

How regularization affects the critical points in linear networks

Authored by:

Amirhossein Taghvaei
Jin W. Kim
Prashant Mehta

Abstract

This paper is concerned with the problem of representing and learning a linear transformation using a linear neural network. In recent years, there is a growing interest in the study of such networks, in part due to the successes of deep learning. The main question of this body of research (and also of our paper) is related to the existence and optimality properties of the critical points of the mean-squared loss function. An additional primary concern of our paper pertains to the robustness of these critical points in the face of (a small amount of) regularization. An optimal control model is introduced for this purpose and a learning algorithm (backprop with weight decay) derived for the same using the Hamilton's formulation of optimal control. The formulation is used to provide a complete characterization of the critical points in terms of the solutions of a nonlinear matrix-valued equation, referred to as the characteristic equation. Analytical and numerical tools from bifurcation theory are used to compute the critical points via the solutions of the characteristic equation.

1 Paper Body

This paper is concerned with the problem of representing and learning a linear transformation with a linear neural network. Although a classical problem (Baldi and Hornik [1989, 1995]), there has been a renewed interest in such networks (Saxe et al. [2013], Kawaguchi [2016], Hardt and Ma [2016], Gunasekar et al. [2017]) because of the successes of deep learning. The motivation for studying linear networks is to gain insight into the optimization problem for the more general nonlinear networks. A ?

Financial support from the NSF CMMI grant 1462773 is gratefully acknowledged.

31st Conference on Neural Information Processing Systems (NIPS 2017), Long Beach, CA, USA.

focus of the recent research on these (and also nonlinear) networks has been on the analysis of the critical points of the non-convex loss function (Dauphin et al. [2014], Choromanska et al. [2015a,b], Soudry and Carmon [2016], Bhojanapalli et al. [2016]). This is also the focus of our paper. Problem: The input-output model is assumed to be of the following linear form: $Z = RX_0 + \eta$

(1)

where $X_0 \in \mathbb{R}^{d_0}$ is the input, $Z \in \mathbb{R}^{d_1}$ is the output, and $\eta \in \mathbb{R}^{d_1}$ is the noise. The input X_0 is modeled as a random variable whose distribution is denoted as p_0 . Its second moment is denoted as $\Sigma_0 := E[X_0 X_0^T]$ and assumed to be finite. The noise η is assumed to be independent of X_0 , with zero mean and finite variance. The linear transformation $R \in \mathbb{R}^{d_1 \times d_0}$ is assumed to satisfy a property (P1) introduced in Sec. 3 ($\mathbb{R}^{d_1 \times d_0}$ denotes the set of $d_1 \times d_0$ matrices). The problem is to learn the weights of a linear neural network from i.i.d. input-output samples $\{(X_0^k, Z^k)\}_{k=1}^K$. Solution architecture: is a continuous-time linear feedforward neural network model: $\dot{X} = A X + B U$

(2)

where $A \in \mathbb{R}^{d_1 \times d_1}$ are the network weights indexed by continuous-time (surrogate for layer) $t \in [0, T]$, and X_0 is the initial condition at time $t = 0$ (same as the input data). The parameter T denotes the network depth. The optimization problem is to choose the weights A over the time-horizon $[0, T]$ to minimize the mean-squared loss function: $E[\|X(T) - Z\|^2]$

(3)

This problem is referred to as the $[\eta = 0]$ problem. Backprop is a stochastic gradient descent algorithm for learning the weights A . In general, one obtains (asymptotic) convergence of the learning algorithm to a (local) minimum of the optimization problem Lee et al. [2016], Ge et al. [2015]. This has spurred investigation of the critical points of the loss function (3) and the optimality properties (local vs. global minima, saddle points) of these points. For linear multilayer (discrete) neural networks (MNN), strong conclusions have been obtained under rather mild conditions: every local minimum is a global minimum and every critical point that is not a local minimum is a saddle point Kawaguchi [2016], Baldi and Hornik [1989]. For the discrete counterpart of the $[\eta = 0]$ problem (referred to as the linear residual network in Hardt and Ma [2016]), an even stronger conclusion is possible: all critical points of the $[\eta = 0]$ problem are global minimum. In experiments, some of these properties are also empirically observed in deep nonlinear networks; cf., Choromanska et al. [2015b], Dauphin et al. [2014], Saxe et al. [2013]. In this paper, we consider the following regularized form of the optimization problem: $\min_{A, \lambda} J[A] = E[\|X(T) - Z\|^2] + \lambda \int_0^T \|A(t)\|^2 dt$ Subject to: $\dot{X} = A X + B U$, $X_0 \sim p_0$

(4)

where $\lambda \in \mathbb{R}_+ := \{x \in \mathbb{R} : x \geq 0\}$ is a regularization parameter. In literature, this form of regularization is referred to as weight decay [Goodfellow et al., 2016, Sec. 7.1.1]. Eq. (4) is an example of an optimal control problem and is referred to as such. The limit $\lambda \rightarrow 0$ is referred to as $[\eta = 0^+]$ problem.

