

---

# On the Optimization of Deep Networks: Implicit Acceleration by Overparameterization

---

Sanjeev Arora<sup>1 2</sup> Nadav Cohen<sup>2</sup> Elad Hazan<sup>1 3</sup>

## Abstract

Conventional wisdom in deep learning states that increasing depth improves expressiveness but complicates optimization. This paper suggests that, sometimes, increasing depth can speed up optimization. The effect of depth on optimization is decoupled from expressiveness by focusing on settings where additional layers amount to overparameterization – linear neural networks, a well-studied model. Theoretical analysis, as well as experiments, show that here depth acts as a preconditioner which may accelerate convergence. Even on simple convex problems such as linear regression with  $\ell_p$  loss,  $p > 2$ , gradient descent can benefit from transitioning to a non-convex overparameterized objective, more than it would from some common acceleration schemes. We also prove that it is mathematically impossible to obtain the acceleration effect of overparameterization via gradients of any regularizer.

## 1. Introduction

How does depth help? This central question of deep learning still eludes full theoretical understanding. The general consensus is that there is a trade-off: increasing depth improves expressiveness, but complicates optimization. Superior expressiveness of deeper networks, long suspected, is now confirmed by theory, albeit for fairly limited learning problems (Eldan & Shamir, 2015; Raghu et al., 2016; Lee et al., 2017; Cohen et al., 2017; Daniely, 2017; Arora et al., 2018). Difficulties in optimizing deeper networks have also been long clear – the signal held by a gradient gets buried as it propagates through many layers. This is known as the “vanishing/exploding gradient problem”. Modern techniques such as batch normalization (Ioffe & Szegedy, 2015) and residual connections (He et al., 2015) have somewhat alleviated these difficulties in practice.

---

<sup>1</sup>Department of Computer Science, Princeton University, Princeton, NJ, USA <sup>2</sup>School of Mathematics, Institute for Advanced Study, Princeton, NJ, USA <sup>3</sup>Google Brain, USA. Correspondence to: Nadav Cohen <cohennadav@ias.edu>.

Given the longstanding consensus on expressiveness vs. optimization trade-offs, this paper conveys a rather counter-intuitive message: increasing depth can *accelerate* optimization. The effect is shown, via first-cut theoretical and empirical analyses, to resemble a combination of two well-known tools in the field of optimization: *momentum*, which led to provable acceleration bounds (Nesterov, 1983); and *adaptive regularization*, a more recent technique proven to accelerate by Duchi et al. (2011) in their proposal of the AdaGrad algorithm. Explicit mergers of both techniques are quite popular in deep learning (Kingma & Ba, 2014; Tieleman & Hinton, 2012). It is thus intriguing that merely introducing depth, with no other modification, can have a similar effect, but *implicitly*.

There is an obvious hurdle in isolating the effect of depth on optimization: if increasing depth leads to faster training on a given dataset, how can one tell whether the improvement arose from a true acceleration phenomenon, or simply due to better representational power (the shallower network was unable to attain the same training loss)? We respond to this hurdle by focusing on *linear neural networks* (cf. Saxe et al. (2013); Goodfellow et al. (2016); Hardt & Ma (2016); Kawaguchi (2016)). With these models, adding layers does not alter expressiveness; it manifests itself only in the replacement of a matrix parameter by a product of matrices – an *overparameterization*.

We provide a new analysis of linear neural network optimization via direct treatment of the differential equations associated with gradient descent when training arbitrarily deep networks on arbitrary loss functions. We find that the overparameterization introduced by depth leads gradient descent to operate as if it were training a shallow (single layer) network, while employing a particular preconditioning scheme. The preconditioning promotes movement along directions already taken by the optimization, and can be seen as an acceleration procedure that combines momentum with adaptive learning rates. Even on simple convex problems such as linear regression with  $\ell_p$  loss,  $p > 2$ , overparameterization via depth can significantly speed up training. Surprisingly, in some of our experiments, not only did overparameterization outperform naïve gradient descent, but it was also faster than two well-known acceleration methods – AdaGrad (Duchi et al., 2011) and AdaDelta (Zeiler, 2012).

In addition to purely linear networks, we also demonstrate (empirically) the implicit acceleration of overparameterization on a non-linear model, by replacing hidden layers with depth-2 linear networks. The implicit acceleration of overparameterization is different from standard regularization – we prove its effect cannot be attained via gradients of *any* fixed regularizer.

Both our theoretical analysis and our empirical evaluation indicate that acceleration via overparameterization need not be computationally expensive. From an optimization perspective, overparameterizing using wide or narrow networks has the same effect – it is only the depth that matters.

## 2. Related Work

Theoretical study of optimization in deep learning is a highly active area of research. Works along this line typically analyze critical points (local minima, saddles) in the landscape of the training objective, either for linear networks (see for example Kawaguchi (2016); Hardt & Ma (2016) or Baldi & Hornik (1989) for a classic account), or for specific non-linear networks under different restrictive assumptions (cf. Choromanska et al. (2015); Haefele & Vidal (2015); Soudry & Carmon (2016); Safran & Shamir (2017)). Other works characterize other aspects of objective landscapes, for example Safran & Shamir (2016) showed that under certain conditions a monotonically descending path from initialization to global optimum exists (in compliance with the empirical observations of Goodfellow et al. (2014)).

The dynamics of optimization was studied in Fukumizu (1998) and Saxe et al. (2013), for linear networks. Like ours, these works analyze gradient descent through its corresponding differential equations. Fukumizu (1998) focuses on linear regression with  $\ell_2$  loss, and does not consider the effect of varying depth – only a two (single hidden) layer network is analyzed. Saxe et al. (2013) also focuses on  $\ell_2$  regression, but considers any depth beyond two (inclusive), ultimately concluding that increasing depth can *slow down* optimization, albeit by a modest amount. In contrast to these two works, our analysis applies to a general loss function, and any depth including one. Intriguingly, we find that for  $\ell_p$  regression, acceleration by depth is revealed only when  $p > 2$ . This explains why the conclusion reached in Saxe et al. (2013) differs from ours.

Turning to general optimization, accelerated gradient (momentum) methods were introduced in Nesterov (1983), and later studied in numerous works (see Wibisono et al. (2016) for a short review). Such methods effectively accumulate gradients throughout the entire optimization path, using the collected history to determine the step at a current point in time. Use of preconditioners to speed up optimization is also a well-known technique. Indeed, the classic Newton’s method can be seen as preconditioning based on second

derivatives. Adaptive preconditioning with only first-order (gradient) information was popularized by the BFGS algorithm and its variants (cf. Nocedal (1980)). Relevant theoretical guarantees, in the context of regret minimization, were given in Hazan et al. (2007); Duchi et al. (2011). In terms of combining momentum and adaptive preconditioning, Adam (Kingma & Ba, 2014) is a popular approach, particularly for optimization of deep networks.

## 3. Warmup: $\ell_p$ Regression

We begin with a simple yet striking example of the effect being studied. For linear regression with  $\ell_p$  loss, we will see how even the slightest overparameterization can have an immense effect on optimization. Specifically, we will see that simple gradient descent on an objective overparameterized by a single scalar, corresponds to a form of accelerated gradient descent on the original objective.

