta} = (\vartheta_1, \vartheta_2, \dots, \vartheta_m)$  with  $\vartheta_i \in \mathbb{R}^D$ , we have

$$L^m(\boldsymbol{\theta}) = L^m(\mathcal{P}_\pi \boldsymbol{\theta}).$$

<sup>1</sup>When the units are 1-dimensional, symmetric losses are symmetric functions (Kung et al., 2009; Sagan, 2013).

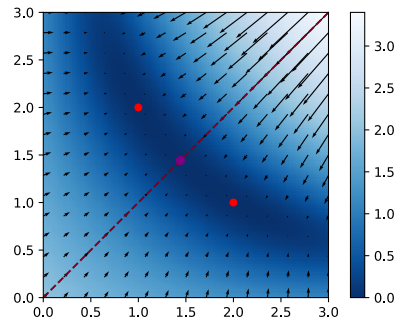


Figure 2. No gradient pointing outside of a symmetry subspace. The gradient flow of a permutation-symmetric loss  $L(w_1, w_2) = \log(\frac{1}{2}((w_1 + w_2 - 3)^2 + (w_1 w_2 - 2)^2) + 1)$ . Red: permutation-symmetric global minima, purple: saddle, dashed line: the symmetry subspace.

The term *unit* may refer to a Gaussian vector in the context of Gaussian mixture models, to a factor in the context of matrix factorization, or to a neuron in the context of neural networks. The symmetry subspaces are defined by the constraint that at least two units are identical:

**Definition 2.2.** Let  $i_1, \dots, i_k \in [m]$  be distinct indices. The **symmetry subspace**  $\mathcal{H}_{i_1, \dots, i_k}$  is defined as

$$\mathcal{H}_{i_1, \dots, i_k} := \{(\vartheta_1, \dots, \vartheta_m) \in \mathbb{R}^{Dm} : \vartheta_{i_1} = \dots = \vartheta_{i_k}\}.$$

As each constraint  $\vartheta_i = \vartheta_j$  suppresses  $D$  degrees of freedom, we have  $\dim(\mathcal{H}_{i_1, \dots, i_k}) = D(m - k + 1)$ . The largest symmetry subspaces are  $\mathcal{H}_{i,j}$ 's: any other symmetry subspace is the intersection of such subspaces.

Let  $\boldsymbol{\rho} : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^{Dm}$  denote the gradient flow under a symmetric loss

$$\dot{\boldsymbol{\rho}}(t) = -\nabla L^m(\boldsymbol{\rho}(t)) \quad (1)$$

for  $t \geq 0$  and for a given initialization  $\boldsymbol{\rho}(0)$ . In Figure 2, we observe that the gradient on the symmetry subspace is tangent to it. In general, the gradient components of a symmetry subspace pointing to neighbor regions cancel out due to permutation symmetry.

**Lemma 2.1.** Let  $L^m : \mathbb{R}^{Dm} \rightarrow \mathbb{R}$  be a symmetric loss on  $m$  units thus a  $C^1$  function and let  $\boldsymbol{\rho} : \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^{Dm}$  be its gradient flow. If  $\boldsymbol{\rho}(0) \in \mathcal{H}_{i_1, \dots, i_k}$ , the gradient flow stays inside the symmetry subspace, i.e.  $\boldsymbol{\rho}(t) \in \mathcal{H}_{i_1, \dots, i_k}$  for all  $t > 0$ . If  $\boldsymbol{\rho}(0) \notin \mathcal{H}_{i,j}$  for all  $i \neq j \in [m]$ , that is outside of all symmetry subspaces, the gradient flow does not visit any symmetry subspace in finite time.

**Remark 2.1.** Lemma 2.1 does not exclude the following scenario: if there is a critical point on the symmetry subspace that is attractive in some directions orthogonal to the symmetry subspace, the gradient flow can reach it in infinite time (i.e. at convergence).

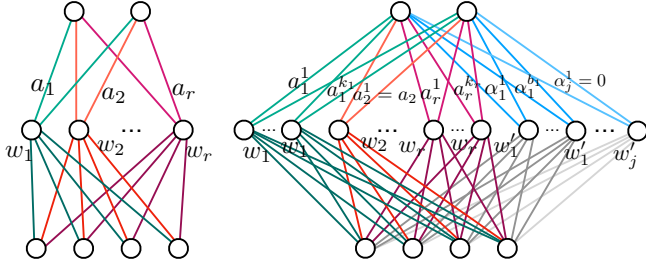


Figure 3. **Left:** Parameters  $\theta^r$  of an irreducible point in a network of  $r$  neurons with  $w_i \neq w_j$  for all  $i \neq j$  and  $a_i \neq 0$  for all  $i$ . **Right:** example of a reducible point in  $\Gamma_s(\theta^r)$  in an expanded network of  $m > r$  neurons. The incoming weight vector of the first neuron is replicated  $k_1$  times, the second one only once, etc.

### 3. Foundations: Invariances in 2-Layer ANNs

In this section, we discuss the implications of the permutation symmetry for the ANN landscapes and identify further invariances in network function parameterization. This approach will allow us to describe the precise geometry of the global minima manifold (Subsection 4.1) and the symmetry-induced critical points (Subsection 4.2) in overparameterized ANNs.

Let  $f^{(2)} : \mathbb{R}^{d_{\text{in}}} \rightarrow \mathbb{R}^{d_{\text{out}}}$  be a two-layer ANN of width  $m$

$$f^{(2)}(x|\theta) = \sum_{i=1}^m a_i \sigma(w_i \cdot x)$$

where  $\theta = (w_1, \dots, w_m, a_1, \dots, a_m)$  is an  $m$ -neuron point in the parameter space  $\mathbb{R}^{D^m}$  with  $w_i \in \mathbb{R}^{d_{\text{in}}}$  and  $a_i \in \mathbb{R}^{d_{\text{out}}}$  so that  $D = d_{\text{in}} + d_{\text{out}}$  and  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  is a  $C^1$  activation function with  $\sigma(x) \neq 0$  for all  $x \in \mathbb{R}$ .<sup>2</sup> Sometimes, we will write  $\theta^m := \theta$  to emphasize the number of neurons.

