

# Tractable Landscapes for Nonconvex Optimization

Shihua Zhang

November 7, 2024

# Overview

- 1 Challenges in Nonconvex Landscapes
- 2 Cases With a Unique Global Minimum
- 3 Symmetry, Saddle Points and Locally Optimizable Functions
- 4 Case Study: Top Eigenvector of a Matrix
- 5 Mode Connectivity of Neural Networks
- 6 Summary

# Background

- Deep learning relies on optimizing a **nonconvex** loss.
- Even simple algorithms such as gradient descent often **optimize the objective value to zero or near-zero**.
- **Goal**: How to optimize the nonconvex landscapes efficiently and identify their properties (for machine learning models)?
- Only apply to simpler nonconvex problems than deep learning.
- How to analyze deep learning with such landscape analysis is still **open**.

# Global and Local Minimum

## Definition 1 (Global/Local minimum)

1. For an objective function  $f(w) : \mathbb{R}^d \rightarrow \mathbb{R}$ , a point  $w^*$  is a **global minimum** if for every  $w$  we have  $f(w^*) \leq f(w)$ .
2. A point  $w$  is a **local minimum/maximum** if there exists a radius  $\epsilon > 0$  such that for every  $\|w' - w\|_2 \leq \epsilon$ , we have  $f(w) \leq f(w')$  ( $f(w) \geq f(w')$  for local maximum).
3. A point  $w$  with  $\nabla f(w) = 0$  is called a **critical point**, and for smooth functions all local minimum/maximum are critical points.

- Here we work with functions whose global minimum exists, and use  $f(w^*)$  to denote its optimal value.

# Spurious Local Minimum

## Definition 2 (Spurious local minimum)

For an objective function  $f(w) : \mathbb{R}^d \rightarrow \mathbb{R}$ , a point  $w$  is a **spurious local minimum** if it is a local minimum, but  $f(w) > f(w^*)$ .

- Many optimization algorithms are based on the idea of **local search**, thus cannot escape from a spurious local minimum.
- Many nonconvex objectives do not have spurious local minima.

# Saddle Points

## Definition 3 (Saddle point)

For an objective function  $f(w) : \mathbb{R}^d \rightarrow \mathbb{R}$ , a point  $w$  is a **saddle point** if  $\nabla f(w) = 0$ , and for every radius  $\epsilon > 0$ , there exists  $w^+, w^-$  within distance  $\epsilon$  of  $w$  such that  $f(w^-) < f(w) < f(w^+)$ .

- This definition covers all cases but makes it **very hard to verify** whether a point is a saddle point.
- In most cases, it is possible to tell whether a point is a saddle point, local minimum or local maximum based on **its Hessian**.

# Second Order Sufficient Condition

## Theorem 4

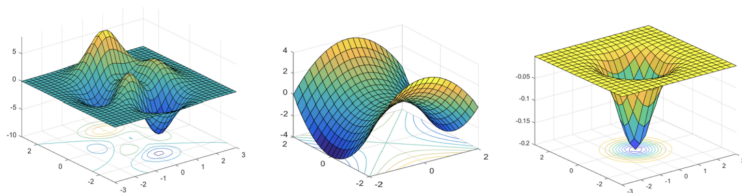
*For an objective function  $f(w) : \mathbb{R}^d \rightarrow \mathbb{R}$  and a critical point  $w$  ( $\nabla f(w) = 0$ ), we know:*

- If  $\nabla^2 f(w) \succ 0$ ,  $w$  is a local minimum.*
- If  $\nabla^2 f(w) \prec 0$ ,  $w$  is a local maximum.*
- If  $\nabla^2 f(w)$  has both a positive and a negative eigenvalue,  $w$  is a saddle point.*

- **Proof Hint:** looking at the second-order Taylor expansion.
- The three cases do not cover all the possible Hessian matrices.

# Flat Regions

- **Challenge:** Even if a function does not have any spurious local minima or saddle point, it can still be hard to optimize.
- **Difficulty:** even if the norm  $\|\nabla f(\mathbf{w})\|_2$  is small, unlike convex functions, one cannot conclude that  $f(\mathbf{w})$  is close to  $f(\mathbf{w}^*)$ .