Consider the objective for a scalar linear regression problem with  $\ell_p$  loss ( $p$  – even positive integer):

$$L(\mathbf{w}) = \mathbb{E}_{(\mathbf{x}, y) \sim S} \left[ \frac{1}{p} (\mathbf{x}^\top \mathbf{w} - y)^p \right]$$

$\mathbf{x} \in \mathbb{R}^d$  here are instances,  $y \in \mathbb{R}$  are continuous labels,  $S$  is a finite collection of labeled instances (training set), and  $\mathbf{w} \in \mathbb{R}^d$  is a learned parameter vector. Suppose now that we apply a simple overparameterization, replacing the parameter vector  $\mathbf{w}$  by a vector  $\mathbf{w}_1 \in \mathbb{R}^d$  times a scalar  $w_2 \in \mathbb{R}$ :

$$L(\mathbf{w}_1, w_2) = \mathbb{E}_{(\mathbf{x}, y) \sim S} \left[ \frac{1}{p} (\mathbf{x}^\top \mathbf{w}_1 w_2 - y)^p \right]$$

Obviously the overparameterization does not affect the expressiveness of the linear model. How does it affect optimization? What happens to gradient descent on this non-convex objective?

**Observation 1.** *Gradient descent over  $L(\mathbf{w}_1, w_2)$ , with fixed small learning rate and near-zero initialization, is equivalent to gradient descent over  $L(\mathbf{w})$  with particular adaptive learning rate and momentum terms.*

To see this, consider the gradients of  $L(\mathbf{w})$  and  $L(\mathbf{w}_1, w_2)$ :

$$\begin{aligned} \nabla_{\mathbf{w}} &:= \mathbb{E}_{(\mathbf{x}, y) \sim S} \left[ (\mathbf{x}^\top \mathbf{w} - y)^{p-1} \mathbf{x} \right] \\ \nabla_{\mathbf{w}_1} &:= \mathbb{E}_{(\mathbf{x}, y) \sim S} \left[ (\mathbf{x}^\top \mathbf{w}_1 w_2 - y)^{p-1} w_2 \mathbf{x} \right] \\ \nabla_{w_2} &:= \mathbb{E}_{(\mathbf{x}, y) \sim S} \left[ (\mathbf{x}^\top \mathbf{w}_1 w_2 - y)^{p-1} \mathbf{w}_1^\top \mathbf{x} \right] \end{aligned}$$

Gradient descent over  $L(\mathbf{w}_1, w_2)$  with learning rate  $\eta > 0$ :

$$\mathbf{w}_1^{(t+1)} \leftarrow \mathbf{w}_1^{(t)} - \eta \nabla_{\mathbf{w}_1}^{(t)}, \quad w_2^{(t+1)} \leftarrow w_2^{(t)} - \eta \nabla_{w_2}^{(t)}$$

The dynamics of the underlying parameter  $\mathbf{w} = \mathbf{w}_1 w_2$  are:

$$\begin{aligned} \mathbf{w}^{(t+1)} &= \mathbf{w}_1^{(t+1)} w_2^{(t+1)} \\ &\leftarrow (\mathbf{w}_1^{(t)} - \eta \nabla_{\mathbf{w}_1}^{(t)}) (w_2^{(t)} - \eta \nabla_{w_2}^{(t)}) \\ &= \mathbf{w}_1^{(t)} w_2^{(t)} - \eta w_2^{(t)} \nabla_{\mathbf{w}_1}^{(t)} - \eta \nabla_{w_2}^{(t)} \mathbf{w}_1^{(t)} + \mathcal{O}(\eta^2) \\ &= \mathbf{w}^{(t)} - \eta (w_2^{(t)})^2 \nabla_{\mathbf{w}^{(t)}} - \eta (w_2^{(t)})^{-1} \nabla_{w_2^{(t)}} \mathbf{w}^{(t)} + \mathcal{O}(\eta^2) \end{aligned}$$

$\eta$  is assumed to be small, thus we neglect  $\mathcal{O}(\eta^2)$ . Denoting  $\rho^{(t)} := \eta (w_2^{(t)})^2 \in \mathbb{R}$  and  $\gamma^{(t)} := \eta (w_2^{(t)})^{-1} \nabla_{w_2^{(t)}} \in \mathbb{R}$ , this gives:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \rho^{(t)} \nabla_{\mathbf{w}^{(t)}} - \gamma^{(t)} \mathbf{w}^{(t)}$$

Since by assumption  $\mathbf{w}_1$  and  $w_2$  are initialized near zero,  $\mathbf{w}$  will initialize near zero as well. This implies that at every iteration  $t$ ,  $\mathbf{w}^{(t)}$  is a weighted combination of past gradients. There thus exist  $\mu^{(t,\tau)} \in \mathbb{R}$  such that:

$$\mathbf{w}^{(t+1)} \leftarrow \mathbf{w}^{(t)} - \rho^{(t)} \nabla_{\mathbf{w}^{(t)}} - \sum_{\tau=1}^{t-1} \mu^{(t,\tau)} \nabla_{\mathbf{w}^{(\tau)}}$$

We conclude that the dynamics governing the underlying parameter  $\mathbf{w}$  correspond to gradient descent with a momentum term, where both the learning rate ( $\rho^{(t)}$ ) and momentum coefficients ( $\mu^{(t,\tau)}$ ) are time-varying and adaptive.

#### 4. Linear Neural Networks

Let  $\mathcal{X} := \mathbb{R}^d$  be a space of objects (*e.g.* images or word embeddings) that we would like to infer something about, and let  $\mathcal{Y} := \mathbb{R}^k$  be the space of possible inferences. Suppose we are given a training set  $\{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^m \subset \mathcal{X} \times \mathcal{Y}$ , along with a (point-wise) loss function  $l : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}_{\geq 0}$ . For example,  $\mathbf{y}^{(i)}$  could hold continuous values with  $l(\cdot)$  being the  $\ell_2$  loss:  $l(\hat{\mathbf{y}}, \mathbf{y}) = \frac{1}{2} \|\hat{\mathbf{y}} - \mathbf{y}\|_2^2$ ; or it could hold one-hot vectors representing categories with  $l(\cdot)$  being the softmax-cross-entropy loss:  $l(\hat{\mathbf{y}}, \mathbf{y}) = -\sum_{r=1}^k y_r \log(e^{\hat{y}_r} / \sum_{r'=1}^k e^{\hat{y}_{r'}})$ , where  $y_r$  and  $\hat{y}_r$  stand for coordinate  $r$  of  $\mathbf{y}$  and  $\hat{\mathbf{y}}$  respectively. For a predictor  $\phi$ , *i.e.* a mapping from  $\mathcal{X}$  to  $\mathcal{Y}$ , the overall training loss is  $L(\phi) := \frac{1}{m} \sum_{i=1}^m l(\phi(\mathbf{x}^{(i)}), \mathbf{y}^{(i)})$ . If  $\phi$  comes from some parametric family  $\Phi := \{\phi_\theta : \mathcal{X} \rightarrow \mathcal{Y} | \theta \in \Theta\}$ , we view the corresponding training loss as a function of the parameters, *i.e.* we consider  $L^\Phi(\theta) := \frac{1}{m} \sum_{i=1}^m l(\phi_\theta(\mathbf{x}^{(i)}), \mathbf{y}^{(i)})$ . For example, if the parametric family in question is the class of (directly parameterized) linear predictors:

$$\Phi^{lin} := \{\mathbf{x} \mapsto W\mathbf{x} | W \in \mathbb{R}^{k,d}\} \quad (1)$$

the respective training loss is a function from  $\mathbb{R}^{k,d}$  to  $\mathbb{R}_{\geq 0}$ .