The training dataset of size  $N$  is denoted by  $\text{Trn} = \{(x_k, y_k)\}_{k=1}^N$  where  $x_k \in \mathbb{R}^{d_{\text{in}}}$ ,  $y_k \in \mathbb{R}^{d_{\text{out}}}$ . The training loss  $L^m : \mathbb{R}^{D^m} \rightarrow \mathbb{R}$  is

$$L^m(\theta) = \frac{1}{N} \sum_{(x,y) \in \text{Trn}} c(f^{(2)}(x|\theta), y) \quad (2)$$

where  $c : \mathbb{R}^{d_{\text{out}}} \times \mathbb{R}^{d_{\text{out}}} \rightarrow [0, +\infty)$  is a single-sample loss that is  $C^1$  in its first component and  $c(\hat{y}, y) = 0$  if and only if  $\hat{y} = y$ , such as the least-squares loss or the logistic loss.

Since  $f^{(2)}(x|\theta)$  is invariant under the permutation of neurons  $\vartheta_i := [w_i, a_i] \in \mathbb{R}^D$ , the concatenation of the incoming and outgoing weight vectors, and both  $\sigma$  and  $c$  are  $C^1$ ,  $L^m$  is a symmetric loss (Def. 2.1). Therefore the symmetry subspaces  $\vartheta_i = \vartheta_j$  are invariant under the gradient flow (Lemma 2.1). ANN functions exhibit further invariances:

<sup>2</sup>We exclude homogenous activation functions, such as ReLU and linear function (for linear networks), where the scaling invariance should also be considered.

**Definition 3.1.** We call an  $m$ -neuron point  $\theta^m$  **irreducible** if it has  $m$  distinct incoming weight vectors  $w_i$ , and no zero outgoing weight vector, i.e.  $a_i \neq 0$  for all  $i \in [m]$ . Otherwise we say that  $\theta^m$  is **reducible**.

Any reducible point  $\theta^m$  is equivalent to a point  $\theta^{m-1}$  with  $(m-1)$ -neurons in that they produce the same function  $f^{(2)}(x|\theta^m) = f^{(2)}(x|\theta^{m-1})$  where  $\theta^{m-1}$  is

1.  $(w_2, w_3, \dots, w_m, a_1 + a_2, a_3, \dots, a_m)$  if  $w_1 = w_2$ ,
2.  $(w_2, w_3, \dots, w_m, a_2, a_3, \dots, a_m)$  if  $a_1 = 0$ .

Note that because of permutation symmetry, the above reductions hold whenever two incoming weight vectors are equal, i.e.  $w_i = w_j$ , or any one of the outgoing vectors is zero  $a_i = 0$ . Moreover, if  $\theta^{m-1}$  is also reducible, we can continue dropping neurons as above until we find an irreducible point  $\theta^r$ . Equivalently (going in the opposite direction), an irreducible  $r$ -neuron point

$$\theta^r = (w_1, \dots, w_r, a_1, \dots, a_r)$$

yields an affine subspace of equal loss points in a network with width  $m \geq r$  (see Figure 3):

**Definition 3.2.** For  $r \geq 1$ ,  $j \geq 0$  with  $r + j \leq m$ , let  $s = (k_1, \dots, k_r, b_1, \dots, b_j)$  be an  $(r+j)$ -tuple of integers such that  $\text{sum}(s) := k_1 + \dots + k_r + b_1 + \dots + b_j = m$  with  $k_i \geq 1$  and  $b_i \geq 0$ . The **affine subspace**  $\Gamma_s(\theta^r)$  of an irreducible point  $\theta^r$  is

$$\begin{aligned} & \{ \underbrace{(w_1, \dots, w_1)}_{k_1}, \dots, \underbrace{(w_r, \dots, w_r)}_{k_r}, \underbrace{(w'_1, \dots, w'_1)}_{b_1}, \dots, \underbrace{(w'_j, \dots, w'_j)}_{b_j}, \\ & a_1^1, \dots, a_1^{k_1}, \dots, a_r^1, \dots, a_r^{k_r}, \alpha_1^1, \dots, \alpha_1^{b_1}, \dots, \alpha_j^1, \dots, \alpha_j^{b_j} \} : \\ & \text{where } \sum_{i=1}^{k_t} a_t^i = a_t \text{ for } t \in [r] \text{ and } \sum_{i=1}^{b_t} \alpha_t^i = 0 \text{ for } t \in [j]. \end{aligned} \quad (3)$$

Note that all  $\theta^m \in \Gamma_s(\theta^r)$  implement the same function:

$$\begin{aligned} f^{(2)}(x|\theta^m) &= \sum_{t=1}^r \sum_{i=1}^{k_t} a_t^i \sigma(w_t \cdot x) + \sum_{t=1}^j \sum_{i=1}^{b_t} \alpha_t^i \sigma(w'_t \cdot x) \\ &= f^{(2)}(x|\theta^r). \end{aligned}$$

Neurons with incoming weight vectors  $w'$  and outgoing weight vectors adding up to zero are called in the following **‘zero-type’ neurons**. Moreover, the network function remains invariant under any permutation of neurons in Definition 3.2. Each permutation defines another affine subspace

$$\mathcal{P}_\pi \Gamma_s(\theta^r) := \{ \mathcal{P}_\pi \theta^m : \theta^m \in \Gamma_s(\theta^r) \text{ and } \pi \in S_m \}$$