**Figure:** Obstacles for nonconvex optimization. From left to right: local minimum, saddle point and flat region.



# Overview

- 1 Challenges in Nonconvex Landscapes
- 2 Cases With a Unique Global Minimum
- 3 Symmetry, Saddle Points and Locally Optimizable Functions
- 4 Case Study: Top Eigenvector of a Matrix
- 5 Mode Connectivity of Neural Networks
- 6 Summary

# Cases With a Unique Global Minimum

We first consider the case that is similar to convex objectives.

- The objective functions we look at **have no spurious local minima or saddle points**.
- **Obstacle**: points with small gradients may not be near-optimal.
- **Main Idea**: identify properties of the objective function, such that it keeps decreasing during the optimization process.

# Cases With a Unique Global Minimum

## Definition 5

Let  $f(w)$  be an objective function with a unique global minimum  $w^*$ , then:

**Polyak-Lojasiewicz:**  $f$  satisfies Polyak-Lojasiewicz if there exists a value  $\mu > 0$  such that for every  $w$ ,  $\|\nabla f(w)\|_2^2 \geq \mu (f(w) - f(w^*))$

**Weakly-quasi-convex:**  $f$  is weakly-quasi-convex if there exists a value  $\mu > 0$  such that for every  $w$ ,  
 $\langle \nabla f(w), w - w^* \rangle \geq \mu (f(w) - f(w^*))$

**Restricted Secant Inequality (RSI):**  $f$  satisfies RSI if there exists a value  $\mu$  such that for every  $w$ ,  $\langle \nabla f(w), w - w^* \rangle \geq \mu \|w - w^*\|_2^2$

# Cases With a Unique Global Minimum

Any one of these three properties can imply fast convergence together with some smoothness of  $f$ .

## Theorem 6

*If an objective function  $f$  satisfies one of Polyak-Lojasiewicz, weakly-quasi-convex or RSI, and  $f$  is smooth, then gradient descent converges to global minimum with a geometric rate.*

- Polyak-Lojasiewicz and RSI requires standard smoothness, weakly-quasi-convex requires a special smoothness property detailed in [Hardt et al., 2016].
- We will use **generalized linear model (GLM)** as an example to show how some of these properties can be used.

# Generalized Linear Model

In GLM ([Kalai and Sastry, 2009], [Kakade et al., 2011]), the input consists of samples  $\{x^{(i)}, y^{(i)}\}$  that are drawn from a distribution  $\mathcal{D}$ , where  $(x, y) \sim \mathcal{D}$  satisfies

$$y = \sigma(w_*^\top x) + \epsilon$$

# Generalized Linear Model

In GLM ([Kalai and Sastry, 2009], [Kakade et al., 2011]), the input consists of samples  $\{x^{(i)}, y^{(i)}\}$  that are drawn from a distribution  $\mathcal{D}$ , where  $(x, y) \sim \mathcal{D}$  satisfies

$$y = \sigma(w_*^\top x) + \epsilon$$

- $\sigma: \mathbb{R} \rightarrow \mathbb{R}$  is a known monotone function,  $\epsilon$  is a noise that satisfies  $\mathbb{E}[\epsilon \mid x] = 0$ .
- Consider Expected loss:  $L(w) = \frac{1}{2} \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[ (y - \sigma(w^\top x))^2 \right]$ .

# Generalized Linear Model

In GLM ([Kalai and Sastry, 2009], [Kakade et al., 2011]), the input consists of samples  $\{x^{(i)}, y^{(i)}\}$  that are drawn from a distribution  $\mathcal{D}$ , where  $(x, y) \sim \mathcal{D}$  satisfies

$$y = \sigma(w_*^\top x) + \epsilon$$

- $\sigma: \mathbb{R} \rightarrow \mathbb{R}$  is a known monotone function,  $\epsilon$  is a noise that satisfies  $\mathbb{E}[\epsilon | x] = 0$ .
- Consider Expected loss:  $L(w) = \frac{1}{2} \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[ (y - \sigma(w^\top x))^2 \right]$ .
- GLM: **learn a single neuron** where  $\sigma$  is its nonlinearity.

# Generalized Linear Model

How to prove prop. (e.g., weakly-quasi-convex or RSI) for GLM?

