

---

# Unique Properties of Flat Minima in Deep Networks: Supplementary Material

---

Rotem Mulayoff<sup>1</sup> Tomer Michaeli<sup>1</sup>

This document contains supplementary material for the article ‘Unique Properties of Flat Minima in Deep Networks’, and includes the following parts:

- I. Stability of minima
- II. Proof of Lemma 1
- III. Scalar networks
- IV. Proof of Lemma 3
- V. The missing parts of the proof of Theorem 1
- VI. Proof of Theorem 2
- VII. Proof of Theorem 3
- VIII. Proof of Theorem 4
- IX. More details about our experiments

## I. Stability of Minima

It is well known that for a  $\beta$ -smooth function that is bounded from below, GD with constant step size  $0 < \eta < 2/\beta$  converges to a stationary point. A twice continuously differentiable function  $f$  is  $\beta$ -smooth *if and only if*  $\lambda_{\max}(\nabla^2 f(\mathbf{w})) \leq \beta$  for every point  $\mathbf{w} \in \mathbb{R}^N$ . Hence, convergence to a stationary point is guaranteed if  $\lambda_{\max}(\nabla^2 f(\mathbf{w})) \leq 2/\eta$  for all  $\mathbf{w} \in \mathbb{R}^N$ . This seemingly stringent global requirement can in fact also be replaced by a local one, as shown by Wu et al. (2018). Specifically, they use the following.

**Definition S1.** Let  $\mathbf{w}^*$  be a stationary point of  $f$ . Consider the linearized dynamical system of GD, namely

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \nabla^2 f(\mathbf{w}^*)(\mathbf{w}_k - \mathbf{w}^*). \quad (\text{S1})$$

Then  $\mathbf{w}^*$  is said to be linearly stable if there exists a constant  $C \in \mathbb{R}$ , such that  $\|\mathbf{w}_k\| \leq C \|\mathbf{w}_0\|$  for all  $k > 0$ .

In other words,  $\mathbf{w}^*$  is linearly stable if once we have arrived near this critical point, we stay around it. In their paper, Wu et al. (2018) show that  $\mathbf{w}^*$  is a linearly stable minimizer if

$$\lambda_{\max} \left( (\mathbf{I} - \eta \nabla^2 f(\mathbf{w}^*))^2 \right) \leq 1. \quad (\text{S2})$$

Note that for all  $i \in \{1, \dots, N\}$

$$\begin{aligned} \lambda_i \left( (\mathbf{I} - \eta \nabla^2 f(\mathbf{w}^*))^2 \right) &= \lambda_i^2 (\mathbf{I} - \eta \nabla^2 f(\mathbf{w}^*)) \\ &= (1 - \eta \lambda_i (\nabla^2 f(\mathbf{w}^*)))^2 \\ &= 1 - \eta \lambda_i (\nabla^2 f(\mathbf{w}^*)) (2 - \eta \lambda_i (\nabla^2 f(\mathbf{w}^*))), \end{aligned} \quad (\text{S3})$$

---

<sup>1</sup>Department of Electrical Engineering, Technion – Israel Institute of Technology, Haifa, Israel. Correspondence to: Rotem Mulayoff <rotem.mulayof@gmail.com>.

where  $\lambda_i$  is the  $i$ th largest eigenvalue. Since  $\eta$  and  $\{\lambda_i\}$  are all nonnegative, it follows that (S2) is equivalent to

$$\lambda_{\max}(\nabla^2 f(\mathbf{w}^*)) \leq \frac{2}{\eta}. \quad (\text{S4})$$

This result asserts that flat minima are stable solutions for GD. In their paper, they also provide a similar result for stochastic GD (SGD), which shows that the sharpness of a minimum should increase to ensure stability for SGD as well.

## II. Proof of Lemma 1

Nar & Sastry (2018) showed that under the Lemma's conditions, the weight matrices converge at a linear rate to  $\mathbf{W}_i = \mathbf{T}^{\frac{1}{m}}$  for all  $i$ . According to our Theorem 4, this solution is a flattest minimum, thus demonstrating (ii).

Arora et al. (2018) showed that gradient flow (GF) satisfies

$$\mathbf{W}_i(t)\mathbf{W}_i^T(t) = \mathbf{W}_{i+1}^T(t)\mathbf{W}_{i+1}(t), \quad \forall t \geq 0 \quad (\text{S5})$$

in our setting. Denoting the SVD of  $\mathbf{W}_i(t)$  by  $\mathbf{U}_i(t)\mathbf{S}_i(t)\mathbf{V}_i^T(t)$ , we thus have that  $\mathbf{U}_i(t)\mathbf{S}_i^2(t)\mathbf{U}_i^T(t) = \mathbf{V}_{i+1}^T(t)\mathbf{S}_{i+1}^2(t)\mathbf{V}_{i+1}(t)$ , which implies that<sup>1</sup>

$$\mathbf{U}_i(t) = \mathbf{V}_{i+1}(t), \quad \mathbf{S}_i(t) = \mathbf{S}_{i+1}(t), \quad \forall t \geq 0. \quad (\text{S6})$$

Assume that GF converges to a global minimum and let  $\mathbf{U}\mathbf{S}\mathbf{V}^T$  denote the SVD of  $\mathbf{T}$ . Since  $\{\mathbf{S}_i(t)\}_{i=1}^m$  are identical, they converge to the same limit,  $\bar{\mathbf{S}}$ . Let  $\mathbf{W}_i = \mathbf{U}_i\bar{\mathbf{S}}\mathbf{V}_i^T$  denote the limit of  $\mathbf{W}_i(t)$ . Then, from (S6), we have that  $\mathbf{V}_{i+1}^T\mathbf{U}_i = \mathbf{I}$  for all  $i$ . Consequently,

$$\mathbf{W}_m\mathbf{W}_{m-1}\cdots\mathbf{W}_1 = \mathbf{U}_m\bar{\mathbf{S}}^m\mathbf{V}_1^T. \quad (\text{S7})$$

But since the left hand side equals  $\mathbf{T}$  by assumption, the right hand side must coincide with the SVD of  $\mathbf{T}$ . This means that  $\bar{\mathbf{S}} = \mathbf{S}^{\frac{1}{m}}$ . Again, by Theorem 4, this is a flattest minimum, thus demonstrating (i).

## III. Scalar Networks

### III.1. The Set of Flattest Minima

As mentioned in the main text, in the scalar case, the end-to-end function  $f_{\mathbf{w}}(x)$  implemented by the network is given by

$$f_{\mathbf{w}}(x) = \prod_{j=1}^m w_j x, \quad (\text{S8})$$

where  $\mathbf{w} = [w_1, w_2, \dots, w_m]^T$ . In our analysis we consider a quadratic loss function, *i.e.*

$$\ell(\mathbf{w}) = \hat{\mathbb{E}} \left[ (y - f_{\mathbf{w}}(x))^2 \right]. \quad (\text{S9})$$

Our goal is to characterize the set of flattest minima of the loss w.r.t.  $\mathbf{w}$ . It is well known that the optimal coefficient for linear estimation is given by

$$\tau = \frac{\hat{\sigma}_{xy}}{\hat{\sigma}_x^2}, \quad (\text{S10})$$

where  $\hat{\sigma}_x^2 = \mathbb{E}[x^2]$  is the empirical second-order moment of  $x$ , and  $\hat{\sigma}_{xy} = \mathbb{E}[xy]$  is the empirical cross second-order moment between  $x$  and  $y$ . Therefore, at any global minimum of  $\ell(\mathbf{w})$ , we have

$$\prod_{j=1}^m w_j = \tau. \quad (\text{S11})$$

<sup>1</sup>The SVD can be non-unique, however there necessarily exists a decomposition satisfying  $\mathbf{U}_i(t) = \mathbf{V}_{i+1}(t)$ .