In our context, a depth- $N$  ( $N \geq 2$ ) linear neural network, with hidden widths  $n_1, n_2, \dots, n_{N-1} \in \mathbb{N}$ , is the following parametric family of linear predictors:  $\Phi^{n_1 \dots n_{N-1}} := \{\mathbf{x} \mapsto W_N W_{N-1} \dots W_1 \mathbf{x} | W_j \in \mathbb{R}^{n_j, n_{j-1}}, j=1 \dots N\}$ , where by definition  $n_0 := d$  and  $n_N := k$ . As customary, we refer to each  $W_j, j=1 \dots N$ , as the weight matrix of layer  $j$ . For simplicity of presentation, we hereinafter omit from our notation the hidden widths  $n_1 \dots n_{N-1}$ , and simply write  $\Phi^N$  instead of  $\Phi^{n_1 \dots n_{N-1}}$  ( $n_1 \dots n_{N-1}$  will be specified explicitly if not clear by context). That is, we denote:

$$\Phi^N := \quad (2)$$

$$\{\mathbf{x} \mapsto W_N W_{N-1} \dots W_1 \mathbf{x} | W_j \in \mathbb{R}^{n_j, n_{j-1}}, j=1 \dots N\}$$

For completeness, we regard a depth-1 network as the family of directly parameterized linear predictors, *i.e.* we set  $\Phi^1 := \Phi^{lin}$  (see Equation 1).

The training loss that corresponds to a depth- $N$  linear network –  $L^{\Phi^N}(W_1, \dots, W_N)$ , is a function from  $\mathbb{R}^{n_1, n_0} \times \dots \times \mathbb{R}^{n_N, n_{N-1}}$  to  $\mathbb{R}_{\geq 0}$ . For brevity, we will denote this function by  $L^N(\cdot)$ . Our focus lies on the behavior

of gradient descent when minimizing  $L^N(\cdot)$ . More specifically, we are interested in the dependence of this behavior on  $N$ , and in particular, in the possibility of increasing  $N$  leading to acceleration. Notice that for any  $N \geq 2$  we have:

$$L^N(W_1, \dots, W_N) = L^1(W_N W_{N-1} \dots W_1) \quad (3)$$

and so the sole difference between the training loss of a depth- $N$  network and that of a depth-1 network (classic linear model) lies in the replacement of a matrix parameter by a product of  $N$  matrices. This implies that if increasing  $N$  can indeed accelerate convergence, it is not an outcome of any phenomenon other than favorable properties of depth-induced overparameterization for optimization.

#### 5. Implicit Dynamics of Gradient Descent

In this section we present a new result for linear neural networks, tying the dynamics of gradient descent on  $L^N(\cdot)$  – the training loss corresponding to a depth- $N$  network, to those on  $L^1(\cdot)$  – training loss of a depth-1 network (classic linear model). Specifically, we show that gradient descent on  $L^N(\cdot)$ , a complicated and seemingly pointless overparameterization, can be directly rewritten as a particular preconditioning scheme over gradient descent on  $L^1(\cdot)$ .

When applied to  $L^N(\cdot)$ , gradient descent takes on the following form:

$$W_j^{(t+1)} \leftarrow (1 - \eta\lambda)W_j^{(t)} - \eta \frac{\partial L^N}{\partial W_j}(W_1^{(t)}, \dots, W_N^{(t)}) \quad (4)$$

$$, j = 1 \dots N$$

$\eta > 0$  here is a learning rate, and  $\lambda \geq 0$  is an optional weight decay coefficient. For simplicity, we regard both  $\eta$  and  $\lambda$  as fixed (no dependence on  $t$ ). Define the underlying *end-to-end weight matrix*:

$$W_e := W_N W_{N-1} \dots W_1 \quad (5)$$

Given that  $L^N(W_1, \dots, W_N) = L^1(W_e)$  (Equation 3), we view  $W_e$  as an optimized weight matrix for  $L^1(\cdot)$ , whose dynamics are governed by Equation 4. Our interest then boils down to the study of these dynamics for different choices of  $N$ . For  $N = 1$  they are (trivially) equivalent to standard gradient descent over  $L^1(\cdot)$ . We will characterize the dynamics for  $N \geq 2$ .

To be able to derive, in our general setting, an explicit update rule for the end-to-end weight matrix  $W_e$  (Equation 5), we introduce an assumption by which the learning rate is small, *i.e.*  $\eta^2 \approx 0$ . Formally, this amounts to translating Equation 4 to the following set of differential equations:

$$\dot{W}_j(t) = -\eta\lambda W_j(t) - \eta \frac{\partial L^N}{\partial W_j}(W_1(t), \dots, W_N(t)) \quad (6)$$

$$, j = 1 \dots N$$

where  $t$  is now a continuous time index, and  $\dot{W}_j(t)$  stands for the derivative of  $W_j$  with respect to time. The use of differential equations, for both theoretical analysis and algorithm design, has a long and rich history in optimization

research (see Helmke & Moore (2012) for an overview). When step sizes (learning rates) are taken to be small, trajectories of discrete optimization algorithms converge to smooth curves modeled by continuous-time differential equations, paving way to the well-established theory of the latter (cf. Boyce et al. (1969)). This approach has led to numerous interesting findings, including recent results in the context of acceleration methods (e.g. Su et al. (2014); Wibisono et al. (2016)).

With the continuous formulation in place, we turn to express the dynamics of the end-to-end matrix  $W_e$ :

**Theorem 1.** *Assume the weight matrices  $W_1 \dots W_N$  follow the dynamics of continuous gradient descent (Equation 6). Assume also that their initial values (time  $t_0$ ) satisfy, for  $j = 1 \dots N - 1$ :*

$$W_{j+1}^\top(t_0)W_{j+1}(t_0) = W_j(t_0)W_j^\top(t_0) \quad (7)$$

Then, the end-to-end weight matrix  $W_e$  (Equation 5) is governed by the following differential equation:

$$\begin{aligned} \dot{W}_e(t) = & -\eta\lambda N \cdot W_e(t) \\ & -\eta \sum_{j=1}^N [W_e(t)W_e^\top(t)]^{\frac{j-1}{N}} \cdot \\ & \frac{dL^1}{dW}(W_e(t)) \cdot [W_e^\top(t)W_e(t)]^{\frac{N-j}{N}} \end{aligned} \quad (8)$$

where  $[\cdot]^{\frac{j-1}{N}}$  and  $[\cdot]^{\frac{N-j}{N}}$ ,  $j = 1 \dots N$ , are fractional power operators defined over positive semidefinite matrices.