The objective is rewritten as:

$$\begin{aligned} L(w) &= \frac{1}{2} \mathbb{E}_{(x,y) \sim \mathcal{D}} \left[ (y - \sigma(w^\top x))^2 \right] \\ &= \frac{1}{2} \mathbb{E}_{(x,\epsilon)} \left[ (\epsilon + \sigma(w_*^\top x) - \sigma(w^\top x))^2 \right] \\ &= \frac{1}{2} \mathbb{E}_{\epsilon} [\epsilon^2] + \frac{1}{2} \mathbb{E}_x \left[ (\sigma(w_*^\top x) - \sigma(w^\top x))^2 \right]. \end{aligned} \tag{1}$$

- This decomposition is helpful as  $\frac{1}{2} \mathbb{E}_{\epsilon} [\epsilon^2]$  is just a constant.



# Generalized Linear Model

Consider the **derivative** of the objective:

$$\nabla L(w) = \mathbb{E}_x \left[ \left( \sigma(w^\top x) - \sigma(w_*^\top x) \right) \sigma'(w^\top x) x \right]. \quad (2)$$

Then we have:

$$\begin{aligned} & \langle \nabla L(w), w - w_* \rangle \\ &= \mathbb{E}_x \left[ \left( \sigma(w^\top x) - \sigma(w_*^\top x) \right) \sigma'(w^\top x) (w^\top x - w_*^\top x) \right] \\ &= \mathbb{E}_x \left[ \sigma'(\xi) \sigma'(w^\top x) (w^\top x - w_*^\top x)^2 \right]. \end{aligned} \quad (3)$$

By making more assumptions on  $\sigma$  and the distribution of  $x$ , it is possible to lowerbound  $\langle \nabla L(w), w - w_* \rangle$  in the form required by either weakly-quasi-convex or RSI.

# Overview

- 1 Challenges in Nonconvex Landscapes
- 2 Cases With a Unique Global Minimum
- 3 Symmetry, Saddle Points and Locally Optimizable Functions
- 4 Case Study: Top Eigenvector of a Matrix
- 5 Mode Connectivity of Neural Networks
- 6 Summary

# Permutation Symmetry for Neural Networks

Consider a two-layer neural network  $h_{\theta}(x) : \mathbb{R}^d \rightarrow \mathbb{R}$ . The parameters  $\theta$  is  $(w_1, w_2, \dots, w_k)$ .

- The function can be evaluated as  $h_{\theta}(x) = \sum_{i=1}^k \sigma(\langle w_i, x \rangle)$ .
- Given a dataset  $(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)}) \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}$ .
- The objective  $f(\theta) = L(h_{\theta}) = \mathbb{E}_{(x,y) \sim \mathcal{D}} [\ell((x, y), h_{\theta})]$  has **permutation symmetry**.
- That is, for any permutation  $\pi(\theta)$  that permutes the weights of the neurons,  $f(\theta) = f(\pi(\theta))$ .

# Permutation Symmetry for Neural Networks

The **permutation symmetry** has many implications:

- If the global minimum  $\theta^*$  is a point where not all neurons have the same weight, then there must be **equivalent global minimum**  $f(\pi(\theta^*))$  for every permutation  $\pi$ .
- An objective with this symmetry must also be nonconvex, because if it were convex, the point  $\bar{\theta} = \frac{1}{k!} \sum_{\pi} \pi(\theta^*)$  must be a global minimum.
- However, for  $\bar{\theta}$  the weight vectors of the neurons are all equal to  $\frac{1}{k} \sum_{i=1}^k w_i$ , so  $h_{\bar{\theta}}(x) = k\sigma\left(\left\langle \frac{1}{k} \sum_{i=1}^k w_i, x \right\rangle\right)$  is equivalent to a neural network with a single neuron.

# Permutation Symmetry for Neural Networks

The **permutation symmetry** has many implications:

- $f$  must be **nonconvex**.
- It is also possible to show that functions with symmetry must **have saddle points**.
- To optimize  $f$ , the algorithm needs to be able to either **avoid or escape from saddle points**.
- More concretely, one would like to find **a second order stationary point**.

# Second order stationary point (SOSP)

## Definition 7 (Second order stationary point (SOSP))

For an objective function  $f(w) : \mathbb{R}^d \rightarrow \mathbb{R}$ , a point  $w$  is a second order stationary point if  $\nabla f(w) = 0$  and  $\nabla^2 f(w) \succeq 0$

- The conditions for SOSP are known as the **second order necessary conditions** for a local minimum.
- The optimization algorithms can be used to find **an approximate** second order stationary point.