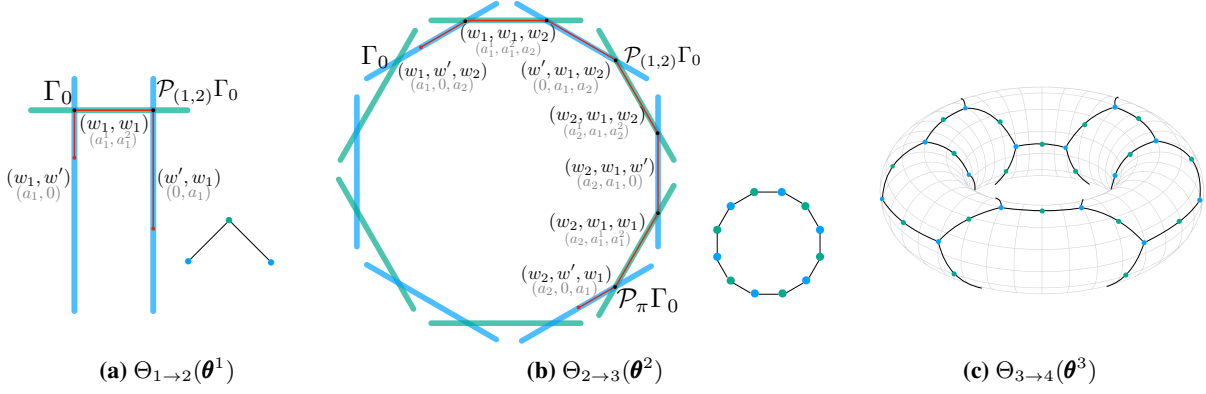


Figure 4. The geometry of the expansion manifold  $\Theta_{r \rightarrow m}$  with  $m = r + 1$  and the connectivity graph of the affine subspaces. The arrangement of the subspaces is demonstrated geometrically only in (a)-(b), but their connectivity graph is shown in all three cases. Blue subspaces have one vanishing output weight, green subspaces have two identical incoming weight vectors. (a) Case of a network with two hidden neurons with parameters  $(w_1, w', a_1, 0)$  that is reducible to a network with a single hidden neuron. The base subspace  $\Gamma_0$  is connected to a neighbor subspace  $\mathcal{P}_{(1,2)}\Gamma_0$  via three line segments: we first shift  $w'$  towards  $w_1$  while keeping the other parameters fixed and then move  $a_1^1$  from  $a_1$  to 0 while keeping  $a_1^1 + a_1^2 = a_1$ . The connectivity graph (bottom right) shows each subspace as an appropriately colored dot. (b) Case of a network with three hidden neurons with parameters  $(w_1, w', w_2, a_1, 0, a_2)$  that is reducible to a network with two hidden neurons.  $\Gamma_0$  is connected to any other subspace  $\mathcal{P}_\pi\Gamma_0$  through transitions from one neighbor to the next. Note that there are  $T(2, 3) = 12$  subspaces. (c) The connectivity graph of subspaces for the expansion  $3 \rightarrow 4$ , there are  $T(3, 4) = 60$  subspaces (24 blue and 36 green), where each blue subspace is connected to three green subspaces and each green subspace is connected to two blue subspaces.

where  $\mathcal{P}_\pi$  permutes the neurons  $\vartheta_i = [w_i, a_i]$  of  $\theta^m$ . We call the union of these affine subspaces the expansion manifold of  $\theta^r$ :

**Definition 3.3.** For  $r \leq m$ , the **expansion manifold**  $\Theta_{r \rightarrow m}(\theta^r) \subset \mathbb{R}^{D^m}$  of an irreducible  $r$ -neuron point  $\theta^r$  is defined by

$$\Theta_{r \rightarrow m}(\theta^r) := \bigcup_{\substack{s=(k_1, \dots, k_r, b_1, \dots, b_j) \\ \pi \in S_m}} \mathcal{P}_\pi \Gamma_s(\theta^r),$$

where  $s$  is a tuple with  $k_i \geq 1$ ,  $b_i \geq 0$  such that  $\text{sum}(s) = m$ .

Since  $\Theta_{r \rightarrow m}(\theta^r)$  is an equal-loss manifold, the gradient flow can cross it at most for once. Therefore  $\Theta_{r \rightarrow m}(\theta^r)$  is not an invariant manifold like the symmetry subspaces. Next, we describe the precise geometry of the expansion manifolds

**Theorem 3.1.** For  $m \geq r$ , the expansion manifold  $\Theta_{r \rightarrow m}(\theta^r)$  of an irreducible point  $\theta^r$  consists of exactly<sup>3</sup>

$$T(r, m) := \sum_{j=0}^{m-r} \sum_{\substack{\text{sum}(s)=m \\ k_i \geq 1, b_i \geq 1}} \binom{m}{k_1, \dots, k_r, b_1, \dots, b_j} \frac{1}{c_1! \dots c_{m-r}!}$$

distinct affine subspaces (none is including another one) of dimension at least  $\min(d_{\text{in}}, d_{\text{out}})(m - r)$ , where  $c_i$  is the number of occurrences of  $i$  among  $(b_1, \dots, b_j)$ .

<sup>3</sup>  $\binom{n_1 + \dots + n_r}{n_1, \dots, n_r}$  denotes the coefficient  $\frac{(n_1 + \dots + n_r)!}{n_1! \dots n_r!}$ .

For  $m > r$ ,  $\Theta_{r \rightarrow m}(\theta^r)$  is connected: any pair of distinct points  $\theta, \theta' \in \Theta_{r \rightarrow m}(\theta^r)$  is connected via a union of line segments  $\gamma : [0, 1] \rightarrow \Theta_{r \rightarrow m}(\theta^r)$  such that  $\gamma(0) = \theta$  and  $\gamma(1) = \theta'$ .

*Proof (Sketch).* The number of affine subspaces  $T$  is equal to the distinct permutations of the incoming weight vectors  $(w_1, \dots, w_r, w'_1, \dots, w'_j)$  for all possible tuples  $s$  where  $w_i$ 's are distinct and  $w'_i$ 's are dummy variables representing zero-type neurons (the neurons that do not contribute to the network function since their outgoing weight vectors sum to zero). The normalization factor  $1/c_1!c_2! \dots c_{m-r}!$  cancels the repetitions coming from the zero-type neurons  $(w'_1, \dots, w'_j)$ . For example for the standard case  $m = r$ , there is no room for zero-type neurons. As a result we have

$$T(r, r) = \sum_{\substack{k_1 + \dots + k_r = r \\ k_i \geq 1}} \binom{r}{k_1, \dots, k_r} = \binom{r}{1, \dots, 1} = r!$$

distinct subspaces of dimension  $\min(d_{\text{in}}, d_{\text{out}})(m - r) = 0$ .

