

On the interplay of network structure and gradient convergence in deep learning

Vamsi K. Ithapu[†] and Sathya N. Ravi[†] and Vikas Singh^{‡,† *†}

October 6, 2016

Abstract

The regularization and output consistency behavior of dropout and layer-wise pretraining for learning deep networks have been fairly well studied. However, our understanding of how the asymptotic convergence of backpropagation in deep architectures is related to the structural properties of the network and other design choices (like denoising and dropout rate) is less clear at this time. An interesting question one may ask is whether the network architecture and input data statistics may guide the choices of learning parameters and vice versa. In this work, we explore the association between such structural, distributional and learnability aspects vis-à-vis their interaction with parameter convergence rates. We present a framework to address these questions based on convergence of backpropagation for general nonconvex objectives using first-order information. This analysis suggests an interesting relationship between feature denoising and dropout. Building upon these results, we obtain a setup that provides systematic guidance regarding the choice of learning parameters and network sizes that achieve a certain level of convergence (in the optimization sense) often mediated by statistical attributes of the inputs. Our results are supported by a set of experimental evaluations as well as independent empirical observations reported by other groups.

1 Introduction

The successful deployment of deep learning in a broad spectrum of applications including localizing objects in images, analyzing particle accelerator data, converting speech into text, predicting the activity of drug molecules, and designing clinical trials [1, 2, 3, 4, 5] provides compelling evidence that they can learn very complex concepts with fairly minimal feature engineering or preprocessing. This success is attributed to the idea of composing simple but non-linear modules that each transform the lower levels (raw or normalized data) into a representation at a more abstract level [6, 7, 8]. While this high level representation learning procedure is quite general, the problem at hand may at least partly govern the choice of the architecture and require certain modifications to the algorithm. Motivated by various experimental considerations, several variants of deep architectures and corresponding regularization schemes have been developed [9, 10, 11, 12]. Complementary to such algorithmic and empirical developments, there is also a growing recent interest in better understanding the mathematical properties of these algorithms. Over the last few years, several interesting results have been presented [13, 14, 15, 16, 17, 18, 19, 20]. While some of these address the hypothesis spaces and the sets of functions learnable by deep networks, the others analyze the parameter space learned via backpropagation. A more detailed discussion about these works is included in Section 1.1.

Our work is related to the foregoing recent results dealing with a theoretical characterization of the properties of deep networks. At the high level, we study the relationship between the learnability of the network

^{*†}Computer Sciences, University of Wisconsin-Madison

^{†‡}Biostatistics and Medical Informatics, University of Wisconsin-Madison

(convergence of parameter estimation in an optimization sense) and the structure of deep networks. In particular, we are interested in a set of questions that can better elucidate the interplay between certain concepts that intuitively seem related, but for the most part, have been studied separately. For instance, what is the influence of the structural aspects of the architecture (number of layers, activation types) on practical considerations such as dropout rate, stepsizes, and so on? Can we *guide* the choice of some of these parameters based on their relationship to the others using information about the statistical moments of the data? More generally, can we evaluate which network converges faster, for a given learning problem? Or more specifically, what is the best network structure and activation functions for a given learning task? Understanding such an interplay of the *architecture* and *data statistics* vis-à-vis distinct *learning schemes* is necessary to complement and even facilitate the continuing empirical successes of deep networks, especially in regimes where deep networks have not yet been thoroughly tested (like learning with small datasets). Further, the interplay allows for relating the input data statistics to deep network learnability. This may enable constructing nonlinearities beyond convolution and max-pooling (which remain important for computer vision datasets), for instance, to tackle structured data like brain images or genetic data where the sample size is a bottleneck.

Overview: The most commonly used procedure for parameter estimation in deep networks is mini-batch stochastic gradients [21]. This involves deriving the average gradient from the error computed on a few training instances, adjusting the weights/parameters accordingly, and iterating until convergence. Our general goal is to tie this procedure (to the extent possible) to the network structure. In contrast to the other recent results [19], we directly work with stochastic gradients with no strong assumptions on the data/network structure [16], which allows obtaining results that are quite generally applicable. The starting point of our analysis is a recent work by [22] dealing with the convergence of stochastic gradients for arbitrary nonconvex problems using a first-order oracle. We build upon and adapt this analysis by first addressing single-layer networks and unsupervised pretraining, and then, the more general case of multi-layer nets with dropout. More importantly, apart from addressing the interplay of network structure and convergence, the algorithms we present, with minor tweaks are easily deployable to the standard training pipeline. Further, our bounds natively take into account the standard regularization schemes like dropout and layer-wise pretraining, making them potentially more useful in practice.

A brief description of our **contributions** are: **(a)** Based on the idea of randomly stopping gradient updates from [22], we present a framework for analyzing mini-batch stochastic gradients on multi-layer deep networks, and prove gradient convergence of deep networks learned via dropout (with/without layer-wise pretraining). **(b)** Later, we derive explicit relationships between the network structure, learning parameters and input data statistics. Our results corroborate many empirical studies, and may *guide* the choices of network/learning hyper-parameters.

1.1 Related Work

The body of literature addressing backpropagation in neural networks, and deep networks in general, is very large and dates back at least to the early 1970s. Hence, we restrict the discussion only to those works that fall immediately within the context of this paper. While there are a number of seminal papers addressing variants of backpropagation, and stochastic gradients, focusing on the convergence of neural networks training [23, 24], a number of recent works [25, 26, 14, 19, 20] provide a fresh treatment of these problems and analyze efficient learning schemes in the context of deep networks specifically. The solution space of backpropagation in this setting has also been addressed in recent results [27, 13], providing new insights into the types of functions learnable by deep networks [15, 16]. [20] have shown that better generalization can be achieved by ensuring smaller training times, while [14] describes in detail the non-convex landscape of deep learning objectives and the goodness of local optima. Beyond these optimization related works, several authors have independently addressed the regularization and learnability aspects of deep networks.

For example, [28, 11] extensively analyzed dropout, [18] develops a probabilistic theory of deep learning, and [17] studies the existence of deep representations. [29] presents a generative model for ReLU-type deep networks, and very recently, [30] studied the equivalence of arbitrary deep networks. Complimentary to these, [19] uses a tensor decomposition perspective.

