



CENTER FOR  
**Brains  
Minds+  
Machines**

CBMM Memo No. 067

September 15, 2017

## Theory of Deep Learning III: Generalization Properties of SGD

by

Chiyuan Zhang<sup>1</sup> Qianli Liao<sup>1</sup> Alexander Rakhlin<sup>2</sup> Brando Miranda<sup>1</sup> Noah Golowich<sup>1</sup> Tomaso Poggio<sup>1</sup>

<sup>1</sup>Center for Brains, Minds, and Machines, McGovern Institute for Brain Research,  
Massachusetts Institute of Technology, Cambridge, MA, 02139.

<sup>2</sup> University of Pennsylvania

**Abstract:** In Theory III we characterize with a mix of theory and experiments the consistency and generalization properties of deep convolutional networks trained with Stochastic Gradient Descent in classification tasks. A present perceived puzzle is that deep networks show good predictive performance when overparametrization relative to the number of training data suggests overfitting. We describe an explanation of these empirical results in terms of the following new results on SGD:

1. SGD concentrates in probability - like the classical Langevin equation – on large volume, “flat” minima, selecting flat minimizers which are with very high probability also *global minimizers*.
2. Minimization by GD or SGD on flat minima can be approximated well by minimization on a linear function of the weights suggesting pseudoinverse solutions.
3. Pseudoinverse solutions are known to be intrinsically regularized with a regularization parameter  $\lambda$  which decreases as  $\frac{1}{T}$  where  $T$  is the number of iterations. This can qualitatively explain all the generalization properties empirically observed for deep networks.
4. GD and SGD are connected closely to robust optimization. This provides an alternative way to show that GD and SGD perform implicit regularization.

These results explain the puzzling findings about fitting randomly labeled data while performing well on natural labeled data. They also explain while overparametrization does not result in overfitting. Quantitative, non-vacuous bounds are still missing, as it has almost always been the case for most practical applications of machine learning. We describe in the appendix an alternative approach that explains more directly, with tools of linear algebra the same qualitative properties and puzzles of generalization in deep polynomial networks.



**This work was supported by the Center for Brains, Minds and Machines (CBMM), funded by NSF STC award CCF - 123 1216.**

# 1 Introduction

## change Appendix theory 4

In the last few years, deep learning has been tremendously successful in many important applications of machine learning. However, our understanding of deep learning is still far from complete. A satisfactory characterization of deep learning should at least cover the following parts: 1) representation power — what types of functions can neural networks (DNNs) represent well and what are the advantages and disadvantages over using shallow models? 2) optimization of the empirical loss — can we characterize the convergence of stochastic gradient descent (SGD) on the non-convex empirical loss encountered in deep learning? 3) why do the deep learning models, despite being highly over-parameterized, still predict well?

The first two questions are addressed in Theory I<sup>[1]</sup> and Theory II<sup>[2]</sup> respectively. In this paper, we focus on the last question, and try to provide a hybrid theoretical and empirical view of it. More specifically, we study learning, especially in over-parameterized DNNs, by optimizing the empirical loss with SGD.

We remind that *generalization* is defined as the following property: the gap between empirical and expected error goes to zero with increasing size of the training set. *Consistency* is a different property: the empirical error converges to the optimal Bayes error with increasing size of the training set. Both properties can be distribution dependent or independent and both have strong and weak forms.

In the case of *one-pass* SGD, where each training point is only visited at most once, the algorithm is optimizing the expected loss directly. Therefore, there is no need to define the empirical loss and to measure generalization: consistency is key here and it holds under rather general conditions. However, in practice, unless one has access to infinite training samples, one-pass SGD is rarely used. Instead, it is almost always better to run many passes of SGD over the same training set. In this case, the algorithm is optimizing the empirical loss, and the deviation between the empirical loss and the expected loss (i.e. the generalization error) must be controlled.

In statistical learning theory, the deviation is usually controlled by restricting the complexity of the hypothesis space. For example, in binary classification, for a hypothesis space with VC-dimension  $d$  and  $N$  i.i.d. training samples, the generalization error could be upper bounded by  $\mathcal{O}(\sqrt{d/N})$ . In the *distribution-free setting*, the VC dimension also provide a lower bound for the generalization error. On the other hand, for overparametrization ( $d \gg N$ ), there is a data distribution under which the generalization error could be arbitrarily bad<sup>[3]</sup>. As we will discuss later this worst case behavior was recently demonstrated by a randomization test on large deep neural networks that have the full capability of shattering the whole training set<sup>[4]</sup>. In those experiments, zero-error minimizers for the empirical loss are found by SGD. Since the test performance must be at chance level, the worst possible generalization error is observed. On the other hand, those same networks are found to achieve low expected error on image classification datasets, without obvious regularization-like mechanisms such as weight decay or data augmentation. This creates an apparent puzzle that we will discuss in section 6.1 as the classical characterization of distribution-independent generalization no longer readily applies. We observe that the puzzle is misleading because consistency may hold in the absence of generalization. A very classical example is 1-Nearest Neighbor, which is an algorithm with zero empirical error for all  $N$  and good asymptotic performance in terms of expected error. More in general, k-NN algorithms for all fixed  $k$  do not have generalization. In fact, if the Bayes error is zero, k-NN algorithms are consistent for all  $k$ !

In this paper, we speak *very loosely* about two regimes for training deep convolutional networks (in the specific case of CIFAR). In the underconstrained regime – which is the standard one for deep learning applications – in which  $n \leq W$ , where  $n$  is the size of the training set and  $W$  is the number of weights– the empirical error can be zero, when, as described in Theory I and II, the regression function underlying the problem is contained in the class of (compositional) functions that are well approximated by the convolutional network and the weights in the network are sufficient for it (this exclude too many “dummy” weights) . On the other hand, the generalization error - the difference between empirical and expected error – may not go to zero for  $n \rightarrow \infty$ . We show that the expected error in the underconstrained regime can be low because SGD implicitly regularizes. In the over-constrained regime –  $n \geq W$  – generalization follows from classical bounds on neural networks with a finite number of weights (assuming the weights themselves are bounded). As side-effects of these basic results we show close relations between SGDL, robust optimization wrt perturbations of the weights and regularization. Finally, we discuss how recent puzzling results on generalization by deep (and shallow) networks can be explained.

In the rest of the paper, we try to address this set of issues at the level of formal rigor of a physicist (not a mathematician: this will come later). With a mix of theoretical and empirical results, we show that isotropic flatness around the global minimizers play a key role in the generalization of over-parameterized deep neural networks.

Notice that in all computer simulations reported in this paper, we turn off all the “tricks” used to improve performance such as data augmentation, weight decay etc. in order to study the basic properties of the SGD algorithm. As a consequence, performance is not state of the art but maximum performance is not our goal here.