The symbol $\text{tr}(\cdot)$ and superscript \top are used to denote matrix trace and matrix transpose, respectively. The regularized problem is important because of the following reasons: 2

(i) The learning algorithms are believed to converge to the critical points of the regularized $[\gamma = 0+]$ problem, a phenomenon known as implicit regularization Neyshabur et al. [2014], Zhang et al. [2016], Gunasekar et al. [2017]. (ii) It is shown in the paper that the stochastic gradient descent (for the functional J) yields the following learning algorithm for the weights $A_t : (k+1)$

$$\begin{aligned} & A_t^{(k)} \\ & = A_t^{(k)} \\ & + \eta_k (\nabla_{A_t} J \\ & + \text{backprop update}) \end{aligned} \quad (5)$$

for $k = 1, 2, \dots$, where η_k is the learning rate parameter. Thus, the parameter η models dissipation (or weight decay) in backprop. In an implementation of backprop, one would expect to obtain critical points of the $[\gamma = 0+]$ problem. The outline of the remainder of this paper is as follows: The Hamilton's formulation is introduced for the optimal control problem (4) in Sec. 2; cf., LeCun et al. [1988], Farotimi et al. [1991] for related constructions. The Hamilton's equations are used to obtain a formula for the gradient of J , and subsequently derive the stochastic gradient descent learning algorithm of the form (5). The equations for the critical points of J are obtained by applying the Maximum Principle of optimal control (Prop. 1). Remarkably, the Hamilton's equations for the critical points can be solved in closed-form to obtain a characterization of the critical points in terms of the solutions of a nonlinear matrix-valued equation, referred to as the characteristic equation (Prop. 2). For a certain special case, where the matrix R is normal, analytical results are obtained based on the use of the implicit function theorem (Thm. 2). Numerical continuation is employed to compute the solutions for this and the more general non-normal cases (Examples 1 and 2).

2

Hamilton's formulation and the learning algorithm

Definition 1. The control Hamiltonian is the function $H = \text{tr}(B^\top \dot{x} B)$ (6) 2 where $x \in \mathbb{R}^d$ is the state, $y \in \mathbb{R}^d$ is the co-state, and $B \in \mathbb{M}_d(\mathbb{R})$ is the weight matrix. The ∂_x ∂_y partial derivatives are denoted as $\partial_x H(x, y, B) := B^\top y$, $\partial_y H(x, y, B) := Bx$, and $\partial_B H(x, y, B) := \dot{x} y x^\top - \dot{y} B$. $H(x, y, B) = y^\top Bx$

Pontryagin's Maximum Principle (MP) is used to obtain the Hamilton's equations for the solution of the optimal control problem (4). The MP represents a necessary condition satisfied by any minimizer. Conversely, a solution of the Hamilton's equation is a critical point of the functional J . The proof of the following proposition appears in the supplementary material. Proposition 1. Consider the terminal cost optimal control A_t is the minimizer and X_t is the corresponding trajectory. $Y : [0, T] \rightarrow \mathbb{R}^d$ such that $dX_t = \partial_x H(X_t, Y_t, A_t)$

$\dot{H} = \frac{d}{dt} H(X_t, Y_t, A_t) = \text{tr}(A_t^T \dot{Y}_t) + \frac{d}{dt} \text{tr}(A_t^T X_t) + \frac{d}{dt} \text{tr}(A_t^T Y_t) = \text{tr}(A_t^T \dot{Y}_t) + \text{tr}(\dot{A}_t^T X_t) + \text{tr}(\dot{A}_t^T Y_t)$ and A_t maximizes the expected value of the Hamiltonian $A_t = \arg \max_{A \in \mathbb{R}^{d \times d}}$

$$E[H(X_t, Y_t, A)]$$

$$B \in \mathbb{R}^{d \times d}$$

problem (4) with $\epsilon = 0$. Suppose Then there exists a random process

$$X_0 \sim p_0$$

$$(7)$$

$$Y_T = Z + X_T$$

$$(8)$$

$$E[Y_t | X_t] =$$

$$(9)$$

$$(\epsilon_0)$$

$$=$$

Conversely, if there exists A_t and the pair (X_t, Y_t) such that equations (7)-(8)-(9) are satisfied, then A_t is a critical point of the optimization problem (4). \square

