
Depth Creates No Bad Local Minima

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 In deep learning, *depth*, as well as *nonlinearity*, create non-convex loss surfaces.
2 Then, does depth alone create bad local minima? In this paper, we prove that
3 without nonlinearity, depth alone does not create bad local minima, although it
4 induces non-convex loss surface. Using this insight, we greatly simplify a recently
5 proposed proof to show that all of the local minima of feedforward deep linear
6 neural networks are global minima. Our theoretical results generalize previous
7 results with fewer assumptions, and this analysis provides a method to show similar
8 results beyond square loss in deep linear models.

9 1 Introduction

10 Deep learning has recently had a profound impact on the machine learning, computer vision, and
11 artificial intelligence communities. In addition to its practical successes, previous studies have
12 revealed several reasons why deep learning has been successful from the viewpoint of its *model*
13 *classes*. An (over-)simplified explanation is the harmony of its great expressivity and *big data*:
14 because of its great expressivity, deep learning can have less *bias*, while a large training dataset leads
15 to less *variance*. The great expressivity can be seen from an aspect of representation learning as well:
16 whereas traditional machine learning makes use of features designed by human users or experts as a
17 type of prior, deep learning tries to learn features from the data as well. More accurately, a key aspect
18 of the model classes in deep learning is the *generalization* property; despite its great expressivity,
19 deep learning model classes can maintain great generalization properties (Livni et al., 2014; Mhaskar
20 et al., 2016; Poggio et al., 2016). This would distinguish deep learning from other possibly too
21 flexible methods, such as shallow neural networks with too many hidden units, and traditional kernel
22 methods with a too powerful kernel. Therefore, the practical success of deep learning seems to be
23 supported by the great quality of its model classes.

24 However, having a great model class is not so useful if we cannot find a good model in the model
25 class via training. Training a deep model is typically framed as non-convex optimization. Because of
26 its non-convexity and high dimensionality, it has been unclear whether we can *efficiently* train a deep
27 model. Note that the difficulty comes from the combination of non-convexity and high dimensionality
28 in weight parameters. If we can reformulate the training problem into several decoupled training
29 problems, with each having a small number of weight parameters, we can effectively train a model
30 via non-convex optimization as theoretically shown in Bayesian optimization and global optimization
31 literatures (Kawaguchi et al., 2015; Wang et al., 2016; Kawaguchi et al., 2016). As a result of
32 non-convexity and high-dimensionality, it was shown that training a general neural network model is
33 NP-hard (Blum & Rivest, 1992). However, such a hardness-result in a worst case analysis would not
34 tightly capture what is going on in practice, as we seem to be able to efficiently train deep models in
35 practice.

36 To understand its practical success beyond worst case analysis, theoretical and practical investi-
37 gations on the training of deep models have recently become an active research area (Saxe et al.,
38 2014; Dauphin et al., 2014; Choromanska et al., 2015; Haeffele & Vidal, 2015; Shamir, 2016;

39 Kawaguchi, 2016; Swirszcz et al., 2016; Arora et al., 2016; Freeman & Bruna, 2016; Soudry &
40 Hoffer, 2017).

41 An important property of a deep model is that the non-convexity comes from *depth*, as well as
42 *nonlinearity*: indeed, depth by itself creates highly non-convex optimization problems. One way to
43 see a property of the non-convexity induced by depth is the non-uniqueness owing to *weight-space*
44 *symmetries* (Krkova & Kainen, 1994): the model represents the same function mapping from the input
45 to the output with different distinct settings in the weight space. Accordingly, there are many distinct
46 globally optimal points and many distinct points with the same loss values due to weight-space
47 symmetries, which would result in a non-convex epigraph (i.e., non-convex function) as well as
48 non-convex sublevel sets (i.e., non-quasiconvex function). Thus, it has been unclear whether *depth* by
49 itself can create a difficult non-convex loss surface. The recent work (Kawaguchi, 2016) indirectly
50 showed, as a consequence of its main theoretical results, that depth does not create bad local minima
51 of deep linear model with Frobenius norm although it creates potentially bad saddle points.

52 In this paper, we directly prove that all local minima of deep linear model corresponds to local minima
53 of shallow model. Building upon this new theoretical insight, we propose a simpler proof for one
54 of the main results in the recent work (Kawaguchi, 2016); all of the local minima of feedforward
55 deep linear neural networks with Frobenius norm are global minima. The power of this proof can go
56 beyond Frobenius norm: as long as the loss function satisfies Theorem 3.2, all local minima of deep
57 linear model corresponds to local minimum of shallow model.

58 2 Main Result

59 To examine the effect of depth alone, we consider the following optimization problem of feedforward
60 deep linear neural networks with the square error loss:

$$\underset{W}{\text{minimize}} \quad L(W) = \frac{1}{2} \|W_H W_{H-1} \cdots W_1 X - Y\|_F^2, \quad (1)$$

61 where $W_i \in \mathbb{R}^{d_i \times d_{i-1}}$ is the weight matrix, $X \in \mathbb{R}^{d_0 \times m}$ is the input training data, and $Y \in \mathbb{R}^{d_H \times m}$
62 is the target training data. Let $p = \arg \min_{0 \leq i \leq H} d_i$ be the index corresponding to the smallest width.
63 Note that for any W , we have $\text{rank}(W_H W_{H-1} \cdots W_1) \leq d_p$. To analyze optimization problem
64 (1), we also consider the following optimization problem with a “shallow” linear model, which is
65 equivalent to problem (1) in terms of the global minimum value:

$$\underset{R}{\text{minimize}} \quad F(R) = \|RX - Y\|_F^2 \quad \text{s.t.} \quad \text{rank}(R) \leq d_p, \quad (2)$$

66 where $R \in \mathbb{R}^{d_H \times d_0}$. Note that problem (2) is non-convex, unless $d_p = \min(d_H, d_0)$, whereas
67 problem (1) is non-convex, even when $d_p \geq \min(d_H, d_0)$ with $H > 1$. In other words, deep
68 parameterization creates a non-convex loss surface even without nonlinearity.