## 1.1 Related work

Deep Learning references start with Hinton’s back-propagation and with LeCun’s convolutional networks (see<sup>[5]</sup> for a nice review). Of course, multi-layer convolutional networks have been around at least as far back as the optical processing era of the 70s. The Neocognitron<sup>[6]</sup> was a convolutional neural network that was trained to recognize characters, inspired by the hierarchical processing postulated by Hubel and Wiesel<sup>[7]</sup> from their recordings of simple and complex cells in visual cortex. The property of *compositionality* was a main motivation for hierarchical models of visual cortex such as HMAX which was described as a pyramid of AND and OR layers<sup>[8]</sup>, that is a sequence of

conjunctions and disjunctions. In Theory I<sup>[1]</sup> we have provided formal conditions under which deep networks can use the compositionality to avoid the curse of dimensionality.

A number of recent papers have mentioned some of the ideas explored in this paper. For instance Soatto et al.<sup>[9]</sup> remark that almost-flat regions of the energy landscape are robust to data perturbations, noise in the activations as well as perturbations of the parameters — which are as we show later directly related to good generalization. They also note that wide valleys should result in better generalization and that optimization algorithms in deep learning seem to discover exactly that.

Recent work by Keskar et al.<sup>[10]</sup> is even more relevant. The authors estimate the loss in a neighborhood of the weights to argue that small batch size in SGD (i.e., larger gradient estimation noise, see later) generalizes better than large mini-batches and also results in significantly flatter minima. In particular, they note that the stochastic gradient descent method used to train deep nets, operate in a small-batch regime wherein a fraction of the training data, usually between 32 and 512 data points, is sampled to compute an approximation to the gradient. They discuss the common observation that when using a larger batch there is a significant degradation in the quality of the model, as measured by its ability to generalize. We provide theoretical arguments for the cause for this generalization drop in the large-batch regime, supporting the numerical evidence of Keskar et al.. The latter shows that large-batch methods tend to converge to sharp minimizers of the training and testing functions — and that sharp minima lead to poor generalization. In contrast, small-batch methods consistently converge to minimizers that generalize better, and our theoretical arguments support a commonly held view that this is due to the inherent noise in the gradient estimation.

On the other hand, as shown in<sup>[11]</sup>, sharp minimizers do not necessarily lead to bad generalization performance. Due to the parameterization redundancy in deep neural networks, given any (flat) minimizers, one can artificially transform the weights to land in a sharp but equivalent minimizer because the function defined by the weights are the same. However, the argument in<sup>[11]</sup> does not conflict with our argument that flat minimizers generalize well. Moreover, it is not clear whether SGD will select the somewhat artificial sharp minimizers created by the equivalent transformations. Notice that isotropic flat minima in the loss wrt *all weights* – in which we are interested – cannot be transformed in sharp minima.

Another recent paper (after the first version of this memo)<sup>[12]</sup> is relevant in terms of the intuition about margin and low expected error (but not generalization as claimed!). In<sup>[12]</sup>, the authors prove nonvacuous bounds for the expected error of *stochastic neural networks*, that is, neural networks for which the weights  $w$  are chosen randomly from some distribution  $\mathcal{D}$ . It remains an open question, however, whether the fact that the expected error is small in expectation over the weights  $w \sim \mathcal{D}$ , implies that for a specific weight vector, suppose  $w_0 := \mathbb{E}_{w \sim \mathcal{D}}[w]$ , leads to expected error that is nearly as low.

In another line of work<sup>[13, 14]</sup>, the authors prove generalization error bounds on neural networks that are expressed primarily in terms of various norms of the weights of the network. Such bounds are obtained by controlling the Rademacher complexity of certain classes of norm-bounded neural networks; these bounds generalize some of the results in<sup>[15]</sup>. The idea that generalization error bounds depending on the size of the parameters rather than the number of parameters dates back to<sup>[16]</sup>, where similar, though weaker, generalization bounds were shown, based on the fat-shattering dimension.

More recently,<sup>[17]</sup> has proven another norm-based bound for fully-connected networks, which depends on the product of the spectral norms of the matrices defining each layer of the network, as well as the sum of the vectorized  $\ell_1$  norms of each layer. This work emphasized the importance of *margin* in ensuring generalization, similar to earlier works on boosting<sup>[18]</sup>. In the binary classification case, where the neural network computes a function  $f : \mathbb{R}^d \rightarrow \mathbb{R}$ , and the class of  $x \in \mathbb{R}^d$  is decided based off of  $\text{sign}(f(x))$ , the margin of a training point  $x_i$  can be defined as  $f(x_i)y_i$ . The results of<sup>[17]</sup> (as well as<sup>[18]</sup>, and related works such as<sup>[15]</sup>) show that if  $f$  has large margin on many of the training data points, then tighter generalization bounds can be proven.

The puzzles mentioned in the abstract and text refer mainly to the recent work of one of us<sup>[4]</sup>, which is one of the motivations of this paper. Even more recently,<sup>[19]</sup> describe similar puzzles for kernel machines trained without regularization terms by SGD. Our results should be helpful in explaining them.

## 2 Background Facts

We describe here several previous results and observations that are key for our results about generalization.

### 2.1 Landscape of the Empirical Risk

In Theory II<sup>[2]</sup> we have described some features of the landscape of the empirical risk, for the case of deep networks of the compositional type (with weight sharing, though the proofs do not need the weight sharing assumption). We assumed over-parametrization, that is more parameters than training data points, as most of the successful deep networks. Under these conditions, setting the empirical error to zero yields a system of  $\epsilon$ -approximating polynomial equations that have an infinite number of solutions (for the network weights). Alternatively, one can replace the RELUs in the network with an approximating univariate polynomial and verify empirically that the network behavior is essentially unchanged. The associated system of equations allows for a large number of solutions – when is not inconsistent – which are *degenerate, that is flat* in several of the dimensions (in CIFAR there are about  $10^6$  unknown parameters for  $6 \times 10^4$  equations. Notice that

solutions with zero empirical error are global minimizers. No other solution with zero-error exists with a deeper minimum or less generic degeneracy. Empirically we observe (see Theory II) that zero-minimizers correspond to isotropically flat regions – not just valleys. In this paper, we will use the word flatness in two distinct meanings (the context makes clear which meaning): one to refer to degeneracy of the empirical minimizers, the other to refer to isotropical flatness around the zero of the empirical loss.

