

# Characterization of Gradient Dominance and Regularity Conditions for Neural Networks

**Yi Zhou**  
Ohio State University  
**Yingbin Liang**  
Ohio State University

zhou.1172@osu.edu

liang.889@osu.edu

## Abstract

The past decade has witnessed a successful application of deep learning to solving many challenging problems in machine learning and artificial intelligence. However, the loss functions of neural networks are still far from being well understood from a theoretical aspect. In this paper, we enrich the current understanding of the landscape of the square loss functions for three types of neural networks, i.e., linear networks, linear residual networks, and one-hidden-layer nonlinear networks. Specifically, when the parameter matrices are square, we establish two quadratic types of landscape properties for the square loss of these neural networks: the gradient dominance condition within the neighborhood of their full rank global minimizers and the regularity condition along certain directions and within the neighborhood of their global minimizers.

## 1 Introduction

The significant success of deep learning (see, e.g., [2]) has influenced many fields such as machine learning, artificial intelligence, computer vision, natural language processing, etc. Consequently, there is a rising interest in understanding the fundamental properties of deep neural networks. Among them, the landscape (also referred to as geometry) of the loss functions of neural networks is an important aspect, since it is central to determine the performance of optimization algorithms that are designed to minimize these loss functions.

This paper focuses on two important landscape properties for nonconvex optimization. The first property is referred to as *gradient dominance condition* as we describe below. Consider a global minimizer  $\mathbf{x}^*$  of a generic function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ , and a neighborhood  $\mathcal{B}_{\mathbf{x}^*}(\delta)$  around  $\mathbf{x}^*$ . The local gradient dominance condition with regard to  $\mathbf{x}^*$  is given by, for some  $\lambda > 0$

$$\forall \mathbf{x} \in \mathcal{B}_{\mathbf{x}^*}(\delta), f(\mathbf{x}) - f(\mathbf{x}^*) \leq \lambda \|\nabla f(\mathbf{x})\|_2^2,$$

This condition has been shown to hold for a variety of nonconvex machine learning problems, e.g., phase retrieval [11] and blind deconvolution [5]. If the gradient descent algorithm iterates in the neighborhood  $\mathcal{B}_{\mathbf{x}^*}(\delta)$ , then the gradient dominance condition, together with a Lipschitz property of the gradient of the objective function, guarantees a linear convergence of the function value residual  $f(\mathbf{x}) - f(\mathbf{x}^*)$  [4, 7]. The second property is referred to as local *regularity condition*, which is given by, for some  $\alpha, \beta > 0$

$$\forall \mathbf{x} \in \mathcal{B}_{\mathbf{x}^*}(\delta), \langle \mathbf{x} - \mathbf{x}^*, \nabla f(\mathbf{x}) \rangle \geq \alpha \|\nabla f(\mathbf{x})\|_2^2 + \beta \|\mathbf{x} - \mathbf{x}^*\|_2^2,$$

This condition can be viewed as a restricted version of the strong convexity, and it has been shown to guarantee a linear convergence of the iterate residual  $\|\mathbf{x} - \mathbf{x}^*\|$  of the gradient descent algorithm [6, 1]. Problems such as phase retrieval [1], affine rank minimization [9, 8] and matrix completion [10] have been shown to satisfy the local regularity condition. This paper studies these two landscape properties for the square loss function of linear, linear residual, and one-hidden-layer nonlinear neural networks under the setting where all parameter matrices are square.