For the general case  $m > r$ , the proof for connectivity follows from the following observations. We start from a base subspace  $\Gamma_0 = \Gamma_s(\theta^r)$ , where there is a zero-type neuron with outgoing weight vector exactly zero<sup>4</sup> at position  $i^*$ . The neighbor subspaces  $\mathcal{P}_{(i^*, i)}\Gamma_0$ , where  $(i^*, i) \in S_m$  is

<sup>4</sup>If all zero-type neurons are part of a group with more than one neuron, we can choose the first neuron in a group and set its outgoing weight vector to zero while respecting the condition in Eq. 3.

a transposition that permutes two neurons only, are connected to the base subspace via three line segments (Figure 4-a). Since any permutation is a composition of transpositions, permuted subspaces  $\mathcal{P}_\pi \Gamma_0$  can be reached via a union of line segments by going from one neighbor to the next (Figure 4-b). ■

#### 4. Overparameterized ANN Landscapes

In this section, we study the geometry of the global minima manifold and the critical subspaces, i.e. affine subspaces containing only critical points, in two-layer overparameterized neural networks. In particular, we show how the affine subspaces that form the global minima manifold are connected to one another (Subsection 4.1). We then find a hierarchy of saddles induced by permutation symmetries, which we call symmetry-induced critical points (Subsection 4.2). Finally, we compare the number of affine subspaces that form the global minima manifold with the number of those that contain symmetry-induced critical points (Subsection 4.3). Generalizations to multi-layer networks are discussed in Section 5.

We assume that there is a minimal width  $r^*$  such that  $\theta_*$  achieves zero loss, i.e.  $L^{r^*}(\theta_*) = 0$ , that the point  $\theta_*$  is unique up to permutation, and that any network with width  $r^* - 1$  has loss  $> 0$  at every point. We call the wider networks with width  $m > r^*$  **overparameterized** and the narrower networks with width  $r < r^*$  **underparameterized**. Note that  $\theta_*$  is irreducible by minimality of  $r^*$ .

##### 4.1. The global minima manifold

Applying Theorem 3.1 to the expansion manifold of a global minimum  $\theta_*$  of the minimal-width network, we obtain a connected manifold of global minima in an overparameterized network of width  $m$ :

**Corollary 4.1.** *In an overparameterized network with width  $m > r^*$ , the expansion manifold of global minima  $\Theta_{r^* \rightarrow m}(\theta_*)$  is connected.*

We have found a connected manifold  $\Theta_{r^* \rightarrow m}(\theta_*)$  of global minima. Furthermore, since  $\Theta_{r^* \rightarrow m}(\theta_*)$  is an expansion manifold, its geometry is precisely as described in Theorem 3.1, and illustrated in Figure 4. The next question is whether  $\Theta_{r^* \rightarrow m}(\theta_*)$  contains all the zero-loss points.

In the remaining part of this subsection, we give a positive answer to this question in a specific setting. We consider a modified loss function:

$$L_\mu^m(\theta) = \int_{\mathbb{R}^{d_{\text{in}}}} c(f^{(2)}(x|\theta), f^*(x)) \mu(dx),$$

where  $\mu$  is an input data distribution with support  $\mathbb{R}^{d_{\text{in}}}$ , and  $f^* : \mathbb{R}^{d_{\text{in}}} \rightarrow \mathbb{R}^{d_{\text{out}}}$  is a true data-generating function. The

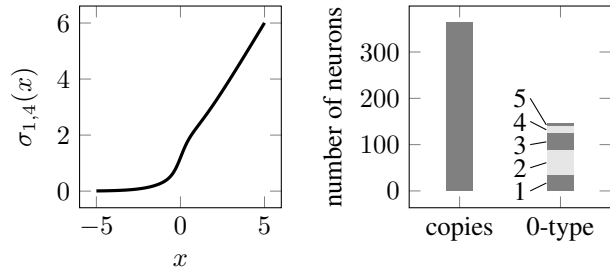


Figure 5. **Left:** The function  $\sigma_{\alpha,\gamma}(x) = \sigma_{\text{soft}}(x) + \alpha\sigma_{\text{sig}}(\gamma x)$  satisfies the technical condition of Theorem 4.2. With this activation function, data is generated by a teacher network of width 4. All 50 student networks with width 10 find a global minimum by reaching loss values below  $10^{-16}$ . **Right:** The 500 = 50 × 10 hidden neurons of all the 50 student networks are classified as copies of teacher neurons or zero-type neurons with vanishing sum of output weights. The zero-type neurons are further classified according to group size: there are 34 neurons with vanishing output weight (group size 1), 54 neurons that have a partner neuron with the same input weights and the sum of output weights equal to 0 (group size 2) etc. All zero-type neurons and replications of weight vectors can be pruned.

assumption on the activation  $\sigma$  in Theorem 4.2 below is only required for this theorem but not in Subsections 4.2 or 4.3. We find that there is no global minimum point outside of the expansion manifold  $\Theta_{r^* \rightarrow m}(\theta_*)$  for the modified loss  $L_\mu^m$  and for a certain class of activation functions (see Figure 5 for an example):

**Theorem 4.2.** *Suppose that the activation function  $\sigma$  is  $C^\infty$ , that  $\sigma(0) \neq 0$ , and that  $\sigma^{(n)}(0) \neq 0$  for infinitely many even and odd values of  $n$  (where  $\sigma^{(n)}$  denotes the  $n$ -th derivative of  $\sigma$ ). For  $m > r^*$ , let  $\theta$  be an  $m$ -neuron point, and  $\theta_*$  be a unique  $r^*$ -neuron global minimum up to permutation, i.e.  $L_\mu^{r^*}(\theta_*) = 0$ . If  $L_\mu^m(\theta) = 0$ , then  $\theta \in \Theta_{r^* \rightarrow m}(\theta_*)$ . (See Appendix-B.3 for the proof.)*

**Remark 4.1.** *The function  $\sigma_{\alpha,\gamma}(x) = \sigma_{\text{soft}}(x) + \alpha\sigma_{\text{sig}}(\gamma x)$  with  $\alpha, \gamma > 0$  (Figure 5) satisfies the conditions of Theorem 4.2, but the standard softplus  $\sigma_{\text{soft}}(x) = \ln[1 + \exp(x)]$  or sigmoidal  $\sigma_{\text{sig}}(x) = 1/[1 + \exp(-x)]$  functions do not. For these, the analysis must include additional invariances.*

**Remark 4.2.** *If a global minimum is found by gradient descent in overparameterized networks, then the final set of parameters can be classified into groups of replicated weight vectors according to Definition 3.2 (Figure 5). The classification can be exploited for pruning the network.*

**Remark 4.3.** *Kuditipudi et al. (2019) construct an example of a finite-size dataset (in contrast with our infinite dataset framework) for two-layer overparameterized ReLU networks where they find discrete global minima points.*

## 4.2. Symmetry-induced critical points

In this subsection, we consider an overparameterized network with a fixed width  $m > r^*$  and study critical points in an expansion manifold  $\Theta_{r \rightarrow m}(\theta_*^r)$  where we assume that  $\theta_*^r$  is an irreducible critical point of an underparameterized network with width  $r < r^*$ . Observe that  $\theta_*^r$  is not a zero-loss point since  $r^*$  is the minimal width to achieve zero loss. We consider only those points without zero-type neurons in  $\Theta_{r \rightarrow m}(\theta_*^r)$ , we show that these have zero gradient, and therefore are critical points of  $L^m$ .