## 2.2 SGD: Basic Setting

Let  $Z$  be a probability space with an unknown measure  $\rho$ . A training set  $S_n$  is a set of i.i.d. samples  $z_i, i = 1, \dots, n$  from  $\rho$ . Assume a hypothesis  $\mathcal{H}$  is chosen in advance of training. Here we assume  $\mathcal{H}$  is a  $p$ -dimensional Hilbert space, and identify elements of  $\mathcal{H}$  with  $p$ -dimensional vectors in  $\mathbb{R}^p$ . A loss function is a map  $V : \mathcal{H} \times Z \rightarrow \mathbb{R}_+$ . Moreover, we assume the expected loss

$$I(f) = \mathbb{E}_z V(f, z) \tag{1}$$

exists for all  $f \in \mathcal{H}$ . We consider the problem of finding a minimizer of  $I(f)$  in a closed subset  $K \subset \mathcal{H}$ . We denote this minimizer by  $f_K$  so that

$$I(f_K) = \min_{f \in K} I(f) \tag{2}$$

In general, the existence and uniqueness of a minimizer is not guaranteed unless some further assumptions are specified.

Since  $\rho$  is unknown, we are not able evaluate  $I(f)$ . Instead, we try to minimize the empirical loss

$$I_{S_n}(f) = \hat{\mathbb{E}}_{z \sim S_n} V(f, z) = \frac{1}{n} \sum_{i=1}^n V(f, z_i) \tag{3}$$

as a proxy. In deep learning, the most commonly used algorithm is SGD and its variants. The basic version of SGD is defined by the following iterations:

$$f_{t+1} = \Pi_K (f_t - \gamma_t \nabla V(f_t, z_t)) \tag{4}$$

where  $z_t$  is a sampled from the training set  $S_n$  uniformly at random, and  $\nabla V(f_t, z_t)$  is an unbiased estimator of the full gradient of the empirical loss at  $f_t$ :

$$\mathbb{E}_{z_t \sim S_n} [\nabla V(f_t, z_t)] = \nabla \hat{I}(f_t)$$

$\gamma_t$  is a decreasing sequence of non-negative numbers, usually called the *learning rates* or *step sizes*.  $\Pi_K : \mathcal{H} \rightarrow K$  is the projection map into  $K$ , and when  $K = \mathcal{H}$ , it becomes the identity map. It is interesting that the following equation, labeled SGDL, and studied by several authors, including <sup>[20]</sup>, seem to work as well as or better than the usual repeat SGD used to train deep networks, as discussed in Section 5:

$$f_{t+1} = f_t - \gamma_n \nabla V(f_t, z_t) + \gamma'_t W_t. \tag{5}$$

Here  $W_t$  is a standard Gaussian vector in  $\mathbb{R}^p$  and  $\gamma'_t$  is a sequence going to zero.

We consider a situation in which the expected cost function  $I(f)$  can have, possibly multiple, *global* minima. As argued by <sup>[21]</sup> there are two ways to prove convergence of SGD. The first method consists of partitioning the parameter space into several attraction basins, assume that after a few iterations the algorithm confines the parameters in a single attraction basin, and proceed as in the convex case. A simpler method, instead of proving that the function  $f$  converges, proves that the cost function  $I(f)$  and its gradient  $\mathbb{E}_z \nabla V(f, z)$  converge.

Existing results show that when the learning rates decrease with an appropriate rate, and subject to relatively mild assumptions, stochastic gradient descent converges almost surely to a global minimum when the objective function is convex or pseudoconvex<sup>1</sup>, and otherwise converges almost surely to a local minimum. This direct optimization shortcuts the usual discussion for batch ERM about differences between optimizing the empirical risk on  $S_n$  and the expected risk.

Often extra-assumptions are made to ensure convergence and generalization by SGD. Here we observe that simulations on standard databases remain essentially unchanged if the domain of the weights is assumed to be a torus which is compact, because the weights remain bounded in most cases.

## 3 Consistency and Generalization

“Folk theorems” as well as classical bounds (see 8.1) suggest that generalization requires more data than the number  $W$  of “effective” parameters – though of course the number of parameters is not a general measure of the network complexity that is relevant for generalization. This is typically not the case in the current use of deep networks. Curiously, the recent wave of papers related to prediction performance of deep learning has mostly neglected the important fact that good accuracy on new data (consistency) can take place in the absence of generalization – defined as the convergence of empirical risk to expected risk as the number of data  $N$  increases. As we mentioned

<sup>1</sup>In convex analysis, a pseudoconvex function is a function that behaves like a convex function with respect to finding its local minima, but need not actually be convex. Informally, a differentiable function is pseudoconvex if it is increasing in any direction where it has a positive directional derivative.

already, a classical example is the NN algorithm which can be proven to perform quite well in general but does not generalize in a distribution-independent way.

Loosely speaking we distinguish two phases (in principle both  $n$  and  $W$  grow to infinity):

1.  $n \geq W$  where  $N$  is the number of training data and  $W$  is the number of “effective” parameters: this is the “classical” situation for generalization
2.  $n \leq W$ : this is the “new” regime in which SGD performs well in terms of consistency but may not have generalization.

Figure 1 shows that at least for this experiment on CIFAR there seem to be generalization for  $N > W$  as expected and, though there is no generalization for  $N < W$ , the test error – as a proxy for the expected error – is good. The random label case in Figure 2 makes it very clear. The expected error is at chance level throughout; the training error is zero for  $n < W$ , but starts to increase to reach the expected error for  $n > W$ . The behavior in the “classical” regime – on the right side of the figures – can be accounted for by several theoretical approaches that we describe in this section. The behavior in the second regime seems to require a novel approach. We show however that it can be explained with rather classical and even linear results (in particular see Appendix 9.10).