2 Preliminaries

2.1 Notation

Let $\mathbf{x} \in \mathbb{R}^{d_x}$ and $\mathbf{y} \in \mathbb{R}^{d_y}$ denote the input feature vector and the corresponding output (or label) respectively. Given multiple $\{\mathbf{x}, \mathbf{y}\} \in \mathcal{X}$, the unknown input-to-output mapping is modeled by a L -layered neural network (L -NN), which comprises of the input (or visible) unit \mathbf{x} , followed by $L - 1$ hidden representations $\mathbf{h}^1, \dots, \mathbf{h}^{L-1}$ and the output (or final) unit \mathbf{y} [7]. The lengths of these $L + 1$ representations are $d_0 = d_x, d_1, \dots, d_{L-1}, d_L = d_y$ respectively. Each layer transforms the representations from the previous layer by first applying an affine transform, followed by a non-linearity (in general, non-convex and not necessarily point-wise) [7, 8]. The layer-wise transformations are denoted by $\mathbf{W}_l \in \mathbb{R}^{d_l \times d_{l-1}}$ for $l = 1, \dots, L$. Using these, the hidden representations are given by $\mathbf{h}^l = \sigma(\mathbf{W}_l, \mathbf{h}^{l-1})$ for $l = 1, \dots, L - 1$ ($\mathbf{h}^0 = \mathbf{x}$), and the output layer is $\mathbf{y} = \sigma(\mathbf{W}_L, \mathbf{h}^{L-1})$, where $\sigma(\cdot)$ represents the non-linear function/mapping between layers. For a 1-NN with no hidden layers, we simply have $\mathbf{y} = \sigma(\mathbf{W}, \mathbf{x})$ where \mathbf{W} 's are the unknowns. Note that the bias in the affine transformation is handled by augmenting features with 1 whenever necessary. We will restrict ourselves to point-wise sigmoid non-linearity i.e., for any $\mathbf{v} \in \mathbb{R}^d$, $\sigma(\mathbf{v}) = \{\frac{1}{1+\exp(-v_i)}\}$ ($i = 1, \dots, d$). The distributional hyper-parameters of interest are $\mu_{\mathbf{x}} = \frac{1}{d_x} \sum_j \mathbb{E} x_j$ and $\tau_{\mathbf{x}} = \frac{1}{d_x} \sum_j \mathbb{E}^2 x_j$, which correspond to the average first moment and average squared first moment of the inputs respectively (the average is across the d_x dimensions). For simplicity we assume $\mathbf{x} \in [0, 1]^{d_x}$ and $\mathbf{y} \in [0, 1]^{d_y}$, and so $\mu_{\mathbf{x}} \in [0, 1]$ and $\tau_{\mathbf{x}} \in [0, 1]$.

Consider the following minimization performed via mini-batch stochastic gradients [21],

$$\min_{\mathbf{W}} f(\mathbf{W}) := \mathbb{E}_{\mathbf{x}, \mathbf{y}} \mathcal{L}(\mathbf{x}, \mathbf{y}; \mathbf{W}) \quad (1)$$

where $\mathcal{L}(\cdot)$ denotes some loss function parameterized by \mathbf{W} and applied to data instances $\{\mathbf{x}, \mathbf{y}\}$. Denote $\eta := \{\mathbf{x}, \mathbf{y}\} \sim \mathcal{X}$. The mini-batch stochastic gradient update using B samples η^1, \dots, η^B is $\mathbf{W} \leftarrow \mathbf{W} - \gamma G(\eta; \mathbf{W})$ where γ is the stepsize. $G(\eta; \mathbf{W})$ computed at \mathbf{W} using the sample set η^1, \dots, η^B is given by

$$G(\eta; \mathbf{W}) = \frac{1}{B} \sum_{i=1}^B \nabla_{\mathbf{W}} \mathcal{L}(\eta^i; \mathbf{W}) \quad (2)$$

Depending on $\mathcal{L}(\cdot)$, (1) corresponds to the backpropagation learning of different classes of neural networks using stochastic gradients. To address many such broad families, we exhaustively develop the analysis for three interesting and general classes of deep networks, starting with single-layer networks, followed by unsupervised pretraining via box-constrained denoising autoencoders [10], and finally multi-layer deep networks with dropout [31]. Due to space restrictions, extensions to the more complex convolutional and recurrent networks including, fancier non-linearities like rectified linear units or max-outs [32], are deferred to the longer version of this paper. For each of these settings, the loss function $\mathcal{L}(\cdot)$ is defined as follows. Due to space constraints we restrict the description of these classes to minimum, and refer the readers to [7, 10, 31], and the references there in, for further and extensive details.

- 1-NN: Single-layer Network

$$\mathcal{L}(\mathbf{x}, \mathbf{y}; \mathbf{W}) = \|\mathbf{y} - \sigma(\mathbf{W}\mathbf{x})\|^2 \quad (3)$$

- **DA: Box-constrained Denoising Autoencoder**

$$\mathcal{L}(\mathbf{x}, \mathbf{y}; \mathbf{W}) = \|\mathbf{x} - \sigma(\mathbf{W}^T \sigma(\mathbf{W}(\mathbf{x} * \mathbf{z})))\|^2 \quad (4)$$

where $\mathbf{W} \in [-w_m, w_m]^{d_h \times d_v}$. $*$ denotes element-wise product and \mathbf{z} is a binary vector of length d_x . Given the denoising rate ζ , the binary scalars $z_1, \dots, z_{d_x} \sim \text{Bernoulli}(\zeta)$ are noise-indicators, such that, whenever $z_i = 0$, the i^{th} element of the input \mathbf{x} , x_i is nullified. We denote $\mathbf{x} * \mathbf{z}$ as $\tilde{\mathbf{x}}$. $[-w_m, w_m]$ is the box-constraint on the unknowns \mathbf{W} , which forces the learned representations to *not* saturate around 0 and/or 1, and has been widely used in variants of both autoencoder design and backpropagation itself [10, 33]. The nature of this constraint is similar in part to other regularization schemes like [26] and [25, 31].

- **L-NN: Multi-layer Network**

$$\begin{aligned} \mathbf{h}^0 &= \mathbf{x}; \quad \mathbf{h}^l = \sigma(\mathbf{W}_l(\mathbf{h}^{l-1} * \mathbf{z}^l)) \\ \mathcal{L}(\mathbf{x}, \mathbf{y}; \mathbf{W}) &= \|\mathbf{y} - \sigma(\mathbf{W}_L(\mathbf{h}^{L-1} * \mathbf{z}^L))\|^2 \end{aligned} \quad (5)$$

where $l = 1, \dots, L-1$. Recall the dropout scheme in training deep networks which address the overfitting problem [31]. In each gradient update iteration, a random fraction of the hidden and/or input units are dropped based on the (given) dropout rates ζ_1, \dots, ζ_L [31]. Similar to DA, the dropped out units for each layer are denoted by a set of binary vectors $\mathbf{z}^1, \dots, \mathbf{z}^L$ such that $z_i^l \sim \text{Bernoulli}(\zeta_l)$ for $l = 1, \dots, L$. Within each iteration, this results in randomly sampling a smaller sub-network with approximately $\prod_{l=1}^L \zeta_l$ fraction of all the transformation parameters. Only these $\prod_{l=1}^L \zeta_l$ fraction of weights are updated in the current iteration, while the rest are not, and the re-sampling and updating process is repeated. For simplicity of analysis we use the same dropout rate, denoted by ζ , for all layers. With some abuse of notation, we use \mathbf{W}_l^k to represent the k^{th} gradient update of the l^{th} layer transformation \mathbf{W}_l .