*Proof.* (sketch – full details in Appendix A.1) If  $\lambda = 0$  (no weight decay) then one can easily show that  $W_{j+1}^\top(t)\dot{W}_{j+1}(t) = \dot{W}_j(t)W_j^\top(t)$  throughout optimization. Taking the transpose of this equation and adding to itself, followed by integration over time, imply that the difference between  $W_{j+1}^\top(t)W_{j+1}(t)$  and  $W_j(t)W_j^\top(t)$  is constant. This difference is zero at initialization (Equation 7), thus will remain zero throughout, i.e.:

$$W_{j+1}^\top(t)W_{j+1}(t) = W_j(t)W_j^\top(t) \quad , \quad \forall t \geq t_0 \quad (9)$$

A slightly more delicate treatment shows that this is true even if  $\lambda > 0$ , i.e. with weight decay included.

Equation 9 implies alignment of the (left and right) singular spaces of  $W_j(t)$  and  $W_{j+1}(t)$ , simplifying the product  $W_{j+1}(t)W_j(t)$ . Successive application of this simplification allows a clean computation for the product of all layers (that is,  $W_e$ ), leading to the explicit form presented in theorem statement (Equation 8).  $\square$

Translating the continuous dynamics of Equation 8 back to discrete time, we obtain the sought-after update rule for the end-to-end weight matrix:

$$\begin{aligned} W_e^{(t+1)} \leftarrow & (1 - \eta\lambda N)W_e^{(t)} \\ & -\eta \sum_{j=1}^N [W_e^{(t)}(W_e^{(t)})^\top]^{\frac{j-1}{N}} \cdot \\ & \frac{dL^1}{dW}(W_e^{(t)}) \cdot [(W_e^{(t)})^\top W_e^{(t)}]^{\frac{N-j}{N}} \end{aligned} \quad (10)$$

This update rule relies on two assumptions: first, that the

learning rate  $\eta$  is small enough for discrete updates to approximate continuous ones; and second, that weights are initialized on par with Equation 7, which will approximately be the case if initialization values are close enough to zero. It is customary in deep learning for both learning rate and weight initializations to be small, but nonetheless above assumptions are only met to a certain extent. We support their applicability by showing empirically (Section 8) that the end-to-end update rule (Equation 10) indeed provides an accurate description for the dynamics of  $W_e$ .

A close look at Equation 10 reveals that the dynamics of the end-to-end weight matrix  $W_e$  are similar to gradient descent over  $L^1(\cdot)$  – training loss corresponding to a depth-1 network (classic linear model). The only difference (besides the scaling by  $N$  of the weight decay coefficient  $\lambda$ ) is that the gradient  $\frac{dL^1}{dW}(W_e)$  is subject to a transformation before being used. Namely, for  $j = 1 \dots N$ , it is multiplied from the left by  $[W_e W_e^\top]^{\frac{j-1}{N}}$  and from the right by  $[W_e^\top W_e]^{\frac{N-j}{N}}$ , followed by summation over  $j$ . Clearly, when  $N = 1$  (depth-1 network) this transformation reduces to identity, and as expected,  $W_e$  precisely adheres to gradient descent over  $L^1(\cdot)$ . When  $N \geq 2$  the dynamics of  $W_e$  are less interpretable. We arrange it as a vector to gain more insight:

**Claim 1.** *For an arbitrary matrix  $A$ , denote by  $\text{vec}(A)$  its arrangement as a vector in column-first order. Then, the end-to-end update rule in Equation 10 can be written as:*

$$\begin{aligned} \text{vec}(W_e^{(t+1)}) \leftarrow & (1 - \eta\lambda N) \cdot \text{vec}(W_e^{(t)}) \\ & -\eta \cdot P_{W_e^{(t)}} \text{vec}\left(\frac{dL^1}{dW}(W_e^{(t)})\right) \end{aligned} \quad (11)$$

where  $P_{W_e^{(t)}}$  is a positive semidefinite preconditioning matrix that depends on  $W_e^{(t)}$ . Namely, if we denote the singular values of  $W_e^{(t)} \in \mathbb{R}^{k,d}$  by  $\sigma_1 \dots \sigma_{\max\{k,d\}} \in \mathbb{R}_{\geq 0}$  (by definition  $\sigma_r = 0$  if  $r > \min\{k,d\}$ ), and corresponding left and right singular vectors by  $\mathbf{u}_1 \dots \mathbf{u}_k \in \mathbb{R}^k$  and  $\mathbf{v}_1 \dots \mathbf{v}_d \in \mathbb{R}^d$  respectively, the eigenvectors of  $P_{W_e^{(t)}}$  are:

$$\text{vec}(\mathbf{u}_r \mathbf{v}_{r'}^\top) \quad , \quad r = 1 \dots k \quad , \quad r' = 1 \dots d$$

with corresponding eigenvalues:

$$\sum_{j=1}^N \sigma_r^{2\frac{N-j}{N}} \sigma_{r'}^{2\frac{j-1}{N}} \quad , \quad r = 1 \dots k \quad , \quad r' = 1 \dots d$$

*Proof.* The result readily follows from the properties of the Kronecker product – see Appendix A.2 for details.  $\square$

Claim 1 implies that in the end-to-end update rule of Equation 10, the transformation applied to the gradient  $\frac{dL^1}{dW}(W_e)$  is essentially a preconditioning, whose eigendirections and eigenvalues depend on the singular value decomposition of  $W_e$ . The eigendirections are the rank-1 matrices  $\mathbf{u}_r \mathbf{v}_{r'}^\top$ , where  $\mathbf{u}_r$  and  $\mathbf{v}_{r'}$  are left and right (respectively) singular vectors of  $W_e$ . The eigenvalue of  $\mathbf{u}_r \mathbf{v}_{r'}^\top$  is  $\sum_{j=1}^N \sigma_r^{2(N-j)/N} \sigma_{r'}^{2(j-1)/N}$ , where  $\sigma_r$  and  $\sigma_{r'}$  are the singular values of  $W_e$  corresponding to  $\mathbf{u}_r$  and  $\mathbf{v}_{r'}$  (respectively). When  $N \geq 2$ , an increase in  $\sigma_r$  or  $\sigma_{r'}$  leads to



an increase in the eigenvalue corresponding to the eigendirection  $\mathbf{u}_r \mathbf{v}_r^\top$ . Qualitatively, this implies that the preconditioning favors directions that correspond to singular vectors whose presence in  $W_e$  is stronger. We conclude that the effect of overparameterization, *i.e.* of replacing a classic linear model (depth-1 network) by a depth- $N$  linear network, boils down to modifying gradient descent by promoting movement along directions that fall in line with the current location in parameter space. A-priori, such a preference may seem peculiar – why should an optimization algorithm be sensitive to its location in parameter space? Indeed, we generally expect sensible algorithms to be translation invariant, *i.e.* be oblivious to parameter value. However, if one takes into account the common practice in deep learning of initializing weights near zero, the location in parameter space can also be regarded as the overall movement made by the algorithm. We thus interpret our findings as indicating that overparameterization promotes movement along directions already taken by the optimization, and therefore can be seen as a form of acceleration. This intuitive interpretation will become more concrete in the subsection that follows.

A final point to make, is that the end-to-end update rule (Equation 10 or 11), which obviously depends on  $N$  – number of layers in the deep linear network, does *not* depend on the hidden widths  $n_1 \dots n_{N-1}$  (see Section 4). This implies that from an optimization perspective, overparameterizing using wide or narrow networks has the same effect – it is only the depth that matters. Consequently, the acceleration of overparameterization can be attained at a minimal computational price, as we demonstrate empirically in Section 8.