69 Though we only consider the Frobenius loss here, the proof holds for general cases. As long as
70 the loss function satisfies Theorem 3.2, all local minima of deep linear model corresponds to local
71 minimum of shallow model.

72 Our first main result states that even though deep parameterization creates a non-convex loss surface,
73 it does not create new bad local minima. In other words, every local minimum in problem (1)
74 corresponds to a local minimum in problem (2).

75 **Theorem 2.1.** (Depth creates no new bad local minima) *Assume that X and Y have full row rank. If
76 $\bar{W} = \{\bar{W}_1, \dots, \bar{W}_H\}$ is a local minimum of problem (1), then $\bar{R} = \bar{W}_H \bar{W}_{H-1} \cdots \bar{W}_1$ achieves the
77 value of a local minimum of problem (2).*

78 Therefore, we can deduce the property of the local minima in problem (1) from those in problem (2).
79 Accordingly, we first analyze the local minima in problem (2), and obtain the following statement.
80

81 **Theorem 2.2.** (No bad local minima for rank restricted shallow model) *If X has full row rank, all
82 local minima of optimization problem (2) are global minima.*

83 By combining Theorems 2.1 and 2.2, we conclude that every local minimum is a global minimum for
 84 feedforward deep linear networks with a square error loss.

85 **Theorem 2.3.** (No bad local minima for deep linear neural networks) *If X and Y have full row rank,*
 86 *then all local minima of problem (1) are global minima.*

87 Theorem 2.3 generalizes one of the main results in (Kawaguchi, 2016) with fewer assumptions.
 88 Following the theoretical work with a random matrix theory (Dauphin et al., 2014; Choromanska
 89 et al., 2015), the recent work (Kawaguchi, 2016) showed that under some strong assumptions, all of
 90 the local minima are global minima for a class of nonlinear deep networks. Furthermore, the recent
 91 work (Kawaguchi, 2016) proved the following properties for a class of general deep linear networks
 92 with arbitrary depth and width: 1) the objective function is non-convex and non-concave; 2) all of the
 93 local minima are global minima; 3) every other critical point is a saddle point; and 4) there is no saddle
 94 point with the Hessian having no negative eigenvalue for shallow networks with one hidden layer,
 95 whereas such saddle points exist for deeper networks. Theorem 2.3 generalizes the second statement
 96 with fewer assumptions; the previous papers (Baldi, 1989; Kawaguchi, 2016) assume that the data
 97 matrix $YX^T(XX^T)^{-1}XY^T$ has distinct eigenvalues, whereas we do not assume that.

98 3 Proof

99 In this section, we provide the proofs of Theorems 2.1, 2.2, and 2.3.

100 3.1 Proof of Theorem 2.1

101 In order to deduce the proof of Theorem 2.1, we need some fundamental facts in linear algebra. The
 102 next two lemmas recall some basic facts of perturbation theory for singular value decomposition
 103 (SVD).

104 Let M and \bar{M} be two $m \times n$ ($m \geq n$) matrices with SVDs

$$B = U\Sigma V^T = (U_1, U_2) \begin{pmatrix} \Sigma_1 & \\ & \Sigma_2 \end{pmatrix} \begin{pmatrix} V_1^T \\ V_2^T \end{pmatrix}$$

$$\bar{B} = \bar{U}\bar{\Sigma}\bar{V}^T = (\bar{U}_1, \bar{U}_2) \begin{pmatrix} \bar{\Sigma}_1 & \\ & \bar{\Sigma}_2 \end{pmatrix} \begin{pmatrix} \bar{V}_1^T \\ \bar{V}_2^T \end{pmatrix},$$

105 where $\Sigma_1 = \text{diag}(\sigma_1, \dots, \sigma_k)$, $\Sigma_2 = \text{diag}(\sigma_{k+1}, \dots, \sigma_n)$, $\bar{\Sigma}_1 = \text{diag}(\bar{\sigma}_1, \dots, \bar{\sigma}_k)$, $\bar{\Sigma}_2 =$
 106 $\text{diag}(\bar{\sigma}_{k+1}, \dots, \bar{\sigma}_n)$, U, V, \bar{U} and \bar{V} are orthogonal matrices.

107 **Lemma 3.1. Continuity of Singular Value** *The singular value σ_i of a matrix is a continuous map*
 108 *of entries of the matrix.*

109 **Lemma 3.2. (Wedin, 1972) Continuity of Singular Space**

If

$$\rho := \min \left\{ \min_{1 \leq i \leq k, 1 \leq j \leq n-k} |\sigma_i - \bar{\sigma}_{k+j}|, \min_{1 \leq i \leq k} \sigma_i \right\} > 0,$$

110 *then:*

$$\sqrt{\|\sin(U_1, \bar{U}_1)\|_F^2 + \|\sin(V_1, \bar{V}_1)\|_F^2} \leq \frac{\sqrt{\|(\bar{M} - M) V_1\|_F^2 + \|(\bar{M}^* - M^*) U_1\|_F^2}}{\rho}.$$

111 For a fixed matrix B , we say “matrix A is a perturbation of matrix B ” if $\|A - B\|_\infty$ is $o(1)$, which
 112 means that the difference between A and B is much smaller than any non-zero number in matrix
 113 B .

