
Bottleneck Structure in Learned Features: Low-Dimension vs Regularity Tradeoff

Anonymous Author(s)

Affiliation

Address

email

1 Corrections

2 1.1 First Correction

3 **Theorem 1** (Theorem 3 from the main). *For all inputs x where $\text{Rank} Jf(x) = R^{(0)}(f; \Omega)$, we have*
 4 $R^{(1)}(f) \geq 2 \log |Jf(x)|_+$, furthermore:

- 5 1. *If $R^{(0)}(f \circ g) = R^{(0)}(f) = R^{(0)}(g)$, then $R^{(1)}(f \circ g) \leq R^{(1)}(f) + R^{(1)}(g)$.*
- 6 2. *If $R^{(0)}(f + g) = R^{(0)}(f) + R^{(0)}(g)$, then $R^{(1)}(f + g) \leq R^{(1)}(f) + R^{(1)}(g)$.*
- 7 3. *If $P_{\text{Im}A^T} \Omega$ and $A\Omega$ are $k = \text{Rank} A$ dimensional and completely positive (i.e. they can be*
 8 *embedded with an isometric linear map into \mathbb{R}_+^m for some m), then $R^{(1)}(x \mapsto Ax; \Omega) =$*
 9 $2 \log |A|_+$.

10 *Proof.* For the first bound, we remember that $R(f; \Omega, L) \geq L \|Jf\|_{2/L}^{2/L}$, therefore

$$R^{(1)}(f; \Omega) = \lim_{L \rightarrow \infty} R(f; \Omega, L) - LR^{(0)}(f; \Omega) \geq \lim_{L \rightarrow \infty} L \sum_{i=1}^{\text{Rank} Jf(x)} s_i (Jf(x))^{\frac{2}{L}} - 1 \geq \sum_{i=1}^{\text{Rank} Jf(x)} 2 \log s_i (Jf(x))$$

11 where we used $s^{\frac{2}{L}} - 1 = e^{\frac{2}{L} \log s} - 1 \geq \frac{2}{L} \log s$.

12 (1) Since $R(f \circ g; \Omega, L_1 + L_2) \leq R(f; L_1) + R(g; L_2)$, we have

$$\begin{aligned} R^{(1)}(f \circ g; \Omega) &= \lim_{L_1 + L_2 \rightarrow \infty} R(f \circ g; \Omega, L_1 + L_2) - (L_1 + L_2)R^{(0)}(f \circ g; \Omega) \\ &\leq \lim_{L_1 \rightarrow \infty} R(f; \Omega, L_1) - L_1 R^{(0)}(f; \Omega) + \lim_{L_2 \rightarrow \infty} R(g; \Omega, L_2) - L_2 R^{(0)}(g; \Omega) \\ &= R^{(1)}(f; \Omega) + R^{(1)}(g; \Omega). \end{aligned}$$

13 (2) Since $R(f + g; \Omega, L) \leq R(f; \Omega, L) + R(g; \Omega, L)$, we have

$$\begin{aligned} R^{(1)}(f + g; \Omega) &= \lim_{L \rightarrow \infty} R(f + g; \Omega, L) - LR^{(0)}(f + g; \Omega) \\ &\leq \lim_{L \rightarrow \infty} R(f; \Omega, L) - LR^{(0)}(f; \Omega) + \lim_{L \rightarrow \infty} R(g; \Omega, L) - LR^{(0)}(g; \Omega) \\ &= R^{(1)}(f; \Omega) + R^{(1)}(g; \Omega). \end{aligned}$$

14 (3) By the first bound, we know that $R^{(1)}(x \mapsto Ax; \Omega) \geq 2 \log |A|_+$, we now need to show
 15 $R^{(1)}(x \mapsto Ax; \Omega) \leq 2 \log |A|_+$. Let us define the set of completely positive representations as the

16 set of bilinear kernels $K(x, y) = x^T B^T B y$ such that Bx has non-negative entries for all $x \in \Omega$ (we
 17 say that a kernel K is completely positive over Ω if it can be represented in this way for some choice
 18 of B). The set of completely positive representations is convex, since for $K(x, y) = x^T B^T B y$ and
 19 $\tilde{K}(x, y) = x^T \tilde{B}^T \tilde{B} y$, we have

$$\frac{K(x, y) + \tilde{K}(x, y)}{2} = x^T \begin{pmatrix} \frac{1}{\sqrt{2}} B \\ \frac{1}{\sqrt{2}} \tilde{B} \end{pmatrix}^T \begin{pmatrix} \frac{1}{\sqrt{2}} B \\ \frac{1}{\sqrt{2}} \tilde{B} \end{pmatrix} y.$$

20 The conditions that there are O_{in} and O_{out} with $O_{in}^T O_{in} = P_{\text{Im}A^T}$ and $O_{out}^T O_{out} = P_{\text{Im}A}$ such that
 21 $O_{in} \Omega \in \mathbb{R}_+^{k_1}$ and $O_{out} A \Omega \in \mathbb{R}_+^{k_2}$ is equivalent to saying that the kernels $K_{in}(x, y) = x^T P_{\text{Im}A^T} y$
 22 and $K_{out}(x, y) = x^T A^T A x$ are completely positive over Ω .