## 2 Preliminaries of Neural Networks

Throughout,  $(\mathbf{X}, \mathbf{Y})$  denotes the input and output data matrix pair. We assume that  $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{d \times m}$ , and  $\Sigma := \Sigma_{\mathbf{X}\mathbf{Y}}^\top \Sigma_{\mathbf{X}\mathbf{X}}^{-1} \Sigma_{\mathbf{X}\mathbf{Y}}$  is full rank with distinct eigenvalues. The kronecker product is denoted as “ $\otimes$ ”. For a matrix  $\mathbf{X}$ , its spectral norm is denoted by  $\|\mathbf{X}\|$ . The smallest nonzero singular value is denoted by  $\eta_{\min}(\mathbf{X})$ . For a collection of matrix variables  $\mathbf{W} := \{\mathbf{W}_1, \dots, \mathbf{W}_l\}$ , we denote  $\text{vec}(\mathbf{W}) := [\text{vec}(\mathbf{W}_1)^\top \dots \text{vec}(\mathbf{W}_l)^\top]^\top$ ,

where  $\text{vec}(\mathbf{W}_k)$  denotes the vertical stack of the columns of  $\mathbf{W}_k$ . We also denote a collection of natural numbers as  $[n] := \{1, \dots, n\}$ .

Consider a feed forward linear neural network with  $l-1$  hidden layers, where each layer  $k \in [l]$  is parameterized by a matrix  $\mathbf{W}_k \in \mathbb{R}^{d \times d}$ . We consider the square loss of the linear network:

$$h(\mathbf{W}) := \frac{1}{2} \|\mathbf{W}_l \mathbf{W}_{l-1} \dots \mathbf{W}_1 \mathbf{X} - \mathbf{Y}\|_F^2. \quad (1)$$

The linear residual network further introduces the residual structure to the linear network. That is, one adds a shortcut (identity map) for every, say,  $r$  hidden layers. Assuming we have in total  $l$  residual units. The  $k$ -th residual unit is parameterized by the parameters  $\mathbf{A}_{kq} \in \mathbb{R}^{d \times d}, \forall q \in [r]$ , and we denote  $\mathbf{B}_k := \mathbf{I} + \mathbf{A}_{kr} \dots \mathbf{A}_{k1}$ . We consider the square loss of a linear residual network:

$$f(\mathbf{A}) := \frac{1}{2} \|(\mathbf{I} + \mathbf{A}_{lr} \dots \mathbf{A}_{l1}) \dots (\mathbf{I} + \mathbf{A}_{kr} \dots \mathbf{A}_{k1}) \dots (\mathbf{I} + \mathbf{A}_{1r} \dots \mathbf{A}_{11}) \mathbf{X} - \mathbf{Y}\|_F^2. \quad (2)$$

Consider a nonlinear neural network with one hidden layer, where the layer parameters satisfy  $\mathbf{V}_1, \mathbf{V}_2 \in \mathbb{R}^{d \times d}$  and the hidden neurons adopt a differentiable nonlinear activation function  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ . We consider the square loss of the nonlinear network with one hidden layer:

$$g(\mathbf{V}) := \frac{1}{2} \|\mathbf{V}_2 \sigma(\mathbf{V}_1 \mathbf{X}) - \mathbf{Y}\|_F^2, \quad (3)$$

where  $\sigma$  acts on  $\mathbf{V}_1 \mathbf{X}$  *entrywise*. In particular, we consider a class of activation functions that satisfy the condition  $\text{range}(\sigma) = \mathbb{R}$ . A typical example of such activation function is the class of parametric ReLU activation functions, i.e.,  $\sigma(x) = \max\{x, ax\}$ , where  $0 < a < 1$ . The following theorem characterizes a useful property of the global minimizers of these networks.

**Theorem 1.** *Consider the global minimizers  $\mathbf{W}^*, \mathbf{A}^*, \mathbf{V}^*$  of  $h(\mathbf{W}), f(\mathbf{A}), g(\mathbf{V})$ , respectively. Then for all  $k \in [l, 1)$   $\mathbf{W}_k^*$  is full rank; 2)  $\mathbf{B}_k^*$  is full rank; and 3)  $\sigma(\mathbf{W}_1^* \mathbf{X})$  is full rank.*

### 3 Gradient Dominance Condition

We first establish the local gradient dominance condition for the loss  $h(\mathbf{W})$  of linear networks.