114 Lemma 3.2 implies that any SVD for a perturbed matrix is a perturbation of some SVD for the
 115 original matrix under full rank condition. More formally:

116 **Lemma 3.3.** Let \bar{M} be a full-rank matrix with singular value decomposition $\bar{M} = \bar{U}\bar{\Sigma}\bar{V}^T$. M is a
 117 perturbation of \bar{M} . Then, there exists one SVD of M , $M = U\Sigma V^T$, such that U is a perturbation of
 118 \bar{U} , Σ is a perturbation of $\bar{\Sigma}$ and V is a perturbation of \bar{V} . (Notice that SVD of a matrix may not be
 119 unique due to rotation of the eigen-space corresponding to the same eigenvalue)

120 **Proof:** With the small perturbation of matrix \bar{M} , Lemma 3.1 shows that the singular values does not
 121 change much. Thus, if $\|\bar{M} - M\|_\infty$ is small enough, $|\sigma_i - \bar{\sigma}_i|$ is also small for all i . Remember that
 122 all singular values of \bar{M} are positive. By letting Σ_1 contain only the singular value σ_i (which may be
 123 multiple, and hence U_1 and V_1 are the singular spaces corresponding to the singular value σ_i), we
 124 have $\rho > 0$ in Lemma 3.2, thus Lemma 3.2 implies that the singular space of the perturbed matrix
 125 corresponding to singular value σ_i in the initial matrix does not change much. The statement of the
 126 lemma follows by combining this result for the different singular values together (i.e., consider each
 127 index i for different σ_i in the above argument). \square

128 We say that W satisfies the rank condition, if $\text{rank}(W_H \cdots W_1) = d_p$. Any perturbation of the
 129 products of matrices is the product of the perturbed matrices, when the original matrix satisfies the
 130 rank constraint. More formally:

131 **Theorem 3.1.** Let $\bar{R} = \bar{W}_H \bar{W}_{H-1} \cdots \bar{W}_1$ with $\text{rank}(\bar{R}) = d_p$. Then, for any R , such that R is a
 132 perturbation of \bar{R} and $\text{rank}(R) \leq d_p$, there exists $\{W_1, W_2, \dots, W_H\}$, such that W_i is perturbation
 133 of \bar{W}_i for all $i \in \{1, \dots, H\}$ and $R = W_H W_{H-1} \cdots W_1$.

134 We will prove the theorem by induction. When $H = 2$, we can easily show that the perturbation of the
 135 product of two matrices is the product of one matrix and the perturbation of the other matrix. When
 136 $H = k \geq 3$, we let M be the product of two specific matrices, and by induction the perturbation of
 137 the product (R) is the product of a perturbation of M and perturbations of the other $H - 2$ matrix.
 138 And a perturbation of M is also the product of perturbations of those two specific matrices, which
 139 proves the statement when $H = k$.

140 **Proof:** The case with $H = 1$ holds by setting $W_1 = R$. We prove the lemma with $H \geq 2$ by
 141 induction.

142 We first consider the base case where $H = 2$ with $\bar{R} = \bar{W}_2 \bar{W}_1$.

143 Let $\bar{R} = \bar{U}\bar{\Sigma}\bar{V}^T$ be the SVD of \bar{R} . It follows Lemma 3.3 that there exists an SVD of R , $R = U\Sigma V^T$,
 144 such that U is a perturbation of \bar{U} , Σ is a perturbation of $\bar{\Sigma}$ and V is a perturbation of \bar{V} .
 145 Because $\text{rank}(\bar{R}) = d_p$, with a small perturbation, the positive singular values remain strictly positive,
 146 whereby, $\text{rank}(R) \geq d_p$. Together with the assumption $\text{rank}(R) \leq d_p$, we have $\text{rank}(R) = d_p$. Let
 147 $\bar{S}_2 = \bar{U}^T \bar{W}_2$ and $\bar{S}_1 = \bar{W}_1 \bar{V}$. Note that $\bar{U}\bar{\Sigma}\bar{V}^T = \bar{R} = \bar{W}_2 \bar{W}_1$. Hence, $\bar{S}_2 \bar{S}_1 = \bar{\Sigma}$ is a diagonal
 148 matrix. Remember Σ is a perturbation of $\bar{\Sigma}$, thus there is an S_2 , which is a perturbation of \bar{S}_2 (each
 149 row of S_2 is a scale of the corresponding row of \bar{S}_2), such that $S_2 \bar{S}_1 = \Sigma$. Let $W_2 = U S_2$ and
 150 $W_1 = \bar{S}_1 V$. Then, W_1 is a perturbation of \bar{W}_1 , W_2 is a perturbation of \bar{W}_2 , and $W_1 W_2 = R$, which
 151 proves the case when $H = 2$.