## Definition 8 (Approximate second order stationary point)

For an objective function  $f(w) : \mathbb{R}^d \rightarrow \mathbb{R}$ , a point  $w$  is a  $(\epsilon, \gamma)$ -second order stationary point ( $(\epsilon, \gamma)$ -SOSP) if  $\|\nabla f(w)\|_2 \leq \epsilon$  and  $\lambda_{\min}(\nabla^2 f(w)) \geq -\gamma$

# Locally Optimizable

Define a class of functions that can be optimized efficiently and allow symmetry and saddle points.

## Definition 9 (Locally optimizable functions)

An objective function  $f(w)$  is **locally optimizable**, if for every  $\tau > 0$ , there exists  $\epsilon, \gamma = \text{poly}(\tau)$  such that every  $(\epsilon, \gamma)$ -SOSP  $w$  of  $f$  satisfies  $f(w) \leq f(w_*) + \tau$ .

- Roughly speaking, **an objective function is locally optimizable** if every local minimum of the function is also a global minimum, and the Hessian of every saddle point has a negative eigenvalue.

# Locally Optimizable Functions

Locally optimizable objective functions:

- Matrix sensing [Hardt et al., 2016]
- Matrix completion [Ge et al., 2016]
- Dictionary learning [Sun et al., 2016]
- Tensor decomposition [Ge et al., 2015]
- Certain objective for two-layer neural network [Ge et al., 2017]



# Overview

- 1 Challenges in Nonconvex Landscapes
- 2 Cases With a Unique Global Minimum
- 3 Symmetry, Saddle Points and Locally Optimizable Functions
- 4 Case Study: Top Eigenvector of a Matrix
- 5 Mode Connectivity of Neural Networks
- 6 Summary

# Top Eigenvector of a Matrix

Here we look at a simple example of a **locally optimizable function**.

- Given a symmetric PSD matrix  $M \in \mathbb{R}^{d \times d}$ , the goal is to find its top eigenvector.
- More precisely, using SVD we can write  $M$  as

$$M = \sum_{i=1}^d \lambda_i v_i v_i^\top$$

Here  $v_i$ 's are orthonormal vectors that are eigenvectors of  $M$ , and  $\lambda_i$ 's are the eigenvalues.

- For simplicity, we assume  $\lambda_1 > \lambda_2 \geq \lambda_3 \geq \dots \geq \lambda_d > 0$

# Top Eigenvector of a Matrix

There are many objective functions whose global optima give the top eigenvector.

- For PSD matrix  $M$ , the global optima of

$$\max_{\|x\|_2=1} x^\top Mx$$

is the top eigenvector of  $M$ . However, this formulation requires a constraint.

- We instead work with an unconstrained version whose correctness follows from [Eckart-Young Theorem](#):

$$\min_{x \in \mathbb{R}^d} f(x) := \frac{1}{4} \|M - xx^\top\|_F^2$$

# Top Eigenvector of a Matrix

Consider the following unconstrained optimization problem:

$$\min_{x \in \mathbb{R}^d} f(x) := \frac{1}{4} \|M - xx^\top\|_F^2$$

- This function does have **a symmetry** in the sense that  $f(x) = f(-x)$ .
- Under the assumptions, the only global minima of it are  $x = \pm\sqrt{\lambda_1}v_1$ . They are the only  $2^{nd}$ -order stationary points.
- Two **proof** strategies:
  - Characterizing all critical points
  - Finding directions of improvements

# Characterizing All Critical Points

The first idea is simple:

- Solve the Eq.  $\nabla f(x) = 0$  to get the position of all critical points.

# Characterizing All Critical Points

The first idea is simple:

- Solve the Eq.  $\nabla f(x) = 0$  to get the position of all critical points.
- For the critical points that are not the desired global minimum, try to prove that they are **local maximum or saddle points**.