### 5.1. Single Output Case

To facilitate a straightforward presentation of our findings, we hereinafter focus on the special case where the optimized models have a single output, *i.e.* where  $k = 1$ . This corresponds, for example, to a binary (two-class) classification problem, or to the prediction of a numeric scalar property (regression). It admits a particularly simple form for the end-to-end update rule of Equation 10:

**Claim 2.** Assume  $k = 1$ , *i.e.*  $W_e \in \mathbb{R}^{1,d}$ . Then, the end-to-end update rule in Equation 10 can be written as follows:

$$W_e^{(t+1)} \leftarrow (1 - \eta\lambda N) \cdot W_e^{(t)} - \eta \|W_e^{(t)}\|_2^{2-\frac{2}{N}} \cdot \left( \frac{dL^1}{dW}(W_e^{(t)}) + (N-1) \cdot Pr_{W_e^{(t)}} \left\{ \frac{dL^1}{dW}(W_e^{(t)}) \right\} \right) \quad (12)$$

where  $\|\cdot\|_2^{2-\frac{2}{N}}$  stands for Euclidean norm raised to the power of  $2 - \frac{2}{N}$ , and  $Pr_W\{\cdot\}$ ,  $W \in \mathbb{R}^{1,d}$ , is defined to be the projection operator onto the direction of  $W$ :

$$Pr_W : \mathbb{R}^{1,d} \rightarrow \mathbb{R}^{1,d} \quad (13)$$

$$Pr_W\{V\} := \begin{cases} \frac{W}{\|W\|_2} V^\top \cdot \frac{W}{\|W\|_2} & , W \neq 0 \\ 0 & , W = 0 \end{cases}$$

*Proof.* The result follows from the definition of a fractional power operator over matrices – see Appendix A.3.  $\square$

Claim 2 implies that in the single output case, the effect of overparameterization (replacing classic linear model by depth- $N$  linear network) on gradient descent is twofold: first, it leads to an *adaptive learning rate* schedule, by introducing the multiplicative factor  $\|W_e\|_2^{2-2/N}$ ; and second, it amplifies (by  $N$ ) the projection of the gradient on the direction of  $W_e$ . Recall that we view  $W_e$  not only as the optimized parameter, but also as the overall movement made in optimization (initialization is assumed to be near zero). Accordingly, the adaptive learning rate schedule can be seen as gaining confidence (increasing step sizes) when optimization moves farther away from initialization, and the gradient projection amplification can be thought of as a certain type of *momentum* that favors movement along the azimuth taken so far. These effects bear potential to accelerate convergence, as we illustrate qualitatively in Section 7, and demonstrate empirically in Section 8.

## 6. Overparameterization Effects Cannot Be Attained via Regularization

Adding a regularizer to the objective is a standard approach for improving optimization (though lately the term regularization is typically associated with generalization). For example, AdaGrad was originally invented to compete with the best regularizer from a particular family. The next theorem shows (for single output case) that the effects of overparameterization cannot be attained by adding a regularization term to the original training loss, or via any similar modification. This is not obvious a-priori, as unlike many acceleration methods that explicitly maintain memory of past gradients, updates under overparameterization are by definition the gradients of *something*. The assumptions in the theorem are minimal and also necessary, as one must rule-out the trivial counter-example of a constant training loss.

**Theorem 2.** Assume  $\frac{dL^1}{dW}$  does not vanish at  $W = 0$ , and is continuous on some neighborhood around this point. For a given  $N \in \mathbb{N}$ ,  $N > 2$ ,<sup>1</sup> define:

$$F(W) := \quad (14)$$

$$\|W\|_2^{2-\frac{2}{N}} \cdot \left( \frac{dL^1}{dW}(W) + (N-1) \cdot Pr_W \left\{ \frac{dL^1}{dW}(W) \right\} \right)$$

where  $Pr_W\{\cdot\}$  is the projection given in Equation 13. Then, there exists no function (of  $W$ ) whose gradient field is  $F(\cdot)$ .

*Proof.* (sketch – full details in Appendix A.4) The proof uses the fundamental theorem for line integrals, which states that the integral of  $\nabla g$  for any differentiable function  $g$  amounts to 0 along every closed curve.

Overparameterization changes gradient descent’s behavior: instead of following the original gradient  $\frac{dL^1}{dW}$ , it follows some other direction  $F(\cdot)$  (see Equations 12 and 14) that

<sup>1</sup> For the result to hold with  $N = 2$ , additional assumptions on  $L^1(\cdot)$  are required; otherwise any non-zero linear function  $L^1(W) = WU^\top$  serves as a counter-example – it leads to a vector field  $F(\cdot)$  that is the gradient of  $W \mapsto \|W\|_2 \cdot WU^\top$ .

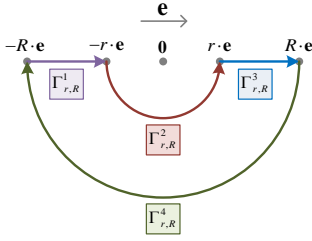


Figure 1. Curve  $\Gamma_{r,R}$  over which line integral is non-zero.

is a *function* of the original gradient as well as the current point  $W$ . We think of this change as a transformation that maps one *vector field*  $\phi(\cdot)$  to another  $-F_\phi(\cdot)$ :

$$F_\phi(W) = \begin{cases} \left\| \|W\|^{2-\frac{2}{N}} \left( \phi(W) + (N-1) \left\langle \phi(W), \frac{W}{\|W\|} \right\rangle \frac{W}{\|W\|} \right) \right\|, & W \neq 0 \\ 0, & W = 0 \end{cases}$$

Notice that for  $\phi = \frac{dL^1}{dW}$ , we get exactly the vector field  $F(\cdot)$  defined in theorem statement. The mapping  $\phi \mapsto F_\phi$  is linear. Moreover, because of the linearity of line integrals, for any curve  $\Gamma$ , the functional  $\phi \mapsto \int_\Gamma F_\phi$  – a mapping of vector fields to scalars, is linear as well.

We show that  $F(\cdot)$  contradicts the fundamental theorem for line integrals. To do so, we construct a closed curve  $\Gamma = \Gamma_{r,R}$  for which the linear functional  $\phi \mapsto \oint_\Gamma F_\phi$  does not vanish at  $\phi = \frac{dL^1}{dW}$ . Let  $\mathbf{e} := \frac{dL^1}{dW}(W=0) / \left\| \frac{dL^1}{dW}(W=0) \right\|$ , which is well-defined since by assumption  $\frac{dL^1}{dW}(W=0) \neq 0$ . For  $r < R$  we define  $\Gamma_{r,R} := \Gamma_{r,R}^1 \rightarrow \Gamma_{r,R}^2 \rightarrow \Gamma_{r,R}^3 \rightarrow \Gamma_{r,R}^4$  as illustrated in Figure 1. With the definition of  $\Gamma_{r,R}$  in place, we decompose  $\frac{dL^1}{dW}$  into a constant vector field  $\kappa \equiv \frac{dL^1}{dW}(W=0)$  plus a residual  $\xi$ . We explicitly compute the line integrals along  $\Gamma_{r,R}^1 \dots \Gamma_{r,R}^4$  for  $F_\kappa$ , and derive bounds for  $F_\xi$ . This, along with the linearity of the functional  $\phi \mapsto \int_\Gamma F_\phi$ , provides a lower bound on the line integral of  $F(\cdot)$  over  $\Gamma_{r,R}$ . We show the lower bound is positive as  $r, R \rightarrow 0$ , thus  $F(\cdot)$  indeed contradicts the fundamental theorem for line integrals.  $\square$

## 7. Illustration of Acceleration

To this end, we showed that overparameterization (use of depth- $N$  linear network in place of classic linear model) induces on gradient descent a particular preconditioning scheme (Equation 10 in general and 12 in the single output case), which can be interpreted as introducing some forms of momentum and adaptive learning rate. We now illustrate qualitatively, on a very simple hypothetical learning problem, the potential of these to accelerate optimization.