152 For the inductive step, given that the lemma holds for the case with $H = k \geq 2$, let us consider
 153 the case when $H = k + 1 \geq 3$ with $\bar{R} = \bar{W}_{k+1} \bar{W}_k \cdots \bar{W}_1$. Let \mathcal{I} be an index set defined as
 154 $\mathcal{I} = \{p, p - 1\}$ if $p \geq 2$, $\mathcal{I} = \{p + 2, p + 1\}$ if $p = 0$ or $p = 1$. We denote the i -th element of a
 155 set \mathcal{I} by \mathcal{I}_i . Then, $\bar{M} = \bar{W}_{\mathcal{I}_2} \bar{W}_{\mathcal{I}_1}$ exists as $k + 1 \geq 3$. Note that \bar{R} can be written as a product of
 156 k matrices with \bar{M} (for example, $\bar{R} = \bar{W}_H \cdots \bar{W}_{\mathcal{I}_1+1} \bar{M} \bar{W}_{\mathcal{I}_2-1} \cdots \bar{W}_1$). Thus, from the inductive
 157 hypothesis, for any R , such that R is a perturbation of \bar{R} and $\text{rank}(R) \leq d_p$, there exists a set of
 158 desired k matrices M and W_i for $i \in \{1, \dots, k + 1\} \setminus \mathcal{I}$, such that W_i is perturbation of \bar{W}_i for all
 159 $i \in \{1, \dots, k + 1\} \setminus \mathcal{I}$, M is perturbation of \bar{M} , and the product is equal to R . Meanwhile, because
 160 \bar{M} is either a d_p by d_{p-2} matrix or a d_{p+2} by d_p matrix, we have $\text{rank}(\bar{M}) \leq d_p$ and $\text{rank}(M) \leq d_p$,
 161 and it follows $\text{rank}(\bar{R}) = d_p$ that $\text{rank}(\bar{M}) = d_p$. Thus, by setting $\bar{R} \leftarrow \bar{M}$ and $R \leftarrow M$ (note that
 162 d_p in $\bar{R} = \bar{W}_{k+1} \bar{W}_k \cdots \bar{W}_1$ is equal to d_p in $\bar{M} = \bar{W}_{\mathcal{I}_2} \bar{W}_{\mathcal{I}_1}$), we can apply the proof for the case of
 163 $H = 2$ to conclude: there exists $\{W_{\mathcal{I}_2}, W_{\mathcal{I}_1}\}$, such that W_i is perturbation of \bar{W}_i for all $i \in \mathcal{I}$, and
 164 $M = W_{\mathcal{I}_2} W_{\mathcal{I}_1}$. Combined with the above statement from the inductive hypothesis, this implies the
 165 lemma with $H = k + 1$, whereby we finish the proof by induction. \square

166 The next two theorems show that, for any local minimum of $L(\cdot)$, there is another local minimum of
 167 $L(\cdot)$, whose function value is the same as the original and it satisfies the rank constraint.

168 **Theorem 3.2.** Let $W = \{W_1, \dots, W_H\}$ be a local minimum of problem (1) and $R \triangleq$
169 $W_H W_{H-1} \dots W_1$. If W_i is not of full rank, then there exists a \bar{W}_i , such that \bar{W}_i is of full rank,
170 \bar{W}_i is a perturbation of W_i , $\bar{W} = \{W_1, \dots, W_{i-1}, \bar{W}_i, W_{i+1}, \dots, W_H\}$ is a local minimum of
171 problem (1), and $L(W) = L(\bar{W})$.

172 The idea of the proof is that if we just change one weight W_i and keep all other weights, it becomes a
173 convex least square problem. Then we are able to perturb \bar{W}_i to maintain the objective value as well
174 as the perturbation is full rank.

175 **Proof of Theorem 3.2** For notational convenience, let $A = W_{i-1} \dots W_1 X$ and $B = W_{i+1} \dots W_H$,
176 and let $L_i(W_i) = \frac{1}{2} \|B^T W_i A - Y\|_F^2$. Because W is a local minimum of L , W_i is a local minimum
177 of L_i . Let $A = U_1^T D_1 V_1$ and $B = U_2^T D_2 V_2$ are the SVDs of A and B , respectively, where D_i is a
178 diagonal matrix with the first s_i terms being strictly positive, $i = 1, 2$. Minimizing L_i over W_i is a
179 least square problem, and the normal equation is

$$BB^T W_i A A^T = B Y A^T, \quad (3)$$

180 hence

$$\begin{aligned} W_i &\in (BB^T)^+ B Y A^T (A A^T)^+ + \{M | BB^T M A A^T = 0\} \\ &= U_2 D_2^+ V_2^T Y V_1 D_1^+ U_1^T + \{U_2 K U_1^T | K_{1:s_2, 1:s_1} = 0\}, \end{aligned}$$

181 where $(\cdot)^+$ is a Moore–Penrose pseudo-inverse and K is a matrix with suitable dimension with the
182 entries in the top left $s_2 \times s_1$ rectangular being 0.

183 Since $V_2^T Y V_1$ is of full rank,

$$\text{rank}(D_2^+ V_2^T Y V_1 D_1^+) \geq \max\{0, s_2 + s_1 - \max\{d_i, d_{i-1}\}\}$$

184 Thus, we can choose a proper K (which contains $d_i + d_{i-1} - s_2 - s_1$ 1s at proper positions with all other
185 terms being 0s) such that $D_2^+ V_2^T Y V_1 D_1^+ + K$ is of full rank, whereby $U_2 (D_2^+ V_2^T Y V_1 D_1^+ + K) U_1^T$
186 is of full rank. Therefore, there is a full rank \hat{W}_i that satisfies the normal equation (3).

187 Let $\bar{W}_i(\mu) = W_i + \mu (\hat{W}_i - W_i)$. Then, $\bar{W}_i(\mu)$ also satisfies the normal equation, and $L(\bar{W}(\mu)) =$
188 $L_i(\bar{W}_i(\mu)) = L_i(W_i) = L(W)$, for any $\mu > 0$.