To compute the Hessian matrix of  $\ell(\mathbf{w})$ , we first calculate the partial derivative w.r.t.  $w_k$ ,

$$\begin{aligned}\frac{\partial}{\partial w_k} \ell(\mathbf{w}) &= \frac{\partial}{\partial w_k} \hat{\mathbb{E}} \left[ (y - f_{\mathbf{w}}(x))^2 \right] = -2 \hat{\mathbb{E}} \left[ (y - f_{\mathbf{w}}(x)) \frac{\partial}{\partial w_k} f_{\mathbf{w}}(x) \right] \\ &= -2 \hat{\mathbb{E}} \left[ \left( y - \prod_{j=1}^m w_j x \right) \prod_{j \neq k} w_j x \right] = 2 \left( \hat{\sigma}_x^2 \prod_{j=1}^m w_j - \hat{\sigma}_{xy} \right) \prod_{j \neq k} w_j\end{aligned}\quad (\text{S12})$$

We now complete the derivation by differentiating (S12) w.r.t.  $w_q$ ,

$$\begin{aligned}\frac{\partial^2}{\partial w_q \partial w_k} \ell(\mathbf{w}) &= \frac{\partial}{\partial w_q} \left[ 2 \left( \hat{\sigma}_x^2 \prod_{j=1}^m w_j - \hat{\sigma}_{xy} \right) \prod_{j \neq k} w_j \right] \\ &= 2 \prod_{j \neq k} w_j \frac{\partial}{\partial w_q} \left( \hat{\sigma}_x^2 \prod_{j=1}^m w_j - \hat{\sigma}_{xy} \right) + 2 \left( \hat{\sigma}_x^2 \prod_{j=1}^m w_j - \hat{\sigma}_{xy} \right) \frac{\partial}{\partial w_q} \prod_{j \neq k} w_j.\end{aligned}\quad (\text{S13})$$

Eq. (S11) asserts that  $\hat{\sigma}_x^2 \prod_{j=1}^m w_j - \hat{\sigma}_{xy} = 0$  at global minima, therefore the second term in (S13) vanishes, and we obtain

$$\frac{\partial^2}{\partial w_q \partial w_k} L(\mathbf{w}) = 2 \prod_{j \neq k} w_j \frac{\partial}{\partial w_q} \left( \hat{\sigma}_x^2 \prod_{j=1}^m w_j - \hat{\sigma}_{xy} \right) = 2 \hat{\sigma}_x^2 \left( \prod_{j \neq k} w_j \right) \left( \prod_{j \neq q} w_j \right).\quad (\text{S14})$$

Hence, using (S10), we can express the elements of the Hessian matrix  $\mathbf{H}_{\mathbf{w}}$  as

$$(\mathbf{H}_{\mathbf{w}})_{k,q} = 2 \hat{\sigma}_x^2 \tau^2 \frac{1}{w_k w_q}.\quad (\text{S15})$$

Let us define the vector  $\mathbf{z} = [w_1^{-1}, w_2^{-1}, \dots, w_m^{-1}]^T$ , then the Hessian matrix can be equivalently written as

$$\mathbf{H}_{\mathbf{w}} = 2 \hat{\sigma}_x^2 \tau^2 \mathbf{z} \mathbf{z}^T.\quad (\text{S16})$$

This shows that the Hessian is a rank one matrix, which implies that it has only one nonzero eigenvalue, with a corresponding eigenvector  $\mathbf{z}$ . Therefore,

$$\lambda_{\max}(\mathbf{H}_{\mathbf{w}}) \mathbf{z} = \mathbf{H}_{\mathbf{w}} \mathbf{z} = 2 \hat{\sigma}_x^2 (\tau)^2 \mathbf{z} \mathbf{z}^T \mathbf{z} = 2 \hat{\sigma}_x^2 (\tau)^2 \|\mathbf{z}\|^2 \mathbf{z},\quad (\text{S17})$$

so that the eigenvalue is given by

$$\lambda_{\max}(\mathbf{H}_{\mathbf{w}}) = 2 \hat{\sigma}_x^2 \tau^2 \|\mathbf{z}\|^2 = 2 \hat{\sigma}_x^2 \tau^2 \sum_{j=1}^m \frac{1}{w_j^2}.\quad (\text{S18})$$

To determine the sharpness of the flattest minima, we need to solve the problem

$$\min_{\mathbf{w} \in \mathbb{R}^m} \lambda_{\max}(\mathbf{H}_{\mathbf{w}}) \quad \text{s.t.} \quad \prod_{j=1}^m w_j = \tau.\quad (\text{S19})$$

By the inequality of the arithmetic and geometric means, we have that for any feasible point  $\mathbf{w}$

$$\sum_{j=1}^m \frac{1}{w_j^2} \geq m \times \left( \prod_{j=1}^m \frac{1}{w_j^2} \right)^{\frac{1}{m}} = m \times \tau^{-\frac{2}{m}}.\quad (\text{S20})$$

Therefore, for all feasible points,

$$\lambda_{\max}(\mathbf{H}_{\mathbf{w}}) \geq 2m \hat{\sigma}_x^2 \tau^{2(1-\frac{1}{m})}.\quad (\text{S21})$$

On the other hand, this inequality can be achieved by setting  $|w_1| = |w_2| = \dots = |w_m|$ . This shows that the right-hand-side is precisely the sharpness of the flattest minimum, so that

$$\Omega_0 = \left\{ \mathbf{w} \in \mathbb{R}^m : \prod_{j=1}^m \text{sgn}(w_j) = \text{sgn}(\tau) \quad \text{and} \quad |w_j| = \sqrt[m]{|\tau|} \quad \forall j \right\}.\quad (\text{S22})$$

### III.2. Proof of Lemma 2

In this section we examine the behavior of the loss function on a line connecting two minima.

**Claim S1.** Assume that  $\tau > 0$  and let  $\mathbf{w}^{(1)}$  and  $\mathbf{w}^{(2)}$  be minimizers of the objective (S9) in  $\mathbb{R}_+^m$ . Then, along the line connecting  $\mathbf{w}^{(1)}$  and  $\mathbf{w}^{(2)}$ , the loss function will appear sharper around  $\mathbf{w}^{(1)}$  than around  $\mathbf{w}^{(2)}$  if

$$\sum_{i=1}^m \frac{w_i^{(2)}}{w_i^{(1)}} > \sum_{i=1}^m \frac{w_i^{(1)}}{w_i^{(2)}}. \quad (\text{S23})$$

*Proof.* The direction vector of the connecting line is  $\boldsymbol{\alpha} = (\mathbf{w}^{(1)} - \mathbf{w}^{(2)}) / \|\mathbf{w}^{(1)} - \mathbf{w}^{(2)}\|$ . Along this direction, the behavior of the loss function around  $\mathbf{w}^{(i)}$  is given by

$$\ell(\mathbf{w}^{(i)} + \eta\boldsymbol{\alpha}) \approx \ell(\mathbf{w}^{(i)}) + \eta\boldsymbol{\alpha}^T \nabla \ell(\mathbf{w}^{(i)}) + \frac{\eta^2}{2} \boldsymbol{\alpha}^T \mathbf{H}(\mathbf{w}^{(i)}) \boldsymbol{\alpha}. \quad (\text{S24})$$

Since  $\nabla \ell(\mathbf{w}^{(1)}) = \nabla \ell(\mathbf{w}^{(2)}) = \mathbf{0}$  and  $\ell(\mathbf{w}^{(1)}) = \ell(\mathbf{w}^{(2)})$ , the loss function will appear sharper around  $\mathbf{w}^{(1)}$  than around  $\mathbf{w}^{(2)}$ , if  $\boldsymbol{\alpha}^T \mathbf{H}(\mathbf{w}^{(1)}) \boldsymbol{\alpha} > \boldsymbol{\alpha}^T \mathbf{H}(\mathbf{w}^{(2)}) \boldsymbol{\alpha}$ . From (S16), this condition is equivalent to  $\|\boldsymbol{\alpha}^T \mathbf{z}^{(1)}\|^2 > \|\boldsymbol{\alpha}^T \mathbf{z}^{(2)}\|^2$ , or more explicitly,

$$\left( (\mathbf{w}^{(1)} - \mathbf{w}^{(2)})^T \mathbf{z}^{(1)} \right)^2 > \left( (\mathbf{w}^{(1)} - \mathbf{w}^{(2)})^T \mathbf{z}^{(2)} \right)^2. \quad (\text{S25})$$

Since  $\mathbf{z} = [w_1^{-1}, w_2^{-1}, \dots, w_m^{-1}]^T$ , this inequality can be written as

$$\left| \sum_{i=1}^m \frac{w_i^{(2)}}{w_i^{(1)}} - m \right| > \left| \sum_{i=1}^m \frac{w_i^{(1)}}{w_i^{(2)}} - m \right|. \quad (\text{S26})$$

Note that

$$\sum_{i=1}^m \frac{w_i^{(2)}}{w_i^{(1)}} \geq m \sqrt[m]{\prod_{i=1}^m \frac{w_i^{(2)}}{w_i^{(1)}}} = m \sqrt[m]{\frac{\prod_{i=1}^m w_i^{(2)}}{\prod_{j=1}^m w_j^{(1)}}} = m \sqrt[m]{\frac{\tau}{\tau}} = m. \quad (\text{S27})$$

Similarly,  $\sum_{i=1}^m \frac{w_i^{(1)}}{w_i^{(2)}} \geq m$ . Therefore, (S26) can be reduced to

$$\sum_{i=1}^m \frac{w_i^{(2)}}{w_i^{(1)}} > \sum_{i=1}^m \frac{w_i^{(1)}}{w_i^{(2)}}. \quad (\text{S28})$$