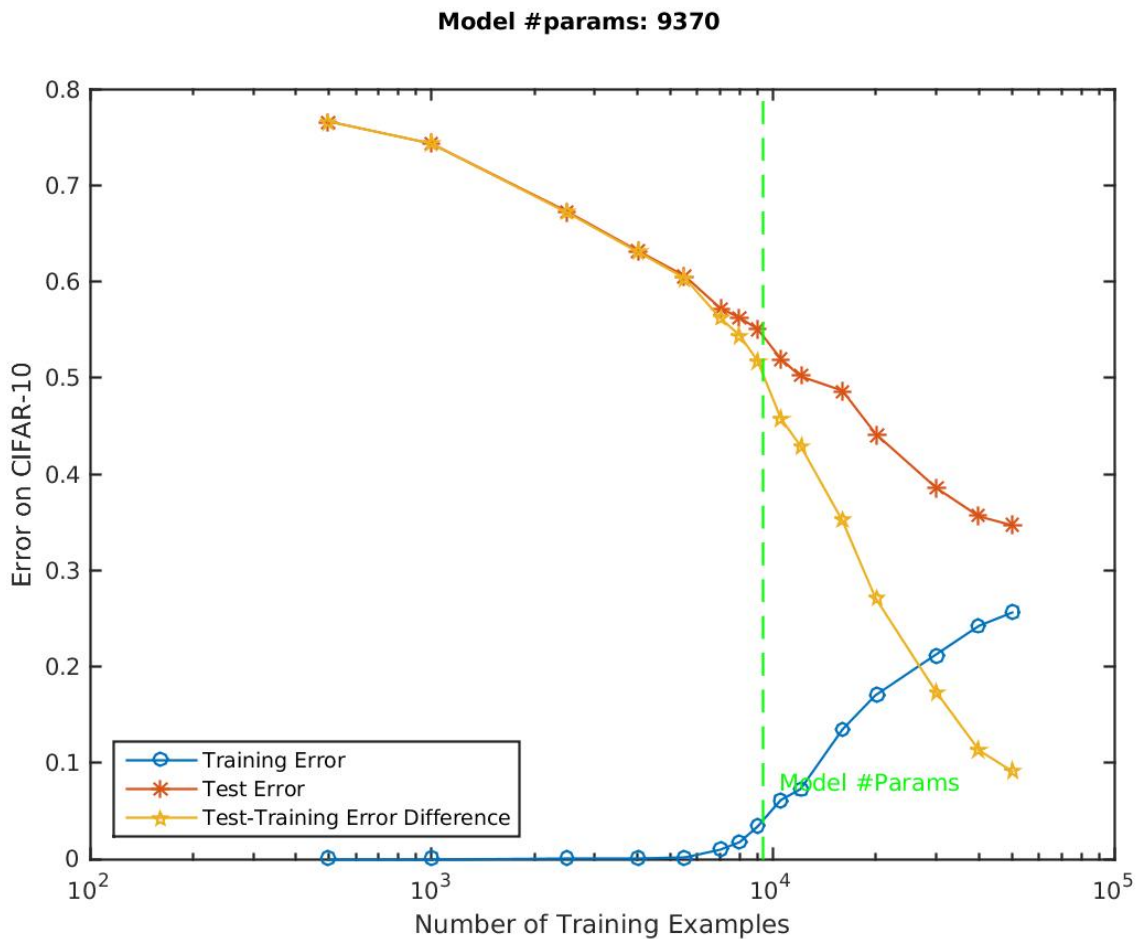


Figure 1: The figure shows the behavior of a deep network trained on subsets of the CIFAR database. The network is a 5-layer all convolutional network (i.e., no pooling) with 16 channels per hidden layer, resulting in only  $W \approx 10000$  weights instead of the typical 300,000. Neither data augmentation nor regularization is performed. The two regimes described in the text are for  $N > W$  and  $N < W$  respectively.

The classical theory characterizes generalization behavior as a function of  $n$  the number of training examples. But a more practical question is which architecture to choose for a fixed, given number of examples  $N$ . In fact, Figures 1 and 2 prompt some additional questions:

- for a fixed  $N$  is it better wrt expected error to train a deep network with  $N \geq W$  or with  $N \leq W$ ?
- is the mechanism explaining generalization the same that explain good expected error performance for  $N \leq W$ ?

Model #params: 9370

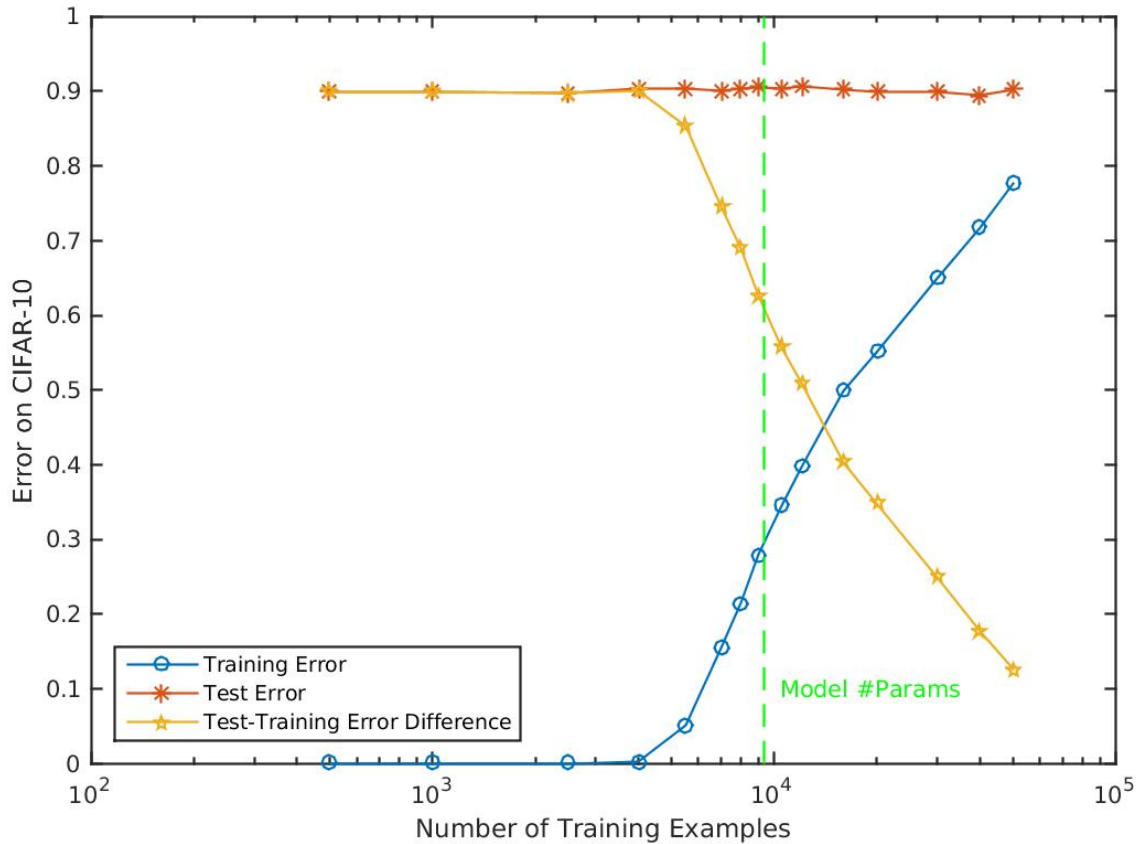


Figure 2: The figure shows the behavior of a deep network trained on subsets of the CIFAR database in which the labels have been randomly scrambled. The network is a 5-layer all convolutional network (i.e., no pooling) with 16 channels per hidden layer, resulting in only  $W \approx 10000$  weights instead of the typical 300,000. Neither data augmentation nor regularization is performed. The two regimes described in the text are for  $N > W$  and  $N < W$  respectively.