Note that W is a local minimum of $L(W)$. Thus, there exists a $\delta > 0$, such that for any W^0 satisfying
 $\|W^0 - W\|_\infty \leq \delta$, we have $L(W^0) \geq L(W)$. It follows from \hat{W}_i being full rank that there exists a
small enough μ , such that $\bar{W}_i(\mu)$ is full rank and $\|\bar{W}_i(\mu) - W_i\|_\infty$ is arbitrarily small (in particular,
 $\|\bar{W}_i(\mu) - W_i\|_\infty \leq \frac{\delta}{2}$), because the non-full-rank matrices are discrete on the line of $\bar{W}_i(\mu)$ with
parameter $\mu > 0$ by considering the determine of $W_i^T(\mu) W_i(\mu)$ or $W_i(\mu) W_i^T(\mu)$ as a polynomial
of λ . Therefore, for any W^0 , such that $\|W^0 - \bar{W}(\mu)\|_\infty \leq \frac{\delta}{2}$, we have

$$\|W^0 - W\|_\infty \leq \|W^0 - \bar{W}(\mu)\|_\infty + \|\bar{W}_i(\mu) - W_i\|_\infty \leq \delta,$$

whereby

$$L(W^0) \geq L(W) = L(\bar{W}(\mu)).$$

189 This shows that $\bar{W}(\mu) = \{W_1, \dots, W_{i-1}, \bar{W}_i(\mu), W_{i+1}, \dots, W_H\}$ is also a local minimum of
190 problem (1) for some small enough μ . \square

191 **Lemma 3.4.** Let $R = AB$ for two given matrices $A \in R^{d_1 \times d_2}$ and $B \in R^{d_2 \times d_3}$. If $d_1 \leq d_2$,
192 $d_1 \leq d_3$ and $\text{rank}(A) = d_1$, then any perturbation of R is the product of A and perturbation of B .

193 **Proof:** Let $A = U D V^T$ be the SVD of A , then, $R = U D V^T B$. Let \bar{R} be a perturbation of R and
194 let $\bar{B} = B + V D^+ U^T (\bar{R} - R)$. Then, \bar{B} is a perturbation of B and $A \bar{B} = \bar{R}$ by noticing $DD^+ = I$,
195 as A has full row rank. \square

196 **Theorem 3.3.** If $\bar{W} = \{\bar{W}_1, \dots, \bar{W}_H\}$ is a local minimum with \bar{W}_i being full rank, then, there exists
197 $\hat{W} = \{\hat{W}_1, \dots, \hat{W}_H\}$, such that \hat{W}_i is a perturbation of \bar{W}_i for all $i \in \{1, \dots, H\}$, \hat{W} is a local
198 minimum, $L(\hat{W}) = L(\bar{W})$, and $\text{rank}(\hat{W}_H \hat{W}_{H-1} \dots \hat{W}_1) = d_p$.

199 In the proof of Theorem 3.3, we will use Theorem 3.2 and Lemma 3.4 to show that we can perturb
200 $\bar{W}_{p-1}, \bar{W}_{p-2}, \dots, \bar{W}_1$ in sequence to make sure the perturbed weight is still the optimal solution
201 and $\text{rank}(\hat{W}_p \hat{W}_{p-1}) = d_p$. Similar strategy can make sure $\text{rank}(\hat{W}_H \hat{W}_{H-1} \cdots \hat{W}_{p+1}) = d_p$, which
202 then proves the whole theorem.

Proof of Theorem 3.3 : If $p \neq 1$, consider

$$L_1(T) := \|\bar{W}_H \cdots \bar{W}_{p+1} T \bar{W}_{p-2} \cdots \bar{W}_1 X - Y\|_F^2.$$

203 Then, it follows from Lemma 3.4 and \bar{W} is a local minimum of $L(W)$ that \bar{T} is a local minimum of L_1 ,
204 where $\bar{T} = \bar{W}_p \bar{W}_{p-1}$. It follows from Theorem 3.2 that there exists \hat{T} , such that \hat{T} is close enough to
205 \bar{T} , \hat{T} is a local minimum of $L_1(T)$, $L_1(\hat{T}) = L_1(\bar{T})$, and $\text{rank}(\hat{T}) = d_p$. Note \hat{T} is a perturbation
206 of \bar{T} , whereby, from Lemma 3.4, there exists \hat{W}_p, \hat{W}_{p-1} , which are perturbations of \bar{W}_p and \bar{W}_{p-1} ,
207 respectively, such that $\hat{W}_p \hat{W}_{p-1} = \hat{T}$. Thus, $\hat{W}^0 = (\bar{W}_H, \dots, \bar{W}_{p+1}, \hat{W}_p, \hat{W}_{p-1}, \bar{W}_{p-2}, \dots, \bar{W}_1)$
208 is a local minimum of $L(W)$, $L(\hat{W}) = L(\bar{W})$ and $\text{rank}(\hat{W}_p \hat{W}_{p-1}) = d_p$.

209 By that analogy, we can find $\hat{W}_p \cdots \hat{W}_1$, such that $\hat{W}^1 = (\bar{W}_H, \dots, \bar{W}_{p+1}, \hat{W}_p, \hat{W}_{p-1}, \dots, \hat{W}_1)$
210 is a local minimum of $L(W)$, \hat{W}_i is a perturbation of \bar{W}_i for $i = 1, \dots, p$, $L(\hat{W}^1) = L(\bar{W})$ and
211 $\text{rank}(\hat{W}_p \hat{W}_{p-1} \cdots \hat{W}_1) = d_p$.

212 Similarly, we can find $\hat{W}_H \cdots \hat{W}_{p+1}$, such that $\hat{W}^2 = (\hat{W}_H, \dots, \hat{W}_{p+1}, \hat{W}_p, \hat{W}_{p-1}, \dots, \hat{W}_1)$ is a
213 local minimum of $L(W)$, \hat{W}_i is a perturbation of \bar{W}_i for $i = p+1, \dots, H$, $L(\hat{W}^2) = L(\hat{W}^1) =$
214 $L(\bar{W})$ and $\text{rank}(\hat{W}_H \hat{W}_{H-1} \cdots \hat{W}_{p+1}) = d_p$.