■

Notice that the loss function is symmetric in a sense that if we flip the sign of two scalar layers, then it remains the same. Therefore, without loss of generality, we can restrict our analysis to a single orthant. Let  $\tau > 0$ , and  $\mathbf{w}^{(1)}$  be the flattest minimum in  $\mathbb{R}_+^m$ , i.e.  $w_i^{(1)} = \tau^{1/m}$  for all  $i \in \{1, \dots, m\}$ . Given a second minimum  $\mathbf{w}^{(2)} \in \mathbb{R}_+^m$  for which the connecting line between  $\mathbf{w}^{(1)}$  and  $\mathbf{w}^{(2)}$  is loyal to the true sharpness, we can construct a third solution  $\mathbf{w}^{(3)}$  that will appear deceptively flatter than  $\mathbf{w}^{(1)}$  over their connecting line. Specifically, let us set

$$w_i^{(3)} = \frac{(w_i^{(1)})^2}{w_i^{(2)}}. \quad (\text{S29})$$

Clearly,  $\mathbf{w}^{(3)}$  is a global minimum as  $\prod_{i=1}^m w_i^{(3)} = \tau$ . Since  $\mathbf{w}^{(2)}$  appears sharper than  $\mathbf{w}^{(1)}$  along their connecting line, then according to Claim S1 we have

$$\sum_{i=1}^m \frac{w_i^{(1)}}{w_i^{(2)}} > \sum_{i=1}^m \frac{w_i^{(2)}}{w_i^{(1)}}. \quad (\text{S30})$$

Thus,

$$\sum_{i=1}^m \frac{w_i^{(3)}}{w_i^{(1)}} = \sum_{i=1}^m \frac{w_i^{(1)}}{w_i^{(2)}} > \sum_{i=1}^m \frac{w_i^{(2)}}{w_i^{(1)}} = \sum_{i=1}^m \frac{w_i^{(1)}}{w_i^{(3)}}. \quad (\text{S31})$$

Therefore, by Claim S1,  $\mathbf{w}^{(1)}$  appears sharper than  $\mathbf{w}^{(3)}$  along their connecting line.

In the special case of two layer networks ( $m = 2$ ), we have that for any minimizer,  $w_2^{(i)} = \tau/w_1^{(i)}$ . Hence,

$$\frac{w_1^{(1)}}{w_1^{(2)}} + \frac{w_2^{(1)}}{w_2^{(2)}} = \frac{w_2^{(2)}}{w_2^{(1)}} + \frac{w_1^{(2)}}{w_1^{(1)}}. \quad (\text{S32})$$

This means that the minima will appear equally sharp.

#### IV. Proof of Lemma 3

In this section we derive the Hessian matrix defined in (18) at a global minimum point, *i.e.* for  $\mathbf{w} \in \Omega$ . Throughout this section we will be using the following properties of the Kronecker product. For any matrices  $\mathbf{M}_1, \mathbf{M}_2, \mathbf{M}_3, \mathbf{M}_4$ ,

$$\text{vec}(\mathbf{M}_1 \mathbf{M}_2 \mathbf{M}_3) = (\mathbf{M}_3^T \otimes \mathbf{M}_1) \text{vec}(\mathbf{M}_2), \quad (\text{P1})$$

$$(\mathbf{M}_1 \otimes \mathbf{M}_2)^T = (\mathbf{M}_1^T \otimes \mathbf{M}_2^T), \quad (\text{P2})$$

$$(\mathbf{M}_1 \otimes \mathbf{M}_2)(\mathbf{M}_3 \otimes \mathbf{M}_4) = (\mathbf{M}_1 \mathbf{M}_3) \otimes (\mathbf{M}_2 \mathbf{M}_4). \quad (\text{P3})$$

Let us start the computation of  $\mathbf{H}_{\mathbf{w}}$  by rearranging the loss function so as to simplify the differentiation w.r.t.  $\mathbf{w}_k$ . Specifically, we have that

$$\begin{aligned} \ell(\mathbf{w}) &= \hat{\mathbb{E}} \left[ \left\| y - \prod_{j=1}^m \mathbf{W}_j x \right\|^2 \right] \\ &= \hat{\mathbb{E}} \left[ \left\| y - \left( \prod_{i=k+1}^m \mathbf{W}_i \right) \mathbf{w}_k \left( \prod_{j=1}^{k-1} \mathbf{W}_j x \right) \right\|^2 \right] \\ &= \hat{\mathbb{E}} \left[ \left\| y - \left( \prod_{j=1}^{k-1} \mathbf{W}_j x \right)^T \otimes \left( \prod_{i=k+1}^m \mathbf{W}_i \right) \mathbf{w}_k \right\|^2 \right] \\ &= \hat{\mathbb{E}} \left[ \left\| y - \left[ x^T \left( \prod_{j=1}^{k-1} \mathbf{W}_j \right)^T \right] \otimes \left[ \mathbf{I} \left( \prod_{i=k+1}^m \mathbf{W}_i \right) \right] \mathbf{w}_k \right\|^2 \right] \\ &= \hat{\mathbb{E}} \left[ \left\| y - \left[ x \otimes \mathbf{I} \right]^T \left[ \left( \prod_{j=1}^{k-1} \mathbf{W}_j \right)^T \otimes \left( \prod_{i=k+1}^m \mathbf{W}_i \right) \right] \mathbf{w}_k \right\|^2 \right], \end{aligned} \quad (\text{S33})$$

where in the third equality we used property (P1), and in the last we used properties (P2) and (P3). To simplify expressions, we define the following matrices

$$\mathbf{U}_k \triangleq \left( \prod_{j=1}^{k-1} \mathbf{W}_j \right)^T \otimes \left( \prod_{i=k+1}^m \mathbf{W}_i \right) \quad \text{and} \quad \mathbf{X} \triangleq x \otimes \mathbf{I}. \quad (\text{S34})$$

Thus, the loss function (2) is given by

$$\ell(\mathbf{w}) = \hat{\mathbb{E}} \left[ \left\| y - \mathbf{X}^T \mathbf{U}_k \mathbf{w}_k \right\|^2 \right]. \quad (\text{S35})$$

Now we are ready to calculate the partial derivative of  $\ell(\mathbf{w})$  w.r.t  $\mathbf{w}_k$ . Notice that  $\mathbf{U}_k$  is not a function of  $\mathbf{w}_k$ , therefore

$$\frac{\partial}{\partial \mathbf{w}_k} \hat{\mathbb{E}} \left[ \left\| y - \mathbf{X}^T \mathbf{U}_k \mathbf{w}_k \right\|^2 \right] = -2 \hat{\mathbb{E}} \left[ \mathbf{U}_k^T \mathbf{X} (y - \mathbf{X}^T \mathbf{U}_k \mathbf{w}_k) \right] = 2 \mathbf{U}_k^T \left( \hat{\mathbb{E}} \left[ \mathbf{X} \mathbf{X}^T \right] \mathbf{U}_k \mathbf{w}_k - \hat{\mathbb{E}} \left[ \mathbf{X} y \right] \right). \quad (\text{S36})$$

Furthermore,

$$\hat{\mathbb{E}}[\mathbf{X}\mathbf{X}^T] = \hat{\mathbb{E}}[(x \otimes \mathbf{I})(x \otimes \mathbf{I})^T] = \hat{\mathbb{E}}[(xx^T \otimes \mathbf{I})] = (\hat{\mathbb{E}}[xx^T]) \otimes \mathbf{I} = \hat{\Sigma}_x \otimes \mathbf{I}, \quad (\text{S37})$$

where in the second equality we used properties (P2) and (P3), and in the third equality we used the linearity of the Kronecker product. Additionally,

$$\hat{\mathbb{E}}[\mathbf{X}y] = \hat{\mathbb{E}}[(x \otimes \mathbf{I})y] = \hat{\mathbb{E}}[\text{vec}(yx^T)] = \text{vec}(\hat{\Sigma}_{yx}), \quad (\text{S38})$$

where in the second step we used (P1). Overall we have that

$$\frac{\partial}{\partial \mathbf{w}_k} \ell(\mathbf{w}) = 2\mathbf{U}_k^T \left[ (\hat{\Sigma}_x \otimes \mathbf{I}) \mathbf{U}_k \mathbf{w}_k - \text{vec}(\hat{\Sigma}_{yx}) \right]. \quad (\text{S39})$$

Next we prepare Eq. (S39) for differentiation w.r.t  $\mathbf{w}_q$ . First, for all  $k$

$$\mathbf{U}_k \mathbf{w}_k = \left[ \left( \prod_{j=1}^{k-1} \mathbf{W}_j \right)^T \otimes \left( \prod_{i=k+1}^m \mathbf{W}_i \right) \right] \text{vec}(\mathbf{W}_k) = \text{vec} \left( \left( \prod_{i=k+1}^m \mathbf{W}_i \right) \mathbf{W}_k \left( \prod_{j=1}^{k-1} \mathbf{W}_j \right) \right) = \text{vec} \left( \prod_{i=1}^m \mathbf{W}_i \right). \quad (\text{S40})$$