# Computing Gradient and Hessian

Expand  $f(x + \delta)$  as follows:

$$\begin{aligned} f(x + \delta) &= \frac{1}{4} \|M - (x + \delta)(x + \delta)^\top\|_F^2 \\ &= \frac{1}{4} \|M - xx^\top - (x\delta^\top + \delta x^\top) - \delta\delta^\top\|_F^2 \\ &= \frac{1}{4} \|M - xx^\top\|_F^2 - \frac{1}{2} \langle M - xx^\top, x\delta + \delta x^\top \rangle \\ &\quad + \left[ \frac{1}{4} \|x\delta^\top + \delta x^\top\|_F^2 - \frac{1}{2} \langle M - xx^\top, \delta\delta^\top \rangle \right] + o(\|\delta\|_2^2) \end{aligned}$$

Thus we have:

$$\nabla f(x) = (xx^\top - M)x, \quad \nabla^2 f(x) = \|x\|_2^2 I + 2xx^\top - M$$

# Characterizing Critical Points

Set  $\nabla f(x) = 0$ , we have:

$$Mx = xx^\top x = \|x\|_2^2 x$$

- The only solutions to  $Mx = \lambda x$  are if  $\lambda$  is an eigenvalue and  $x$  is (a scaled version) of the corresponding eigenvector.
- $x = \pm\sqrt{\lambda_i}v_i$  or  $x = 0$ . And  $x = \pm\sqrt{\lambda_1}v_1$  are intended solutions.



# Characterizing Critical Points

Set  $\nabla f(x) = 0$ , we have:

$$Mx = xx^\top x = \|x\|_2^2 x$$

- The only solutions to  $Mx = \lambda x$  are if  $\lambda$  is an eigenvalue and  $x$  is (a scaled version) of the corresponding eigenvector.
- $x = \pm\sqrt{\lambda_i}v_i$  or  $x = 0$ . And  $x = \pm\sqrt{\lambda_1}v_1$  are intended solutions.
- Next we need to show for every other critical point, its Hessian has a negative direction, i.e., there exists a  $\delta$  such that  $\delta^\top [\nabla^2 f(x)] \delta < 0$ .
- Key: The main step of the proof involves guessing what is this direction  $\delta$ . In this case we will choose  $\delta = v_1$ .

# Characterizing Critical Points

When  $x = \pm\sqrt{\lambda_i}v_i$ , and  $\delta = v_1$ , we have:

$$\delta^\top [\nabla^2 f(x)] \delta = v_1^\top \left[ \left\| \sqrt{\lambda_i} v_i \right\|_2^2 I + 2\lambda_i v_i v_i^\top - M \right] v_1 = \lambda_i - \lambda_1 < 0$$

The proof for  $x = 0$  is very similar.

Combining all the steps above, we proved the theorem:

## Theorem 10 (Properties of critical points)

*The only critical points of  $f(x)$  are of the form  $x = \pm\sqrt{\lambda_i}v_i$  or  $x = 0$ . For all critical points except  $x = \pm\sqrt{\lambda_1}v_1$ ,  $\nabla^2 f(x)$  has a negative eigenvalue.*

# Characterizing Critical Points

When  $x = \pm\sqrt{\lambda_i}v_i$ , and  $\delta = v_1$ , we have:

$$\delta^\top [\nabla^2 f(x)] \delta = v_1^\top \left[ \left\| \sqrt{\lambda_i} v_i \right\|_2^2 I + 2\lambda_i v_i v_i^\top - M \right] v_1 = \lambda_i - \lambda_1 < 0$$

The proof for  $x = 0$  is very similar.

Combining all the steps above, we proved the theorem:

## Theorem 10 (Properties of critical points)

*The only critical points of  $f(x)$  are of the form  $x = \pm\sqrt{\lambda_i}v_i$  or  $x = 0$ . For all critical points except  $x = \pm\sqrt{\lambda_1}v_1$ ,  $\nabla^2 f(x)$  has a negative eigenvalue.*

- The only  $2^{\text{nd}}$ -order stationary points are  $x = \pm\sqrt{\lambda_1}v_1$ , so all  $2^{\text{nd}}$ -order stationary points are also global minima.

# Finding Directions of Improvements

It is often infeasible to **enumerate all the solutions for  $\nabla f(x) = 0$** .