Remark 1. The Maximum Principle can also be used to derive analogous (difference) equations in discrete-time as well as nonlinear settings. It is equivalent to the method of Lagrange multipliers that is used to derive the backprop algorithm in MNN, e.g., LeCun et al. [1988]. The continuous-time limit is considered here because the computations are simpler and the results are more insightful. Similar considerations have also motivated the study of continuous-time limit of other types of optimization algorithms, e.g., Su et al. [2014], Wibisono et al. [2016]. The Hamiltonian is also used to express the first order variation in the functional J . For this purpose, define the Hilbert space of matrix-valued functions $L^2([0, T]; \mathbb{R}^{d \times d}) := \{A : [0, T] \rightarrow \mathbb{R}^{d \times d} \mid \int_0^T \text{tr}(A_t^T A_t) dt < \infty\}$ with the inner product $\langle A, V \rangle_{L^2} := \int_0^T \text{tr}(A_t^T V_t) dt$. For any $A \in L^2$, the gradient of the functional J evaluated at A is denoted as $\nabla J[A] \in L^2$. It is defined using the directional derivative formula: $J(A + V) - J(A) \approx \langle \nabla J[A], V \rangle_{L^2} := \lim_{h \rightarrow 0} \frac{J(A + hV) - J(A)}{h}$

where $V \in L^2$ prescribes the direction (variation) along which the derivative is being computed. The explicit formula for ∇J is given by

$\nabla J[A] := E \left[\frac{d}{dt} H(X_t, Y_t, A_t) = \text{tr}(A_t^T \dot{Y}_t) + \text{tr}(\dot{A}_t^T X_t) + \text{tr}(\dot{A}_t^T Y_t) \right] B$ where X_t and Y_t are the obtained by solving the Hamilton's equations (7)-(8) with the prescribed (not necessarily optimal) weight matrix $A \in L^2$. The significance of the formula is that the steepest descent in the objective function J is obtained by moving in the direction of the steepest (for each fixed $t \in [0, T]$) ascent in the Hamiltonian H . Consequently, a stochastic gradient descent algorithm to learn the weights is as follows: (k+1)

$$A_t^{(k)}$$

$$(k)$$

$$= A_t^{(k)}$$

$$(k)$$

$$\sim p_k(A_t^{(k)} | X_t)$$

$$(k)$$

$$\sim Y_t$$

The derivation of the bound (14) is equally straightforward and appears as part of the supplementary material. Although the result is attractive, the conclusion is somewhat misleading because (as we will demonstrate with examples) even a small amount of regularization can lead to local (but not global) minimum as well as saddle point solutions. Assumption: The following assumption is made throughout the remainder of this paper: (i) Property P1: The matrix R has no eigenvalues on $\mathbb{R}^+ := \{x \in \mathbb{R} : x \geq 0\}$. The matrix R is non-derogatory. That is, no eigenvalue of R appears in more than one Jordan block. For the scalar ($d = 1$) case, this property means R is strictly positive. For the scalar case, the $\text{RT} \int_0^t A dt$ fundamental solution is given by the closed form formula $\Phi(t, 0) = e^{At}$. Thus, the positivity of R is seen to be necessary to obtain a meaningful solution. For the vector case, this property represents a sufficient condition such that $\log(R)$ can be defined as a real-valued matrix. That is, under property (P1), there exists a (not necessarily unique) matrix $\log(R) \in \mathbb{M}_d(\mathbb{R})$ whose matrix exponential $e^{\log(R)} = R$; cf., Culver [1966], Higham [2014]. The logarithm is trivially a minimum for the $\|\cdot\|_F = 0$ problem. Indeed, $\text{At} \rightarrow T \log(R)$ gives $\log(R)$