Particularly, this means that the value of the term  $\mathbf{U}_k \mathbf{w}_k$  is the same for all  $k$ . Hence,  $\mathbf{U}_k \mathbf{w}_k = \mathbf{U}_q \mathbf{w}_q$  and therefore

$$\frac{\partial}{\partial \mathbf{w}_k} \ell(\mathbf{w}) = 2\mathbf{U}_k^T \left[ (\hat{\Sigma}_x \otimes \mathbf{I}) \mathbf{U}_q \mathbf{w}_q - \text{vec}(\hat{\Sigma}_{yx}) \right]. \quad (\text{S41})$$

Now, let us differentiate the vector  $\frac{\partial}{\partial \mathbf{w}_k} \ell(\mathbf{w})$  w.r.t. the scalar  $w_{q,l}$ , which is the  $l$ th element in the vector  $\mathbf{w}_q$ . Notice that  $\mathbf{U}_k$  and  $\mathbf{w}_q$  itself are the only terms which depend on  $\mathbf{w}_q$ . Therefore, by the product rule of differentiation using the denominator-layout notation<sup>2</sup>

$$\begin{aligned} \frac{\partial^2}{\partial w_{q,l} \partial \mathbf{w}_k} \ell(\mathbf{w}) &= 2 \frac{\partial}{\partial w_{q,l}} \left( \mathbf{U}_k^T \left[ (\hat{\Sigma}_x \otimes \mathbf{I}) \mathbf{U}_q \mathbf{w}_q - \text{vec}(\hat{\Sigma}_{yx}) \right] \right) \\ &= 2 \left[ (\hat{\Sigma}_x \otimes \mathbf{I}) \mathbf{U}_q \mathbf{w}_q - \text{vec}(\hat{\Sigma}_{yx}) \right]^T \left( \frac{\partial}{\partial w_{q,l}} \mathbf{U}_k^T \right) + 2 \left( \frac{\partial}{\partial w_{q,l}} \mathbf{w}_q \right) \mathbf{U}_q^T (\hat{\Sigma}_x \otimes \mathbf{I}) \mathbf{U}_k. \end{aligned} \quad (\text{S42})$$

However, at a global minimum

$$\left( \hat{\Sigma}_x \otimes \mathbf{I} \right) \mathbf{U}_q \mathbf{w}_q = \left( \hat{\Sigma}_x \otimes \mathbf{I} \right) \text{vec} \left( \prod_{i=1}^m \mathbf{W}_i \right) = \text{vec} \left( \prod_{i=1}^m \mathbf{W}_i \hat{\Sigma}_x \right) = \text{vec} \left( \hat{\Sigma}_{yx} \hat{\Sigma}_x^{-1} \hat{\Sigma}_x \right) = \text{vec}(\hat{\Sigma}_{yx}), \quad (\text{S43})$$

where in the first equality we used (S40), in the second equality we used (P1) and in the third equality we used the assumption that  $\mathbf{w} \in \Omega$ . Hence, for all  $1 \leq q \leq m$  we have that

$$\left( \hat{\Sigma}_x \otimes \mathbf{I} \right) \mathbf{U}_q \mathbf{w}_q - \text{vec}(\hat{\Sigma}_{yx}) = \mathbf{0}. \quad (\text{S44})$$

Therefore, (S42) is reduced to

$$\frac{\partial^2}{\partial w_{q,l} \partial \mathbf{w}_k} \ell(\mathbf{w}) = 2 \left( \frac{\partial}{\partial w_{q,l}} \mathbf{w}_q \right) \mathbf{U}_q^T (\hat{\Sigma}_x \otimes \mathbf{I}) \mathbf{U}_k, \quad (\text{S45})$$

for all  $1 \leq l \leq m$ . Hence,

$$\frac{\partial^2}{\partial w_q \partial \mathbf{w}_k} \ell(\mathbf{w}) = 2 \left( \frac{\partial}{\partial \mathbf{w}_q} \mathbf{w}_q \right) \mathbf{U}_q^T (\hat{\Sigma}_x \otimes \mathbf{I}) \mathbf{U}_k = 2\mathbf{U}_q^T (\hat{\Sigma}_x \otimes \mathbf{I}) \mathbf{U}_k = 2\mathbf{U}_q^T \left( \hat{\Sigma}_x^{\frac{1}{2}} \otimes \mathbf{I} \right)^T \left( \hat{\Sigma}_x^{\frac{1}{2}} \otimes \mathbf{I} \right) \mathbf{U}_k, \quad (\text{S46})$$

<sup>2</sup>Where the derivative  $\frac{\partial \mathbf{A}}{\partial z}$  of a matrix  $\mathbf{A}$  w.r.t. a scalar  $z$  is laid out according to  $\mathbf{A}^T$ .

where  $\hat{\Sigma}_x^{\frac{1}{2}}$  is the symmetric square root matrix of  $\hat{\Sigma}_x$ . Let us define the matrices  $\{\Phi_k\}_{k=1}^m$  as

$$\Phi_k = U_k^T \left( \hat{\Sigma}_x^{\frac{1}{2}} \otimes I \right)^T = \left[ \left( \prod_{j=1}^{k-1} W_j \right) \otimes \left( \prod_{i=k+1}^m W_i \right)^T \right] \left( \hat{\Sigma}_x^{\frac{1}{2}} \otimes I \right) = \left( \prod_{j=1}^{k-1} W_j \hat{\Sigma}_x^{\frac{1}{2}} \right) \otimes \left( \prod_{i=k+1}^m W_i \right)^T. \quad (\text{S47})$$

Thus,

$$\frac{\partial^2}{\partial w_q \partial w_k} \ell(\mathbf{w}) = 2\Phi_q \Phi_k^T. \quad (\text{S48})$$

Finally, the Hessian matrix of the loss function  $\ell(\mathbf{w})$  is given by

$$\mathbf{H}_w = 2\Phi\Phi^T, \quad (\text{S49})$$

where  $\Phi = [\Phi_1^T, \Phi_2^T, \dots, \Phi_m^T]^T$ .

## V. The Missing Parts of the Proof of Theorem 1

### V.1. Proof of Lemma 4

The proof is straightforward. We have

$$\sum_{k=1}^m \|\Psi_k\|_F^2 \geq \sum_{k=1}^m \|\Psi_k\|_2^2 \geq m \left[ \prod_{k=1}^m \|\Psi_k\|_2 \right]^{\frac{2}{m}} \geq m \left[ \left\| \prod_{k=1}^m \Psi_k \right\|_2 \right]^{\frac{2}{m}}, \quad (\text{S50})$$

where in the first inequality we used the fact that  $\|\Psi\|_F \geq \|\Psi\|_2$  for any matrix  $\Psi \in \mathbb{R}^{d_1 \times d_2}$ . The second inequality is due to the inequality of arithmetic and geometric means. In the final inequality we used the fact that  $\|\cdot\|_2$  is a sub-multiplicative matrix norm, meaning  $\|\Psi\|_2 \|\Phi\|_2 \geq \|\Psi\Phi\|_2$  for any pair of matrices  $\Psi \in \mathbb{R}^{d_1 \times d_2}$ ,  $\Phi \in \mathbb{R}^{d_2 \times d_3}$ .

### V.2. Maximal Value of $\nu$

On the one hand, for any  $B \in \mathbb{R}^{d_y \times d_x}$  such that  $\|B\|_F = 1$ ,

$$\left\| (BT^T)^{m-1} B \right\|_2 \leq \|B\|_2^m \|T\|_2^{m-1} \leq (\sigma_{\max}(T))^{m-1}, \quad (\text{S51})$$

where in the second inequality we used  $\|B\|_2 \leq \|B\|_F = 1$ . On the other hand, this upper bound is achieved by  $B = \mathbf{u}\mathbf{v}^T$ , as

$$\begin{aligned} \left\| (\mathbf{u}\mathbf{v}^T T^T)^{m-1} \mathbf{u}\mathbf{v}^T \right\|_2 &= \left\| \mathbf{u} (\mathbf{v}^T T^T \mathbf{u})^{m-1} \mathbf{v}^T \right\|_2 = (\mathbf{v}^T T^T \mathbf{u})^{m-1} \|\mathbf{u}\mathbf{v}^T\|_2 \\ &= (\sigma_{\max}(T))^{m-1} \|\mathbf{u}\| \|\mathbf{v}\| = (\sigma_{\max}(T))^{m-1}. \end{aligned} \quad (\text{S52})$$

Therefore,

$$\max_{\|B\|_F=1} \nu(B) = 2m \times (\sigma_{\max}(T))^{2(1-\frac{1}{m})}. \quad (\text{S53})$$

### V.3. Maximal Eigenvalue at the Canonical Solution (27)

In (20) we have

$$\lambda_{\max}(\mathbf{H}_w) = \max_{\|b\|=1} 2\|\Phi b\|^2. \quad (\text{S54})$$