- A L-NN may be pretrained layer-wise before supervised tuning [7, 10, 33]. Here, the $L-1$ hidden layers are first pretrained, for instance, using the box-constrained DA from (4). This gives the estimates of the $L-1$ transformations $\mathbf{W}_1, \dots, \mathbf{W}_{L-1}$, which (along with \mathbf{y} s) are then used to *initialize* the L-NN. This ‘pretrained’ L-NN is then learned using dropout as in (5). This is the classical deep learning regime [7, 33], and we discuss this case separately from the fully-supervised version which has *no* layer-wise pretraining. Several studies have already shown interesting empirical relationships between dropout and the DA [28], and we complement this body of work by providing explicit relationships between them.

2.2 Roadmap

The gradient update in (2) is central to the ideas described in this paper. By *tracking* the behavior of these gradients as they propagate across multiple layers of the network, with minimal assumptions on the loss function $\mathcal{L}(\eta; \mathbf{W})$, we can assess the convergence of the overall parameter estimation scheme while taking into account the influence (and structure) of each of the layers involved. Since the updates are stochastic, ideally, we are interested in the “expected” gradients over a certain number of iterations, fixed ahead of time. Motivated by this intuition, our main idea adapted from [22], is to randomly sample the number of gradient update iterations. Specifically, let N denote the maximum possible number of iterations that can be performed keeping in mind the memory and time constraints (in general, N is very large). The stopping distribution $\mathbb{P}_R(\cdot)$ gives the probability that k^{th} iteration ($k = 1, \dots, N$) is the *last or stopping* iteration. We denote this randomly sampled stopping iteration as $R \in \{1, \dots, N\}$, and so,

$$\mathbb{P}_R(\cdot) := \frac{p_R^k}{\sum_{k=1}^N p_R^k} \quad \text{where} \quad p_R^k = \text{Pr}(R = k) \quad (6)$$

$\mathbb{P}_R(\cdot)$ can either be fixed ahead of time or learned from a hyper-training procedure. An alternate way to interpret the stopping distribution is by observing that p_R^k represents the probability that the estimates (\mathbf{W} 's) at the k^{th} iteration are desired final solutions returned by the learning procedure.

By proceeding with this *random stopping mini-batch stochastic gradients*, and using some Lipschitz properties of the objective and certain distributional characteristics of the input data, we can analyze the three loss functions (3), (4) and (5). The stopping iteration R is random and the loss function in (1) includes an expectation over data instances $\eta = \{\mathbf{x}, \mathbf{y}\}$. Therefore, we are interested in the *expectation* of the gradients $\nabla_{\mathbf{W}} f(\mathbf{W}^k)$ computed over $R \sim \mathbb{P}_R(\cdot)$ and $\eta \sim \mathcal{X}$. Few of the hyper-parameters that play significant role are, the variance of $G(\eta; \mathbf{W})$ in (2) denoted by e^s , e^{da} and e^m for single-layer, pretraining and multi-layer cases respectively, the distributional hyper-parameters μ_x and τ_x , and the denoising/dropout rate ζ along with the box-constraint. Apart from the convergence bounds, we also derive sample sizes required for large deviation estimates of \mathbf{W} 's which marginalizes the influence of the random stopping iteration. This leads to bounds on training time and the minimum required training dataset size, which in turn can be used to ensure certain level of generalization using existing results [20]. The treatment for the simple case of single-layer networks serves as basis for the more general settings. Due to space constraints, the proofs for all the technical results are included in an adjoining technical report at <http://pages.cs.wisc.edu/~7Evamsi/files/techrep-dlinterplay.pdf>

Why gradient norm? We may ask if characterizing the behavior of the gradients is ideal for the goal of characterizing the interplay. There are more than a few reasons why this strategy is at least sensible. First, note that it is NP-hard to check local minima even for simple non-convex problems [34] — so, an analysis using the norms of the gradients is an attractive alternative, especially, if it leads to a similar main result. Second, a direct way of approaching our central question of the architecture and convergence interplay is by analyzing the gradients themselves. Clearly, learnability of the network would be governed by the goodness of the estimates, which in turn depend on the convergence of stochastic gradients. In most cases, this naturally requires an asymptotic analysis of the norm, which, for nonconvex objectives with minimal assumptions (like Lipschitz continuity) was unavailable until very recently [22]. Third, working with the gradient norm directly allows for assessing faster training times, which, as argued in [20], is vital for better generalization. Alternatively, using the gradient convergence as a surrogate for this training time allows for modulating the network structure to improve generalization.

3 Single-layer Networks

Consider a 1-NN with the corresponding loss function from (3). We learn it via the random stopping mini-batch stochastic gradients (batch size B and $R \sim \mathbb{P}_R(k)$, $k = 1, \dots, N$). This learning procedure is summarized in Alg. 1 in the supplemental technical report (see the link from Section 2); referred from here on as single-layer randomized stochastic gradients, *single-layer RSG*. \mathbf{W}^1 and \mathbf{W}^R are the initial and final estimates respectively, and γ^k is the stepsize at the k^{th} iteration. With no prior information about how to setup the stopping distribution, one would simply choose $\mathbb{P}_R(k) := \text{Unif}[1, N]$. The first result summarizes the decay of the expected gradients in single-layer RSG for this setting. $D_f = f(\mathbf{W}^1) - f^*$ is the initial deviation of the objective from unknown optimum \mathbf{W}^* .

Theorem 3.1 (1-NN, Constant Stepsize). *Consider a single-layer RSG with no dropout and constant stepsize $\gamma^k = \gamma \forall k$. Let $e_\gamma^s = (1 - \frac{13}{16}\gamma)$, $e^s = \frac{13d_x d_y}{256}$ and $R \sim \text{Unif}[1, N]$. The expected gradients are given by*

$$\mathbb{E}_{R,\eta}(\|\nabla_{\mathbf{W}} f(\mathbf{W}^R)\|^2) \leq \frac{1}{e_\gamma^s} \left(\frac{D_f}{N\gamma} + \frac{e^s \gamma}{B} \right) \quad (7)$$

and the optimal constant stepsize is $\gamma_o = \sqrt{\frac{Bf(\mathbf{W}^1)}{e^s N}}$

Remarks: We should point out that the *asymptotic* behavior of gradients in backpropagation, including variants with adaptive stepsizes, momentum etc. have been well studied [24, 23, 27, 13]. This aspect is *not* novel to our work specifically, as also described in Section 1. However, to our knowledge, relatively few ‘explicit’ results about the convergence *rates* are known; although imposing restrictions on the objective does lead to improved guarantees [35]. This may be because of the lack of such results in the numerical optimization literature in the context of general nonconvex objectives. The recent result in [22] presents one of the first such results addressing general nonconvex objectives with a first-order oracle. Consequently, we believe that Theorem 3.1 gives non-trivial information and, as we will show in later sections, leads to interesting new results in the context of neural networks. We now explain (7) briefly.