Consider the task of linear regression, assigning to vectors in  $\mathbb{R}^2$  labels in  $\mathbb{R}$ . Suppose that our training set consists of two points in  $\mathbb{R}^2 \times \mathbb{R}$ :  $([1, 0]^\top, y_1)$  and  $([0, 1]^\top, y_2)$ . Assume also that the loss function of interest is  $\ell_p$ ,  $p \in 2\mathbb{N}$ :  $\ell_p(\hat{y}, y) = \frac{1}{p}(\hat{y} - y)^p$ . Denoting the learned parameter by  $\mathbf{w} = [w_1, w_2]^\top$ , the overall training loss can be written as:<sup>2</sup>

<sup>2</sup> We omit the averaging constant  $\frac{1}{2}$  for conciseness.

$$L(w_1, w_2) = \frac{1}{p}(w_1 - y_1)^p + \frac{1}{p}(w_2 - y_2)^p$$

With fixed learning rate  $\eta > 0$  (weight decay omitted for simplicity), gradient descent over  $L(\cdot)$  gives:

$$w_i^{(t+1)} \leftarrow w_i^{(t)} - \eta(w_i^{(t)} - y_i)^{p-1}, \quad i = 1, 2$$

Changing variables per  $\Delta_i = w_i - y_i$ , we have:

$$\Delta_i^{(t+1)} \leftarrow \Delta_i^{(t)} (1 - \eta(\Delta_i^{(t)})^{p-2}), \quad i = 1, 2 \quad (15)$$

Assuming the original weights  $w_1$  and  $w_2$  are initialized near zero,  $\Delta_1$  and  $\Delta_2$  start off at  $-y_1$  and  $-y_2$  respectively, and will eventually reach the optimum  $\Delta_1^* = \Delta_2^* = 0$  if the learning rate is small enough to prevent divergence:

$$\eta < \frac{2}{y_i^{p-2}}, \quad i = 1, 2$$

Suppose now that the problem is ill-conditioned, in the sense that  $y_1 \gg y_2$ . If  $p = 2$  this has no effect on the bound for  $\eta$ .<sup>3</sup> If  $p > 2$  the learning rate is determined by  $y_1$ , leading  $\Delta_2$  to converge very slowly. In a sense,  $\Delta_2$  will suffer from the fact that there is no ‘‘communication’’ between the coordinates (this will actually be the case not just with gradient descent, but with most algorithms typically used in large-scale settings – AdaGrad, Adam, *etc.*).

Now consider the scenario where we optimize  $L(\cdot)$  via overparameterization, *i.e.* with the update rule in Equation 12 (single output). In this case the coordinates are coupled, and as  $\Delta_1$  gets small ( $w_1$  gets close to  $y_1$ ), the learning rate is effectively scaled by  $y_1^{2-\frac{2}{N}}$  (in addition to a scaling by  $N$  in coordinate 1 only), allowing (if  $y_1 > 1$ ) faster convergence of  $\Delta_2$ . We thus have the luxury of temporarily slowing down  $\Delta_2$  to ensure that  $\Delta_1$  does not diverge, with the latter speeding up the former as it reaches safe grounds. In Appendix B we consider a special case and formalize this intuition, deriving a concrete bound for the acceleration of overparameterization.

## 8. Experiments

Our analysis (Section 5) suggests that overparameterization – replacement of a classic linear model by a deep linear network, induces on gradient descent a certain preconditioning scheme. We qualitatively argued (Section 7) that in some cases, this preconditioning may accelerate convergence. In this section we put these claims to the test, through a series of empirical evaluations based on TensorFlow toolbox (Abadi et al. (2016)). For conciseness, many of the details behind our implementation are deferred to Appendix C.

We begin by evaluating our analytically-derived preconditioning scheme – the end-to-end update rule in Equation 10. Our objective in this experiment is to ensure that our analysis, continuous in nature and based on a particular assumption on weight initialization (Equation 7), is indeed applicable to practical scenarios. We focus on the single output

<sup>3</sup> Optimal learning rate for gradient descent on quadratic objective does not depend on current parameter value (*cf.* Goh (2017)).

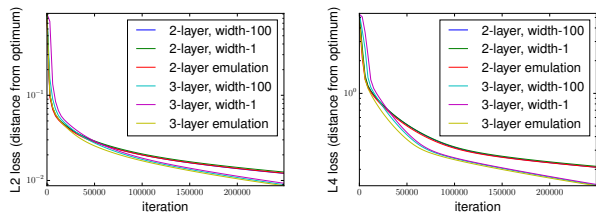


Figure 2. (to be viewed in color) Gradient descent optimization of deep linear networks (depths 2, 3) vs. the analytically-derived preconditioning schemes (over single layer model; Equation 12). Both plots show training objective (left –  $\ell_2$  loss; right –  $\ell_4$  loss) per iteration, on a numeric regression dataset from UCI Machine Learning Repository (details in text). Notice the emulation of preconditioning schemes. Notice also the negligible effect of network width – for a given depth, setting size of hidden layers to 1 (scalars) or 100 yielded similar convergence (on par with our analysis).

case, where the update-rule takes on a particularly simple (and efficiently implementable) form – Equation 12. The dataset chosen was UCI Machine Learning Repository’s “Gas Sensor Array Drift at Different Concentrations” (Ver-gara et al., 2012; Rodriguez-Lujan et al., 2014). Specifically, we used the dataset’s “Ethanol” problem – a scalar regression task with 2565 examples, each comprising 128 features (one of the largest numeric regression tasks in the repository). As training objectives, we tried both  $\ell_2$  and  $\ell_4$  losses. Figure 2 shows convergence (training objective per iteration) of gradient descent optimizing depth-2 and depth-3 linear networks, against optimization of a single layer model using the respective preconditioning schemes (Equation 12 with  $N = 2, 3$ ). As can be seen, the preconditioning schemes reliably emulate deep network optimization, suggesting that, at least in some cases, our analysis indeed captures practical dynamics.

Alongside the validity of the end-to-end update rule, Figure 2 also demonstrates the negligible effect of network width on convergence, in accordance with our analysis (see Section 5). Specifically, it shows that in the evaluated setting, hidden layers of size 1 (scalars) suffice in order for the essence of overparameterization to fully emerge. Unless otherwise indicated, all results reported hereinafter are based on this configuration, *i.e.* on scalar hidden layers. The computational toll associated with overparameterization will thus be virtually non-existent.