**Theorem 2.** *Consider  $h(\mathbf{W})$  of the linear neural network with  $m = d$ . Consider a global minimizer  $\mathbf{W}^*$  and let  $\tau = \frac{1}{2} \min_{k \in [l]} \eta_{\min}(\mathbf{W}_k^*)$ . Then any point in the neighborhood of  $\mathbf{W}^*$  defined as  $\{\mathbf{W} : \|\mathbf{W}_k - \mathbf{W}_k^*\| < \tau, \forall k \in [l]\}$  satisfies*

$$h(\mathbf{W}) - h(\mathbf{W}^*) \leq \lambda_h \|\nabla_{\text{vec}(\mathbf{W})} h(\mathbf{W})\|_2^2, \text{ where } \lambda_h = (2l\tau^{2(l-1)}\eta_{\min}^2(\mathbf{X}))^{-1}. \quad (4)$$

We note that Theorem 1 guarantees that any global minimizer  $\mathbf{W}^*$  of  $h(\mathbf{W})$  is full rank, and hence the parameter  $\tau$  defined in Theorem 2 is strictly positive. The gradient dominance condition implies a linear convergence of the function value to the global minimum via a gradient descent algorithm if the iterations stay in this  $\tau$  neighborhood. In particular, a larger parameter  $\tau$  (a larger minimum singular value) implies a smaller  $\lambda_h$ , which yields a faster convergence of the function value to the global minimum via the gradient descent algorithm. Next, we establish the local gradient dominance condition for the loss  $f(\mathbf{A})$  of linear residual networks.

**Theorem 3.** *Consider  $f(\mathbf{A})$  of the linear residual neural network with  $m = d$ . Consider a full-rank global minimizer  $\mathbf{A}^*$ , and let  $\tau = \frac{1}{2} \min_{k \in [l]} \eta_{\min}(\mathbf{B}_k^*)$ ,  $\tilde{\tau} = \frac{1}{2} \min_{k \in [l], q \in [r]} \eta_{\min}(\mathbf{A}_{kq}^*)$ , and pick  $\hat{\tau}$  sufficiently small such that any point in the neighborhood of  $\mathbf{A}^*$  defined as  $\{\mathbf{A} : \|\mathbf{A}_{kq} - \mathbf{A}_{kq}^*\| < \hat{\tau}, \forall k \in [l], q \in [r]\}$  satisfies  $\|\mathbf{B}_k - \mathbf{B}_k^*\| < \tau$  for all  $k \in [l]$ . Then any point in the neighborhood of  $\mathbf{A}^*$  defined as  $\{\mathbf{A} : \|\mathbf{A}_{kq} - \mathbf{A}_{kq}^*\| < \min\{\hat{\tau}, \tilde{\tau}\}, \forall k \in [l], q \in [r]\}$  satisfies*

$$f(\mathbf{A}) - f(\mathbf{A}^*) \leq \lambda_f \|\nabla_{\text{vec}(\mathbf{A})} f(\mathbf{A})\|_2^2, \text{ where } \lambda_f = \left(2lr\tilde{\tau}^{2(r-1)}\tau^{2(l-1)}\eta_{\min}^2(\mathbf{X})\right)^{-1}. \quad (5)$$

Compare to the gradient dominance condition obtained in [3], which is applicable to the neighborhood of the origin for the residual network with  $r = 1$ , the above result characterizes the gradient dominance condition in the neighborhood of any full-rank minimizer  $\mathbf{A}^*$  and to more general residual networks with  $r > 1$ .

We note that the parameter  $\lambda_f$  in Theorem 3 depends on both  $\tau$  and  $\tilde{\tau}$ , where  $\tau$  captures the overall property of each residual unit and  $\tilde{\tau}$  captures the property of individual linear unit in each residual unit. Hence, in general,

the  $\lambda_f$  in Theorem 3 for linear residual networks is very different from the  $\lambda_h$  in the gradient dominance condition in Theorem 2 for linear networks. When the shortcut depth  $r$  becomes large, the parameter  $\lambda_f$  involves  $\tilde{\tau}$  that depends on more unparameterized variables in  $\mathbf{A}^*$ , and hence becomes more similar to the parameter  $\lambda_h$  of linear networks.