While the first term corresponds to the goodness of fit of the network (showing the influence of D_f), the second term encodes the network degrees of freedom (e^s), and is larger for big (or fatter) networks. Clearly, the bounds decreases as N (and/or B) increase, and it is interesting to see that the effect of network size ($d_x d_y$) is negligible for large batch sizes. The second term in (7) induces a kind of ‘bias’, that depends on e^s and B . As B increases, (2) will be much closer to a full-batch gradient update, there by reducing this bias (a classical property of mini-batch stochastic gradients [36, 26]). Further, this term indicates that larger batch sizes are necessary for fatter networks with large d_x or d_y . One caveat of the generality of (7) is that it might be loose in the worst case where $D_f \sim 0$ (i.e., when \mathbf{W}^1 is already a good estimate of the stationary point). e_γ^s increases the overall bound as γ increases. The optimal stepsize γ_o in Theorem 3.1 is calculated by balancing the two terms in (7), using which we see that the expected gradients are $\mathcal{O}\left(\frac{1}{\sqrt{N}}\right)$ for a given network and $\mathcal{O}(\sqrt{d_x d_y})$ for a given N .

We see from the proof of Theorem 3.1 that (see technical report), allowing constant γ^k s directly corresponds to choosing uniform $\mathbb{P}_R(\cdot)$ i.e., they are *equivalent* in some sense. Nevertheless, using alternate $\mathbb{P}_R(\cdot)$ in tandem with constant γ^k s results in different sets of constants in (7), while the overall (high-level) dependence on N and network size $d_x d_y$ would remain the same. Beyond constant γ^k s, the most common setting is to use decaying γ^k s. The corresponding result is in Theorem 3.2 with $\gamma^k = \frac{\gamma}{k^\rho}$ (for some $\rho > 0$), and $\mathcal{H}_N(\theta) = \sum_{i=1}^N \frac{1}{i^\theta}$ is the generalized harmonic number.

Corollary 3.2 (1-NN, Decreasing Stepsize). *Consider a single-layer RSG with no dropout and stepsizes $\gamma^k = \frac{\gamma}{k^\rho}$. Let $e^s = \frac{13d_x d_y}{256}$ and the probability of stopping at k^{th} iteration $p_R^k = \gamma^k(1 - \frac{13}{16}\gamma^k)$. The expected gradients are given by*

$$\mathbb{E}_{R,\eta}(\|\nabla_{\mathbf{W}} f(\mathbf{W}^R)\|^2) \leq \frac{16}{3\mathcal{H}_N(\rho)} \left(\frac{D_f}{\gamma} + \frac{e^s \gamma \mathcal{H}_N(2\rho)}{B} \right) \quad (8)$$

Remarks: Whenever $\rho \approx 0$, the stepsizes are approximately constant, and we have $\mathcal{H}_N(\rho) \approx N$ (assuming that $\gamma < 1$). Here the bound in (8) is at least as large as (7) making the two results consistent with each other. The dependence of ρ , and its interaction with N and other hyper-parameters, on the decay of the expected gradients is reasonably complex. Nevertheless, broadly, as ρ increases, the bound first decreases and eventually increases becoming more looser. This trend is expected. To see this, first observe that large ρ implies strong decay of the stepsizes. For sufficiently small γ , this results in stopping iteration probability p_R^k that decreases as k increases (refer to its definition from Theorem 3.2), i.e., large ρ results in $R \ll N$. The bound is implying the trivial fact that the expected gradients are going to be large whenever the gradient updating is stopped early.

Overall, these observations clearly imply that (7) and (8) are capturing all the intuitive trends one would expect to see from using stochastic gradients, in turn, making the analysis and results all the more useful. Lastly, it may not be appropriate to bias R to be far smaller than N . One can relax this and instead use arbitrary $\mathbb{P}_R(\cdot)$ that is monotonically *increasing*, i.e., $p_R^k \leq p_R^{k+1} \forall k$, there by pushing R towards N with high probability. A similar decay of gradients (like Corollary 3.2) result can be obtained for this case, but

we omit it from this shorter version of the paper. Theorem 3.1 and Corollary 3.2 describe convergence of 1-NN for *one run* of stochastic gradient. In practice, one is more interested in a large deviation bound over multiple runs, especially because of the randomization over η and R (see (7)). We define such a large deviation estimate, using $\mathbf{W}^{R_1}, \dots, \mathbf{W}^{R_T}$ computed from $T > 1$ *independent* runs of single-layer RSG, and compute the minimum N required to achieve such an estimate.

Definition 1 ((ϵ, δ) -solution). Given $\epsilon > 0$ and $0 < \delta \ll 1$, an (ϵ, δ) -solution of a single-layer network is given by $\arg \min_t \|\nabla_{\mathbf{W}} f(\mathbf{W}^{R_t})\|^2$ such that $Pr(\min_t \|\nabla_{\mathbf{W}} f(\mathbf{W}^{R_t})\|^2 \leq \epsilon) \geq 1 - \delta$.

Corollary 3.3 (Computational Complexity). To compute a (ϵ, δ) -solution for single-layer RSG with no dropout, optimal constant stepsizes and $\mathbb{P}_R := Unif[1, N]$, we need

$$N(\epsilon, \delta) \geq \frac{4f(\mathbf{W}^1)e^s}{B\delta^{2/T}\epsilon^2}(1 + \bar{\delta})^2, \quad \bar{\delta} = \frac{13B\epsilon\delta^{1/T}}{32e^s} \quad (9)$$

Remarks: To illustrate the practical usefulness of this result, consider an example network with $d_x = 100$, $d_y = 5$ and say $f(\mathbf{W}^1) \sim d_y$. For a $(0.05, 0.05)$ -solution with $T = 10$ runs and a batchsize $B = 50$, (9) gives $N > 7.8 \times 10^3$. If the number of epochs is $C = 200$, then under the (reasonable) assumption that $SC \approx BN$ (where S is sample size), this leads to ~ 2000 instances. Clearly, (9) implies that N increases as the network size ($d_x d_y$) and $f(\mathbf{W}^1)$ increase, and for obtaining good large deviation estimates (i.e., $\epsilon, \delta \approx 0$), the number of required iterations is large. Corollary 3.3 does not use any information about the input data statistics (e.g., moments), and so expectedly, (9) may overestimate N and S . Later, we make use of such data statistics to analyze networks with hidden layers, while the rest of the recipe extends from the 1-NN setting.