**Key:** For every point  $x$  that is not a global minimum, we define **its direction of improvements** as below:

## Definition 11

- For an objective function  $f$  and a point  $x$ , we say  $\delta$  is a direction of improvement (of  $f$  at  $x$ ) if  $|\langle \nabla f(x), \delta \rangle| > 0$  or  $\delta^\top [\nabla^2 f(x)] \delta < 0$ .
- We say  $\delta$  is an  $(\epsilon, \gamma)$ -direction of improvement (of  $f$  at  $x$ ) if  $|\langle \nabla f(x), \delta \rangle| > \epsilon \|\delta\|_2$  or  $\delta^\top [\nabla^2 f(x)] \delta < -\gamma \|\delta\|_2^2$ .

- If  $\delta$  is a direction of improvement for  $f$  at  $x$ , then moving along one of  $\delta$  or  $-\delta$  for a small enough step can decrease the objective function.

# Finding Directions of Improvements

For simplicity, consider an simpler version of the top eigenvector problem, where  $M = zz^\top$  is a rank-1 matrix, and  $z$  is a unit vector. Then

$$\min_x f(x) = \frac{1}{4} \|zz^\top - xx^\top\|_F^2 \quad (4)$$

- Which direction should we move to decrease the objective?
- One only have the optimal direction  $z$  and the current direction  $x$ , so the natural guesses would be  $z$ ,  $x$  or  $z - x$ .

# Finding Directions of Improvements

## Lemma 12

*For objective function  $f$ , there exists a universal constant  $c > 0$  such that for any  $\tau < 1$ , if neither  $x$  or  $z$  is an  $(c\tau, 1/4)$ -direction of improvement for the point  $x$ , then  $f(x) \leq \tau$ .*

- The proof of this lemma involves some detailed calculation.
- To get some intuition, we first think about **what happens if neither  $x$  or  $z$  is a direction of improvement.**

# Finding Directions of Improvements

## Lemma 13

*For objective function  $f$ , if neither  $x$  or  $z$  is a direction of improvement of  $f$  at  $x$ , then  $f(x) = 0$ .*

## Proof.

If  $x$  is not a direction of improvement, we must have:

$$\langle \nabla f(x), x \rangle = 0 \implies \|x\|_2^4 = \langle x, z \rangle^2$$

If  $z$  is not a direction of improvement, we know  $z^\top [\nabla^2 f(x)] z \geq 0$  which means

$$\|x\|^2 + 2\langle x, z \rangle^2 - 1 \geq 0 \implies \|x\|^2 \geq 1/3$$



# Finding Directions of Improvements

## Proof.

Consider the fact that  $\langle x, z \rangle^2 \leq \|x\|_2^2 \|z\|_2^2 = \|x\|_2^2$ , thus we have  $\langle x, z \rangle^2 = \|x\|_2^4 \geq 1/9$ .

Finally, since  $z$  is not a direction of improvement, we know  $\langle \nabla f(x), z \rangle = 0$ , which implies  $\langle x, z \rangle (\|x\|_2^2 - 1) = 0$ . We have already proved  $\langle x, z \rangle^2 \geq 1/9 > 0$ , thus  $\|x\|_2^2 = 1$ .

Again we know  $\langle x, z \rangle^2 = \|x\|_2^4 = 1$ . The only two vectors with  $\langle x, z \rangle^2 = 1$  and  $\|x\|_2^2 = 1$  are  $x = \pm z$ . □



# Finding Directions of Improvements

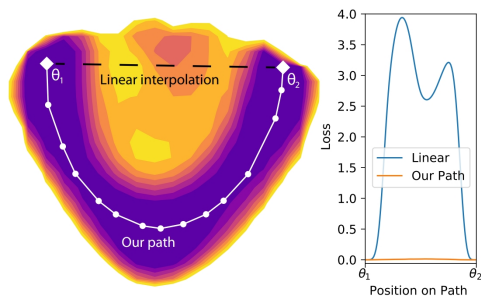
- The proof of Lemma 12 is very similar to Lemma 13, except we need to **allow slacks in every equation and inequality** we use.
- Lemma 12 and Lemma 13 **both use directions  $x$  and  $z$** . It is also possible to use the direction  $x - z$  when  $\langle x, z \rangle \geq 0$  (and  $x + z$  when  $\langle x, z \rangle < 0$ ).
- Both ideas can **be generalized** to handle the case when  $M = ZZ^T$  where  $Z \in \mathbb{R}^{d \times r}$ , so  $M$  is a rank- $r$  matrix.