The empirical answer to these questions is shown in Figures 3 and 4 and it is surprising: it appears that the increase in the number of parameters does not yield any overfitting. We will first deal with the classical regime  $N \ll W$  and then take on the challenge of explaining the non-overfitting behavior shown on the right of Figure 3.

#### 4 Classical generalization bounds

This section though necessary for completeness is too boring. As a consequence, it landed in the Appendix which the reader may consult because some of the definitions necessary for the following are there. The Appendix reviews a number of approaches used to prove generalization – that is convergence to zero of the gap between empirical and expected error for  $n \rightarrow \infty$ . Though some of them do not explicitly require  $n \gg W$ , they “seem” empirically vacuous in our experiments with CIFAR for  $n \ll W$ .

#### 5 The puzzle: SGD for the underdetermined case

The puzzle seems outside classical learning theory: why there is no overfitting for increasing overparametrization? Why is the expected error so low in the non-classical regime of  $N < W$  in which we cannot speak about generalization? Our conjecture is that

- Conjecture 1.**
- *SGD, while minimizing the empirical loss also maximizes the volume, that is “flatness”, of the minima, effectively linearizing the problem wrt weights*
  - *in the linear case SGD and GD converge to the pseudoinverse-based solution*

Training data size: 50000

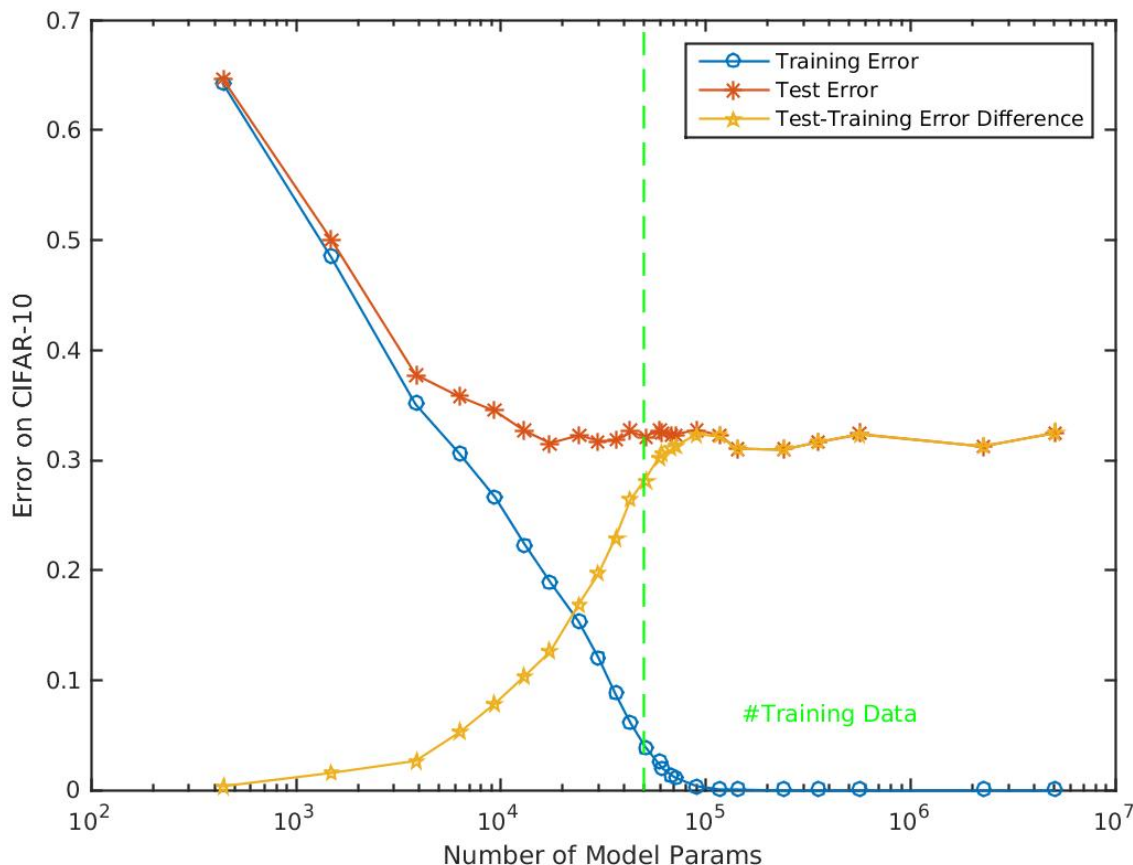


Figure 3: The previous figures show dependence on  $N$  – number of training examples – for a fixed architecture with  $W$  parameters. This figure shows dependence on  $W$  for a fixed training set with  $N$  examples. The network is again a 5-layer all convolutional network (i.e., no pooling). All hidden layers have the same number of channels. Neither data augmentation nor regularization is performed. The classical theory explains the generalization behavior on the left; the challenge is to explain the lack of overfitting for  $W > N$ . As expected from Theory II<sup>[2]</sup>, there is zero error as soon as  $W = N$  and for  $W \geq N$ .

Our argument can be loosely described as follows. In the overparametrized regime, SGD selects with high probability large volume minima, that mostly correspond to flat minima, because those correspond to large regions of zero loss. Furthermore for flat minima SGD and GD converge to the pseudoinverse-based solution for the unknown weights. This means that the solution found is the minimum norm one among the infinite degenerate solutions that correspond to zero empirical error. In this sense, those minimum norm solutions also correspond to large margin solutions: the intuition is that they are robust to perturbations of the weights and thus to corresponding perturbations of the input vector. The last step in the argument is that large margin solutions generalize well in the case of classification. Thus increasing the number of parameters does not change this regularizing behavior of SGD, exactly as in the linear case. As a side remark, we will later show that SGD is similar to robust optimization, providing implicit regularization.

A similar argument underlies the following intuition that links Theory II with this memo: the unique zero minimizers for  $n < W$  but close to  $W$  become degenerate, that is flat, for  $n \leq W$ . Among them SGD chooses the one with the minimum  $L_2$  norm and thus with the best generalization. Notice that overparametrization *does not necessarily improve the test error* which is optimal around the the boundary between over- and under-parametrization.

We consider the steps of our argument in turn, *starting with properties of SGD that have been so far unrecognized* from the machine learning point of view, to the best of our knowledge.