Choosing $\mathbb{P}_R(\cdot)$: An important issue in the results presented is the choice of $\mathbb{P}_R(\cdot)$. Firstly, Theorem 3.1 and Corollary 3.2 cover the typical choices of stopping criteria. Secondly, from a practical stand-point, once can use multiple choices of $\mathbb{P}_R(\cdot)$ from a dictionary of distributions \mathcal{P} , and choose the best one via some cross-validation. For instance, the best $\mathbb{P}_R(\cdot) \in \mathcal{P}$ can be selected based on a validation dataset either by directly computing the empirical average of the gradients, or using alternate measures like generalization performance. This is valid because several recent results justify using number of training iterations as a ‘surrogate’ for the generalization performance [20]. Further, the analysis also allows for probing the nature of the R after fixing $\mathbb{P}_R(\cdot)$ i.e., depending on the Hessian at R^{th} iteration, gradient updates may be continued if needed, and the technical results presented here would still hold for this post-hoc increase in iterations. The longer version of the paper will have extensive analysis about the influence and appropriate choices of $\mathbb{P}_R(\cdot)$

4 Unsupervised Pretraining

Building upon the results from 3.1–3.3, we now consider single-layer networks that perform unsupervised learning via box-constrained DA [10] (see the loss function from (4)). Unlike the single-layer setting, because of the constraint on \mathbf{W} , we will now be interested in the expected *projected* gradients $\nabla_{\mathbf{W}} \tilde{f}(\mathbf{W})$, which simply are the Euclidean projections of $\nabla_{\mathbf{W}} f(\mathbf{W})$ on $\mathbf{W} \in [-w_m, w_m]^{d_h \times d_x}$. Alg. 2 in the supplemental technical report summarizes this DA RSG procedure. To ensure broader discussion of the interplay (of network structure and learning) we restrict ourselves to the case of constant γ^k s with $R \sim Unif[1, N]$. We point out that the trends inferred from our results would still be appropriate for alternative stepsizes and $\mathbb{P}_R(\cdot)$. The following result bounds $\nabla_{\mathbf{W}} \tilde{f}(\mathbf{W})$ for DA RSG. e^{da} and e_γ^{da} encode the network structure and learning constants.

Theorem 4.1 (DA, constant stepsize). Consider a DA RSG with $\mathbf{W} \in [-w_m, w_m]^{d_h \times d_x}$. Let

$$e^{da} = \frac{d_x d_h}{16} \left[1 + \frac{\zeta d_x w_m}{4} \mu_{\mathbf{x}} + \left(\frac{5\zeta}{16} - \frac{\zeta^2}{4} \right) (\zeta d_x w_m)^2 \tau_{\mathbf{x}} \right] \quad (10)$$

and \mathcal{U}_{da} denote the Lipschitz constant of $\nabla_{\mathbf{W}} f(\mathbf{W})$ with loss function from (4). Using constant stepsize $\gamma < \frac{2}{\mathcal{U}_{da}}$ and denoting $e_{\gamma}^{da} = 1 - \frac{\mathcal{U}_{da}}{2} \gamma$, the expected projected gradients are

$$\mathbb{E}(\|\nabla_{\mathbf{W}} \tilde{f}(\mathbf{W}^R)\|^2) \leq \frac{D_f}{N \gamma e_{\gamma}^{da}} + \frac{e^{da}}{B} \left(1 + \frac{1}{e_{\gamma}^{da}} \right) \quad (11)$$

and the optimal stepsize is $\gamma_o \approx \sqrt{\frac{2BD_f}{\mathcal{U}_{da} e^{da} N}}$

Remarks: Similar to (7) from Theorem 3.1 for 1-NN, the above result in (11) combines the contributions from the output goodness of fit (D_f), the number of free parameters ($d_h d_x$) and the stepsize choice (γ , e_{γ}^{da}). All the remarks from Theorem 3.1 would still apply here – D_f is balanced out by the number of iterations N , the second term involving the variance of gradients (e^{da}) and the batchsize controls the ‘bias’ from the network’s degrees of freedom. However, unlike 1-NN, here the dependence on the network structure is much more involved. The input and hidden layer lengths do not contribute equally (see e^{da} from (10)) which can be partly explained by the asymmetric structure of the loss (refer to (4)). For smaller constraint sets i.e., small w_m , and hence small e^{da} , we expect the projected gradients to typically have small magnitude, which is clearly implied by (11). In practice, w_m is reasonably small [37] and, in general, W_{ij} ’s have been shown to emulate a ‘fat’ Gaussian with mean centered around zero [38].

Denoising Rates versus Network structure: (10) and its resulting structure in (11) seem to imply a non-trivial interplay between the network size $d_h d_x$, the data statistics (via $\mu_{\mathbf{x}}$ and $\tau_{\mathbf{x}}$) and the denoising rate ζ . Here we analyze this for few commonly encountered cases. A trivial observation from (11) is that it is always beneficial to use smaller (or thinner) networks, resulting in faster decay of expected gradients as iterations increase. This trend, in tandem with observations from [20] about generalization imply the superiority of thinner networks – [39] and others have already shown some empirical evidence for this behavior. As was seen with 1-NNs, fatter DAs will need large batchsize.

Small data moments; ($\mu_{\mathbf{x}} \approx 0$, $\tau_{\mathbf{x}} \approx 0$) In this trivial case, there is nothing much to learn. e^{da} will be as small as possible, resulting in faster convergence, and the influence of ζ on (11) is nullified as well. Hence, independent of the complexity of the task, there is no necessity for using large ζ s whenever the input data averages are small. The only contributing factor is the network size, and clearly, smaller networks are better.

Small denoising rate; ($0 \ll \zeta < 1$) Here $\mathbf{x} * \mathbf{z} \approx \mathbf{x}$, and the noise in gradients $\nabla_{\mathbf{W}} \mathcal{L}(\eta; \mathbf{W})$ is almost entirely from data statistics. Within this setting,

- Small $\mu_{\mathbf{x}}$ and $\tau_{\mathbf{x}}$ leads to faster convergence. As they increase, large batches and N are needed.
- For large $\mu_{\mathbf{x}}$ and $\tau_{\mathbf{x}}$, smaller stepsizes and input length (d_x) are required to control (11). In the pathological case where $\mu_{\mathbf{x}}, \tau_{\mathbf{x}} \sim 1$ and $\zeta \approx 1$, the bounds may be too loose to be relevant.

Independent of how small (11) is, large ζ always leads to overfitting and poor hidden representations [10]. (11) predicts this from the convergence perspective. To see this, observe that large ζ (and mid-sized network) implies large expected gradients, and hence, the training time N needs to be reasonably large. [20] show that networks with large training times may not generalize well.