To further compare the  $\lambda_f$  in Theorem 3 and the  $\lambda_h$  in Theorem 2, consider a simplified setting of the linear residual network with the shortcut depth  $r = 1$ . Then  $\lambda_f = (2l\tau^{2(l-1)}\eta_{\min}^2(\mathbf{X}))^{-1}$ . Although it takes the same expression as  $\lambda_h$  in Theorem 2 for the linear network, the parameter  $\tau$  is better regularized since  $\mathbf{B}_k, k \in [l]$  are further parameterized by  $\mathbf{I} + \mathbf{A}_k$ . More specifically, when  $\|\mathbf{A}_k\| < 1$ ,  $\eta_{\min}(\mathbf{I} + \mathbf{A}_k)$  (and hence the parameter  $\tau$ ) is regularized away from zero by the identity map, which was also observed by [3]. Consequently, the identity shortcut leads to a smaller  $\lambda_f$  (due to larger  $\tau$ ) compared to a large  $\lambda_h$  when the parameters of linear networks have small spectral norm. Such a smaller  $\lambda_f$  is more desirable for optimization, because the function value approaches closer to the global minimum after one iteration of a gradient descent algorithm. We now establishes the local gradient dominance condition for the loss  $g(\mathbf{V})$  of nonlinear networks.

**Theorem 4.** *Consider the loss function  $g(\mathbf{V})$  of one-hidden-layer nonlinear neural networks with  $m = d$  and with  $\text{range}(\sigma) = \mathbb{R}$ . Consider a global minimizer  $\mathbf{V}^*$ , and let  $\tau = \frac{1}{2}\eta_{\min}(\sigma(\mathbf{V}_1^*\mathbf{X}))$ . Then any point in the neighborhood of  $\mathbf{V}^*$  defined as  $\{\mathbf{V} : \|\sigma(\mathbf{V}_1\mathbf{X}) - \sigma(\mathbf{V}_1^*\mathbf{X})\| \leq \tau\}$  satisfies*

$$g(\mathbf{V}) - g(\mathbf{V}^*) \leq \lambda_g \|\nabla_{\text{vec}(\mathbf{V})}g(\mathbf{V})\|_2^2, \text{ where } \lambda_g = (2\tau^2)^{-1}. \quad (6)$$

We note that Theorem 1 guarantees that  $\sigma(\mathbf{V}_1^*\mathbf{X})$  is full rank, and hence  $\tau$  is well defined. Differently from linear networks, the gradient dominance condition for nonlinear networks holds in a nonlinear  $\tau$  neighborhood that involves the activation function  $\sigma$ . This is naturally due to the nonlinearity of the network. Furthermore, the parameter  $\tau$  depends on the nonlinear term  $\sigma(\mathbf{V}_1\mathbf{X})$ , whereas the  $\tau$  in Theorem 2 of linear networks depends on the individual parameters  $\mathbf{W}_k$ .

## 4 Regularity Condition

For the linear network, define matrix  $\mathbf{G}(\mathbf{W}) := [\mathbf{G}_1(\mathbf{W}), \dots, \mathbf{G}_l(\mathbf{W})]$ , where for all  $k \in [l]$

$$\mathbf{G}_k(\mathbf{W}) := (\mathbf{W}_{k-1} \dots \mathbf{W}_1 \mathbf{X})^\top \otimes (\mathbf{W}_l \dots \mathbf{W}_{k+1}). \quad (7)$$

The following result establishes the local regularity condition for the loss  $h(\mathbf{W})$  of linear networks.