**Definition 4.1.** For  $r \leq m$ , let  $s = (k_1, \dots, k_r)$  be an  $r$ -tuple with  $k_i \geq 1$  and  $\text{sum}(s) = m$ . The **symmetry-induced critical points** are those in the set

$$\bar{\Theta}_{r \rightarrow m}(\theta_*^r) = \bigcup_{\substack{s=(k_1, \dots, k_r) \\ \pi \in S_m}} \mathcal{P}_\pi \bar{\Gamma}_s(\theta_*^r)$$

where the critical (affine) subspace  $\bar{\Gamma}_s(\theta_*^r) \subset \mathbb{R}^{D^m}$  of an irreducible critical point  $\theta_*^r = (w_1^*, \dots, w_r^*, a_1^*, \dots, a_r^*)$  is

$$\left\{ \underbrace{(w_1^*, \dots, w_1^*)}_{k_1}, \dots, \underbrace{(w_r^*, \dots, w_r^*)}_{k_r}, \beta_1^1 a_1^*, \dots, \beta_1^{k_1} a_1^*, \dots, \beta_r^1 a_r^*, \dots, \beta_r^{k_r} a_r^* \right\} : \sum_{i=1}^{k_t} \beta_t^i = 1 \text{ for } t \in [r]. \quad (4)$$

All points in  $\bar{\Theta}_{r \rightarrow m}(\theta_*^r)$  are critical points hence the name symmetry-induced ‘critical points’:

**Proposition 4.3.** For an irreducible critical point  $\theta_*^r$  of  $L^r$ ,  $\bar{\Theta}_{r \rightarrow m}(\theta_*^r)$  is a union of

$$G(r, m) := \sum_{\substack{k_1 + \dots + k_r = m \\ k_i \geq 1}} \binom{m}{k_1, \dots, k_r}$$

distinct non-intersecting affine subspaces of dimension  $m - r$ . All points in  $\bar{\Theta}_{r \rightarrow m}(\theta_*^r)$  are critical points of  $L^m$ .

Proposition 4.3 shows that a critical point of a smaller network  $\theta_*^r$  expands into  $G(r, m)$  critical subspaces in the overparameterized network with width  $m$ . If  $\theta_*^r$  is a strict saddle,  $\bar{\Theta}_{r \rightarrow m}(\theta_*^r)$  contains only strict saddles, since the escape direction is preserved for affine transformations  $\bar{\Gamma}_s$ .

**Proposition 4.4.** For  $C^2$  functions  $c$  and  $\sigma$ , for all  $\theta_*^m \in \bar{\Theta}_{r \rightarrow m}(\theta_*^r)$ , the spectrum of the Hessian  $\nabla^2 L^m(\theta_*^m)$  has  $(m - r)$  zero eigenvalues. Moreover, if  $\theta_*^r$  is a strict saddle, then all points in  $\bar{\Theta}_{r \rightarrow m}(\theta_*^r)$  are also strict saddles, i.e., their Hessian has at least one negative eigenvalue.

If  $\theta_*^r$  is a local minimum, Fukumizu et al. (2019) show that the subspaces for which only one neuron is replicated ( $k_i > 1$ ,  $k_j = 1$  for all  $j \neq i$ ) may contain both local minima and strict saddles depending on the spectrum of a matrix

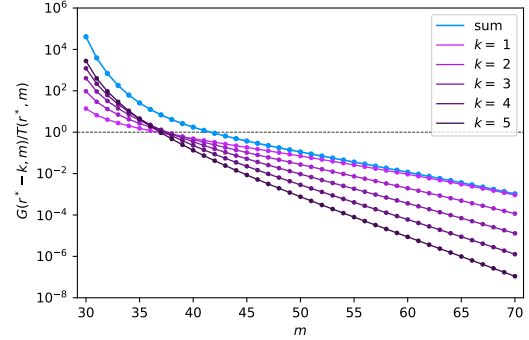


Figure 6. The ratio  $R_k(r^*, m)$  of the multiplier for  $k$ -th level saddle  $G(r^* - k, m)$  to the number of global minima subspaces  $T(r^*, m)$  as the width  $m$  of the overparameterized network increases, plotted for a fixed width  $r^* = 30$  of the minimal network. The ratio of all critical subspaces to the global minima subspaces  $\sum_{k=1}^{r^*-1} a_k G(r^* - k, m) / T(r^*, m)$  is shown in blue assuming  $a_k = 1$  for all  $k$ . Note that for  $m \gg r^*$  the blue curve approaches the curve for  $k = 1$  indicating that only subspaces corresponding to first-level saddles are potentially relevant, yet the global minima subspaces clearly dominate.

of derivatives [see their Theorem 11]. We expect a similar result to hold true for all subspaces in  $\bar{\Theta}_{r \rightarrow m}(\theta_*^r)$ , including arbitrary replications.

**Remark 4.4.** We explore a hierarchy between symmetry-induced critical points in  $\cup_{r < r^*} \bar{\Theta}_{r \rightarrow m}(\theta_*^r)$  in a network of width  $m$ : **first-level saddles** refer to symmetry-induced critical points that are equivalent to a minimum of a network of width  $r^* - 1$ ; more generally,  **$k$ -th level saddles** refer to those equivalent to a minimum of a network of width  $r^* - k$ . Adding neurons enables the network to reach a lower loss minimum thus higher-level symmetry-induced saddles usually attain higher losses. We notice a similarity with Gaussian Process (Bray & Dean, 2007) and spherical spin glass (Auffinger et al., 2013) landscapes, where the higher-order<sup>5</sup> saddles attain higher losses.

Finally, we note that the dimensionality of the global minima subspaces  $\mathcal{P}_\pi \bar{\Gamma}_s(\theta_*^r)$  and the critical subspaces  $\mathcal{P}_\pi \bar{\Gamma}_s(\theta_*^r)$  differ, in particular in the way they depend on  $r$ . What is common is that they are all ‘tiny’ compared to the ambient dimensionality of the parameter space. In the following subsection, we will thus focus on the comparison of the number of critical subspaces and that of global minima subspaces.