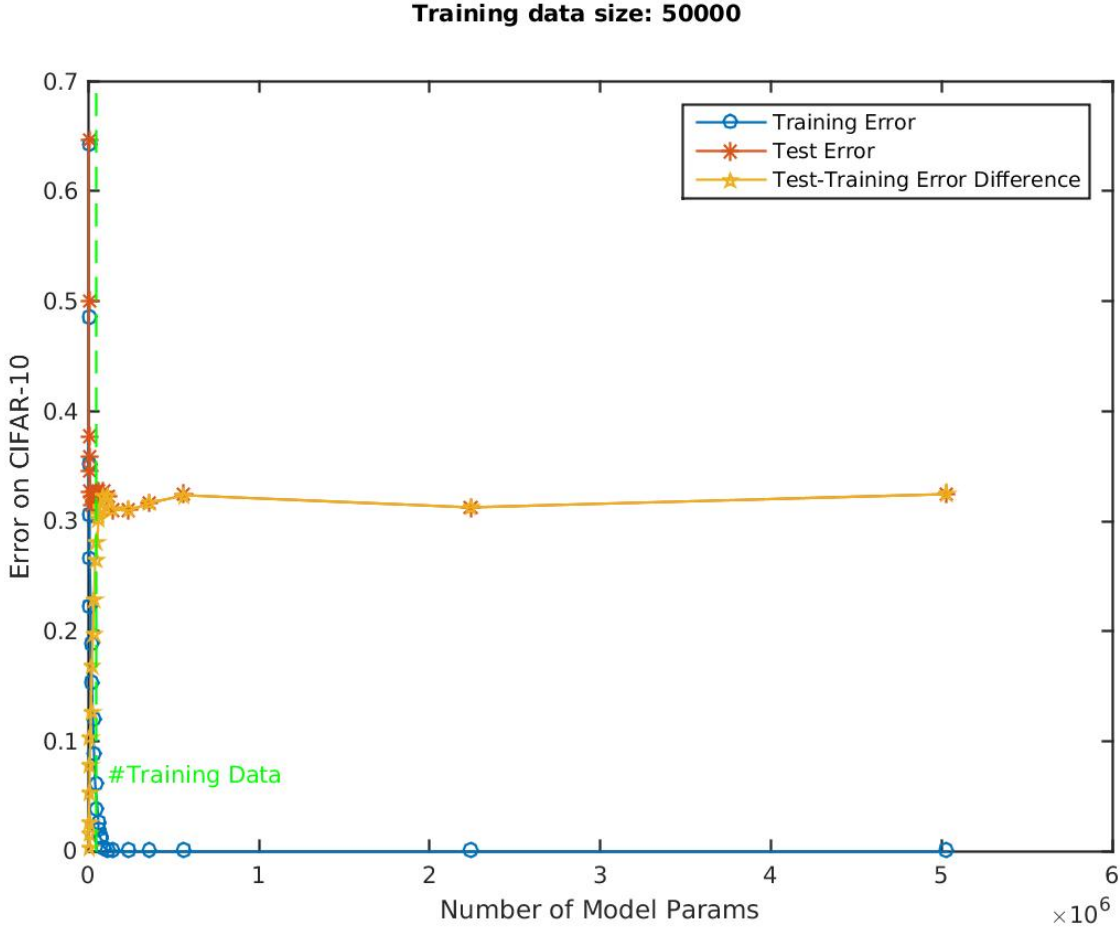


Figure 4: The plot is the same as Figure 3 but with a linear scale for the number of parameters  $W$ .

## 5.1 SGD as an approximate Langevin equation

We consider the usual SGD update defined by the recursion

$$f_{t+1} = f_t - \gamma_t \nabla V(f_t, z_t), \quad (6)$$

where  $z_t$  is fixed,  $\nabla V(f_t, z_t)$  is the gradient of the loss with respect to  $f$  at  $z_t$ , and  $\gamma_t$  is a suitable decreasing sequence. When  $z_t \subset [n]$  is a minibatch, we overload the notation and write  $\nabla V(f_t, z_t) = \frac{1}{|z_t|} \sum_{z \in z_t} \nabla V(f_t, z)$ .

We define a noise “equivalent quantity”

$$\xi_t = \nabla V(f_t, z_t) - \nabla I_{S_n}(f_t), \quad (7)$$

and it is clear that  $\mathbb{E}\xi_t = 0$ .

We write Equation 6 as

$$f_{t+1} = f_t - \gamma_t (\nabla I_{S_n}(f_t) + \xi_t). \quad (8)$$

With typical values used in minibatch (each minibatch corresponding to  $z_t$ ) training of deep nets, it turns out that the vector of gradient updates  $\nabla V(f_t, z_t)$  empirically shows components with approximate Gaussian distributions (see Figure 5). This is expected because of the Central Limit Theorem (each minibatch involves sum over many random choices of datapoints).