# Overview

- 1 Challenges in Nonconvex Landscapes
- 2 Cases With a Unique Global Minimum
- 3 Symmetry, Saddle Points and Locally Optimizable Functions
- 4 Case Study: Top Eigenvector of a Matrix
- 5 Mode Connectivity of Neural Networks
- 6 Summary

# Mode Connectivity

- Although the loss landscape of DNNs is nonconvex with many minima, there are some **tractable** structures.
- Mode connectivity**: Different local minima can be **connected by simple paths** [Garipov et al., 2018, Draxler et al., 2018].



**Figure:** Two minima (found by SGD) are connected by a polygonal chain of low loss, but the loss along the linear path is high [Draxler et al., 2018].

# Linear Mode Connectivity

- Mode connectivity suggests that different local minima are **not isolated**, but essentially **form a connected manifold**.
- **Linear mode connectivity (LMC)**: Connected by a **linear** path.

## Definition 14 (Linear mode connectivity)

Given dataset  $D$  and two modes  $\theta_A, \theta_B$  that  $\text{Err}_D(\theta_A) = \text{Err}_D(\theta_B)$ , two mode  $\theta_A$  and  $\theta_B$  satisfy the linear mode connectivity if

$$\forall \alpha \in [0, 1], \text{Err}_D(\alpha\theta_A + (1 - \alpha)\theta_B) \approx \text{Err}_D(\theta_A)$$

- LMC implies that the minima are **in the same basin**.
- SGD converges to a basin  $\rightarrow$  flat minima.
- **When LMC happens?**

# Spawning Method [Frankle et al., 2020]

- 1 A network is randomly initialized, trained for some epochs.
- 2 **Spawned into two copies** which continue to be **independently trained** by different SGD randomnesses.

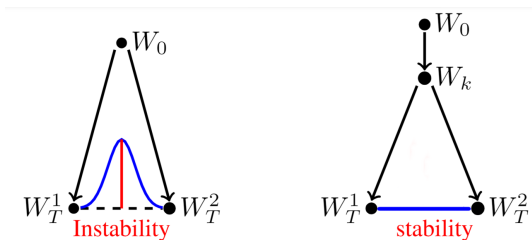


Figure: **Left:** MC; **Right:** LMC after spawning [Frankle et al., 2020].

- **Insight:** The result of optimization is **determined in early stage**.

# Permutation Method [Entezari et al., 2021]

- Recall that DNNs satisfy **permutation symmetry**.
- Conjecture:** SGD solutions are **LMC** after proper permutation.

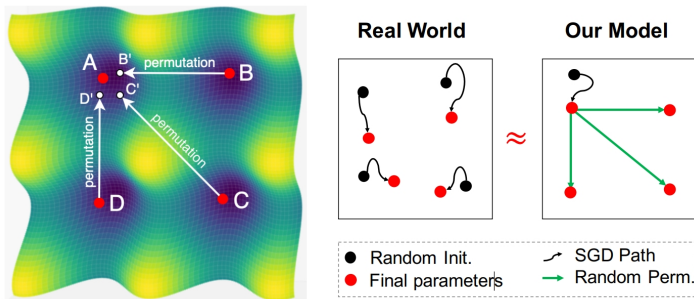


Figure: Permuting minima to the same basin [Entezari et al., 2021].

- Insight:** Permutation symmetry leads to different basins, yet SGD can converge to minima with similar performance.

# Ways to Find Proper Permutation

- **Align** the neurons of independently trained models via permutation [Ainsworth et al., 2022, Qu and Horvath, 2024].
  - **Weight matching:**  $\min_{\pi} \sum_{\ell=1}^L \left\| \mathbf{W}_A^{(\ell)} - \mathbf{P}^{(\ell)} \mathbf{W}_B^{(\ell)} \mathbf{P}^{(\ell-1)\top} \right\|_F^2$
  - **Activation matching:**  $\min_{\pi} \sum_{\ell=1}^L \left\| \mathbf{H}_A^{(\ell)} - \mathbf{P}^{(\ell)} \mathbf{H}_B^{(\ell)} \right\|_F^2$
- LMC allows us to merge models by “teleporting” solutions into a single basin.