**Theorem 5.** *Consider  $h(\mathbf{W})$  of linear neural networks with  $m = d$ . Further consider a global minimizer  $\mathbf{W}^*$ , and let  $\zeta = 2 \max_{k \in [l]} \|\mathbf{W}_k^*\|$ . Then for any  $\delta > 0$ , there exists a sufficiently small  $\epsilon(\delta)$  such that any point  $\mathbf{W}$  that satisfies*

$$\|\mathbf{G}(\mathbf{W}^*)\text{vec}(\mathbf{W} - \mathbf{W}^*)\|_2 \geq \delta \|\text{vec}(\mathbf{W} - \mathbf{W}^*)\|_2 \quad (8)$$

and within the neighborhood of  $\mathbf{W}^*$  defined as  $\{\mathbf{W} : \|\mathbf{W}_k - \mathbf{W}_k^*\|_F < \epsilon(\delta), \forall k \in [l]\}$  satisfies

$$\langle \nabla_{\text{vec}(\mathbf{W})}h(\mathbf{W}), \text{vec}(\mathbf{W} - \mathbf{W}^*) \rangle \geq \alpha \|\nabla_{\text{vec}(\mathbf{W})}h(\mathbf{W})\|_2^2 + \beta \|\text{vec}(\mathbf{W} - \mathbf{W}^*)\|_2^2, \quad (9)$$

where  $\alpha = \gamma/(l\zeta^{2(l-1)}\|\mathbf{X}\|^2)$  and  $\beta = (1 - \gamma)\delta^2/2$  for any  $\gamma \in (0, 1)$ .

We note that the regularity condition has been established for various nonconvex problems [1]. There, the condition was shown to hold within the entire neighborhood of any global minimizer. In comparison, Theorem 5 guarantees the regularity condition for linear networks within a neighborhood of  $\mathbf{W}^*$  along the directions of  $\text{vec}(\mathbf{W} - \mathbf{W}^*)$  that satisfy eq. (8). Furthermore, the parameter  $\delta$  in eq. (8) determines the range of directions that satisfy eq. (8). For example, if we set  $\delta = \eta_{\min}(\mathbf{G}(\mathbf{W}^*))$ , then all  $\mathbf{W}$  such that  $\text{vec}(\mathbf{W} - \mathbf{W}^*) \perp \ker \mathbf{G}(\mathbf{W}^*)$  satisfy eq. (8).

For all  $\mathbf{W}$  that satisfy the regularity condition, it can be shown that one gradient descent iteration yields an update that is closer to the global minimizer  $\mathbf{W}^*$  [1]. Hence,  $\mathbf{W}^*$  serves as an attractive point along the directions of  $\text{vec}(\mathbf{W} - \mathbf{W}^*)$  that satisfy eq. (8). Furthermore, the value of  $\delta$  in eq. (8) affects the parameter  $\beta$  in the regularity condition. Larger  $\delta$  results larger  $\beta$ , which further implies that one gradient descent iteration at the point  $\mathbf{W}$  yields an update that is closer to the global minimizer  $\mathbf{W}^*$  [1].

For the linear residual network, define  $\mathbf{Q}(\mathbf{A}) := [\mathbf{Q}_{11}(\mathbf{A}), \dots, \mathbf{Q}_{lr}(\mathbf{A})]$ , where for all  $k \in [l], q \in [r]$

$$\mathbf{Q}_{kq}(\mathbf{A}) := [(\mathbf{B}_{k-1} \dots \mathbf{B}_1 \mathbf{X})^\top \otimes (\mathbf{B}_l \dots \mathbf{B}_{k+1})] \left[ (\mathbf{A}_{k(q-1)} \dots \mathbf{A}_{k1})^\top \otimes (\mathbf{A}_{kr} \dots \mathbf{A}_{k(q+1)}) \right].$$

We then establish the local regularity condition for the loss  $f(\mathbf{A})$  of linear residual networks.