23 By the convexity of completely positive representations, the interpolation $K_p = pK_{in} + (1-p)K_{out}$
 24 is completely positive for all $p \in [0, 1]$. Now choose for all depths L and all layers $\ell = 1, \dots, L-1$
 25 a matrix $B_{L,\ell}$ such that $K_{p=\frac{\ell}{L}}(x, y) = x^T B_{L,\ell}^T B_{L,\ell} y$ and then choose the weights W_ℓ of the depth
 26 L network as

$$W_\ell = B_{L,\ell} B_{L,\ell-1}^+,$$

27 using the convention $B_{L,0} = I_{d_{in}}$ and $B_{L,L} = I_{out}$. By induction, we show that for any input
 28 $x \in \Omega$ the activation of the ℓ -th hidden layer is $B_{L,\ell} x$. This is true for $\ell = 1$, since $W_1 = B_{L,1}$ and
 29 therefore $p^{(1)}(x) = B_{L,1} x$ which has positive entries so that $q^{(1)}(x) = \sigma(p^{(1)}(x)) = B_{L,1} x$. Then
 30 by induction

$$p^{(\ell)}(x) = W_\ell q^{(\ell-1)}(x) = B_{L,\ell} B_{L,\ell-1}^+ B_{L,\ell-1} x = B_{L,\ell} x,$$

31 which has positive entries, so that again $q^{(\ell)}(x) = \sigma(p^{(\ell)}(x)) = B_{L,\ell} x$. In the end, we get
 32 $p^{(L)}(x) = Ax$ as needed.

33 Let us now compute the Frobenius norms of the weight matrices $\|W_\ell\|_F^2 =$
 34 $\text{Tr} \left[B_{L,\ell}^T B_{L,\ell} \left(B_{L,\ell-1}^T B_{L,\ell-1} \right)^+ \right]$ as $L \rightarrow \infty$, remember that $B_{L,\ell}^T B_{L,\ell} = \frac{\ell}{L} P_{\text{Im}A^T} + (1 - \frac{\ell}{L}) A^T A$,
 35 therefore the matrices $B_{L,\ell}^T B_{L,\ell}$ and $B_{L,\ell-1}^T B_{L,\ell-1}$ converge to each other, so that at first or-
 36 der $B_{L,\ell}^T B_{L,\ell} \left(B_{L,\ell-1}^T B_{L,\ell-1} \right)^+$ converges to $P_{\text{Im}A^T}$, so that $\|W_\ell\|_F^2 \rightarrow \text{Rank}A$, so that
 37 $\sum_{\ell=1}^L \|W_\ell\|_F^2 - L \text{Rank}A$ converges to a finite value as $L \rightarrow \infty$. To obtain this finite limit, we need
 38 to study approximate the next order

$$\begin{aligned} \|W_\ell\|_F^2 - \text{Rank}A &= \sum_{i=1}^{\text{Rank}A} 2 \log s_i(W_i) + O(L^{-2}) \\ &= \log \left| B_{L,\ell}^T B_{L,\ell} \left(B_{L,\ell-1}^T B_{L,\ell-1} \right)^+ \right|_+ + O(L^{-2}) \\ &= \log \left| B_{L,\ell}^T B_{L,\ell} \right|_+ - \log \left| B_{L,\ell-1}^T B_{L,\ell-1} \right|_+ + O(L^{-2}). \end{aligned}$$

39 But as we sum all these second order terms, they cancel out, and we are left with

$$\sum_{\ell=1}^L \|W_\ell\|_F^2 - L \text{Rank}A = 2 \log |A|_+ - 2 \log |I_{\text{Im}A^T}|_+ + O(L^{-1}).$$

40 We have therefore build parameters θ that represent the function $x \mapsto Ax$ with parameter norm
 41 $\|\theta\|^2 = L \text{Rank}A + 2 \log |A|_+ + O(L^{-1})$, which upper bounds the representation cost, thus implying
 42 that $R^{(1)}(x \mapsto Ax; \Omega) \leq 2 \log |A|_+$ as needed. \square

43 1.2 Identity

44 **Proposition 2** (Proposition 4 from the main). *For a domain with $\text{Rank}_J(id; \Omega) =$
 45 $\text{Rank}_{BN}(id; \Omega) = k$, then Ω is k -planar if $R^{(1)}(id; \Omega) = 0$.*

46 *Proof.* We will show that for any two points $x, y \in \Omega$ with k -dim tangent spaces, their tangent
 47 spaces must match if $R^{(1)}(id; \Omega) = 0$.