## 4.3. Width-dependent comparison of the critical subspaces and the global minima subspaces

In the loss landscape of an overparameterized network of width  $m$ , we have the connected global minima mani-

<sup>5</sup>The order of a saddle point is the number of negative eigenvalues of its Hessian.

fold  $\Theta_{r^* \rightarrow m}(\theta_*)$  as well as many subspaces of symmetry-induced critical points in  $\bar{\Theta}_{r \rightarrow m}(\theta_*^r)$ , where  $\theta_*^r$  is an irreducible critical point in a smaller network with some width  $r = r^* - k < r^*$ . In this subsection, we count these subspaces and find

- $T(r^*, m)$  global minima subspaces (Corollary 4.1)
- $G(r^* - k, m)a_k$  critical subspaces for all  $k = 1, \dots, r^* - 1$  (Proposition 4.3)

where  $a_k$  is the number of distinct<sup>6</sup> irreducible critical points in a network with width  $r^* - k$  and where  $G(r^* - k, m)$  is the multiplier.

To compare the number of non-global critical subspaces with the number of global minima subspaces, we give closed-form formulas for  $G$  and  $T$ . This is proven in Appendix-B.5 using Newton’s series for finite differences (Milne-Thomson, 2000) and a counting argument:

**Proposition 4.5.** *For  $r \leq m$ , we have*

$$G(r, m) = \sum_{i=1}^r \binom{r}{i} (-1)^{r-i} i^m$$

$$T(r, m) = G(r, m) + \sum_{u=1}^{m-r} \binom{m}{u} G(r, m-u) g(u)$$

where  $g(u) = \sum_{j=1}^u \frac{1}{j!} G(j, u)$ .

Using Proposition 4.5, we find the following asymptotic behaviors for  $G$  and  $T$ :

**Lemma 4.6.** *For any  $k \geq 0$  fixed, we have,*

$$G(m-k, m) \sim T(m-k, m) \sim \frac{m^k}{2^k k!} m!, \text{ as } m \rightarrow \infty.$$

For any fixed  $r \geq 0$ , we have  $G(r, m) \sim r^m$  as  $m \rightarrow \infty$ .

We are now ready to compare the number of global minima subspaces  $T$  with the number of critical subspaces  $G$  under the assumption that the minimal width  $r^*$  is large. This is realistic since for a real-world dataset the network should be sufficiently wide to achieve zero loss. Applying Lemma 4.6, we find that the symmetry-induced critical points dominate the global minima in mildly overparameterized, and vice versa in vastly overparameterized networks (see Figure 6). A mathematical analysis yields:

**Mildly Overparameterized.** Let  $m = r^* + h$  for fixed  $h$ . We have in the limit  $r^* \rightarrow \infty$  and for fixed  $k$  a ratio:

$$R_k(r^*, m) := \frac{G(r^* - k, m)}{T(r^*, m)} \sim \frac{1}{2^k (h+k) \cdots (h+1)} (r^*)^k. \quad (5)$$

<sup>6</sup>We say two irreducible critical points  $\theta_*^a$  and  $\theta_*^b$  are distinct if  $\theta_*^a \neq \mathcal{P}_\pi \theta_*^b$  for all applicable permutations  $\pi$ .

Thus, for a small amount of overparameterization, the multiplier of the  $k$ -th level saddles  $G(r^* - k, m)$  scales as  $(r^*)^k T(r^*, m)$ , indicating a proliferation of saddles at a rate much larger than that of the global minima. Related to this proliferation, we empirically encounter training failures (i.e. training halts before reaching a global minimum) for typical initializations in this regime (see Figure 1). Moreover, we empirically find traces of approaching a saddle in gradient trajectories in narrow two-layer ANNs trained on MNIST (see Appendix).

**Vastly Overparameterized.** For  $m$  very large, i.e.  $m \gg r^*$ , we have

$$\sum_{k=1}^{r^*-1} R_k(r^*, m) a_k = \frac{\sum_{k=1}^{r^*-1} G(r^* - k, m) a_k}{T(r^*, m)} \leq \left( \frac{r^* - 1}{r^*} \right)^m \quad (6)$$

if  $a_k$ ’s satisfy  $a_k \leq \binom{r^*-1}{k-1}$ . Because the RHS of Eq. (6) decreases down to 0 as  $m \rightarrow \infty$  (at a geometric rate), the global minima dominate *all* symmetry-induced critical points. We note that there could be other critical points in addition to those generated by the symmetries. The calculations above are presented in the Appendix.

## 5. Multi-Layer ANNs

In this section, we introduce the expansion manifold for multi-layer networks that enables obtaining connectivity and counting results on the global minima manifold for multi-layer networks (i.e., generalizing Theorem 3.1 and Corollary 4.1). Finally, we compare the number of affine subspaces of the global minima and symmetry-induced critical points. An ANN with  $L$  layers  $f^{(L)} : \mathbb{R}^{d_{\text{in}}} \rightarrow \mathbb{R}^{d_{\text{out}}}$  with widths  $\mathbf{r} = (r_1, r_2, \dots, r_{L-1})$  is

$$f^{(L)}(x|\theta) = W^{(L)} \sigma(W^{(L-1)} \dots \sigma(W^{(1)} x)) \quad (7)$$

where  $W^{(\ell)} \in \mathbb{R}^{r_\ell \times r_{\ell-1}}$  for  $\ell = 1, \dots, L$  with  $r_0 = d_{\text{in}}$  and  $r_L = d_{\text{out}}$ , the non-linearity  $\sigma$  is applied element-wise, and  $\theta = (W^{(L)}, \dots, W^{(1)}) \in \mathbb{R}^{d(\mathbf{r})}$  is the vector of parameters of dimension  $d(\mathbf{r}) = \sum_{\ell=1}^L r_{\ell-1} r_\ell$ . Observing that any pair of weight matrices  $(W^{(\ell)}, W^{(\ell+1)})$  for  $\ell = 1, \dots, L-1$  forms a two-layer network within the multi-layer network, we say that a multi-layer network is irreducible if all pairs  $(W^{(\ell)}, W^{(\ell+1)})$  are irreducible.