**Theorem 6.** Consider  $f(\mathbf{A})$  of linear residual neural networks with  $m = d$ . Further consider a global minimizer  $\mathbf{A}^*$ , and let  $\zeta = 2 \max_{k \in [l]} \|\mathbf{B}_k^*\|$  and  $\tilde{\zeta} = 2 \max_{k \in [l], q \in [r]} \|\mathbf{A}_{kq}^*\|$ . Then, for any constant  $\delta > 0$ , there exists a sufficiently small  $\epsilon(\delta)$  such that any point  $\mathbf{A}$  that satisfies

$$\|\mathbf{Q}(\mathbf{A}^*) \text{vec}(\mathbf{A} - \mathbf{A}^*)\|_2 \geq \delta \|\text{vec}(\mathbf{A} - \mathbf{A}^*)\|_2$$

and within the neighborhood defined as  $\{\mathbf{A} : \|\mathbf{A}_{kq} - \mathbf{A}_{kq}^*\|_F < \epsilon(\delta), \forall k \in [l], q \in [r]\}$ , satisfies

$$\langle \nabla_{\text{vec}(\mathbf{A})} f(\mathbf{A}), \text{vec}(\mathbf{A} - \mathbf{A}^*) \rangle \geq \alpha \|\nabla_{\text{vec}(\mathbf{A})} f(\mathbf{A})\|_2^2 + \beta \|\text{vec}(\mathbf{A} - \mathbf{A}^*)\|_2^2, \quad (10)$$

where  $\alpha = \gamma / (lr \tilde{\zeta}^{2(r-1)} \zeta^{2(l-1)} \|\mathbf{X}\|^2)$  and  $\beta = (1 - \gamma)\delta^2/2$  with any  $\gamma \in (0, 1)$ .

Similarly to the regularity condition for linear networks, the regularity condition for linear residual networks holds along the directions of  $\text{vec}(\mathbf{A} - \mathbf{A}^*)$  that depends on  $\mathbf{Q}(\mathbf{A}^*)$ . However, the parametrization of  $\mathbf{Q}(\mathbf{A}^*)$  is different from that of  $\mathbf{G}(\mathbf{W}^*)$  of linear networks. To illustrate, consider a simplified setting where the shortcut depth is  $r = 1$ . Then,  $\mathbf{Q}_k(\mathbf{A}^*) = (\mathbf{B}_{k-1}^* \dots \mathbf{B}_1^* \mathbf{X})^\top \otimes (\mathbf{B}_l^* \dots \mathbf{B}_{k+1}^*)$ . Although it takes a similar form as  $\mathbf{G}(\mathbf{W}^*)$  of the linear network, the reparameterization  $\mathbf{B}_k^* = \mathbf{I} + \mathbf{A}_k^*$  keeps  $\eta_{\min}(\mathbf{B}_k^*)$  away from zero when  $\|\mathbf{A}_k^*\|$  is small. This enlarges  $\eta_{\min}(\mathbf{Q}_k(\mathbf{A}^*))$  so that the direction constraint can be satisfied along a wider range of directions. In this way,  $\mathbf{A}^*$  attracts the optimization iteration to converge along a wider range of directions in the neighborhood of the origin.

For the nonlinear network, define matrix

$$\mathbf{H} = [(\mathbf{I} \otimes \mathbf{V}_2^*) \sigma'(\text{diag}(\text{vec}(\mathbf{V}_1^* \mathbf{X})))(\mathbf{X}^\top \otimes \mathbf{I}), \sigma(\mathbf{V}_1^* \mathbf{X})^\top \otimes \mathbf{I}]$$

We now establish the local regularity condition for the loss  $g(\mathbf{V})$  of nonlinear networks.