48 Let $A = J\alpha^{(L-1)}(x)|_{T_x\Omega}$ and $B = J\alpha^{(L-1)}(y)|_{T_y\Omega}$ be the be the Jacobian of the last hidden
 49 activations restricted to the tangent spaces, we know that

$$\begin{aligned} P_{T_x\Omega} &= W_L A \\ P_{T_y\Omega} &= W_L B \end{aligned}$$

50 so that given any weight matrix W_L whose image contains $T_x\Omega$ and $T_y\Omega$, we can write

$$\begin{aligned} A &= W_L^+ P_{T_x\Omega} \\ B &= W_L^+ P_{T_y\Omega}. \end{aligned}$$

51 Without loss of generality, we may assume that the span of $T_x\Omega$ and $T_y\Omega$ is full output space, and
 52 therefore that $W_L W_L^T$ is invertible.

53 Now we now that any parameters that represent the identity on Ω and has $A = J\alpha^{(L-1)}(x)|_{T_x\Omega}$ and
 54 $B = J\alpha^{(L-1)}(y)|_{T_y\Omega}$ must have parameter norm at least

$$\|W_L\|_F^2 + k(L-1) + \max\{2\log|A|_+, 2\log|B|_+\}.$$

55 Subtracting kL and taking $L \rightarrow \infty$, we obtain that

$$R^{(1)}(id; \Omega) \geq \min_{W_L} \|W_L\|_F^2 - k + \max\left\{2\log|W_L^+ P_{T_x\Omega}|_+, 2\log|W_L^+ P_{T_y\Omega}|_+\right\}.$$

56 If we optimize W_L only up to scaling (i.e. optimize aW_L over a) we see that at the optimum, we
 57 always have $\|W_L\|_F^2 = k$. This allows us to rewrite the optimization as

$$R^{(1)}(id; \Omega) \geq \min_{\|W_L\|_F^2=k} \max\left\{2\log|W_L^+ P_{T_x\Omega}|_+, 2\log|W_L^+ P_{T_y\Omega}|_+\right\}.$$

58 The only way to put the first term inside the maximum to 0 is to have $W_L W_L^T = P_{T_x\Omega}$, but this
 59 leads to an exploding second term if $P_{T_x\Omega} \neq P_{T_y\Omega}$. \square

60 Under the assumption of uniform Lipschitzness, one can show a stronger version of the above:

61 **Proposition 3.** For a C -uniformly Lipschitz sequence of ReLU networks representing the function
 62 f , we have

$$R^{(1)}(f) \geq \log|Jf(x)|_+ + \log|Jf(y)|_+ + C^{-2} \|Jf_\theta(x) - Jf_\theta(y)\|_*.$$

63 *Proof.* The decomposition of the difference

$$Jf_\theta(x) - Jf_\theta(y) = \sum_{\ell=1}^{L-1} W_L D_{L-1}(y) \cdots W_{\ell+1} (D_\ell(x) - D_\ell(y)) W_\ell D_{\ell-1}(x) \cdots D_1(x) W_1,$$

64 for the $w_\ell \times w_\ell$ diagonal matrices $D_\ell(x) = \text{diag}(\dot{\sigma}(\tilde{\alpha}_\ell(x)))$, implies the bound

$$\begin{aligned} \|Jf_\theta(x) - Jf_\theta(y)\|_* &\leq \sum_{\ell=1}^{L-1} \|W_L D_{L-1}(y) \cdots D_{\ell+1}(y)\|_{op} \|W_{\ell+1} (D_\ell(x) - D_\ell(y)) W_\ell\|_* \|D_{\ell-1}(x) \cdots D_1(x) W_1\|_{op} \\ &\leq \frac{C^2}{2} \sum_{\ell=1}^{L-1} \left(\|W_{\ell+1} (D_\ell(x) - D_\ell(y))\|_F^2 + \|(D_\ell(x) - D_\ell(y)) W_\ell\|_F^2 \right) \end{aligned}$$

65 since $\|AB\|_* \leq \frac{\|A\|_F^2 + \|B\|_F^2}{2}$.