Large denoising rate, ($0 < \zeta \ll 1$) Here $e^{da} \approx \frac{d_x d_h}{16}$. The influence of data statistics is completely nullified by large corruptions (i.e., small ζ , see (4) and the last two terms in (10)). Unless $\mu_{\mathbf{x}}, \tau_{\mathbf{x}}$ are

unreasonably large the convergence is almost entirely controlled by ζ , which in turn would be faster for thinner networks with large batchsize.

Beyond these prototypical settings, for small stepsizes γ , e_γ^{da} will be as large as possible making the bound tighter. Clearly, the influence of both data moments and denoising rates may be mitigated by increasing B and N . The above described cases are some of the widely used settings, but the interplay from (11) and (10) in Corollary 4.1 is much more involved. The trends and interpretations derived here will be useful in understanding multi-layer networks. Recall the discussion about large deviation estimates from Corollary 3.3. A similar such result relating N to the number of instances S and B can be obtained for DA as well, however, due to space restrictions we omit it here.

5 Multi-layer Networks

We now extend our analysis to multi-layer networks. Using Theorem 4.1 as a starting point, we first consider a L -NN that performs *layer-wise pretraining* using DAs from Section 4 before backpropagation based supervised finetuning. Since the layer-wise pretraining uses box-constraints, we are still interested in the expected projected gradients, accumulated across all the layers. The resulting bound allows us to incorporate feature dropout (during the supervised tuning stage) in both the input and hidden layers later in Section 6. Recall the discussion about L -NN learning mechanisms and the corresponding loss functions (see (5)) from Section 2. A consequence of this general result is a clear and intuitive relation between dropout based supervised learning, layer-wise pretraining and other structural and distributional parameters. Alg. 3 in the technical report presents this multi-layer RSG procedure. For notational convenience we collectively denote the final estimate $\mathbf{W}_1^R, \dots, \mathbf{W}_L^R$ simply as \mathbf{W}^R in the results.

5.1 Layer-wise Pretraining

The following result shows the decay of expected projected gradients for multi-layer RSG. D_f here denotes the initial deviation of the objective *after* the $L - 1$ layers have been pretrained. e_l^m for $l = 1, \dots, L$ encode the structural and learning hyper-parameters of the network, and we assume that all the hidden layers are pretrained to the same degree i.e., each of the $L - 1$ layers are pretrained to a given (α, δ_α) solution.

Corollary 5.1 (Multi-layer Network). *Consider a multi-layer RSG with no dropout learned via layer-wise box-constrained DA pretraining followed by supervised backpropagation with constant stepsizes $\gamma_l \forall l$. Let $e_1^m = \frac{\gamma_1}{4} d_0 d_1 d_2 w_m^l$, $e_l^m = \frac{\gamma_l}{4} d_{l-1} d_l d_{l+1} w_m^l$ and $e_L^m = \frac{13d_{L-1}d_L\gamma_L^2}{256}$. Whenever $\gamma_l < \frac{20}{\alpha d_{l+1} w_m^l}$, and the hidden layers are pretrained for a (α, δ_α) solution (as defined in Definition 1), we have*

$$\mathbb{E} \|\nabla_{\mathbf{W}} \tilde{f}(\mathbf{W}^R)\|^2 \leq \frac{1}{e_\gamma^m} \left(\frac{D_f}{N} + \frac{1}{B} (e_L^m + \alpha \sum_{l=1}^{L-1} e_l^m) \right) \quad (12)$$

where $e_\gamma^m = \min \left\{ \gamma_L - \frac{13}{16}(\gamma_L)^2, \gamma_l - \frac{\alpha d_{l+1} w_m^l}{20}(\gamma_l)^2 \right\}$

Remarks: The structure of (12) is similar to those from (7) and (11). Hence, the trends suggested by the interplay of D_f , L and d_0, \dots, d_L , stepsizes, N and B are similar to those observed from the Theorems 3.1 and 4.1. However, as expected, the interactions are much more complex because of the presence of multiple layers. For fixed network lengths and stepsizes, e_l^m are constants, and encode the variance of gradients within the l^{th} layer and the corresponding free parameters. As network size increases, these constants increase proportionally, there by requiring large N and B . The observations about thinner networks being superior to fatter ones can also be seen from e_l^m s, and (12) suggests that the network depth should *not* be

more than necessary. Such a minimum depth would depend on the trade-off of convergence (from (12)) and generalization [20, 7].

The influence of layer-wise DAs are concealed within α and δ_α and the corresponding D_f . A small α in (12) (which controls the goodness of pretraining) suggests that the influence of the $\mathbf{h}^1, \dots, \mathbf{h}^{L-1}$ on backpropagation tuning is very small. In such a scenario, the tuning is mostly confined to mapping \mathbf{h}^{L-1} to the outputs \mathbf{y} – which is not necessarily a good thing, and one would want to allow for all \mathbf{h}^l s (and \mathbf{W}_l s) to change if needed [10, 40]. On the other hand, a large D_f controls the alternate regime where \mathbf{h}^l s are not “good enough” to abstractly represent the input data. In such a case, aggressive tuning is needed, as suggested by (12), with large N and B . Overall, (12) implies that the goodness (or badness) of pretraining will be passed on to the tuning stage convergence via α and D_f . As suggested earlier, these trade-offs can be related to the training times (proportional to N , derived from generalization [20]). Recently, [41] showed empirical evidence that aggressive pretraining (especially in higher layers) results in \mathbf{h}^l s that may not be transferable across multiple arbitrary tasks. Theorem 5.1 suggests the same – α (and D_f) should not be reasonably small to allow for aggressive tuning to arbitrary learning tasks. It is clear from (12) that we can control the *convergence* of multi-layer nets by preceding the backpropagation tuning with layer-wise pretraining. There is already very strong empirical evidence for the generalization performance of pretraining [33, 10, 40]. Corollary 5.1 complements these studies with guaranteed convergence of gradients.

6 Multi-layer with Dropout

Using ReLUs and dropout, [42, 2, 43] and others have shown that whenever large amounts of labeled data is available, pretraining might not be *necessary* to achieve good generalization, clearly suggesting an underlying relationship between dropout and layer-wise pretraining. The following result summarizes the convergence of expected projected gradients in this general setting where the $L - 1$ layers are not pretrained, and instead dropout is induced in input and all hidden layers, and supervised backpropagation is performed directly with random initializations for $\mathbf{W}_1^1, \dots, \mathbf{W}_L^1$ [31]. Here, ζ denotes the dropout rate for all layers i.e., w.p. $1 - \zeta$ a unit is dropped and all parameters corresponding to this unit are not updated (see Section 2, and (5)). It is reasonable to expect some interaction between α (the pretraining goodness, see Corollary 5.1) and ζ .