Training error: 0.000965

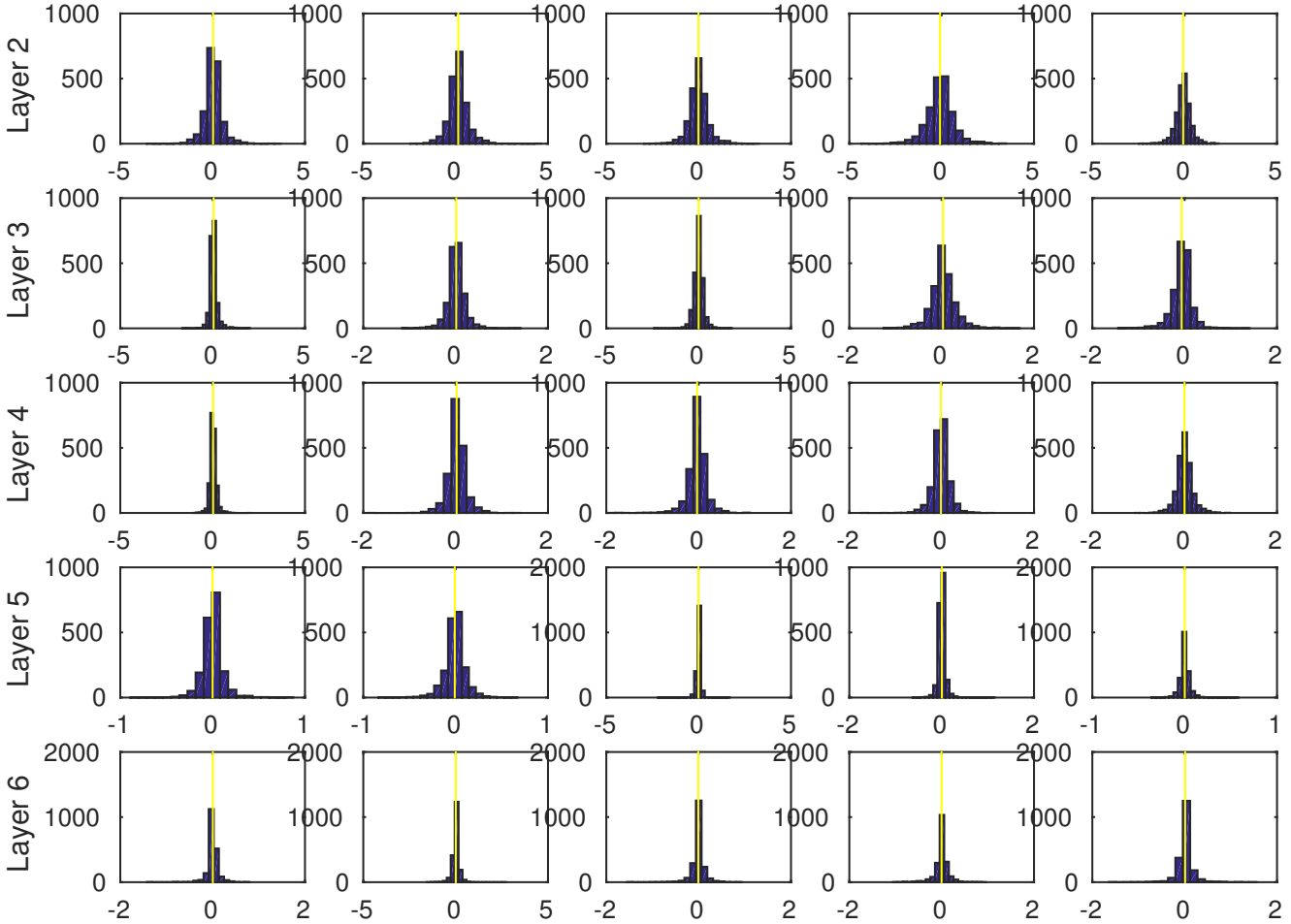


Figure 5: Histograms of some of the components of  $\nabla V(f_t, z_i)$  over  $i$  for fixed  $t$  in the asymptotic regime. Notice that the average corresponds to the gradient of the full loss and that is empirically very small. The histograms look approximately Gaussian as expected (see text) for minibatches that are not too small or too large.

Now we observe that (8) is a discretized Langevin diffusion, albeit with a noise scaled as  $\gamma_n$  rather than  $\sqrt{\gamma_n}$ . In fact, the continuous SGD dynamics corresponds to a stochastic gradient equation using a potential function defined by  $U = I_{S_n}[f] = \frac{1}{n} \sum_{i=1}^n V(f, z_i)$  (see Proposition 3 and section 5 in [22]). If the noise is the derivative of the Brownian motion, it is a Langevin equation – that is a stochastic dynamical system – with an associated Fokker-Planck equation on the probability distributions of  $f_t$ . The asymptotic probability distribution is the Boltzman distribution that is  $\approx e^{-\frac{U}{\gamma K}}$ .

For more details, see for instance section 5 of [23]. Several proofs that adding a white noise term to equation (6) will make it converge to a global minimum are available (see [24]). Notice that the discrete version of the Langevin dynamics is equivalent to a Metropolis-Hastings algorithm for small learning rate (when the rejection step can be neglected).

## 5.2 SGDL concentrates at large volume, “flat” minima

The argument about convergence of SGDL to large volume minima that we call “flat”, is straightforward. The asymptotic distribution reached by a Langevin equation (GDL) – as well as by SGDL – is the Boltzman distribution that is

$$p(f) = \frac{1}{Z} e^{-\frac{U}{T}}, \tag{9}$$

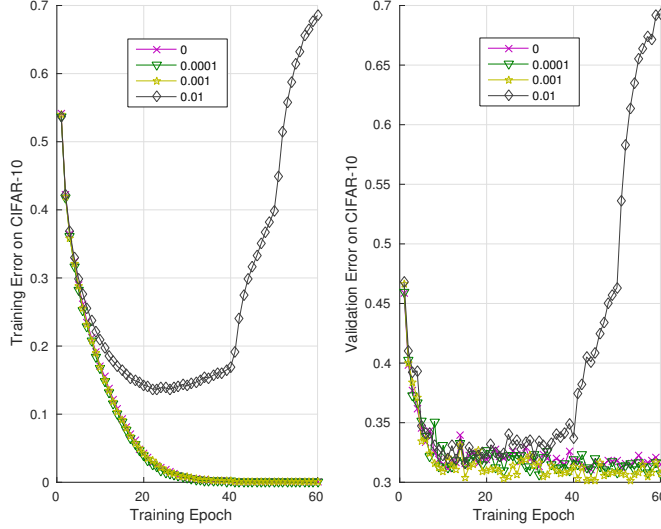


Figure 6: Equation 5 – that is SGD with added Gaussian (with constant power) – behaves in a similar way to standard SGD. Notice that SGDL has slightly better validation performance than SGD.

where  $Z$  is a normalization constant,  $U$  is the loss and  $T$  reflects the noise power. The equation implies, and Figure 8 shows, that SGD prefers degenerate minima relative to non-degenerate ones of the same depth. In addition, among two minimum basins of equal depth, the one with a larger volume, is much more likely in high dimensions (Figure 7). Taken together, these two facts suggest that SGD selects degenerate minimizers and, among those, the ones corresponding to larger isotropic flat regions of the loss. Suppose the landscape of the empirical minima is well-behaved in the sense that deeper minima have broader basin of attraction. Then it is possible to prove that SDGL shows concentration – *because of the high dimensionality* – of its asymptotic distribution Equation 9 – to minima that are the most robust to perturbations of the weights. Notice that these assumptions are satisfied in the zero error case: among zero-minimizer, SGDL selects the ones that are flatter, i.e. have the larger volume<sup>2</sup>.

**Conjecture 2.** *Under regularity assumptions on the landscape of the empirical minima, SGDL corresponds to the following robust optimization*

$$\min_w \max_{(\delta_1, \dots, \delta_n)} \frac{1}{n} \sum_{i=1}^n V(y_i, f_{w+\delta_i}(x_i)). \quad (10)$$

Notice that isotropic and non-isotropic degeneracy is expected to be the key factor for SGDL to converge to zero-minimizers. It is especially important to note that *SGDL and SGD maximize volume and “flatness” of the loss in weight space*. Given a flat minimum, one may ask where SGD will converge to. For situations such as in Figure 9 and for a minimum such as in Figure 10, arguments similar to Appendix 9.6 may hold approximately, suggesting a locally minimum norm solution (see Equation 34). In particular, the weight values found by SGD are expected to be mostly around their average value over the flat region (at least in the case of square loss). We can also consider the polynomial describing  $f(x)$  – the function computed by the network when each RELU is replaced by a univariate polynomial approximant. In this case  $f(x)$  is a linear function in the vector  $X$  comprising all the relevant monomials, implying that SGD converges to the minimum norm solution in  $X$ .