66 Now since

$$\begin{aligned} L \|Jf_\theta(x)\|_{2/L}^{2/L} &\leq \frac{1}{2} \sum_{\ell=1}^L \|W_\ell D_{\ell-1}(x)\|_F^2 \\ L \|Jf_\theta(y)\|_{2/L}^{2/L} &\leq \frac{1}{2} \sum_{\ell=1}^L \|D_\ell(x) W_\ell\|_F^2 \end{aligned}$$

67 with the convention $D_0(x) = I_{d_{in}}$ and $D_L(x) = I_{d_{out}}$. We obtain that

$$\begin{aligned} L \|Jf_\theta(x)\|_{2/L}^{2/L} + L \|Jf_\theta(y)\|_{2/L}^{2/L} &\leq \frac{1}{2} \sum_{\ell=1}^L \|W_\ell D_{\ell-1}(x)\|_F^2 + \|W_\ell D_{\ell-1}(y)\|_F^2 + \|D_\ell(x)W_\ell\|_F^2 + \|D_\ell(y)W_\ell\|_F^2 \\ &\leq \sum_{\ell=1}^L 2 \|W_\ell\|_F^2 - \frac{1}{2} \|W_\ell (D_{\ell-1}(x) - D_{\ell-1}(y))\|_F^2 - \frac{1}{2} \|(D_\ell(x) - D_\ell(y))W_\ell\|_F^2. \end{aligned}$$

68 This implies the bound

$$\|\theta\|^2 \geq \frac{L \|Jf_\theta(x)\|_{2/L}^{2/L} + L \|Jf_\theta(y)\|_{2/L}^{2/L}}{2} + C^{-2} \|Jf_\theta(x) - Jf_\theta(y)\|_*$$

69 and thus

$$R^{(1)}(f) \geq \log |Jf(x)|_+ + \log |Jf(y)|_+ + C^{-2} \|Jf_\theta(x) - Jf_\theta(y)\|_*.$$

70

□

71 1.3 Second Correction

72 **Proposition 4** (Proposition 11 from the main). *If there is a limiting representation as $L \rightarrow 0$ in the*
73 *optimal representation of f , then $R^{(2)}(f) \geq 0$. Furthermore:*

74 1. *If $R^{(0)}(f \circ g) = R^{(0)}(f) = R^{(0)}(g)$ and $R^{(1)}(f \circ g) = R^{(1)}(f) + R^{(1)}(g)$, then*
75 $\sqrt{R^{(2)}(f \circ g)} \leq \sqrt{R^{(2)}(f)} + \sqrt{R^{(2)}(g)}$.

76 2. *If $R^{(0)}(f + g) = R^{(0)}(f) + R^{(0)}(g)$ and $R^{(1)}(f + g) = R^{(1)}(f) + R^{(1)}(g)$, then $R^{(2)}(f +$
77 $g) \leq R^{(2)}(f) + R^{(2)}(g)$.*

78 3. *If $A^p \Omega$ is $k = \text{Rank} A$ -dimensional and completely positive for all $p \in [0, 1]$, where A^p has*
79 *its non-zero singular taken to the p -th power, then $R^{(2)}(x \mapsto Ax; \Omega) = \frac{1}{2} \|\log_+ A^T A\|^2$.*

80 *Proof.* We start from the inequality

$$R(f \circ g; \Omega, L_f + L_g) \leq R(f; g(\Omega), L_f) + R(g; \Omega, L_g).$$

81 We subtract $(L_f + L_g)R^{(0)}(f \circ g) + R^{(1)}(f \circ g)$ divide by $L_f + L_g$ and take the limit of increasing
82 depths L_f, L_g with $\lim_{L_g, L_f \rightarrow \infty} \frac{L_f}{L_f + L_g} = p \in (0, 1)$ to obtain

$$R^{(2)}(f \circ g; \Omega) \leq \frac{1}{1-p} R^{(2)}(f; g(\Omega)) + \frac{1}{p} R^{(2)}(g; \Omega). \quad (1)$$

83 If K_p is the limiting representation at a ratio $p \in (0, 1)$, we have $R^{(2)}(f; \Omega) = \frac{1}{p} R^{(2)}(K_p; \Omega) +$
84 $\frac{1}{1-p} R^{(2)}(K_p \rightarrow f; \Omega)$ and p must minimize the RHS since if it was instead minimized at a different
85 ratio $p' \neq p$, one could find a lower norm representation by mapping to K_p in the first $p'L$ layers and
86 then back to the outputs. Now there are two possibilities, either $R^{(2)}(K_p; \Omega)$ and $R^{(2)}(K_p \rightarrow f; \Omega)$
87 are non-negative in which case the minimum is attained at some $p \in (0, 1)$ and $R^{(2)}(f; \Omega) \geq 0$, or
88 one or both is negative in which case the above is minimized at $p \in \{0, 1\}$ and $R^{(2)}(f; \Omega) = -\infty$.
89 Since we assumed $p \in (0, 1)$, we are in the first case.

90 (1) To prove the first property, we optimize the RHS of 1 over all possible choices of p (and assuming
91 that $R^{(2)}(f; g(\Omega)), R^{(2)}(g; \Omega) \geq 0$) we obtain

$$\sqrt{R^{(2)}(f \circ g; \Omega)} \leq \sqrt{R^{(2)}(f; g(\Omega))} + \sqrt{R^{(2)}(g; \Omega)}.$$

92 (2) This follows from the inequality $R(f + g; \Omega, L) \leq R(f; g(\Omega), L) + R(g; \Omega, L)$ after subtracting
93 the $R^{(0)}$ and $R^{(1)}$ terms, dividing by L and taking $L \rightarrow \infty$.