$X_t = e^{-\int_0^t A ds} X_0$ and thus $X_T = e^{-\int_0^T A ds} X_0 = R X_0$. This shows A can be made arbitrarily small by choosing a large enough depth T of the network. An analogous result for the linear residual MNN appears in [Hardt and Ma, 2016, Thm. 2.1]. The question then is whether the constant solution $A \rightarrow T \log(R)$ is also obtained as a critical point for the $\|\cdot\|_F = 0$ problem? The following proposition provides a complete characterization of the critical points (for the general $\|\cdot\|_F$ problem) in terms of the solutions of a matrix-valued characteristic equation: Proposition 2. The general solution of the Hamilton's equations (7)-(9) is given by

$$X_t = e^{\int_0^t A ds} X_0 \quad Y_t = e^{-\int_0^t A ds} Y_0$$

$$Z_t = e^{\int_0^t A ds} Z_0$$

$$A_t = e^{\int_0^t A ds} A_0 e^{-\int_0^t A ds}$$

$$Z_t = e^{\int_0^t A ds} Z_0$$

$$e^{\int_0^t A ds}$$

$$(T \rightarrow t)C$$

$$C e^{\int_0^t A ds}$$

$$\Phi(t, 0)$$

$$2$$

$$(15) \quad \Phi(t, 0)$$

$$e^{\int_0^t A ds}$$

$$(Z \rightarrow XT)$$

$$(16) \quad (17)$$

Under Property (P1), $\log(R)$ is uniquely defined if and only if all the eigenvalues of R are positive. When not unique there are countably many matrix logarithms, all denoted as $\log(R)$. The principal logarithm of R is the unique such matrix whose eigenvalues lie in the strip $\{z \in \mathbb{C} : -\pi < \text{Im}(z) \leq \pi\}$.

$$5$$

where $C \in \mathbb{M}_d(\mathbb{R})$ is an arbitrary solution of the characteristic equation $\Phi(C) = F \log(R) \Phi(C)$

Remark 4. For the scalar case $\log(\cdot)$ is a single-valued function. Therefore, $A^T C = T \log(R)$ is the unique critical point (minimizer) for the $[\gamma = 0+]$ problem. While the $[\gamma = 0+]$ problem admits a unique minimizer, the $[\gamma = 0]$ problem does not. In fact, any A^T of the form $A^T = T \log(R) + A^T t R T$ where $0 \leq A^T t dt = 0$ is also a minimizer of the $[\gamma = 0]$ problem. So, while there are infinitely many minimizers of the $[\gamma = 0]$ problem, only one of these survives with even a small amount of regularization. A global characterization of critical points as a function of parameters $(\gamma, R, \gamma_0, T) \in \mathbb{R}_+ \times \mathbb{R}_+ \times \mathbb{R}_+ \times \mathbb{R}_+$ is possible and appears as part of the supplementary material. \square Example 1 (Normal matrix case). Consider the characteristic equation (18) with $R = I$ (rotation in the plane by $\gamma/2$), $\gamma_0 = I$ and $T = 1$. For $\gamma = 0$, the normal solutions of the characteristic equation are given by the multi-valued matrix logarithm function: \square Example 1 $\log(R) = (\gamma/2 + 2\pi n) =: C(0; n)$, $n = 0, \pm 1, \pm 2, \dots$. It is easy to verify that $eC(0; n) = R$. $C(0; 0)$ is referred to as the principal logarithm. The software package PyDSTool Clewley et al. [2007] is used to numerically continue the solution $C(\gamma; n)$ as a function of the parameter γ . Fig. 1(a) depicts the solution branches in terms of the $(2, 1)$ entry of the matrix $C(\gamma; n)$ for $n = 0, \pm 1, \pm 2$. The following observations are made concerning these solutions: \square (i) For each fixed $n \neq 0$, there exist a range $(0, \gamma_n)$, there is a qualitative change minimum and a saddle point. At the limit (turning) point $\gamma = \gamma_n$ in the solution from a minimum to a saddle point. \square (ii) As a function of n , γ_n decreases monotonically as $|n|$ increases. For $\gamma \geq \gamma_n$ (iii) Along the branch with a fixed $n \neq 0$, as $\gamma \rightarrow 0$, the saddle point solution “escapes to infinity.” That is as $\gamma \rightarrow 0$, the saddle point solution $C(\gamma; n) \rightarrow (\gamma/2 + (2n \mp 1)\pi)$. The $1 \mp$ associated cost $J[A] \rightarrow 1$ (The cost of global minimizer $J = 0$). (iv) Among the numerically obtained solution branches, the principal branch $C(\gamma; 0)$ has the lowest cost. Fig. 1 (b) depicts the cost for the solutions depicted in Fig. 1 (a). The numerical calculations indicate that while the $[\gamma = 0]$ problem has infinitely many critical points (all global minimizers), only a finitely many critical points persist for any finite positive value of γ . Moreover, there exists both local (but not global) minimum as well as saddle points for this case. Among the solutions computed, the principal branch (continued from the principal logarithm $C(0; 0)$) has the minimum cost. Example 2 (Non-normal \square matrix case). Numerical continuation is used to obtain solutions for non-normal $R = e^{A^T t}$, where γ is a continuation parameter and $T = 1$. Fig. 2(a) depicts a solution branch as a function of parameter γ . The solution is initialized with the normal solution $C(0; 0)$ described in Example 1. By varying γ , the solution is continued to $\gamma = \gamma/2$ (indicated as in Fig. 2(a)). This way, the solution $C = \gamma$ is found for $R = e^{A^T t}$. It is easy to verify that $C(0; 1/2)$ is a solution of the characteristic equation (18) for $\gamma = 0$ and $T = 1$. For this solution, the critical point $\sin(\gamma t) \cos(\gamma t)$ point of the optimal control problem $A^T = \sin(\gamma t) \cos(\gamma t) + \gamma \cos(\gamma t) + \gamma \sin(\gamma t)$