Note that  $\|\Phi b\|^2 = \sum_{k=1}^m \|\Phi_k b\|^2$ . Using the definition of  $\Phi_k$  in (16) we get

$$\|\Phi_k b\|^2 = \left\| \left( \prod_{j=1}^{k-1} W_j \hat{\Sigma}_x^{\frac{1}{2}} \right) \otimes \left( \prod_{i=k+1}^m W_i \right)^T b \right\|^2 = \left\| \left( \prod_{i=k+1}^m W_i \right)^T B \left( \prod_{j=1}^{k-1} W_j \hat{\Sigma}_x^{\frac{1}{2}} \right)^T \right\|_F^2, \quad (\text{S55})$$

where we used (P1) with  $\mathbf{b} = \text{vec}(\mathbf{B})$ . Therefore,

$$\lambda_{\max}(\mathbf{H}_{\mathbf{w}}) = \max_{\|\mathbf{B}\|_{\text{F}}=1} 2 \sum_{k=1}^m \left\| \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{B} \hat{\Sigma}_x^{\frac{1}{2}} \left( \prod_{j=1}^{k-1} \mathbf{W}_j \right)^T \right\|_{\text{F}}^2. \quad (\text{S56})$$

Substituting the canonical solution (27) in this optimization problem, we obtain

$$\sum_{k=1}^m \left\| \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{B} \left( \prod_{j=1}^{k-1} \mathbf{W}_j \right)^T \right\|_{\text{F}}^2 = \sum_{k=1}^m \left\| \left( \prod_{i=k+1}^m \mathbf{S}_i^{\frac{1}{m}} \right)^T \mathbf{U}^T \mathbf{B} \mathbf{V} \left( \prod_{j=1}^{k-1} \mathbf{S}_j^{\frac{1}{m}} \right)^T \right\|_{\text{F}}^2, \quad (\text{S57})$$

where in the first and the last terms of the series ( $k=1$  and  $k=m$ ) we used the fact that  $\mathbf{U}$  and  $\mathbf{V}$  are unitary matrices, so that

$$\left\| \left( \prod_{i=2}^m \mathbf{S}_i^{\frac{1}{m}} \right)^T \mathbf{U}^T \mathbf{B} \right\|_{\text{F}}^2 = \left\| \left( \prod_{i=2}^m \mathbf{S}_i^{\frac{1}{m}} \right)^T \mathbf{U}^T \mathbf{B} \mathbf{V} \right\|_{\text{F}}^2, \quad \left\| \mathbf{B} \mathbf{V} \left( \prod_{j=1}^{m-1} \mathbf{S}_j^{\frac{1}{m}} \right)^T \right\|_{\text{F}}^2 = \left\| \mathbf{U}^T \mathbf{B} \mathbf{V} \left( \prod_{j=1}^{m-1} \mathbf{S}_j^{\frac{1}{m}} \right)^T \right\|_{\text{F}}^2. \quad (\text{S58})$$

Note that  $\prod_{i=k+1}^m \mathbf{S}_i^{\frac{1}{m}}$  is a diagonal  $d_y \times d_k$  matrix, whose  $q$ th diagonal entry is  $(\sigma_q(\mathbf{T}))^{(m-k)/m}$  (where  $\sigma_q(\mathbf{T})$  is the  $q$ th largest singular value of  $\mathbf{T}$ ). Similarly,  $\prod_{j=1}^{k-1} \mathbf{S}_j^{\frac{1}{m}}$  is a diagonal  $d_{k-1} \times d_x$  matrix, whose  $q$ th diagonal entry is  $(\sigma_q(\mathbf{T}))^{(k-1)/m}$ . Therefore, we can write

$$\sum_{k=1}^m \left\| \left( \prod_{i=k+1}^m \mathbf{S}_i^{\frac{1}{m}} \right)^T \mathbf{U}^T \mathbf{B} \mathbf{V} \left( \prod_{j=1}^{k-1} \mathbf{S}_j^{\frac{1}{m}} \right)^T \right\|_{\text{F}}^2 = \sum_{k=1}^m \left\| \left( \mathbf{S}^{\frac{m-k}{m}} \right)^T \mathbf{U}^T \mathbf{B} \mathbf{V} \left( \mathbf{S}^{\frac{k-1}{m}} \right)^T \right\|_{\text{F}}^2, \quad (\text{S59})$$

where  $\mathbf{S}^\alpha$  denotes a  $d_y \times d_x$  diagonal matrix whose  $q$ th diagonal entry is  $(\sigma_q(\mathbf{T}))^\alpha$ . Here, we used the fact that the Frobenius norm is unaffected by zero entries, and thus removed/added zero rows/columns.

Next, we perform the change of variables  $\tilde{\mathbf{B}} = \mathbf{U}^T \mathbf{B} \mathbf{V} \in \mathbb{R}^{d_y \times d_x}$  to obtain the following optimization problem

$$\max_{\tilde{\mathbf{B}} \in \mathbb{R}^{d \times d}} 2 \sum_{k=1}^m \left\| \left( \mathbf{S}^{\frac{m-k}{m}} \right)^T \tilde{\mathbf{B}} \left( \mathbf{S}^{\frac{k-1}{m}} \right)^T \right\|_{\text{F}}^2 \quad \text{s.t.} \quad \|\tilde{\mathbf{B}}\|_{\text{F}}^2 = 1. \quad (\text{S60})$$

Writing the objective in terms of the elements of  $\tilde{\mathbf{B}}$ , which we denote by  $\{\tilde{b}_{i,j}\}$ , gives

$$2 \sum_{k=1}^m \left\| \left( \mathbf{S}^{\frac{m-k}{m}} \right)^T \tilde{\mathbf{B}} \left( \mathbf{S}^{\frac{k-1}{m}} \right)^T \right\|_{\text{F}}^2 = 2 \sum_{k=1}^m \sum_{i=1}^d \sum_{j=1}^d \left[ (\sigma_j(\mathbf{T}))^{\frac{m-k}{m}} (\sigma_i(\mathbf{T}))^{\frac{k-1}{m}} \tilde{b}_{i,j} \right]^2, \quad (\text{S61})$$

where  $d = \min\{d_x, d_y\}$  is the number of singular values of  $\mathbf{T}$ . By changing the order of the summation, we get

$$\max_{\tilde{b}_{1,1}, \dots, \tilde{b}_{d,d} \in \mathbb{R}} 2 \sum_{i,j=1}^d \tilde{b}_{i,j}^2 \sum_{k=1}^m \left[ (\sigma_j(\mathbf{T}))^{\frac{m-k}{m}} (\sigma_i(\mathbf{T}))^{\frac{k-1}{m}} \right]^2 \quad \text{s.t.} \quad \sum_{i,j=1}^d \tilde{b}_{i,j}^2 = 1. \quad (\text{S62})$$

This is a simple linear optimization problem over the unit simplex, whose optimal value is attained at one of the vertices,

$$\max_{i,j \in \{1, \dots, d\}} 2 \sum_{k=1}^m \left[ (\sigma_j(\mathbf{T}))^{\frac{m-k}{m}} (\sigma_i(\mathbf{T}))^{\frac{k-1}{m}} \right]^2. \quad (\text{S63})$$

The maximal value is attained for  $i=j=1$ , thus the value of (21) for the canonical solution is

$$2 \sum_{k=1}^m \left[ (\sigma_1(\mathbf{T}))^{\frac{m-k}{m}} (\sigma_1(\mathbf{T}))^{\frac{k-1}{m}} \right]^2 = 2m \times (\sigma_{\max}(\mathbf{T}))^{2(1-\frac{1}{m})}. \quad (\text{S64})$$

This result shows that the canonical solution (27) is indeed a minimizer of the maximal eigenvalue of the Hessian matrix.

#### V.4. Proof of the Top Eigenvector of $H_w$

On the one hand, according to Section 5, for any flattest minimum point  $w \in \Omega_0$ , the largest eigenvalue satisfies

$$\lambda_{\max}(H_w) = 2m \times (\sigma_{\max}(T))^{2(1-\frac{1}{m})}. \quad (\text{S65})$$

On the other hand, the maximal eigenvalue of the Hessian matrix is the solution to the optimization problem (S56), in which  $b = \text{vec}(B)$  is the eigenvector of  $\hat{H}_w$  (see (S54)). Substituting  $B^* = uv^T$  (i.e.  $b^* = v \otimes u$ ) in the objective function, we get

$$\begin{aligned} 2 \sum_{k=1}^m \left\| \left( \prod_{i=k+1}^m W_i \right)^T B^* \left( \prod_{j=1}^{k-1} W_j \right)^T \right\|_F^2 &\geq 2m \times \left\| \left[ B^* \left( \prod_{i=1}^m W_i \right)^T \right]^{m-1} B^* \right\|_2^{\frac{2}{m}} \\ &= 2m \times \left\| (B^* T^T)^{m-1} B^* \right\|_2^{\frac{2}{m}} \\ &= 2m \times (\sigma_{\max}(T))^{2(1-\frac{1}{m})}, \end{aligned} \quad (\text{S66})$$

where in the second inequality we used Lemma 4 and explicitly unrolled the product, as in (24), and in the last step we used (S52). This proves that  $b^* = v \otimes u$  is an eigenvector of  $\hat{H}_w$  corresponding to the maximal eigenvalue. Now, since  $\hat{H}_w = 2\Phi^T \Phi$  and  $H_w = 2\Phi \Phi^T$ , we have that  $\Phi b^* = \Phi(v \otimes u)$  is the eigenvector of  $H_w$  corresponding to its maximal eigenvalue.