94 (3) If $A = USV^T$, one chooses $W_\ell = U_\ell S^{\frac{1}{2}} U_{\ell-1}^T$ with $U_0 = V$, $U_L = U$ and U_ℓ chosen so that
95 $U_\ell S^{\frac{\ell}{2}} V^T \Omega \in \mathbb{R}_+^{n_\ell}$, choosing large enough widths n_ℓ . This choice of representation of A is optimal,

96 i.e. its parameter norm matches the representation cost $L\text{Tr} \left[S^{\frac{2}{L}} \right] = L\text{Rank}A + 2\log|A|_+ +$
 97 $\frac{1}{2L} \|\log_+ A^T A\|^2 + O(L^{-2})$.

98 We know that

$$\lim_{L \rightarrow \infty} R^{(1)}(\alpha_{\ell_1} \rightarrow \alpha_{\ell_2}; \Omega) = R^{(1)}(f_\theta; \Omega) \lim_{L \rightarrow \infty} \frac{\ell_2 - \ell_1}{L}$$

99

$$\frac{1}{p} R^{(2)}(\alpha; \Omega) + \frac{1}{1-p} R^{(2)}(\alpha \rightarrow f; \Omega) \geq R^{(2)}(f; \Omega)$$

100

$$\frac{1}{p} R^{(2)}(\alpha; \Omega) + \frac{1}{1-p} R^{(2)}(\alpha \rightarrow f; \Omega) \geq R^{(2)}(f; \Omega)$$

101

□

102 2 Local Minima Stability

103 In this section we motivate the assumption that the Jacobian $J\tilde{\alpha}_\ell(x)$ is uniformly bounded in oper-
 104 ator norm as $L \rightarrow \infty$. The idea is that solutions with a blowing up Jacobian $J\tilde{\alpha}_\ell(x)$ correspond to
 105 very narrow local minima.

106 The narrowness of a local minimum is related to the Neural Tangent Kernel (or Fisher matrix). We
 107 have that

$$\text{Tr} \left[\Theta^{(L)}(x, x) \right] = \sum_{\ell=1}^L \|\alpha_{\ell-1}(x)\|^2 \|J(\tilde{\alpha}_\ell \rightarrow f_\theta)(x)\|_F^2.$$

108 If the Jacobian $J\tilde{\alpha}_\ell(x)$ blows up, it is reasonable to expect both $\|\alpha_{\ell-1}(x)\|^2$ and $\|J(\tilde{\alpha}_\ell \rightarrow f_\theta)(x)\|_F^2$
 109 to blow up too:

- 110 1. If $\|J\tilde{\alpha}_{\ell-1}(x)\|_{op}$ is very large then small variation of the inputs x can lead to very large
 111 preactivations $\|\tilde{\alpha}_{\ell-1}(x)\|^2$ and thus activations $\|\alpha_{\ell-1}(x)\|^2$ (if we use the ReLU one needs
 112 to also assume $\tilde{\alpha}_{\ell-1}(x)$ does not have mostly negative entries).
- 113 2. Since $\det J\tilde{\alpha}_\ell(x)$ is bounded, a blowing up top singular value $\|J\tilde{\alpha}_\ell(x)\|_{op}$ implies that the
 114 k -th singular value must implode for $k = \text{Rank}_{BN}(f_\theta; \Omega)$. More precisely, if $\|J\tilde{\alpha}_\ell(x)\|_{op}$
 115 is of order L^γ for $\gamma > 1$, then the k -th singular value is of order $L^{-\frac{\gamma}{k-1}}$ or less, and since

$$Jf_\theta(x) = J(\tilde{\alpha}_\ell \rightarrow f_\theta)(x)J\tilde{\alpha}_\ell(x)$$

116 we need $\|J(\tilde{\alpha}_\ell \rightarrow f_\theta)(x)\|_{op}$ to be at least of order $L^{\frac{\gamma}{k-1}}$ so that the k -th singular value of
 117 $Jf_\theta(x)$ remains of order 1.

118 This suggests that in non-uniformly Lipschitz sequences, the NTK would blow up at a rate faster than
 119 L , whereas uniformly Lipschitz sequence have a NTK of order L . In simpler terms non-uniformly
 120 Lipschitz sequences of local minima are infinitely narrower than their uniformly Lipschitz coun-
 121 terpart, which suggests that finite learning rate GD is naturally biased towards uniformly Lipschitz
 122 local minima.

123 **Theorem 5** (Theorem 5 from the main). *For any point x , we have*

$$\|\partial_{xy}^2 \Theta(x, x)\|_{op} \geq 2L \|Jf_\theta(x)\|_{op}^{2-2/L}$$

124 where $\partial_{xy}^2 \Theta(x, x)$ is understood as a $d_{in}d_{out} \times d_{in}d_{out}$ matrix.