**The global minima manifold.** We define the expansion manifold of an irreducible network with widths  $\mathbf{r}$  into larger widths  $\mathbf{m}$  by taking the sequential expansion manifolds of all pairs  $(W^{(\ell)}, W^{(\ell+1)})$ . More precisely, we define the multi-layer expansion manifold as follows

$$\Theta_{\mathbf{r} \rightarrow \mathbf{m}}(\theta^{\mathbf{r}}) := \{ \phi_1 \in \mathbb{R}^{d(\mathbf{m})} : \phi_{L-1} \in \Theta_{\mathbf{r} \rightarrow \mathbf{m}}^{(L-1)}(\theta^{\mathbf{r}}), \dots, \phi_1 \in \Theta_{\mathbf{r} \rightarrow \mathbf{m}}^{(1)}(\phi_2) \} \quad (8)$$



where  $\Theta_{\mathbf{r} \rightarrow \mathbf{m}}^{(\ell)}(\phi)$  substitutes the pair  $(W^{(\ell)}, W^{(\ell+1)})$  with those of a point in the usual expansion manifold (Def. 3.3). Since each expansion leaves the output of the network unchanged, all points in this expansion have the same loss. Note that the order in which we take these expansions affects the final manifold; expanding from the last layer to the first one gives the largest final manifold. The same final manifold can be obtained via a ‘forward pass’ if one considers expansion up to an equivalence of the incoming weight vectors.

Assume that a minimal  $L$ -layer network achieves a unique (up to permutation) global minimum point  $\theta_*$  with widths  $\mathbf{r}^* = (r_1^*, r_2^*, \dots, r_{L-1}^*)$ . In an overparameterized network of widths  $\mathbf{m} = (m_1, \dots, m_{L-1})$  with  $m_\ell > r_\ell^*$  for all  $\ell \in [L-1]$  (i.e. at least one extra neuron at every hidden layer), we find a connected manifold of global minimum, which is simply the multi-layer expansion manifold  $\Theta_{\mathbf{r}^* \rightarrow \mathbf{m}}(\theta_*)$  of the minimum point  $\theta_*$ . The zero-loss expansion manifold  $\Theta_{\mathbf{r}^* \rightarrow \mathbf{m}}(\theta_*)$  consists of the following number of distinct affine subspaces

$$\prod_{\ell=1}^{L-1} T(r_\ell^*, m_\ell).$$

**Symmetry-induced critical points.** Similarly, we can consider the symmetry-induced critical points for multi-layer networks by applying sequential expansions  $\bar{\Theta}_{\mathbf{r} \rightarrow \mathbf{m}}^{(\ell)}$  to all hidden layers. We note that applying this expansion to a pair  $(W_*^{(\ell)}, W_*^{(\ell+1)})$  of a critical point  $\theta_*^r$  generates a manifold of critical points as in the two-layer case, hence these expansions preserve criticality. The number of affine subspaces in the set of symmetry-induced critical points is

$$\prod_{\ell=1}^{L-1} G(r_\ell, m_\ell).$$

**Application.** Similar to Fig 1-(d,e), we consider the case where a minimal  $L$ -layer network with  $r^*$  neurons at each hidden layer reaches a global minimum point  $\theta_*$ . Let us consider an overparameterization with  $m = r^* + h$  neurons at each hidden layer. The ratio of the number of critical subspaces of  $k$ -th level saddles to the global minima subspaces is

$$R_k(r^*, m)^{L-1} = \left( \frac{G(r^* - k, m)}{T(r^*, m)} \right)^{L-1},$$

which is exponential in depth. Therefore in the mildly overparameterized regime, i.e. when  $h$  is small, we see that the ratio of the number of saddles to that of global minima grows exponentially with depth. In other words, we observe that the dominance of the number of saddles is even more pronounced in the multi-layer case. For the vastly overparameterized regime, i.e. when  $h$  is large, we observe

the opposite effect: the dominance of the number of global minima is stronger in the multi-layer case. Finally, we observe a width-depth trade-off in reaching a dominance of the global minima: one can either increase the width of a two-layer network so that the ratios  $R_k(r^*, m)$  go down to 0; or increase the depth in a network where each layer is just large enough to guarantee  $R_k(r^*, m) < 1$  which eventually decreases the total ratio down to 0.

## 6. Conclusion & Discussion

In this paper, we explicitly characterize the geometry formed by the critical points in overparameterized neural networks. For the global minima, we showed that under mild conditions they live in a manifold consisting of a number of connected affine subspaces. We characterize a certain type of critical points, the so-called symmetry-induced critical points and we showed that they form an explicit number of affine subspaces. From the theoretical point of view, it remains an open question whether there are other critical points in the overparameterized networks in addition to the symmetry-induced ones. We also leave it to future work to study whether all symmetry-induced critical points are strict saddles or not.

Our main result quantifies the scaling of the numbers of global minima subspaces and the subspaces containing symmetry-induced critical points as the width grows. In mildly overparameterized networks, the number of critical subspaces is much greater than that of the global minima subspaces, so that in practice, the gradient trajectories may get influenced by these saddles or even get transiently stuck in their neighborhood for a fraction of typical initializations. However, in vastly overparameterized networks, the number of global minima subspaces dominates that of the critical subspaces so that symmetry-induced saddles play only a marginal role. From a practical point of view, our theoretical results pave the way to applications in optimization of non-convex neural networks loss landscapes via a combination of overparameterization and pruning.

## Acknowledgements

The authors thank the authors of [Lengyel et al. \(2020\)](#) for a discussion about neural network invariances at the very beginning of this project. The authors thank Valentin Schmutz for a discussion, Bernd Illing and Levent Sagun for their detailed feedback on the manuscript. This work is partly supported by Swiss National Science Foundation (no. 200020\_184615) and ERC SG CONSTAMIS. C. Hongler acknowledges support from the Blavatnik Family Foundation, the Latsis Foundation, and the NCCR Swissmap.