Corollary 6.1 (Pretraining vs. Dropout). *Given the input and hidden layer dropout rate ζ , for learning the L -layered network from Corollary 5.1 that is pretrained to a (α, δ_α) solution, we have*

$$\begin{aligned} \mathbb{E}_{R,\eta} \|\nabla_{\mathbf{W}} \tilde{f}(\mathbf{W}^R)\|^2 &\leq \frac{D_f}{N e_\gamma^m \zeta^2} \\ &+ \frac{1}{e_\gamma^m B} \left(\frac{e_L^m}{\zeta} + \alpha e_{L-1}^m + \alpha \zeta \sum_{l=1}^{L-2} e_l^m \right) \end{aligned} \quad (13)$$

Remarks: This is our most general result. Although, the schools of layer-wise pretraining and dropout have been studied independently, from our knowledge, this is the first theoretical result that explicitly combines these two regimes in a systematic way. Recall that e_1^m, \dots, e_L^m encode L -NN’s degrees of freedom. Hence, the first term in (13) corresponds to the outputs’ goodness-of-fit, while the other terms represent the effective ‘reduction’ in the number of free parameters because of pretraining. Independent of ζ , (13) clearly implies that pretraining always improves convergence (since α will reduce). [9, 31] have shown empirically that this is especially true for improving generalization. To better interpret (13), consider the two standard mechanisms — dropout learning with and without layer-wise pretraining. Given L, d_0, \dots, d_L and stepsizes $(\gamma_1, \dots, \gamma_L)$, e_1^m, \dots, e_L^m from Corollary 5.1 are constants. We assume D_f (which depends on pretraining) to be reasonably large because the loss is calculated over predicted vs. observed \mathbf{y} s (see (5)).

6.1 Pretraining + Dropout

- If L -NN is pretrained to a small α , the last two terms in (13) are already as small as they can be, and the rest will dominate them. For a given B and N , the best choice for ζ will then be ≈ 1 i.e., very minimal or *no* dropout. This is essentially “good” layer-wise pretraining followed by supervised fine-tuning, which is known to work well [44, 33]. Hence, (13) presented a succinct and clear theoretical argument for the classical revival of deep networks from the perspective of convergence (and eventually generalization).
- Alternatively, if the pretraining is poor (i.e., large α), and we still operate in the $\zeta \rightarrow 1$ regime, the fine-tuning will update the full network in each iteration. This would result in overfitting – the fundamental argument that necessitates dropout (empirically made clear in [31, 28]). Hence one needs to decrease ζ , resulting in slower convergence because the first two terms will rapidly increase as ζ decreases. This is expected since dropout essentially adds ‘noise’ to the solution path, forcing backpropagation to work with a subset of all activations [11, 28], and overall, involving more work.

6.2 Only Dropout

With no pretraining, there is a complex trade-off between the terms involving ζ in (13). The optimal ζ will need to balance out the variance from the hidden layers (the last term) and the goodness of output approximation (first and second terms). For certain values of D_f and e_l^m ’s, the best ζ will be ~ 0.5 , which was recently reported as the best rate as per dropout dynamics [11]. Large values of B and N will ensure that the bound is small. Large N , in turn, implies larger training dataset size (see the setup from Corollary 3.3 and 5.1). Putting it another way, if large amounts of labeled data is available, one can by-pass pretraining completely and only perform supervised backpropagation forcing all the terms in (13) to reduce with reasonably large B and number of epochs. This was the argument put forth by the recent school of deep learning [2, 43] – fully supervised dropout with very large amounts of training data – where pretraining has been for the most part neglected completely.

Overall, Corollary 6.1 corroborates many existing observations about pretraining and dropout, and provides new tools to *guide* the learning procedure, in general. It guarantees convergence, and allows us to *explicitly* calculate the best possible settings for all hyper-parameters, to achieve a certain level of generalization or training time [20]. A more extensive discussion about the outcomes and insights from the interplay in Corollary 6.1 will be included in the longer version of the paper. The three algorithms presented in this work are only minor modifications over the classical autoencoder and backpropagation pipeline, and so, are straightforward to implement.

7 Experiments

We evaluate the bounds in (11), (12) and (13) using MNIST, CIFAR-10, CALTECH-101 and a brain imaging dataset. Figures 1 summarizes a fraction of the findings on CIFAR-10. The longer version of the paper will have complete set of evaluations. Figure 1(a) shows the expected gradients vs. network sizes. As predicted in (7), the convergence is slower for large networks. Visible and hidden layer lengths have unequal influence on the convergence (see remarks of Corollary 4.1). The influence of denoising rate and w_m is in Figure 1(b). The y -axis shows the expectation of *projected* gradients on $[-w_m, w_m]$, and as suggested by (11), the convergence is faster for small w_m ’s. It is interesting that the choice of ζ ’s has an almost negligible influence (refer Figure 1(b), for a given w_m). Figure 1(c), shows the interaction of network lengths vs ζ , and as observed from Figure 1(a), the networks lengths dominate the convergence with visible layer length being the most influential factor. The plots in Figure 1 correspond to the data moments $\mu_x \sim 0.5$ and $\tau_x \sim 0.3$, implying that the terms in e^{da} (refer (10)) are non-negligible. Figure 1(d) shows the influence of the network

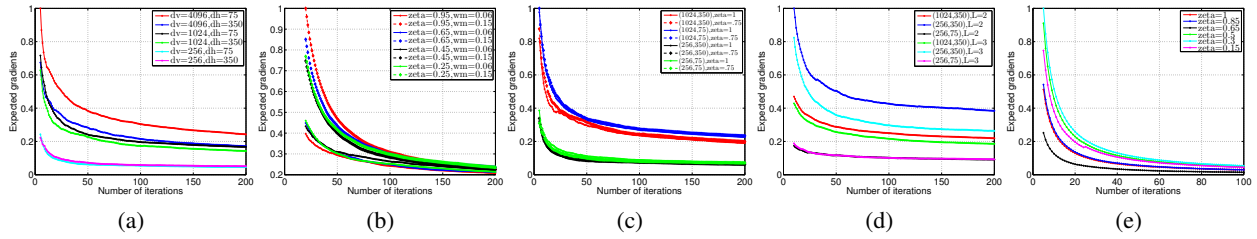


Figure 1: Expected gradients (CIFAR) for pretraining (a,b) and multi-layer net (c,d). (e) Influence of dropout. y -axes are scaled to maximum. $B = 100$.

depth. Clearly, the expected gradients are influenced more by layer lengths than the number of layers itself (12). Figure 1(e) shows the effect of changing the dropout rate in a 3-layered network. Although the overall effect is small, the convergence is slower for very small (red curve) and very large (cyan, magenta curves) ζ s (see (13)).

Acknowledgements: We thank the anonymous reviewers. We are supported by NIH AG040396, NSF CAREER 1252725, NSF CCF 1320755, the UW grants ADRC AG033514, ICTR 1UL1RR025011 and CPCP AI117924.