125 Furthermore, for any two points x, y such that the pre-activations of all neurons of the network
 126 remain constant on the segment $[x, y]$, then either $\|\Theta(x, x)\|_{op}$ or $\|\Theta(y, y)\|_{op}$ is lower bounded by

$$127 \frac{L}{4} \|x - y\|^2 \left\| Jf_\theta(x) \frac{y-x}{\|x-y\|} \right\|_2^{2-2/L}.$$

128 *Proof.* (1) For any point x , we have

$$\begin{aligned} \partial_{x,y} (v^T \Theta(x, x) v) [u, u] &= \sum_{\ell=1}^L u^T W_1^T D_1(x) \cdots D_{\ell-1}(x)^2 \cdots D_1(x) W_1 u v^T W_L D_{L-1}(x) \cdots D_{\ell}(x)^2 \cdots D_{L-1}(x) W_L^T v \\ &= \sum_{\ell=1}^L \|D_{\ell-1}(x) \cdots D_1(x) W_1 u\|_2^2 \|D_{\ell}(x) \cdots D_{L-1}(x) W_L v\|_2^2. \end{aligned}$$

129 On the other hand, we have

$$\begin{aligned} |v^T Jf_{\theta}(x) u| &= |v^T W_L D_{L-1}(x) \cdots D_1(x) W_1 u| \\ &\leq \|D_{\ell}(x) \cdots D_1(x) W_1 u\|_2 \|D_{\ell}(x) \cdots D_{L-1}(x) W_L v\|_2, \end{aligned}$$

130 where we used the fact that $D_{\ell}(x) D_{\ell}(x) = D_{\ell}(x)$. This applies to the case $\ell = L$ and $\ell = 1$ too,
131 using the definition $D_L(x) = I_{d_{out}}$ and $D_0(x) = I_{d_{in}}$. This implies

$$\begin{aligned} \partial_{xy}^2 (v^T \Theta(x, x) v) [u, u] &\geq |v^T Jf_{\theta}(x) u|^2 \sum_{\ell=1}^L \frac{\|D_{\ell-1}(x) \cdots D_1(x) W_1 u\|_2^2}{\|D_{\ell}(x) \cdots D_1(x) W_1 u\|_2^2} \\ &\geq |v^T Jf_{\theta}(x) u|^2 L \left(\frac{\|u\|_2^2}{\|W_L D_{L-1}(x) \cdots D_1(x) W_1 u\|_2^2} \right)^{\frac{1}{L}} \\ &\geq L \frac{|v^T Jf_{\theta}(x) u|^2}{\|Jf_{\theta}(x) u\|_2^{2/L}}. \end{aligned}$$

132 where we used the geometric/arithmetic mean inequality for the second inequality.

133 If u, v are right and left singular vectors of $Jf_{\theta}(x)$ with singular value s , then the above bound
134 equals $Ls^{2-\frac{2}{L}}$.

135 (2) Now let us consider a segment $\gamma(t) = (1-t)x + ty$ between two points x, y with no changes
136 of activations on these paths (i.e. $D_{\ell}(\gamma(t))$ is constant for all $t \in [0, 1]$). Defining $u = \frac{y-x}{\|y-x\|}$ and
137 $v = \frac{Jf_{\theta}(x)u}{\|Jf_{\theta}(x)u\|}$, we have

$$\partial_t v^T \Theta(\gamma(t), \gamma(t)) v = \|x - y\| \partial_x (v^T \Theta(\gamma(t), \gamma(t)) v) [u] + \|x - y\| \partial_y (v^T \Theta(\gamma(t), \gamma(t)) v) [u]$$

138 and since $\partial_{xx} \Theta(\gamma(t), \gamma(t)) = 0$ and $\partial_{yy} \Theta(\gamma(t), \gamma(t)) = 0$ for all $t \in [0, 1]$, we have

$$\partial_t^2 (v^T \Theta(\gamma(t), \gamma(t)) v) = 2 \|x - y\|^2 \partial_{xy}^2 (v^T \Theta(\gamma(t), \gamma(t)) v) [u, u] \geq 2L \|x - y\|^2 \|Jf_{\theta}(x) u\|_2^{2-2/L}.$$

139 Since $v^T \Theta(\gamma(t), \gamma(t)) v \geq 0$ for all $t \in [0, 1]$ then either

$$v^T \Theta(x, x) v \geq \frac{L}{4} \|x - y\|^2 \|Jf_{\theta}(x) u\|_2^{2-2/L}$$

140 or

$$v^T \Theta(y, y) v \geq \frac{L}{4} \|x - y\|^2 \|Jf_{\theta}(x) u\|_2^{2-2/L}.$$