**Theorem 7.** Consider  $g(\mathbf{V})$  of one-hidden-layer nonlinear neural networks with  $m = d$  and  $\text{range}(\sigma) = \mathbb{R}$ . Further consider a global minimizer  $\mathbf{V}^*$  of  $g(\mathbf{V})$ , and let  $\zeta = 2 \max\{\|\sigma(\mathbf{V}_1^* \mathbf{X})\|, \|\mathbf{V}_2^*\|, \|\sigma'(\mathbf{V}_1^* \mathbf{X})\|_\infty\}$ . Then there exists a sufficiently small  $\epsilon(\delta)$  such that any point  $\mathbf{V}$  that satisfies

$$\|\mathbf{H}(\mathbf{V}^*) \text{vec}(\mathbf{V} - \mathbf{V}^*)\|_2 \geq \delta \|\text{vec}(\mathbf{V} - \mathbf{V}^*)\|_2$$

and within the neighborhood of  $\mathbf{V}^*$  defined as  $\{\mathbf{V} : \|\mathbf{V}_k - \mathbf{V}_k^*\|_F < \epsilon(\delta), \forall k \in [2]\}$  satisfies

$$\langle \nabla_{\text{vec}(\mathbf{V})} g(\mathbf{V}), \text{vec}(\mathbf{V} - \mathbf{V}^*) \rangle \geq \alpha \|\nabla_{\text{vec}(\mathbf{V})} g(\mathbf{V})\|_2^2 + \beta \|\text{vec}(\mathbf{V} - \mathbf{V}^*)\|_2^2, \quad (11)$$

where  $\alpha = \gamma / \max\{\|\mathbf{X}\|^2 \zeta^4, \zeta^2\}$  and  $\beta = (1 - \gamma)\delta^2/2$  for any  $\gamma \in (0, 1)$ .

Thus, nonlinear neural networks with one hidden layer also have an amenable landscape near the global minimizers that attracts gradient iterates to converge along the directions restricted by  $\mathbf{H}(\mathbf{V}^*)$ .

## 5 Conclusion

In this paper, we establish the gradient dominance condition and the regularity condition for three types of neural networks in the neighborhood of their global minimizers under certain conditions. It is interesting to exploit these conditions in the convergence analysis of optimization algorithms applied to deep learning networks.

## References

- [1] E. Candès, X. Li, and M. Soltanolkotabi. Phase retrieval via wirtinger flow: Theory and algorithms. *IEEE Transactions on Information Theory*, 61(4):1985–2007, April 2015.
- [2] I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT Press, 2016. <http://www.deeplearningbook.org>.
- [3] M. Hardt and T. Ma. Identity matters in deep learning. *Proc. International Conference on Learning Representations (ICLR)*, 2017.
- [4] H. Karimi, J. Nutini, and M. Schmidt. Linear convergence of gradient and proximal-gradient methods under the Polyak-Łojasiewicz condition. In *Machine Learning and Knowledge Discovery in Databases: European Conference, ECML PKDD*, pages 795–811, 2016.

- [5] X. Li, S. Ling, T. Strohmer, and K. Wei. Rapid, robust, and reliable blind deconvolution via nonconvex optimization. *Arxiv: 1606.04933*, 2016. URL <http://arxiv.org/abs/1606.04933>.
- [6] Y. Nesterov. *Introductory Lectures on Convex Optimization: A Basic Course*. Springer, 2014.
- [7] S. Reddi, A. Hefny, S. Sra, B. Póczós, and A. Smola. Stochastic variance reduction for nonconvex optimization. In *Proc. International Conference on Machine Learning (ICML)*, pages 314–323, 2016.
- [8] S. Tu, R. Boczar, M. Soltanolkotabi, and B. Recht. Low-rank solutions of linear matrix equations via procrustes flow. In *Proc. International Conference on Machine Learning (ICML)*, pages 964–973, 2016.
- [9] Q. Zheng and J. Lafferty. A convergent gradient descent algorithm for rank minimization and semidefinite programming from random linear measurements. In *Proc. International Conference on Neural Information Processing Systems (NIPS)*, pages 109–117, 2015.
- [10] Q. Zheng and J. Lafferty. Convergence analysis for rectangular matrix completion using burer-monteiro factorization and gradient descent. *ArXiv: 1605.07051*, 2016.
- [11] Y. Zhou, H. Zhang, and Y. Liang. Geometrical properties and accelerated gradient solvers of non-convex phase retrieval. In *Proc. Annual Allerton Conference on Communication, Control, and Computing*. 2016.