As a final observation on Figure 2, notice that it exhibits faster convergence with a deeper network. This however does not serve as evidence in favor of acceleration by depth, as we did not set learning rates optimally per model (simply used the common choice of  $10^{-3}$ ). To conduct a fair comparison between the networks, and more importantly, between them and a classic single layer model, multiple learning rates were tried, and the one giving fastest convergence was taken on a per-model basis. Figure 3 shows the results of this experiment. As can be seen, convergence of deeper

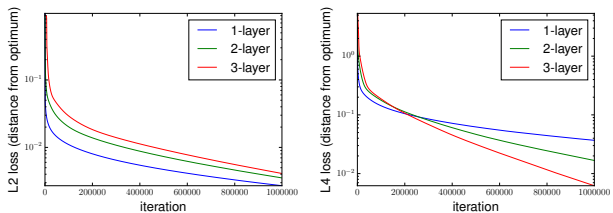


Figure 3. (to be viewed in color) Gradient descent optimization of single layer model vs. linear networks of depth 2 and 3. Setup is identical to that of Figure 2, except that here learning rates were chosen via grid search, individually per model (see Appendix C). Notice that with  $\ell_2$  loss, depth (slightly) hinders optimization, whereas with  $\ell_4$  loss it leads to significant acceleration (on par with our qualitative analysis in Section 7).

networks is (slightly) slower in the case of  $\ell_2$  loss. This falls in line with the findings of Saxe et al. (2013). In stark contrast, and on par with our qualitative analysis in Section 7, is the fact that with  $\ell_4$  loss adding depth significantly accelerated convergence. To the best of our knowledge, this provides first empirical evidence to the fact that depth, even without any gain in expressiveness, and despite introducing non-convexity to a formerly convex problem, can lead to favorable optimization.

In light of the speedup observed with  $\ell_4$  loss, it is natural to ask how the implicit acceleration of depth compares against explicit methods for acceleration and adaptive learning. Figure 4-left shows convergence of a depth-3 network (optimized with gradient descent) against that of a single layer model optimized with AdaGrad (Duchi et al., 2011) and AdaDelta (Zeiler, 2012). The displayed curves correspond to optimal learning rates, chosen individually via grid search. Quite surprisingly, we find that in this specific setting, overparameterizing, thereby turning a convex problem non-convex, is a more effective optimization strategy than carefully designed algorithms tailored for convex problems. We note that this was not observed with all algorithms – for example Adam (Kingma & Ba, 2014) was considerably faster than overparameterization. However, when introducing overparameterization simultaneously with Adam (a setting we did not theoretically analyze), further acceleration is attained – see Figure 4-right. This suggests that at least in some cases, not only plain gradient descent benefits from depth, but also more elaborate algorithms commonly employed in state of the art applications.

An immediate question arises at this point. If depth indeed accelerates convergence, why not add as many layers as one can computationally afford? The reason, which is actually apparent in our analysis, is the so-called *vanishing gradient problem*. When training a very deep network (large  $N$ ), while initializing weights to be small, the end-to-end matrix  $W_e$  (Equation 5) is extremely close to zero, severely attenuating gradients in the preconditioning scheme (Equation 10). A possible approach for alleviating this issue is to initialize weights to be larger, yet small enough such that the

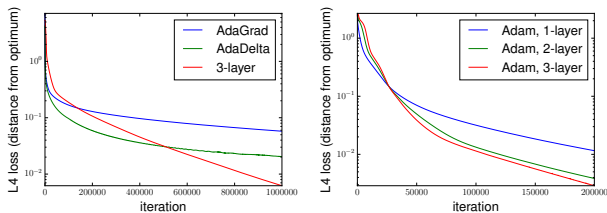


Figure 4. (to be viewed in color) **Left:** Gradient descent optimization of depth-3 linear network vs. AdaGrad and AdaDelta over single layer model. Setup is identical to that of Figure 3-right. Notice that the implicit acceleration of overparameterization outperforms both AdaGrad and AdaDelta (former is actually slower than plain gradient descent). **Right:** Adam optimization of single layer model vs. Adam over linear networks of depth 2 and 3. Same setup, but with learning rates set per Adam’s default in TensorFlow. Notice that depth improves speed, suggesting that the acceleration of overparameterization may be somewhat orthogonal to explicit acceleration methods.

end-to-end matrix does not “explode”. The choice of identity (or near identity) initialization leads to what is known as *linear residual networks* (Hardt & Ma, 2016), akin to the successful residual networks architecture (He et al., 2015) commonly employed in deep learning. Notice that identity initialization satisfies the condition in Equation 7, rendering the end-to-end update rule (Equation 10) applicable. Figure 5-left shows convergence, under gradient descent, of a single layer model against deeper networks than those evaluated before – depths 4 and 8. As can be seen, with standard, near-zero initialization, the depth-4 network starts making visible progress only after about 65K iterations, whereas the depth-8 network seems stuck even after 100K iterations. In contrast, under identity initialization, both networks immediately make progress, and again depth serves as an implicit accelerator.

As a final sanity test, we evaluate the effect of overparameterization on optimization in a non-idealized (yet simple) deep learning setting. Specifically, we experiment with the convolutional network tutorial for MNIST built into TensorFlow,<sup>4</sup> which includes convolution, pooling and dense layers, ReLU non-linearities, stochastic gradient descent with momentum, and dropout (Srivastava et al., 2014). We introduced overparameterization by simply placing two matrices in succession instead of the matrix in each dense layer. Here, as opposed to previous experiments, widths of the newly formed hidden layers were not set to 1, but rather to the minimal values that do not deteriorate expressiveness (see Appendix C). Overall, with an addition of roughly 15% in number of parameters, optimization has accelerated considerably – see Figure 5-right. The displayed results were obtained with the hyperparameter settings hardcoded into the tutorial. We have tried alternative settings (varying learning rates and standard deviations of initializations – see

<sup>4</sup> <https://github.com/tensorflow/models/tree/master/tutorials/image/mnist>

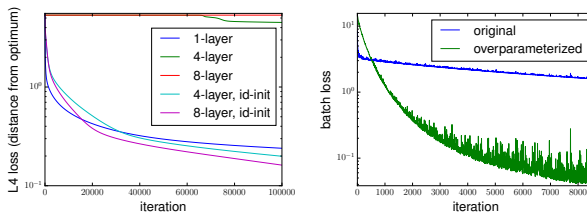


Figure 5. (to be viewed in color) **Left:** Gradient descent optimization of single layer model vs. linear networks deeper than before (depths 4, 8). For deep networks, both near-zero and near-identity initializations were evaluated. Setup identical to that of Figure 3-right. Notice that deep networks suffer from vanishing gradients under near-zero initialization, while near-identity (“residual”) initialization eliminates the problem. **Right:** Stochastic gradient descent optimization in TensorFlow’s convolutional network tutorial for MNIST. Plot shows batch loss per iteration, in original setting vs. overparameterized one (depth-2 linear networks in place of dense layers).

Appendix C), and in all cases observed an outcome similar to that in Figure 5-right – overparameterization led to significant speedup. Nevertheless, as reported above for linear networks, it is likely that for non-linear networks the effect of depth on optimization is mixed – some settings accelerate by it, while others do not. Comprehensive characterization of the cases in which depth accelerates optimization warrants much further study. We hope our work will spur interest in this avenue of research.