215 Noticing that

$$\text{rank}(\hat{W}_H \cdots \hat{W}_1) \geq \text{rank}(\hat{W}_H \hat{W}_{H-1} \cdots \hat{W}_{p+1}) + \text{rank}(\hat{W}_p \hat{W}_{p-1} \cdots \hat{W}_1) - d_p = d_p$$

216 and $\text{rank}(\hat{W}_H \cdots \hat{W}_1) \leq \min_{i=0, \dots, H} d_i = d_p$, we have $\text{rank}(\hat{W}_H \cdots \hat{W}_1) = d_p$, which completes
217 the proof. \square

Proof of Theorem 2.1: It follows from Theorem 3.2 and Theorem 3.3 that there exists another local
minimum $\hat{W} = \hat{W} = \{\hat{W}_1, \dots, \hat{W}_H\}$, such that $L(\hat{W}) = L(\bar{W})$ and $\text{rank}(\hat{W}_H \hat{W}_{H-1} \cdots \hat{W}_1) =$
 d_p . Remember that $\hat{R} = \hat{W}_H \hat{W}_{H-1} \cdots \hat{W}_1$. It then follows from Theorem 3.1 that for any R , such
that R is a perturbation of \hat{R} and $\text{rank}(R) \leq d_p$, we have $R = W_H W_{H-1} \cdots W_1$, where W_i is a
perturbation of \hat{W}_i . Therefore, by noticing \hat{W} is a local minimum of (1), we have

$$F(R) = L(W) \geq L(\hat{W}) = F(\hat{R}),$$

218 which shows that \hat{R} is a local minimum of (2). \square

219 In the proof of Theorem 2.2, we at first show that we just need to consider the case where X is an
220 identity matrix and Y is a diagonal matrix by noticing rotation is invariant under Frobenius norm.
221 Then we show that the local minimum must be a block diagonal and symmetric matrix, and each block
222 term is a projection matrix on the space corresponding to the same eigenvalue of the diagonal matrix
223 Y . Finally, we show that those projection matrices must be onto the eigenspace of Y corresponding
224 to the as large as possible eigenvalues, which then shows that the local minimum shares the same
225 function value.

226 3.2 Proof of Theorem 2.2

227 Let $X = U_1 \Sigma_1 V_1^T$ be the SVD decomposition of X , where Σ_1 is a diagonal matrix with full row
228 rank. Then,

$$\begin{aligned} F(R) &= \|RU_1 \Sigma_1 V_1^T - Y\|_F^2 = \|RU_1 \Sigma_1 - Y V_1\|_F^2 \\ &= \|(RU_1)(\Sigma_1)_{1:d_1, 1:d_1} - (Y V_1)_{1:d_2, 1:d_1}\|_F^2 + \text{Const}, \end{aligned}$$

229 where Const is a constant in R and $(\cdot)_{t_1:t_2, t_3:t_4}$ is a submatrix of (\cdot) , which contains the t_1 to t_2 row
 230 and t_3 to t_4 column of (\cdot) . If R is a local minimum of (2), then $S = RU_1$ is a local minimum of
 231

$$\min_S G(S) = \|S\hat{\Sigma}_1 - \hat{Y}\|_F^2 \quad (4)$$

s.t. $\text{rank}(S) \leq k$,

where $\hat{\Sigma}_1 := (\Sigma_1)_{1:d_1, 1:d_1}$, $\hat{Y} := (YV_1)_{1:d_2, 1:d_1}$ and the difference of objective function values of
 (2) and (4) is a constant. Let $\hat{Y} := U_2 \Sigma_2 V_2^T$ be the SVD of \hat{Y} , then

$$G(S) = \|S\hat{\Sigma}_1 - U_2 \Sigma_2 V_2^T\|_F^2 = \|U_2^T S \hat{\Sigma}_1 V_2 - \hat{\Sigma}_2\|_F^2,$$

232 and if S is a local minimum of $G(S)$, we have $T := U_2^T S \hat{\Sigma}_1 V_2$ is a local minimum of

$$\min_T H(T) = \|T - \Sigma_2\|_F^2 \quad (5)$$

s.t. $\text{rank}(T) \leq k$,

233 and the objective function values of (4) and (5) are the same at corresponding points. Let Σ_2 have
 234 r distinct positive diagonal terms $\lambda_1 > \dots > \lambda_r \geq 0$ with multiplicities m_1, \dots, m_r . Let T^* be a
 235 local minimum of (5), and

$$T^* = U^* \Sigma^* V^{*T} = [U_S^* U_N^*] \begin{bmatrix} \Sigma_S^* & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} V_S^{*T} \\ V_N^{*T} \end{bmatrix}$$

236 be the SVD of T , where Σ_S^* are positive singular values. Let $P_L := U_S^* (U_S^{*T} U_S^*)^{-1} U_S^{*T}$ and
 237 $P_R := V_S^* (V_S^{*T} V_S^*)^{-1} V_S^{*T}$ be the projection matrix to the space spanned by U_S^* and V_S^* , respectively.
 238 Note that $\{T | P_L T = T\} \subseteq \{T | \text{rank}(T) \leq k\}$, thus, T^* is also a local minimum of

$$\min \|T - \Sigma_2\|_F^2 \quad (6)$$

s.t. $P_L T = T$,

239 which is a convex problem, and it can be shown by the first order optimality condition that the only
 240 local minimum of (6) is $T^* = P_L \Sigma_2$. Similarly, we have $T^* = \Sigma_2 P_R$. Then, $D := \Sigma_2 \Sigma_2^T$ is a
 241 diagonal matrix, with r distinct non-zero diagonal terms $\lambda_1^2 > \dots > \lambda_r^2 > 0$ with multiplicities
 242 m_1, \dots, m_r . Therefore,

$$\begin{aligned} P_L D P_L &= P_L \Sigma_2 \Sigma_2^T P_L^T = T^* T^{*T} = \Sigma_2 P_R P_R^T \Sigma_2^T \\ &= \Sigma_2 P_R \Sigma_2^T = \Sigma_2 T^{*T} = \Sigma_2 \Sigma_2^T P_L^T = D P_L. \end{aligned}$$