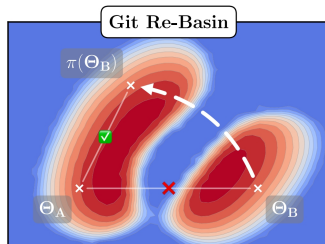
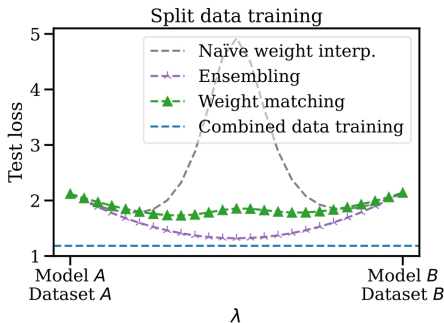


Figure: Git Re-Basin, a **weight averaging method** [Ainsworth et al., 2022].

# LMC Leads to better Model Averaging

- **Goal:** Merge models trained on **disjoint datasets**.
  - Federated Learning
  - Distributed Training
- **Method:** Weight averaging in the same basin via permutation.



**Figure:** Merging ResNets on CIFAR100 **outperforms both input models** while **using half compute required** for ensembling [Ainsworth et al., 2022].



# Overview

- 1 Challenges in Nonconvex Landscapes
- 2 Cases With a Unique Global Minimum
- 3 Symmetry, Saddle Points and Locally Optimizable Functions
- 4 Case Study: Top Eigenvector of a Matrix
- 5 Mode Connectivity of Neural Networks
- 6 Summary

# Summary

- **Cases with a unique global minimum:** identify properties of the objective function, such as PL, weakly-quasi-convex and RSI conditions.
- **The permutation symmetry has many implication:** nonconvex, saddle points and so on.
- **Two strategies for analyzing the landscape:**
  - Characterizing all critical points
  - Finding directions of improvements
- **(Linear) mode connectivity reveals the relationship between local minima in nonconvex landscape:** Connected by simple (linear) paths.

# References I



Ainsworth, S. K., Hayase, J., and Srinivasa, S. (2022).  
Git re-basin: Merging models modulo permutation symmetries.  
*arXiv preprint arXiv:2209.04836*.



Draxler, F., Veschgini, K., Salmhofer, M., and Hamprecht, F. (2018).  
Essentially no barriers in neural network energy landscape.  
*In International conference on machine learning*, pages 1309–1318. PMLR.



Entezari, R., Sedghi, H., Saukh, O., and Neyshabur, B. (2021).  
The role of permutation invariance in linear mode connectivity of neural networks.  
*arXiv preprint arXiv:2110.06296*.



Frankle, J., Dziugaite, G. K., Roy, D., and Carbin, M. (2020).  
Linear mode connectivity and the lottery ticket hypothesis.  
*In International Conference on Machine Learning*, pages 3259–3269. PMLR.



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.



Ge, R., Huang, F., Jin, C., and Yuan, Y. (2015).  
Escaping from saddle points—online stochastic gradient for tensor decomposition.  
*In Conference on learning theory*, pages 797–842. PMLR.

# References II



Ge, R., Lee, J. D., and Ma, T. (2016).

Matrix completion has no spurious local minimum.

In Lee, D., Sugiyama, M., Luxburg, U., Guyon, I., and Garnett, R., editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc.



Ge, R., Lee, J. D., and Ma, T. (2017).

Learning one-hidden-layer neural networks with landscape design.

*arXiv preprint arXiv:1711.00501*.



Hardt, M., Ma, T., and Recht, B. (2016).

Gradient descent learns linear dynamical systems.

*arXiv preprint arXiv:1609.05191*.



Kakade, S. M., Kanade, V., Shamir, O., and Kalai, A. (2011).

Efficient learning of generalized linear and single index models with isotonic regression.

In Shawe-Taylor, J., Zemel, R., Bartlett, P., Pereira, F., and Weinberger, K. Q., editors, *Advances in Neural Information Processing Systems*, volume 24. Curran Associates, Inc.



Kalai, A. T. and Sastry, R. (2009).

The isotron algorithm: High-dimensional isotonic regression.

In *COLT*. Citeseer.

# References III



Qu, X. and Horvath, S. (2024).

Rethink model re-basin and the linear mode connectivity.

*arXiv preprint arXiv:2402.05966.*



Sun, J., Qu, Q., and Wright, J. (2016).

Complete dictionary recovery over the sphere i: Overview and the geometric picture.

*IEEE Transactions on Information Theory*, 63(2):853–884.