References

- [1] Ian Lenz, Honglak Lee, and Ashutosh Saxena. Deep learning for detecting robotic grasps. *The International Journal of Robotics Research*, 34:705–724, 2015.
- [2] A. Krizhevsky, I. Sutskever, and G. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, volume 1, page 4, 2012.
- [3] Geoffrey Hinton, Li Deng, Dong Yu, George E Dahl, et al. Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *Signal Processing Magazine, IEEE*, 29:82–97, 2012.
- [4] Ilya Sutskever, Oriol Vinyals, and Quoc VV Le. Sequence to sequence learning with neural networks. In *NIPS*, pages 3104–3112, 2014.
- [5] Vamsi K Ithapu, Vikas Singh, Ozioma C Okonkwo, et al. Imaging-based enrichment criteria using deep learning algorithms for efficient clinical trials in mild cognitive impairment. *Alzheimer's & Dementia*, 11:1489–1499, 2015.
- [6] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521:436–44, 2015.
- [7] Y. Bengio. Learning deep architectures for AI. *FTML*, 2:1–127, 2009.
- [8] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation learning: A review and new perspectives. *PAMI*, 35:1798–1828, 2013.
- [9] H. Lee, R. Grosse, R. Ranganath, and A. Ng. Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations. In *ICML*, pages 609–616, 2009.
- [10] P. Vincent, H. Larochelle, et al. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *JMLR*, pages 3371–3408, 2010.

- [11] Pierre Baldi and Peter Sadowski. The dropout learning algorithm. *Artificial intelligence*, 210:78–122, 2014.
- [12] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv:1502.03167*, 2015.
- [13] Hongmei Shao and Gaofeng Zheng. Convergence analysis of a back-propagation algorithm with adaptive momentum. *Neurocomputing*, 74:749–752, 2011.
- [14] Yann N Dauphin, Razvan Pascanu, Caglar Gulcehre, Kyunghyun Cho, Surya Ganguli, and Yoshua Bengio. Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. In *NIPS*, pages 2933–2941, 2014.
- [15] R. Livni, S. Shalev-Shwartz, and O. Shamir. On the computational efficiency of training neural networks. In *NIPS*, pages 855–863, 2014.
- [16] Francis Bach. Breaking the curse of dimensionality with convex neural networks. *arXiv:1412.8690*, 2014.
- [17] Sanjeev Arora, Aditya Bhaskara, Rong Ge, and Tengyu Ma. Provable bounds for learning some deep representations. In *ICML*, pages 584–592, 2014.
- [18] Ankit B Patel, Tan Nguyen, and Richard G Baraniuk. A probabilistic theory of deep learning. *arXiv:1504.00641*, 2015.
- [19] Majid Janzamin, Hanie Sedghi, and Anima Anandkumar. Beating the perils of non-convexity: Guaranteed training of neural networks using tensor methods. *CoRR abs/1506.08473*, 2015.
- [20] Moritz Hardt, Benjamin Recht, and Yoram Singer. Train faster, generalize better: Stability of stochastic gradient descent. *arXiv:1509.01240*, 2015.
- [21] Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *COMPSTAT*, pages 177–186. 2010.
- [22] S. Ghadimi and G. Lan. Stochastic first-and zeroth-order methods for nonconvex stochastic programming. *SIAM Optimization*, 23:2341–2368, 2013.
- [23] Sue Becker and Yann Le Cun. Improving the convergence of back-propagation learning with second order methods. In *Proc. connectionist models summer school*, pages 29–37, 1988.
- [24] George D. Magoulas, Michael N. Vrahatis, and George S Androulakis. Improving the convergence of the backpropagation algorithm using learning rate adaptation methods. *Neural Computation*, 11:1769–1796, 1999.
- [25] J. Ngiam, A. Coates, A. Lahiri, B. Prochnow, Q. Le, and A. Ng. On optimization methods for deep learning. In *ICML*, pages 265–272, 2011.
- [26] Yoshua Bengio. Practical recommendations for gradient-based training of deep architectures. In *Neural Networks: Tricks of the Trade*, pages 437–478. 2012.
- [27] Enrique Castillo, Bertha Guijarro-Berdiñas, Oscar Fontenla-Romero, and Amparo Alonso-Betanzos. A very fast learning method for neural networks based on sensitivity analysis. *JMLR*, 7:1159–1182, 2006.

- [28] Stefan Wager, Sida Wang, and Percy S Liang. Dropout training as adaptive regularization. In *NIPS*, pages 351–359, 2013.
- [29] Sanjeev Arora, Yingyu Liang, and Tengyu Ma. Why are deep nets reversible: A simple theory, with implications for training. *arXiv:1511.05653*, 2015.
- [30] Tao Wei, Changhu Wang, Rong Rui, and Chang Wen Chen. Network morphism. *arXiv:1603.01670*, 2016.
- [31] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *JMLR*, 15:1929–1958, 2014.
- [32] Ian J Goodfellow, David Warde-Farley, Mehdi Mirza, Aaron C Courville, and Yoshua Bengio. Maxout networks. *ICML*, 28:1319–1327, 2013.
- [33] D. Erhan, Y. Bengio, A. Courville, et al. Why does unsupervised pre-training help deep learning? *JMLR*, 11:625–660, 2010.
- [34] Katta G Murty and Santosh N Kabadi. Some np-complete problems in quadratic and nonlinear programming. *Math. prog.*, 39:117–129, 1987.
- [35] Subutai Ahmad, Gerald Tesauro, and Yu He. Asymptotic convergence of backpropagation: numerical experiments. In *NIPS*, pages 606–613, 1990.
- [36] Yann A LeCun, Léon Bottou, Genevieve B Orr, and Klaus-Robert Müller. Efficient backprop. In *Neural networks: Tricks of the trade*, pages 9–48. 2012.
- [37] Deepak Pathak, Philipp Krahenbuhl, and Trevor Darrell. Constrained convolutional neural networks for weakly supervised segmentation. In *ICCV*, pages 1796–1804, 2015.
- [38] I Bellido and Emile Fiesler. Do backpropagation trained neural networks have normal weight distributions? In *ICANN*, pages 772–775. 1993.
- [39] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. *arXiv:1412.6550*, 2014.
- [40] A. Saxe, P. Koh, Z. Chen, M. Bhand, B. Suresh, and A. Ng. On random weights and unsupervised feature learning. In *ICML*, pages 1089–1096, 2011.
- [41] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson. How transferable are features in deep neural networks? In *NIPS*, pages 3320–3328, 2014.
- [42] Vinod Nair and Geoffrey E Hinton. Rectified linear units improve restricted boltzmann machines. In *ICML*, pages 807–814, 2010.
- [43] Christian Szegedy, Wei Liu, Yangqing Jia, et al. Going deeper with convolutions. *arXiv:1409.4842*, 2014.
- [44] G. Hinton and R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313:504–507, 2006.