" # ? ?? tan ? ?? sec ? principal logarithm $\log(R) =$, where $? = \sin?1$. The regularization 4 ? sec ? ? tan ? cost for the non-constant solution A_t is strictly smaller than the constant $T1 \log(R)$ solution: $Z 1 Z 1 Z 1 ?2 \downarrow 3.76 = \text{tr}(\log(R) \log(R) \downarrow) dt \text{tr}(A_t A_t \downarrow) dt = \text{tr}(CC \downarrow) dt = t 4 0 0 0$ Next, the parameter $? = ?2$ is fixed, and the solution continued in the parameter $?$. Fig. 2(b) depicts the cost $J[A]$ for the resulting solution branch of critical points (minimum). The cost with the constant $T1 \log(R)$ is also depicted. It is noted that the latter is not a critical point of the optimal control problem for any positive value of $?$. 8

4

Conclusions and directions for future work

In this paper, we studied the optimization problem of learning the weights of a linear neural network with mean-squared loss function. In order to do so, we introduced a novel formulation: (i) The linear network is modeled as a continuous time (surrogate for layer) optimal control problem; (ii) A weight decay type regularization is considered where the interest is in the limit as the regularization parameter $? \rightarrow 0$ (the limit is referred to as the $[? = 0+]$ problem). The Maximum Principle of optimal control theory is used to derive the Hamilton's equations for the critical points. A remarkable result of our paper is that the critical point solutions of the infinite-dimensional problem are completely characterized via the solutions of a finite-dimensional characteristic equation (Eq. (18)). That such a reduction is possible is unexpected because the weight update equation is nonlinear (even in the settings of linear networks). Based on the analysis of the characteristic equation, several conclusions are obtained3 : (i) It has been noted in literature that, for linear networks, all critical points are global minimum. While this is also true here for the $[? = 0]$ and the $[? = 0+]$ problems, even a small amount of regularization alters the picture, e.g., saddle points emerge (Example 1). (ii) The critical points of the regularized $[? = 0+]$ problem is qualitatively very different compared to the non-regularized $[? = 0]$ problem (Remark 4). Several quantitative results on the critical points of the regularized problem are described in Theorem 2 and Examples 1 and 2. (iii) The study of the characteristic equation revealed an unexpected qualitative difference in the critical points between the two cases where $R := E[ZX0;]$ is a normal or non-normal matrix. In the latter (generic) case, the network weights are necessarily non-constant (Prop. 2). We believe that the ideas and tools introduced in this paper will be useful for the researchers working on the analysis of deep learning. In particular, the paper is expected to highlight and spur work on implicit regularization. Some directions for future work are briefly noted next: (i) Non-normal solutions of the characteristic equation: Analysis of the non-normal solutions of the characteristic equation remains an open problem. The non-normal solutions are important because of the following empirical observation (summarized as part of the supplementary material): In numerical experiments with learning, the weights can get stuck at non-normal critical points before eventually converging to a 'good' minimum. (ii) Generalization error: With a finite number of samples $(X0i, Z i)_{i=1}^N$, the characteristic equation (N)

$$?C = F \downarrow (R ? F) ?0$$

$+F(\lambda)Q(N)PNPN(\lambda)$ where $\lambda_0 := \frac{1}{N} \sum_{i=1}^N X_{0i} X_{0i}^T$ and $Q(N) := \frac{1}{N} \sum_{i=1}^N X_{0i} X_{0i}^T$. Sensitivity analysis of the (N) solution of the characteristic equation, with respect to variations in λ_0 and $Q(N)$, can shed light on the generalization error for different critical points. (iii) Second order analysis: The paper does not contain second order analysis of the critical points λ to determine whether they are local minimum or saddle points. Based on certain preliminary results for the scalar case, it is conjectured that the second order analysis is possible in terms of the first order variation for the characteristic equation. 3