141

□

142 Now for the proof that rank-underestimating functions require exploding Jacobians:

143 **Proposition 6** (Proposition 7 from the main). *Let $f^* : \Omega \rightarrow \mathbb{R}^{d_{in}}$ for a bounded domain Ω and with*
144 *Rank $_J(f^*; \Omega) = m$, then with high prob. over the sampling of N i.i.d. inputs $x_1, \dots, x_N \in \Omega$, any*
145 *BN-rank function $\hat{f} = h \circ g$ that fits the data $\hat{f}(x_i) = f^*(x_i)$ must satisfy*

$$\sup_{x \in \Omega} \|Jg(x)\|_{op} \sup_{z \in g(\Omega)} \|Jh(z)\|_{op} = \Omega \left(N^{\frac{1}{m} - \frac{1}{k}} \right).$$

146 *Proof.* Writing $z_i = g(x_i) \in g(\Omega) \subset \mathbb{R}^k$, the shortest path through z_1, \dots, z_N has length at most
147 of order $\text{diam}g(\Omega) N^{1-\frac{1}{k}}$ [2]. But this path is mapped by h to a path through the random points
148 y_1, \dots, y_N which must have length of order at least $N^{1-\frac{1}{m}}$ with high probability [1] since the
149 support of the distribution of the y_i s is m -dimensional. This implies the proposition. □

150 3 Representation Geodesics

151 **Theorem 7** (Theorem 9 from the main). *Consider a sequence $(\theta_L)_L$ of representations of a function*
 152 *$f : \Omega \rightarrow \mathbb{R}^{d_{out}}$ with $R^{(0)}(f; \Omega) = \text{Rank}_J(f; \Omega)$ and $R^{(1)}(f; \Omega) < \infty$. Then any accumulating*
 153 *representation K_p at a ratio p is k -planar, i.e. there are k features ϕ_1, \dots, ϕ_k such that $K_p(x, y) =$*
 154 *$\phi(x)^T \phi(y)$.*

155 *Proof.* Let K_p be an accumulating representation, then it must satisfy $R^{(1)}(K_p \rightarrow K_p; \Omega) = 0$:
 156 for any M we can find M sequences $(\ell_L^{(1)})_L, \dots, (\ell_L^{(M)})_L$ each separated by an infinite amount
 157 of layers as $L \rightarrow \infty$ that all converge to the representation K_p ; if $R^{(1)}(K_p \rightarrow K_p) > 0$ then the
 158 overall first correction would be infinite, since $R^{(1)}(\lim_{L \rightarrow \infty} f_{\theta_L}; \Omega) \geq MR^{(1)}(K_p \rightarrow K_p; \Omega)$ for
 159 all M . Then by Proposition 2 we obtain that K_p must be k -planar. \square

160 4 Technical Results

161 4.1 Regularity Counterexample

162 We give here an example of a simple function whose optimal representation geodesic does not
 163 converge, due to it being not uniformly Lipschitz:

164 **Example 8.** The function $f : \Omega \rightarrow \mathbb{R}^3$ with $\Omega = [0, 1]^3$ defined by

$$f(x, y, z) = \begin{cases} (x, y, z) & \text{if } x \leq y \\ (x, y, z + a(x - y)) & \text{if } x > y \end{cases}$$

165 satisfies $R^{(0)}(f; \Omega) = 3$ and $R^{(1)}(f; \Omega) = 0$. The optimal representations of f are not uniformly
 166 Lipschitz as $L \rightarrow \infty$.

167 *Proof.* While we are not able to identify exactly the optimal representation geodesic for the function
 168 f , we will first show that $R^{(1)}(f; \Omega) = 0$, and then show that the uniform Lipschitzness of the
 169 optimal representations would contradict with Proposition 3.

170 (1) Since the Jacobian takes two values inside \mathbb{R}_+^3 , either the identity I_3 or $\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & -1 & 1 \end{pmatrix}$, we

171 know by Theorem 1 that $R^{(1)}(f; \Omega) \geq 2 \log |I_3|_+ = 0$. We therefore only need to construct a
 172 sequence of parameters of different depths that represent f with a squared parameter norm of or-
 173 der $3L + o(1)$. For simplicity, we only do this construction for even depths (the odd case can be
 174 constructed similarly). We define:

$$\begin{aligned} W_\ell &= \begin{pmatrix} e^\epsilon & 0 & 0 \\ 0 & e^\epsilon & 0 \\ 0 & 0 & e^{-2\epsilon} \end{pmatrix} \text{ for } \ell = 1, \dots, \frac{L}{2} - 1 \\ W_{\frac{L}{2}} &= \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \text{ for } \ell = \frac{L}{2} + 2, \dots, L \\ W_\ell &= \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & e^{-(L-2)\epsilon} \end{pmatrix} \text{ for } \ell = \frac{L}{2} + 2, \dots, L \\ W_\ell &= \begin{pmatrix} e^{-\epsilon} & 0 & 0 \\ 0 & e^{-\epsilon} & 0 \\ 0 & 0 & e^{2\epsilon} \end{pmatrix} \text{ for } \ell = \frac{L}{2} + 2, \dots, L \end{aligned}$$