#### VI. Proof of Theorem 2

Let us start the proof by presenting two lemmas.

**Lemma S1.** *Let  $\hat{\Sigma}_x = I$ . If  $w \in \Omega_0$  then for all  $k \in \{1, 2, \dots, m\}$*

$$\left\| u^T \prod_{i=k+1}^m W_i \right\| \left\| \prod_{j=1}^{k-1} W_j v \right\| = (\sigma_{\max}(T))^{1-\frac{1}{m}}. \quad (\text{S67})$$

*Proof.* First, observe that for  $B^* = uv^T$ , the left-hand side of (S67) can be written as

$$\left\| u^T \prod_{i=k+1}^m W_i \right\| \left\| \prod_{j=1}^{k-1} W_j v \right\| = \left\| \left( \prod_{i=k+1}^m W_i \right)^T B^* \left( \prod_{j=1}^{k-1} W_j \right)^T \right\|_2. \quad (\text{S68})$$

Now, from Theorem 1 and Eq. (21) we can conclude that

$$\lambda_{\max}(H_w) = 2 \sum_{k=1}^m \left\| \left( \prod_{i=k+1}^m W_i \right)^T B^* \left( \prod_{j=1}^{k-1} W_j \right)^T \right\|_2^2 = 2m \times (\sigma_{\max}(T))^{2(1-\frac{1}{m})}. \quad (\text{S69})$$

Note that since  $B^*$  is a rank-1 matrix, the entire expression within the norm is rank-1, which is the reason we could replace the Frobenius norm appearing in (21) by the operator norm (the two norms coincide for rank-1 matrices). Furthermore, by the inequality of arithmetic and geometric means

$$\begin{aligned} 2 \sum_{k=1}^m \left\| \left( \prod_{i=k+1}^m W_i \right)^T B^* \left( \prod_{j=1}^{k-1} W_j \right)^T \right\|_2^2 &\geq 2m \left[ \prod_{k=1}^m \left\| \left( \prod_{i=k+1}^m W_i \right)^T B^* \left( \prod_{j=1}^{k-1} W_j \right)^T \right\|_2^2 \right]^{\frac{1}{m}} \\ &\geq 2m \left[ \left\| \prod_{k=1}^m \left( \prod_{i=k+1}^m W_i \right)^T B^* \left( \prod_{j=1}^{k-1} W_j \right)^T \right\|_2^2 \right]^{\frac{1}{m}} \\ &= 2m \times (\sigma_{\max}(T))^{2(1-\frac{1}{m})}, \end{aligned} \quad (\text{S70})$$

where the second inequality is due to the sub-multiplicativity property of the operator norm, and in the last step we unrolled the product, as in (24), and used (S52). From (S69) and (S70) we obtain that the inequality of arithmetic and geometric means in (S70) is achieved with equality. This happens if and only if all summands in the series are equal. Thus, we conclude that for all  $k \in \{1, 2, \dots, m\}$ ,

$$\left\| \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{B}^* \left( \prod_{j=1}^{k-1} \mathbf{W}_j \right)^T \right\|_2 = (\sigma_{\max}(\mathbf{T}))^{1 - \frac{1}{m}}, \quad (\text{S71})$$

and together with (S68), this implies that

$$\left\| \mathbf{u}^T \prod_{i=k+1}^m \mathbf{W}_i \right\| \left\| \prod_{j=1}^{k-1} \mathbf{W}_j \mathbf{v} \right\| = (\sigma_{\max}(\mathbf{T}))^{1 - \frac{1}{m}}. \quad (\text{S72})$$

■

While Lemma S1 characterizes the norms of the vectors  $\mathbf{u}^T \prod_{i=k+1}^m \mathbf{W}_i$  and  $\prod_{j=1}^{k-1} \mathbf{W}_j \mathbf{v}$ , the next Lemma characterizes their directions.

**Lemma S2.** *Let  $\hat{\Sigma}_x = \mathbf{I}$ . If  $\mathbf{w} \in \Omega_0$  then for all  $k \in \{0, 1, 2, \dots, m\}$*

$$\frac{1}{\left\| \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{u} \right\|} \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{u} = \frac{1}{\left\| \prod_{j=1}^k \mathbf{W}_j \mathbf{v} \right\|} \prod_{j=1}^k \mathbf{W}_j \mathbf{v}. \quad (\text{S73})$$

*Proof.* From Lemma S1 we have

$$\prod_{k=1}^m \left\| \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{u} \right\| \left\| \mathbf{v}^T \left( \prod_{j=1}^{k-1} \mathbf{W}_j \right)^T \right\| = (\sigma_{\max}(\mathbf{T}))^{m-1}. \quad (\text{S74})$$

Recall that in our convention (see Section 4), for  $k = 1, m$  we have

$$\left\| \mathbf{v}^T \left( \prod_{j=1}^0 \mathbf{W}_j \right)^T \right\| = \|\mathbf{v}^T\| = 1, \quad \left\| \left( \prod_{i=m+1}^m \mathbf{W}_i \right)^T \mathbf{u} \right\| = \|\mathbf{u}\| = 1. \quad (\text{S75})$$

Therefore, (S74) can be written as

$$\prod_{k=1}^{m-1} \left\| \mathbf{v}^T \left( \prod_{j=1}^k \mathbf{W}_j \right)^T \right\| \left\| \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{u} \right\| = (\sigma_{\max}(\mathbf{T}))^{m-1}. \quad (\text{S76})$$

On the other hand, by the Cauchy–Schwarz inequality we have

$$\begin{aligned} \prod_{k=1}^{m-1} \left\| \mathbf{v}^T \left( \prod_{j=1}^k \mathbf{W}_j \right)^T \right\| \left\| \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{u} \right\| &\geq \prod_{k=1}^{m-1} \mathbf{v}^T \left( \prod_{j=1}^k \mathbf{W}_j \right)^T \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{u} \\ &= \prod_{k=1}^{m-1} \mathbf{v}^T \mathbf{T}^T \mathbf{u} \\ &= (\sigma_{\max}(\mathbf{T}))^{(m-1)}. \end{aligned} \quad (\text{S77})$$

From (S76) we have that the Cauchy–Schwarz inequalities are achieved with equality. Thus, for all  $k \in \{0, 1, 2, \dots, m\}$

$$\frac{1}{\left\| \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{u} \right\|} \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{u} = \frac{1}{\left\| \prod_{j=1}^k \mathbf{W}_j \mathbf{v} \right\|} \prod_{j=1}^k \mathbf{W}_j \mathbf{v}. \quad (\text{S78})$$

■

Now we are ready to prove Theorem 2. From Lemma S2, we have

$$\frac{1}{\left\| \prod_{j=1}^k \mathbf{W}_j \mathbf{v} \right\|} \prod_{j=1}^k \mathbf{W}_j \mathbf{v} = \frac{1}{\left\| \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{u} \right\|} \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{u}. \quad (\text{S79})$$

Multiplying by  $(\prod_{j=1}^k \mathbf{W}_j)^T$  from the left, (S79) becomes

$$\frac{1}{\left\| \prod_{j=1}^k \mathbf{W}_j \mathbf{v} \right\|} \left( \prod_{j=1}^k \mathbf{W}_j \right)^T \prod_{j=1}^k \mathbf{W}_j \mathbf{v} = \frac{1}{\left\| \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{u} \right\|} \mathbf{T}^T \mathbf{u}. \quad (\text{S80})$$

Note that  $\mathbf{T}^T \mathbf{u} = \sigma_{\max}(\mathbf{T}) \mathbf{v}$ . Therefore,

$$\left[ \left( \prod_{j=1}^k \mathbf{W}_j \right)^T \prod_{j=1}^k \mathbf{W}_j \right] \mathbf{v} = \frac{\left\| \prod_{j=1}^k \mathbf{W}_j \mathbf{v} \right\|}{\left\| \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{u} \right\|} \sigma_{\max}(\mathbf{T}) \mathbf{v}. \quad (\text{S81})$$