### 9. Conclusion

Through theory and experiments, we demonstrated that overparameterizing a neural network by increasing its depth can accelerate optimization, even on very simple problems.

Our analysis of linear neural networks, the subject of various recent studies, yielded a new result: for these models, overparameterization by depth can be understood as a preconditioning scheme with a closed form description (Theorem 1 and the claims thereafter). The preconditioning may be interpreted as a combination between certain forms of adaptive learning rate and momentum. Given that it depends on network depth but not on width, acceleration by overparameterization can be attained at a minimal computational price, as we demonstrate empirically in Section 8.

Clearly, complete theoretical analysis for non-linear networks will be challenging. Empirically however, we showed that the trivial idea of replacing an internal weight matrix by a product of two can significantly accelerate optimization, with absolutely no effect on expressiveness (Figure 5-right).

The fact that gradient descent over classic convex problems such as linear regression with  $\ell_p$  loss,  $p > 2$ , can accelerate from transitioning to a non-convex overparameterized objective, does not coincide with conventional wisdom, and provides food for thought. Can this effect be rigorously quantified, similarly to analyses of explicit acceleration methods such as momentum or adaptive regularization (AdaGrad)?



## Acknowledgments

Sanjeev Arora’s work is supported by NSF, ONR, Simons Foundation, Schmidt Foundation, Mozilla Research, Amazon Research, DARPA and SRC. Elad Hazan’s work is supported by NSF grant 1523815 and Google Brain. Nadav Cohen is a member of the Zuckerman Israeli Postdoctoral Scholars Program, and is supported by Eric and Wendy Schmidt.

## References

- Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., et al. Tensorflow: A system for large-scale machine learning. In *OSDI*, volume 16, pp. 265–283, 2016.
- Arora, R., Basu, A., Mianjy, P., and Mukherjee, A. Understanding deep neural networks with rectified linear units. *International Conference on Learning Representations (ICLR)*, 2018.
- Baldi, P. and Hornik, K. Neural networks and principal component analysis: Learning from examples without local minima. *Neural networks*, 2(1):53–58, 1989.
- Boyce, W. E., DiPrima, R. C., and Haines, C. W. *Elementary differential equations and boundary value problems*, volume 9. Wiley New York, 1969.
- Buck, R. C. *Advanced calculus*. Waveland Press, 2003.
- Choromanska, A., Henaff, M., Mathieu, M., Arous, G. B., and LeCun, Y. The loss surfaces of multilayer networks. In *Artificial Intelligence and Statistics*, pp. 192–204, 2015.
- Cohen, N., Sharir, O., Levine, Y., Tamari, R., Yakira, D., and Shashua, A. Analysis and design of convolutional networks via hierarchical tensor decompositions. *arXiv preprint arXiv:1705.02302*, 2017.
- Daniely, A. Depth separation for neural networks. *arXiv preprint arXiv:1702.08489*, 2017.
- Duchi, J., Hazan, E., and Singer, Y. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research*, 12(Jul):2121–2159, 2011.
- Eldan, R. and Shamir, O. The power of depth for feedforward neural networks. *arXiv preprint arXiv:1512.03965*, 2015.
- Fukumizu, K. Effect of batch learning in multilayer neural networks. *Gen*, 1(04):1E–03, 1998.
- Goh, G. Why momentum really works. *Distill*, 2017. doi: 10.23915/distill.00006. URL <http://distill.pub/2017/momentum>.
- Goodfellow, I., Bengio, Y., Courville, A., and Bengio, Y. *Deep learning*, volume 1. MIT press Cambridge, 2016.
- Goodfellow, I. J., Vinyals, O., and Saxe, A. M. Qualitatively characterizing neural network optimization problems. *arXiv preprint arXiv:1412.6544*, 2014.
- Haefele, B. D. and Vidal, R. Global Optimality in Tensor Factorization, Deep Learning, and Beyond. *CoRR abs/1202.2745*, cs.NA, 2015.
- Hardt, M. and Ma, T. Identity matters in deep learning. *arXiv preprint arXiv:1611.04231*, 2016.
- Hazan, E., Agarwal, A., and Kale, S. Logarithmic regret algorithms for online convex optimization. *Mach. Learn.*, 69(2-3):169–192, December 2007. ISSN 0885-6125.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. *arXiv preprint arXiv:1512.03385*, 2015.
- Helmke, U. and Moore, J. B. *Optimization and dynamical systems*. Springer Science & Business Media, 2012.
- Ioffe, S. and Szegedy, C. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pp. 448–456, 2015.
- Jones, E., Oliphant, T., Peterson, P., et al. SciPy: Open source scientific tools for Python, 2001–. URL <http://www.scipy.org/>. [Online; accessed ;today].
- Kawaguchi, K. Deep learning without poor local minima. In *Advances in Neural Information Processing Systems*, pp. 586–594, 2016.
- Kingma, D. and Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- Lee, H., Ge, R., Risteski, A., Ma, T., and Arora, S. On the ability of neural nets to express distributions. *arXiv preprint arXiv:1702.07028*, 2017.
- Nesterov, Y. A method of solving a convex programming problem with convergence rate  $o(1/k^2)$ . In *Soviet Mathematics Doklady*, volume 27, pp. 372–376, 1983.
- Nocedal, J. Updating quasi-newton matrices with limited storage. *Mathematics of Computation*, 35(151):773–782, 1980.
- Raghu, M., Poole, B., Kleinberg, J., Ganguli, S., and Sohl-Dickstein, J. On the expressive power of deep neural networks. *arXiv preprint arXiv:1606.05336*, 2016.
- Rodriguez-Lujan, I., Fonollosa, J., Vergara, A., Homer, M., and Huerta, R. On the calibration of sensor arrays for pattern recognition using the minimal number of experiments. *Chemometrics and Intelligent Laboratory Systems*, 130:123–134, 2014.
- Safran, I. and Shamir, O. On the quality of the initial basin in overspecified neural networks. In *International Conference on Machine Learning*, pp. 774–782, 2016.
- Safran, I. and Shamir, O. Spurious local minima are common in two-layer relu neural networks. *arXiv preprint arXiv:1712.08968*, 2017.
- Saxe, A. M., McClelland, J. L., and Ganguli, S. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *arXiv preprint arXiv:1312.6120*, 2013.
- Soudry, D. and Carmon, Y. No bad local minima: Data independent training error guarantees for multilayer neural networks. *arXiv preprint arXiv:1605.08361*, 2016.
- Srivastava, N., Hinton, G. E., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(1):1929–1958, 2014.

- Su, W., Boyd, S., and Candes, E. A differential equation for modeling nesterovs accelerated gradient method: Theory and insights. In *Advances in Neural Information Processing Systems*, pp. 2510–2518, 2014.
- Tieleman, T. and Hinton, G. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural networks for machine learning*, 4(2):26–31, 2012.
- Vergara, A., Vembu, S., Ayhan, T., Ryan, M. A., Homer, M. L., and Huerta, R. Chemical gas sensor drift compensation using classifier ensembles. *Sensors and Actuators B: Chemical*, 166: 320–329, 2012.
- Wibisono, A., Wilson, A. C., and Jordan, M. I. A variational perspective on accelerated methods in optimization. *Proceedings of the National Academy of Sciences*, 113(47):E7351–E7358, 2016.
- Zeiler, M. D. Adadelata: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012.