243 Note that the left hand is a symmetric matrix, thus, $D P_L$ is also a symmetric matrix. Meanwhile, P_L
 244 is a symmetric matrix, whereby P_L is a r -block diagonal matrix with each block corresponding to
 245 the same diagonal terms of D . Therefore, $T^* = P_L \Sigma_2$ is also a r -block diagonal matrix.

246 Let

$$T^* = \begin{bmatrix} T_1^* & & & \\ & \ddots & & \\ & & T_r^* & \\ & & & 0 \end{bmatrix},$$

247 where T_i^* is a $m_i \times m_i$ matrix, then $T^* T^{*T} = \Sigma_2 T^{*T}$ implies $T_i^* T_i^{*T} = \lambda_i T_i^{*T}$. Thus, T_i^* is a
 248 symmetric matrix and $\frac{T_i^*}{\lambda_i}$ is a projection matrix. Let $\text{rank}(T_i^*) = d_{p_i}$, then, $\sum_{i=1}^r d_{p_i} \leq p$ and
 249 $\text{tr}(T_i^*) = \lambda_i d_{p_i}$, whereby,

$$\begin{aligned} H(T^*) &= \sum_{i=1}^r \|T_i^* - \lambda_i I_{m_i}\|_F^2 \\ &= \sum_{i=1}^r \text{tr}(T_i^2) - 2\lambda_i \text{tr}(T_i) + m_i \lambda_i^2 \\ &= \sum_{i=1}^r (m_i - d_{p_i}) \lambda_i^2. \end{aligned}$$

250 Let j be the largest number that $\sum_{i=1}^j m_i < d_p$. Then, it is easy to find that the global minima of (6)
 251 satisfy $d_{p_i} = m_i$ for $i \leq j$, $d_{p_{j+1}} = d_p - \sum_{i=1}^j m_i$ and $d_{p_i} = 0$ for $i > j + 1$ which gives all of the
 252 global minima.

253 Now, let us show that all local minima must be global minima. As local minima T^* is a block
 254 diagonal matrix, thus, we can assume without loss of generality that both Σ_2 and T^* are square
 255 matrices, because the all 0 rows and columns in Σ_2 and T do not change anything. Thus, it follows
 256 T_i^* is symmetric that T^* is a symmetric matrix. Remember that $\frac{T_i^*}{\lambda_i}$ is a projection matrix, thus the
 257 eigenvalues of T_i^* are either 0 or λ_i , whereby

$$T^* = \sum_{i=1}^r \sum_{j=1}^{d_{p_i}} \lambda_i u_{ij} u_{ij}^T,$$

258 where u_{ij} is the j th normalized orthogonal eigen-vector of T^* corresponding to eigenvalue λ_i .

259 It is easy to see that, at a local minimum, we have $\sum_{i=1}^r d_{p_i} = d_p$, otherwise, there is a descent
 260 direction by adding a rank 1 matrix to T^* corresponding to one positive eigenvalue. If there exists
 261 i_1, i_2 , such that $i_1 < i_2$, $d_{p_{i_1}} < m_{i_1}$, and $d_{p_{i_2}} \geq 1$, then, there exists \bar{u}_{i_1} , such that $\bar{u}_{i_1} \perp u_{i_1 j}$ for
 262 $j = 1, \dots, d_{p_{i_1}}$. Let

$$T(\theta) := T^* - \lambda_{i_2} u_{i_2 1} u_{i_2 1}^T + (\lambda_{i_1} \sin^2 \theta + \lambda_{i_2} \cos^2 \theta) \\ (u_{i_2 1} \cos \theta + \bar{u}_{i_1} \sin \theta) (u_{i_2 1} \cos \theta + \bar{u}_{i_1} \sin \theta)^T.$$

Then, $\text{rank}(T(\theta)) = \text{rank}(T^*) = d_p$, $T(0) = T^*$ and

$$H(T(\theta)) = H(T^*) + \lambda_1^2 + \lambda_2^2 - (\lambda_1 \sin^2 \theta + \lambda_2 \cos^2 \theta)^2.$$

263 It is easy to check that $H(T(\theta))$ is monotonically decreasing with θ , which gives a descent
 264 at T^* , contradicting with that T^* is a local minimum. Therefore, there is no such i_1 and i_2 , which
 265 shows that T^* is a global minimum. \square

266 3.3 Proof of Theorem 2.3

267 The statement follows from Theorem 2.1 and 2.2.

268 4 Conclusion

269 We have proven that, even though depth creates a non-convex loss surface, it does not create new bad
 270 local minima. Based on this new insight, we have successfully proposed a new simple proof for the
 271 fact that all of the local minima of feedforward deep linear neural networks are global minima as a
 272 corollary.