This shows that  $\mathbf{v}$  is an eigenvector of  $(\prod_{j=1}^k \mathbf{W}_j)^T \prod_{j=1}^k \mathbf{W}_j$ , *i.e.* a singular vector of  $\prod_{j=1}^k \mathbf{W}_j$ . To compute the corresponding singular value, let us multiply this equation by  $\mathbf{v}^T$  from the left to get the following result.

$$\left\| \prod_{j=1}^k \mathbf{W}_j \mathbf{v} \right\| \left\| \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{u} \right\| = \sigma_{\max}(\mathbf{T}). \quad (\text{S82})$$

Recall from Lemma S1 that  $\|\mathbf{u}^T \prod_{i=k+1}^m \mathbf{W}_i\| = (\sigma_{\max}(\mathbf{T}))^{1-1/m} / \|\prod_{j=1}^{k-1} \mathbf{W}_j \mathbf{v}\|$ . Substituting into (S82), we obtain that

$$\left\| \prod_{j=1}^k \mathbf{W}_j \mathbf{v} \right\| = (\sigma_{\max}(\mathbf{T}))^{\frac{1}{m}} \times \left\| \prod_{j=1}^{k-1} \mathbf{W}_j \mathbf{v} \right\|. \quad (\text{S83})$$

By unwrapping this recursive formula with an initial condition for  $k=0$  of  $\|\prod_{j=1}^0 \mathbf{W}_j \mathbf{v}\| = \|\mathbf{v}\| = 1$ , we get

$$\left\| \prod_{j=1}^k \mathbf{W}_j \mathbf{v} \right\| = \sigma_{\max}(\mathbf{T})^{\frac{k}{m}}. \quad (\text{S84})$$

The proof for the left singular vector and its corresponding singular value is the same.

Next, we prove the bound on the intermediate gain. By Theorem 1 and Eq. (21),

$$\max_{\|\mathbf{B}\|_{\text{F}}=1} \sum_{l=1}^m \left\| \left( \prod_{i=l+1}^m \mathbf{W}_i \right)^T \mathbf{B} \left( \prod_{j=1}^{l-1} \mathbf{W}_j \right)^T \right\|_{\text{F}}^2 = m \times (\sigma_{\max}(\mathbf{T}))^{2(1-\frac{1}{m})}. \quad (\text{S85})$$

Now, for any  $k$ , we have that

$$\max_{\|\mathbf{B}\|_{\text{F}}=1} \sum_{l=1}^m \left\| \left( \prod_{i=l+1}^m \mathbf{W}_i \right)^T \mathbf{B} \left( \prod_{j=1}^{l-1} \mathbf{W}_j \right)^T \right\|_{\text{F}}^2 \geq \max_{\|\mathbf{B}\|_{\text{F}}=1} \left\| \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{B} \left( \prod_{j=1}^{k-1} \mathbf{W}_j \right)^T \right\|_{\text{F}}^2. \quad (\text{S86})$$

Furthermore, note that

$$\begin{aligned} \max_{\|\mathbf{B}\|_{\text{F}}=1} \left\| \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{B} \left( \prod_{j=1}^{k-1} \mathbf{W}_j \right)^T \right\|_{\text{F}} &= \max_{\|b\|=1} \left\| \left( \prod_{j=1}^{k-1} \mathbf{W}_j \right) \otimes \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{b} \right\| \\ &= \sigma_{\max} \left( \left( \prod_{j=1}^{k-1} \mathbf{W}_j \right) \otimes \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \right) \\ &= \sigma_{\max} \left( \prod_{i=1}^{k-1} \mathbf{W}_i \right) \times \sigma_{\max} \left( \prod_{i=k+1}^m \mathbf{W}_i \right). \end{aligned} \quad (\text{S87})$$

Therefore, this implies that

$$\sigma_{\max} \left( \prod_{i=1}^{k-1} \mathbf{W}_i \right) \times \sigma_{\max} \left( \prod_{i=k+1}^m \mathbf{W}_i \right) \leq \sqrt{m} \times (\sigma_{\max}(\mathbf{T}))^{(1-\frac{1}{m})}, \quad (\text{S88})$$

or, equivalently, that

$$\sigma_{\max} \left( \prod_{i=1}^{k-1} \mathbf{W}_i \right) \leq \frac{\sqrt{m} \times (\sigma_{\max}(\mathbf{T}))^{(1-\frac{1}{m})}}{\sigma_{\max} \left( \prod_{i=k+1}^m \mathbf{W}_i \right)}. \quad (\text{S89})$$

By the first part of Theorem 2 we have

$$\sigma_{\max} \left( \prod_{i=k+1}^m \mathbf{W}_i \right) \geq \left( \prod_{i=k+1}^m \mathbf{W}_i \mathbf{u} \right) = (\sigma_{\max}(\mathbf{T}))^{(1-\frac{k}{m})}. \quad (\text{S90})$$

Hence, we can further bound (S89) from above as

$$\sigma_{\max} \left( \prod_{i=1}^{k-1} \mathbf{W}_i \right) \leq \frac{\sqrt{m} \times (\sigma_{\max}(\mathbf{T}))^{(1-\frac{1}{m})}}{\sigma_{\max} \left( \prod_{i=k+1}^m \mathbf{W}_i \right)} \leq \sqrt{m} \times (\sigma_{\max}(\mathbf{T}))^{\frac{k-1}{m}}. \quad (\text{S91})$$

The proof of the other direction is similar.

## VII. Proof of Theorem 3

By Lemma S2

$$\frac{1}{\left\| \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{u} \right\|} \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{u} = \frac{1}{\left\| \prod_{j=1}^k \mathbf{W}_j \mathbf{v} \right\|} \prod_{j=1}^k \mathbf{W}_j \mathbf{v}. \quad (\text{S92})$$

Multiplying both sides by  $\mathbf{W}_k^T$  from the left, we obtain

$$\frac{1}{\left\| \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{u} \right\|} \left( \prod_{i=k}^m \mathbf{W}_i \right)^T \mathbf{u} = \frac{1}{\left\| \prod_{j=1}^k \mathbf{W}_j \mathbf{v} \right\|} \left( \mathbf{W}_k^T \mathbf{W}_k \right) \prod_{j=1}^{k-1} \mathbf{W}_j \mathbf{v}. \quad (\text{S93})$$

Using Lemma S2 again for  $k-1$  we have

$$\left( \prod_{i=k}^m \mathbf{W}_i \right)^T \mathbf{u} = \frac{\left\| \left( \prod_{i=k}^m \mathbf{W}_i \right)^T \mathbf{u} \right\|}{\left\| \prod_{j=1}^{k-1} \mathbf{W}_j \mathbf{v} \right\|} \prod_{j=1}^{k-1} \mathbf{W}_j \mathbf{v}. \quad (\text{S94})$$

Plugging this equation in (S93) we get

$$\frac{\left\| \left( \prod_{i=k}^m \mathbf{W}_i \right)^T \mathbf{u} \right\|}{\left\| \prod_{j=1}^{k-1} \mathbf{W}_j \mathbf{v} \right\| \left\| \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{u} \right\|} \prod_{j=1}^{k-1} \mathbf{W}_j \mathbf{v} = \frac{1}{\left\| \prod_{j=1}^k \mathbf{W}_j \mathbf{v} \right\|} \left( \mathbf{W}_k^T \mathbf{W}_k \right) \prod_{j=1}^{k-1} \mathbf{W}_j \mathbf{v}. \quad (\text{S95})$$

From Theorem 2, we know that  $\left\| \left( \prod_{i=k}^m \mathbf{W}_i \right)^T \mathbf{u} \right\| = (\sigma_{\max}(\mathbf{T}))^{(m-k+1)/m}$ ,  $\left\| \prod_{j=1}^k \mathbf{W}_j \mathbf{v} \right\| = (\sigma_{\max}(\mathbf{T}))^{k/m}$ , and  $\left\| \prod_{j=1}^{k-1} \mathbf{W}_j \mathbf{v} \right\| \left\| \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{u} \right\| = (\sigma_{\max}(\mathbf{T}))^{1-1/m}$ . Therefore, (S95) can be reduced to

$$(\sigma_{\max}(\mathbf{T}))^{\frac{2}{m}} \times \mathbf{r}_k = (\mathbf{W}_k^T \mathbf{W}_k) \mathbf{r}_k. \quad (\text{S96})$$

Hence,  $\mathbf{r}_k$  is an eigenvector of  $\mathbf{W}_k^T \mathbf{W}_k$  with a corresponding eigenvalue of  $(\sigma_{\max}(\mathbf{T}))^{2/m}$ . Namely,  $\mathbf{r}_k / \|\mathbf{r}_k\|$  is a singular vector of  $\mathbf{W}_k$  with a corresponding singular value of  $(\sigma_{\max}(\mathbf{T}))^{1/m}$ . Using Lemma S2 we have

$$\frac{1}{\|\mathbf{r}_k\|} \mathbf{W}_k \mathbf{r}_k = (\sigma_{\max}(\mathbf{T}))^{\frac{1}{m}} \times \frac{1}{\|\mathbf{q}_k\|} \mathbf{q}_k. \quad (\text{S97})$$

From this equation we deduce that  $\bar{\mathbf{r}}_k$  and  $\bar{\mathbf{q}}_k$  are pair of singular vectors of  $\mathbf{W}_k$ , with a singular value of  $(\sigma_{\max}(\mathbf{T}))^{1/m}$ . Note that the equality  $\bar{\mathbf{r}}_{k+1} = \bar{\mathbf{q}}_k$  is in fact the result of Lemma S2.