Qualitative aspects of some of the conclusions may be obvious to experts in Deep Learning. The objective here is to obtain quantitative characterization in the (relatively tractable) setting of linear networks.

9

2 References

P. F. Baldi and K. Hornik. Neural networks and principal component analysis: Learning from examples without local minima. *Neural networks*, 2(1):53–58, 1989. P. F. Baldi and K. Hornik. Learning in linear neural networks: A survey. *IEEE Transactions on neural networks*, 6(4):837–858, 1995. S. Bhojanapalli, B. Neyshabur, and N. Srebro. Global optimality of local search for low rank matrix recovery. In *Advances in Neural Information Processing Systems*, pages 3873–3881, 2016. A. Choromanska, M. Henaff, M. Mathieu, G. B. Arous, and Y. LeCun. The loss surfaces of multilayer networks. In *AISTATS*, 2015a. A. Choromanska, Y. LeCun, and G. B. Arous. Open problem: The landscape of the loss surfaces of multilayer networks. In *COLT*, pages 1756–1760, 2015b. R. Clewley, W. E. Sherwood, M. D. LaMar, and J. Guckenheimer. Pydstool, a software environment for dynamical systems modeling, 2007. URL <http://pydstool.sourceforge.net>. W. J. Culver. On the existence and uniqueness of the real logarithm of a matrix. *Proceedings of the American Mathematical Society*, 17(5):1146–1151, 1966. Y. N. Dauphin, R. Pascanu, C. Gulcehre, K. Cho, S. Ganguli, and Y. Bengio. Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. In *Advances in neural information processing systems*, pages 2933–2941, 2014. O. Farotimi, A. Dembo, and T. Kailath. A general weight matrix formulation using optimal control. *IEEE Transactions on neural networks*, 2(3):378–394, 1991. R. Ge, F. Huang, C. Jin, and Y. Yuan. Escaping From Saddle Points ? Online Stochastic Gradient for Tensor Decomposition. *arXiv:1503.02101*, March 2015. I. Goodfellow, Y. Bengio, and A. Courville. *Deep learning*. MIT press, 2016. S. Gunasekar, B. Woodworth, S. Bhojanapalli, B. Neyshabur, and N. Srebro. Implicit regularization in matrix factorization. *arXiv preprint arXiv:1705.09280*, 2017. M. Hardt and T. Ma. Identity matters in deep learning. *arXiv:1611.04231*, November 2016. N. J. Higham. *Functions of matrices*. CRC Press, 2014. K. Kawaguchi. Deep learning without poor local minima. In *Advances In Neural Information Processing Systems*, pages 586–594, 2016. Y. LeCun, D. Touresky, G. Hinton, and T. Sejnowski. A theoretical framework for back-propagation. In

The Connectionist Models Summer School, volume 1, pages 21–28, 1988. J. D. Lee, M. Simchowitz, M. I. Jordan, and B. Recht. Gradient Descent Converges to Minimizers. arXiv:1602.04915, February 2016. B. Neyshabur, R. Tomioka, and N. Srebro. In search of the real inductive bias: On the role of implicit regularization in deep learning. arXiv preprint arXiv:1412.6614, 2014. 10

A. M. Saxe, J. L. McClelland, and S. Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. arXiv:1312.6120, December 2013. D. Soudry and Y. Carmon. No bad local minima: Data independent training error guarantees for multilayer neural networks. arXiv:1605.08361, May 2016. W. Su, S. Boyd, and E. Candes. A differential equation for modeling nesterov’s accelerated gradient method: Theory and insights. In Advances in Neural Information Processing Systems, pages 2510–2518, 2014. A. Wibisono, A. Wilson, and M. Jordan. A variational perspective on accelerated methods in optimization. Proceedings of the National Academy of Sciences, page 201614734, 2016. C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals. Understanding deep learning requires rethinking generalization. arXiv preprint arXiv:1611.03530, 2016.