175 We have for all $x \in \mathbb{R}_+^3$

$$\alpha_{\frac{L}{2}-1}(x) = \begin{pmatrix} e^{\frac{L-2}{2}\epsilon} x_1 \\ e^{\frac{L-2}{2}\epsilon} x_2 \\ e^{-(L-2)\epsilon} x_3 \end{pmatrix}$$

176 and

$$\alpha_{\frac{L}{2}}(x) = \begin{pmatrix} e^{\frac{L-2}{2}\epsilon} x_1 \\ e^{\frac{L-2}{2}\epsilon} x_2 \\ e^{-(L-2)\epsilon} x_3 \\ \sigma(x_1 - x_2) \end{pmatrix}$$

177 and

$$\alpha_{\frac{L}{2}+1}(x) = \begin{pmatrix} e^{\frac{L-2}{2}\epsilon} x_1 \\ e^{\frac{L-2}{2}\epsilon} x_2 \\ e^{-(L-2)\epsilon} (x_3 + \sigma(x_1 - x_2)) \end{pmatrix}$$

178 and

$$f_\theta(x) = \begin{pmatrix} x_1 \\ x_2 \\ x_3 + \sigma(x_1 - x_2) \end{pmatrix}.$$

179 The norm of the parameters is

$$\begin{aligned} & \frac{L-2}{2}(2e^{2\epsilon} + e^{-4\epsilon}) + (3 + 2e^{-(L-2)\epsilon}) + (3 + e^{-2(L-2)\epsilon}) + \frac{L-2}{2}(2e^{-2\epsilon} + e^{4\epsilon}) \\ & = 3L + 2(e^{2\epsilon} - 1) + (e^{-4\epsilon} - 1) + 2e^{-(L-2)\epsilon} + e^{-2(L-2)\epsilon} + 2(e^{-2\epsilon} - 1) + (e^{4\epsilon} - 1) \end{aligned}$$

180 If we take $\epsilon = L^{-\gamma}$ for $\gamma \in (\frac{1}{2}, 1)$, then the terms $2e^{-(L-2)\epsilon}$ and $e^{-2(L-2)\epsilon}$ decay exponentially (at
181 a rate of $e^{L^{1-\gamma}}$), in addition the terms $2(e^{2\epsilon} - 1) + (e^{-4\epsilon} - 1)$ and $2(e^{-2\epsilon} - 1) + (e^{4\epsilon} - 1)$ are of
182 order $L^{-2\gamma}$. This proves that $R^{(1)}(f; \Omega) = 0$.

183 (2) Let us now assume that the optimal representation of f is C -uniform Lipschitz for some constant
184 C , then by Proposition 3, we have that

$$R^{(1)}(f; \Omega) \geq \log |I_3|_+ + \log \left\| \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & -1 & 1 \end{pmatrix} \right\|_+ + C^{-2} \left\| I_3 - \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & -1 & 1 \end{pmatrix} \right\|_* > 0,$$

185 which contradicts with the fact that $R^{(1)}(f; \Omega) = 0$. □

186 4.2 Extension outside FPLFs

187 Since all functions represented by finite depth and width networks are FPLFs, the representation
188 cost of any such function is infinite. But we can define the representation cost of a function f
189 that is the limit of a sequence of FPLF as the infimum over all sequences $f_i \rightarrow f$ converging of
190 $\lim_{i \rightarrow \infty} R(f_i; \Omega)$ (for some choice of convergence type that implies convergence of the Jacobians
191 $Jf_i(x) \rightarrow Jf(x)$). Note that since the representation cost $R(f; \Omega)$ is lower semi-continuous, i.e.
192 $\liminf_{f \rightarrow f_0} R(f; \Omega) \geq R(f_0; \Omega)$, this does not change the definition of the representation cost on
193 the space of FPLFs.

194 The definitions of the decomposition $R^{(0)}, R^{(1)}, R^{(2)}$ can also be similarly extended, and one can
195 check that the properties described in Theorem 1 of [3] and Theorems 1 and [] of this paper all
196 extend as well.

197 References

- 198 [1] Jillian Beardwood, J. H. Halton, and J. M. Hammersley. The shortest path through many points.
199 *Mathematical Proceedings of the Cambridge Philosophical Society*, 55(4):299–327, 1959.
- 200 [2] L. Few. The shortest path and the shortest road through n points. *Mathematika*, 2(2):141–144,
201 1955.
- 202 [3] Arthur Jacot. Implicit bias of large depth networks: a notion of rank for nonlinear functions. In
203 *The Eleventh International Conference on Learning Representations*, 2023.