## VIII. Proof of Theorem 4

On the one hand, according to Theorem 1

$$\min_{\tilde{\mathbf{w}} \in \Omega} \lambda_{\max}(\mathbf{H}\tilde{\mathbf{w}}) = 2m \times \sigma_{\max}(\mathbf{T})^{2(1-\frac{1}{m})}. \quad (\text{S98})$$

On the other hand, given an arbitrary minimum point  $\mathbf{w} \in \Omega$

$$\begin{aligned} \min_{\tilde{\mathbf{w}} \in \Omega} \lambda_{\max}(\mathbf{H}\tilde{\mathbf{w}}) &\leq \lambda_{\max}(\mathbf{H}\mathbf{w}) \\ &= \max_{\mathbf{B} \in \mathbb{R}^{d_y \times d_x}} 2 \sum_{k=1}^m \left\| \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{B} \left( \prod_{j=1}^{k-1} \mathbf{W}_j \right)^T \right\|_{\text{F}}^2 \quad \text{s.t.} \quad \|\mathbf{B}\|_{\text{F}} = 1 \\ &\leq \max_{\mathbf{B}_1, \dots, \mathbf{B}_m \in \mathbb{R}^{d_y \times d_x}} 2 \sum_{k=1}^m \left\| \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{B}_k \left( \prod_{j=1}^{k-1} \mathbf{W}_j \right)^T \right\|_{\text{F}}^2 \quad \text{s.t.} \quad \|\mathbf{B}_1\|_{\text{F}} = \dots = \|\mathbf{B}_m\|_{\text{F}} = 1. \end{aligned} \quad (\text{S99})$$

Here we obtain a separable optimization problem. Let us examine one term from the series

$$\max_{\mathbf{B}_k \in \mathbb{R}^{d_y \times d_x}} \left\| \left( \prod_{i=k+1}^m \mathbf{W}_i \right)^T \mathbf{B}_k \left( \prod_{j=1}^{k-1} \mathbf{W}_j \right)^T \right\|_{\text{F}}^2 \quad \text{s.t.} \quad \|\mathbf{B}_k\|_{\text{F}} = 1. \quad (\text{S100})$$

In (S87), we saw that the value of (S100) is

$$\left( \sigma_{\max} \left( \prod_{i=k+1}^m \mathbf{W}_i \right) \right)^2 \times \left( \sigma_{\max} \left( \prod_{i=1}^{k-1} \mathbf{W}_i \right) \right)^2 \leq \prod_{i \neq k} (\sigma_{\max}(\mathbf{W}_i))^2, \quad (\text{S101})$$

where we used the sub-multiplicity property of the operator norm (the top singular value). If  $\sigma_{\max}(\mathbf{W}_k) = (\sigma_{\max}(\mathbf{T}))^{1/m}$  for all  $k$ , then

$$\prod_{i \neq k} (\sigma_{\max}(\mathbf{W}_i))^2 = \prod_{i \neq k} (\sigma_{\max}(\mathbf{T}))^{\frac{2}{m}} = (\sigma_{\max}(\mathbf{T}))^{2(1-\frac{1}{m})}. \quad (\text{S102})$$

Thus,

$$\lambda_{\max}(\mathbf{H}\mathbf{w}) = 2m \times (\sigma_{\max}(\mathbf{T}))^{2(1-\frac{1}{m})}. \quad (\text{S103})$$

Therefore  $\mathbf{w}$  is a flattest minimum.

## IX. More Details About Our Experiments

### IX.1. Linear Networks

To sample arbitrary global minima, we started with the canonical solution (27), and multiplied the weight matrices by random matrices from the left and right, such that the the left matrix of one layer cancels out the right matrix of the next (thus keeping the end-to-end function unmodified). Specifically, let  $\{\mathbf{A}_i\}_{i=1}^{m-1}$  be Gaussian random matrices with i.i.d. entries, distributed  $\mathcal{N}(0, 1)$ . Then the weights for arbitrary solutions were generated as

$$\mathbf{W}_m = \mathbf{U} \mathbf{S}_m^{\frac{1}{m}} \mathbf{A}_{m-1}, \quad \mathbf{W}_i = \mathbf{A}_i^{-1} \mathbf{S}_i^{\frac{1}{m}} \mathbf{A}_{i-1}, \quad \mathbf{W}_1 = \mathbf{A}_1^{-1} \mathbf{S}_1^{\frac{1}{m}} \mathbf{V}^T. \quad (\text{S104})$$

To obtain flattest minima, we minimized  $\lambda_{\max}(\mathbf{H}\mathbf{w})$  w.r.t. the weights, by taking random steps over the manifold of global minima  $\Omega$ , and greedily progressing towards a flattest solution. In detail, we randomly generated a set of matrices  $\{\mathbf{A}_i^0\}_{i=1}^m$  with i.i.d. normally distributed entries. We then set the initial weights of the network to be  $\mathbf{W}_i^0 = \mathbf{A}_i^0$ , for all  $i \neq j$ , and

$$\mathbf{W}_j^0 = \left( \prod_{i=j+1}^m \mathbf{A}_i \right)^{-1} \mathbf{T} \left( \prod_{i=1}^{j-1} \mathbf{A}_i \right)^{-1}, \quad (\text{S105})$$

where  $j$  was a random integer chosen uniformly over  $\{1, \dots, m\}$ . Next, we iteratively took small random steps over the manifold of global minima according to the following update rule.

$$\begin{aligned} \mathbf{W}_m^{t+1} &= \mathbf{W}_m^t (\mathbf{I} + \varepsilon_t \mathbf{A}_{m-1}^t), \\ \mathbf{W}_i^{t+1} &= (\mathbf{I} + \varepsilon_t \mathbf{A}_i^t)^{-1} \mathbf{W}_i^t (\mathbf{I} + \varepsilon_t \mathbf{A}_{i-1}^t), \\ \mathbf{W}_1^{t+1} &= (\mathbf{I} + \varepsilon_t \mathbf{A}_1^t)^{-1} \mathbf{W}_1^t, \end{aligned} \tag{S106}$$

where  $\varepsilon_t$  is the step size at the  $t$ th iteration, and  $\{\mathbf{A}_i^t\}_{i=1}^m$  are again random matrices with i.i.d. normally distributed entries. We continued to the next iteration only if the spectral norm of the Hessian decreased. Otherwise, we generated an additional set of direction matrices  $\{\mathbf{A}_i^t\}_{i=1}^m$  until we got a decrement. We stopped this process when the objective achieved its minimal value of  $2m \times (\sigma_{\max}(\mathbf{T}))^{2(1-1/m)}$ , up to a minor error.

## IX.2. Nonlinear Networks

The table below summarizes the parameters and the results for the methods we used in the nonlinear setting for Fig. 6.

	Method 1	Method 2
Optimization Algorithm	SGD	Adam
Learning rate	1/2	$3 \times 10^{-4}$
Other parameters	momentum = 0	$\beta_1 = 0.8, \beta_2 = 0.99$
Batch size	100	100
Train loss	$2.69 \times 10^{-2} \pm 2.53 \times 10^{-5}$	$2.74 \times 10^{-2} \pm 4.03 \times 10^{-5}$
Validation loss	$2.70 \times 10^{-2} \pm 1.68 \times 10^{-4}$	$2.80 \times 10^{-2} \pm 1.94 \times 10^{-4}$
$\lambda_{\max}$	$1.76 \pm 9.43 \times 10^{-3}$	$12.9 \pm 2.2$

Table S1. Summary of the two methods we used to train the network.

## References

- Arora, S., Cohen, N., and Hazan, E. E. On the optimization of deep networks: Implicit acceleration by overparameterization. In *35th International Conference on Machine Learning, ICML 2018*, pp. 372–389. International Machine Learning Society (IMLS), 2018.
- Nar, K. and Sastry, S. Step size matters in deep learning. In *Advances in Neural Information Processing Systems*, pp. 3436–3444, 2018.
- Wu, L., Ma, C., and Weinan, E. How SGD selects the global minima in over-parameterized learning: A dynamical stability perspective. In *Advances in Neural Information Processing Systems*, pp. 8279–8288, 2018.