273 The benefits of this new results are not limited to the simplification of the previous proof. For
 274 example, our results apply to problems beyond square loss. Let us consider the shallow prob-
 275 lem (S) minimize $L(R)$ s.t. $\text{rank}(R) \leq d_p$, and and the deep parameterization counterpart (D)
 276 minimize $L(W_H W_{H-1} \cdots W_1)$. Our analysis shows that for any function L , as long as L satisfies
 277 Theorem 3.2, any local minimum of (D) corresponds to a local minimum of (S). This is not limited to
 278 when L is least square loss, and this is why we say depth creates no bad local minima.

279 In addition, our analysis can directly apply to matrix completion unlike previous results. Ge et al.
 280 (2016) show that local minima of the symmetric matrix completion problem are global with high
 281 probability. This should be able to extend to asymmetric case. Denote $f(W) := \sum_{i,j \in \Omega} (Y -$
 282 $W_2 W_1)_{i,j}$, then local minimum of $f(W)$ is global with high probability, where Ω is the observed
 283 entries. Then, our analysis here can directly show that the result can be extended for deep linear
 284 parameterization: for $h(W) := \sum_{i,j \in \Omega} (Y - W_H W_{H-1} \cdots W_1)_{i,j}$, any local minimum of $h(W)$ is
 285 global with high probability.

286 **References**

- 287 Arora, Raman, Basu, Amitabh, Mianjy, Poorya, and Mukherjee, Anirbit. Understanding deep neural
288 networks with rectified linear units. *arXiv preprint arXiv:1611.01491*, 2016.
- 289 Baldi, Pierre. Linear learning: Landscapes and algorithms. In *Advances in neural information*
290 *processing systems*, pp. 65–72, 1989.
- 291 Blum, Avrim L and Rivest, Ronald L. Training a 3-node neural network is np-complete. *Neural*
292 *Networks*, 5(1):117–127, 1992.
- 293 Choromanska, Anna, Henaff, Mikael, Mathieu, Michael, Ben Arous, Gerard, and LeCun, Yann. The
294 loss surfaces of multilayer networks. In *Proceedings of the Eighteenth International Conference*
295 *on Artificial Intelligence and Statistics*, pp. 192–204, 2015.
- 296 Dauphin, Yann N, Pascanu, Razvan, Gulcehre, Caglar, Cho, Kyunghyun, Ganguli, Surya, and Bengio,
297 Yoshua. Identifying and attacking the saddle point problem in high-dimensional non-convex
298 optimization. In *Advances in Neural Information Processing Systems*, pp. 2933–2941, 2014.
- 299 Freeman, C Daniel and Bruna, Joan. Topology and geometry of half-rectified network optimization.
300 *arXiv preprint arXiv:1611.01540*, 2016.
- 301 Ge, Rong, Lee, Jason D, and Ma, Tengyu. Matrix completion has no spurious local minimum. In
302 *Advances in Neural Information Processing Systems*, pp. 2973–2981, 2016.
- 303 Haeffele, Benjamin D and Vidal, René. Global optimality in tensor factorization, deep learning, and
304 beyond. *arXiv preprint arXiv:1506.07540*, 2015.
- 305 Kawaguchi, Kenji. Deep learning without poor local minima. In *Advances in Neural Information*
306 *Processing Systems (NIPS)*, 2016.
- 307 Kawaguchi, Kenji, Kaelbling, Leslie Pack, and Lozano-Pérez, Tomás. Bayesian optimization with
308 exponential convergence. In *In Advances in Neural Information Processing (NIPS)*, 2015.
- 309 Kawaguchi, Kenji, Maruyama, Yu, and Zheng, Xiaoyu. Global continuous optimization with error
310 bound and fast convergence. *Journal of Artificial Intelligence Research*, 56:153–195, 2016.
- 311 Krkova, Vera and Kainen, Paul C. Functionally equivalent feedforward neural networks. *Neural*
312 *Computation*, 6(3):543–558, 1994.
- 313 Livni, Roi, Shalev-Shwartz, Shai, and Shamir, Ohad. On the computational efficiency of training
314 neural networks. In *Advances in Neural Information Processing Systems*, pp. 855–863, 2014.
- 315 Mhaskar, Hrushikesh, Liao, Qianli, and Poggio, Tomaso. Learning functions: When is deep better
316 than shallow. *arXiv preprint arXiv:1603.00988*, 2016.
- 317 Poggio, Tomaso, Mhaskar, Hrushikesh, Rosasco, Lorenzo, Miranda, Brando, and Liao, Qianli. Why
318 and when can deep—but not shallow—networks avoid the curse of dimensionality: a review. *arXiv*
319 *preprint arXiv:1611.00740*, 2016.
- 320 Saxe, Andrew M, McClelland, James L, and Ganguli, Surya. Exact solutions to the nonlinear
321 dynamics of learning in deep linear neural networks. In *International Conference on Learning*
322 *Representations*, 2014.
- 323 Shamir, Ohad. Distribution-specific hardness of learning neural networks. *arXiv preprint*
324 *arXiv:1609.01037*, 2016.
- 325 Soudry, Daniel and Hoffer, Elad. Exponentially vanishing sub-optimal local minima in multilayer
326 neural networks. *arXiv preprint arXiv:1702.05777*, 2017.
- 327 Swirszcz, Grzegorz, Czarnecki, Wojciech Marian, and Pascanu, Razvan. Local minima in training of
328 deep networks. *arXiv preprint arXiv:1611.06310*, 2016.
- 329 Wang, Zi, Zhou, Bolei, and Jegelka, Stefanie. Optimization as estimation with gaussian processes in
330 bandit settings. In *International Conf. on Artificial and Statistics (AISTATS)*, 2016.
- 331 Wedin, Per-Åke. Perturbation bounds in connection with singular value decomposition. *BIT*
332 *Numerical Mathematics*, 12(1):99–111, 1972.