## References

- Arora, S., Du, S. S., Hu, W., Li, Z., and Wang, R. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. *arXiv preprint arXiv:1901.08584*, 2019.
- Auffinger, A., Arous, G. B., and Černý, J. Random matrices and complexity of spin glasses. *Communications on Pure and Applied Mathematics*, 66(2):165–201, 2013.
- Bray, A. J. and Dean, D. S. Statistics of critical points of gaussian fields on large-dimensional spaces. *Physical review letters*, 98(15):150201, 2007.
- Brea, J., Simsek, B., Illing, B., and Gerstner, W. 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*, 2019.
- Chizat, L. and Bach, F. On the global convergence of gradient descent for over-parameterized models using optimal transport. *Advances in neural information processing systems*, 31:3036–3046, 2018.
- Cooper, Y. The critical locus of overparameterized neural networks. *arXiv preprint arXiv:2005.04210*, 2020.
- Dauphin, Y. N., Pascanu, R., Gulcehre, C., Cho, K., Ganguli, S., and Bengio, Y. Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. In *Advances in neural information processing systems*, pp. 2933–2941, 2014.
- Draxler, F., Veschgini, K., Salmhofer, M., and Hamprecht, F. A. Essentially no barriers in neural network energy landscape. *arXiv preprint arXiv:1803.00885*, 2018.
- Du, S., Lee, J., Li, H., Wang, L., and Zhai, X. Gradient descent finds global minima of deep neural networks. In *International Conference on Machine Learning*, pp. 1675–1685. PMLR, 2019.
- Du, S. S., Jin, C., Lee, J. D., Jordan, M. I., Singh, A., and Póczos, B. Gradient descent can take exponential time to escape saddle points. In *Advances in neural information processing systems*, pp. 1067–1077, 2017.
- Du, S. S., Zhai, X., Póczos, B., and Singh, A. Gradient descent provably optimizes over-parameterized neural networks. *arXiv preprint arXiv:1810.02054*, 2018.
- Fort, S. and Jastrzebski, S. Large scale structure of neural network loss landscapes. In *Advances in Neural Information Processing Systems*, pp. 6709–6717, 2019.
- Fort, S., Dziugaite, G. K., Paul, M., Kharaghani, S., Roy, D. M., and Ganguli, S. Deep learning versus kernel learning: an empirical study of loss landscape geometry and the time evolution of the neural tangent kernel. *arXiv preprint arXiv:2010.15110*, 2020.
- Frankle, J. and Carbin, M. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *arXiv preprint arXiv:1803.03635*, 2018.
- Frankle, J., Dziugaite, G. K., Roy, D., and Carbin, M. Linear mode connectivity and the lottery ticket hypothesis. In *International Conference on Machine Learning*, pp. 3259–3269. PMLR, 2020.
- Freeman, C. D. and Bruna, J. Topology and geometry of half-rectified network optimization. *arXiv preprint arXiv:1611.01540*, 2016.
- Fukumizu, K. and Amari, S.-i. Local minima and plateaus in hierarchical structures of multilayer perceptrons. *Neural networks*, 13(3):317–327, 2000.
- Fukumizu, K., Yamaguchi, S., Mototake, Y.-i., and Tanaka, M. Semi-flat minima and saddle points by embedding neural networks to overparameterization. *arXiv preprint arXiv:1906.04868*, 2019.
- Garipov, T., Izmailov, P., Podoprikin, D., Vetrov, D. P., and Wilson, A. G. Loss surfaces, mode connectivity, and fast ensembling of dnns. *Advances in Neural Information Processing Systems*, 31:8789–8798, 2018.
- Geiger, M., Spigler, S., d’Ascoli, S., Sagun, L., Baity-Jesi, M., Biroli, G., and Wyart, M. Jamming transition as a paradigm to understand the loss landscape of deep neural networks. *Physical Review E*, 100(1):012115, 2019.
- Głuch, G. and Urbanke, R. Noether: The more things change, the more stay the same. *arXiv preprint arXiv:2104.05508*, 2021.
- Jacot, A., Gabriel, F., and Hongler, C. Neural tangent kernel: Convergence and generalization in neural networks. In *Advances in neural information processing systems*, pp. 8571–8580, 2018.
- Jin, C., Ge, R., Netrapalli, P., Kakade, S. M., and Jordan, M. I. How to escape saddle points efficiently. *arXiv preprint arXiv:1703.00887*, 2017.
- Kudithipudi, R., Wang, X., Lee, H., Zhang, Y., Li, Z., Hu, W., Ge, R., and Arora, S. Explaining landscape connectivity of low-cost solutions for multilayer nets. In *Advances in Neural Information Processing Systems*, pp. 14601–14610, 2019.
- Kung, J. P., Rota, G.-C., and Yan, C. H. *Combinatorics: the Rota way*. Cambridge University Press, 2009.

- Kunin, D., Sagastuy-Brena, J., Ganguli, S., Yamins, D. L., and Tanaka, H. Neural mechanics: Symmetry and broken conservation laws in deep learning dynamics. *arXiv preprint arXiv:2012.04728*, 2020.
- Lee, J., Xiao, L., Schoenholz, S. S., Bahri, Y., Novak, R., Sohl-Dickstein, J., and Pennington, J. Wide neural networks of any depth evolve as linear models under gradient descent. *arXiv preprint arXiv:1902.06720*, 2019a.
- Lee, J., Schoenholz, S. S., Pennington, J., Adlam, B., Xiao, L., Novak, R., and Sohl-Dickstein, J. Finite versus infinite neural networks: an empirical study. *arXiv preprint arXiv:2007.15801*, 2020.
- Lee, J. D., Panageas, I., Piliouras, G., Simchowitz, M., Jordan, M. I., and Recht, B. First-order methods almost always avoid strict saddle points. *Mathematical programming*, 176(1):311–337, 2019b.
- Lengyel, D., Petangoda, J., Falk, I., Highnam, K., Lazarou, M., Kolbeinsson, A., Deisenroth, M. P., and Jennings, N. R. Genni: Visualising the geometry of equivalences for neural network identifiability. *arXiv preprint arXiv:2011.07407*, 2020.
- Milne-Thomson, L. M. *The calculus of finite differences*. American Mathematical Soc., 2000.
- Nguyen, Q. On connected sublevel sets in deep learning. *arXiv preprint arXiv:1901.07417*, 2019.
- Sagan, B. E. *The symmetric group: representations, combinatorial algorithms, and symmetric functions*, volume 203. Springer Science & Business Media, 2013.
- Sagun, L., Guney, V. U., Arous, G. B., and LeCun, Y. Explorations on high dimensional landscapes. *arXiv preprint arXiv:1412.6615*, 2014.
- Sagun, L., Evci, U., Guney, V. U., Dauphin, Y., and Bottou, L. Empirical analysis of the hessian of over-parametrized neural networks. *arXiv preprint arXiv:1706.04454*, 2017.