### 5.3 SGD finds minimum norm solutions (and maximizes geometrical margin)

Traditionally the functional margin of a correctedly classified data point  $y_i, x_i$  is defined as the quantity  $\gamma = y_i f(x_i)$ <sup>3</sup>. For a linear classifier  $f(x) = Wx$ , the geometrical margin is the distance *in input space* from the classification boundary and is maximized by minimizing  $\|W\|_2$ . In the nonlinear case it is difficult in general to relate the functional margin  $y_i f(x_i)$  to the geometrical margin.

Flatness of the minima around some  $w_0$  values of the weight vector suggests (to an ex-physicist) the following linear approximation of the network function  $f$  realized by the network in the neighborhood of the flat minimum:

$$f(x; w) - f(x; w_0) \approx ([\nabla_w f(x)]_{w_0}, \Delta w) = (\phi(x), \Delta w) \quad (12)$$

<sup>2</sup>Given a zero minimizer there is no other minimizer that has smaller volume AND is deeper.

<sup>3</sup>A perhaps more satisfying approach is the following: propose a new definition of margin that directly capture the idea of measuring robustness to perturbations, show that it is equivalent to the old definition of geometrical margin in the linear case. The new definition we propose is

**Definition 1.** *For a given functional margin  $y_i f(x_i)$ , the geometrical margin is*

$$\gamma = \max_{\|\Delta x\|_2} | \text{sign}(y_i f(x_i)) = \text{sign}(y_i f(x_i + \Delta x)) |. \quad (11)$$

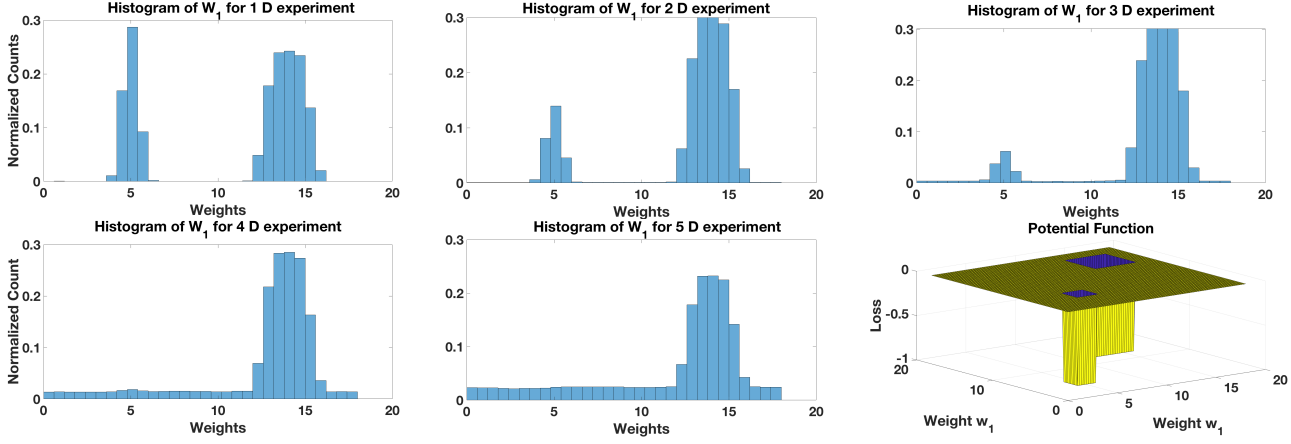


Figure 7: The figure shows the histogram of a one-dimensional slice of the asymptotic distribution obtained by running Langevin Gradient Descent (GDL) on the potential surface on the right. The potential function has two minima: they have the same depth but one has a flat region which is a factor 2 larger in each of the dimensions. The 1D histogram for the first weight coordinate is shown here for dimensionality 1, 2, 3, 4 and 5D. The figures graphically show – as expected from the asymptotic Boltzman distribution – that noisy gradient descent selects with high probability minimizers with larger margin. As expected, higher dimensionality implies higher probability of selecting the flatter minimum.

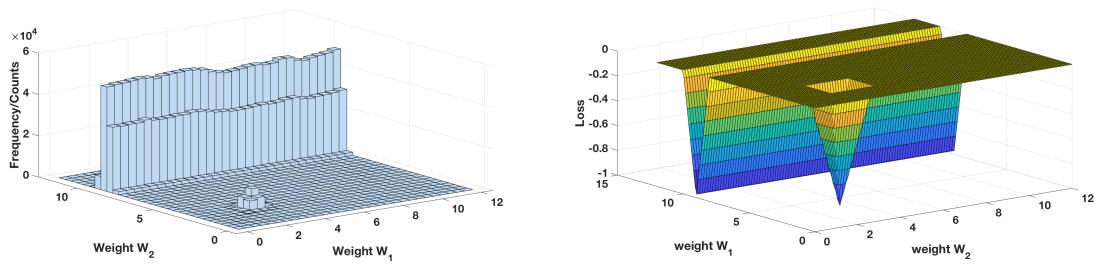


Figure 8: Langevin Gradient Descent (GDL) on the 2D potential function shown above leads to an asymptotic distribution with the histogram shown on the left. As expected from the form of the Boltzman distribution, the Langevin dynamics prefers degenerate minima to non-degenrate minima of the same depth. In high dimensions we expect the asymptotic distribution to concentrate strongly around the degenerate minima as confirmed on figure 9.

where  $(\cdot, \cdot)$  is the scalar product and  $\phi$  is a vector function of the input  $x$  but with the dimensionality of the weight vector, parametrized by the weights (set to  $w_0$ ).

Assume for now that  $w_0 \approx 0$ . Then GD or SGD iterations starting at the minimum of the loss function would converge (see [25]) to the pseudoinverse solution of the linear problem that is

$$\Delta w \approx \Phi = E \tag{13}$$

with solution

$$\Delta w = E\Phi^\dagger \tag{14}$$

where  $\Phi$  is the matrix with columns the  $\phi(x_i)$  vectors evaluated at each of the training data  $x_i$ ,  $E_i = y_i - f(x_i, w_0)$  are the errors. Since the pseudoinverse regularizes and generalizes this would be the solution. The key problem is the assumption  $w_0 = 0$ : there is no obvious guarantee that  $w_i$  or combinations of them corresponding to degenerate directions will have zero or small norm at the minima. The result above only says that such directions will not change. While the Appendix provides a more complete result we state here the conclusion of our argument

**Lemma 1.** Assume that SGD (or GD) reaches a zero-error minimum after  $T$  iterations and the associated weights have small norm. Further iterations are guaranteed to lead to a minimum norm change of parameters. In particular, the  $L_2$  norm of the weights associated with degenerate directions will not increase.

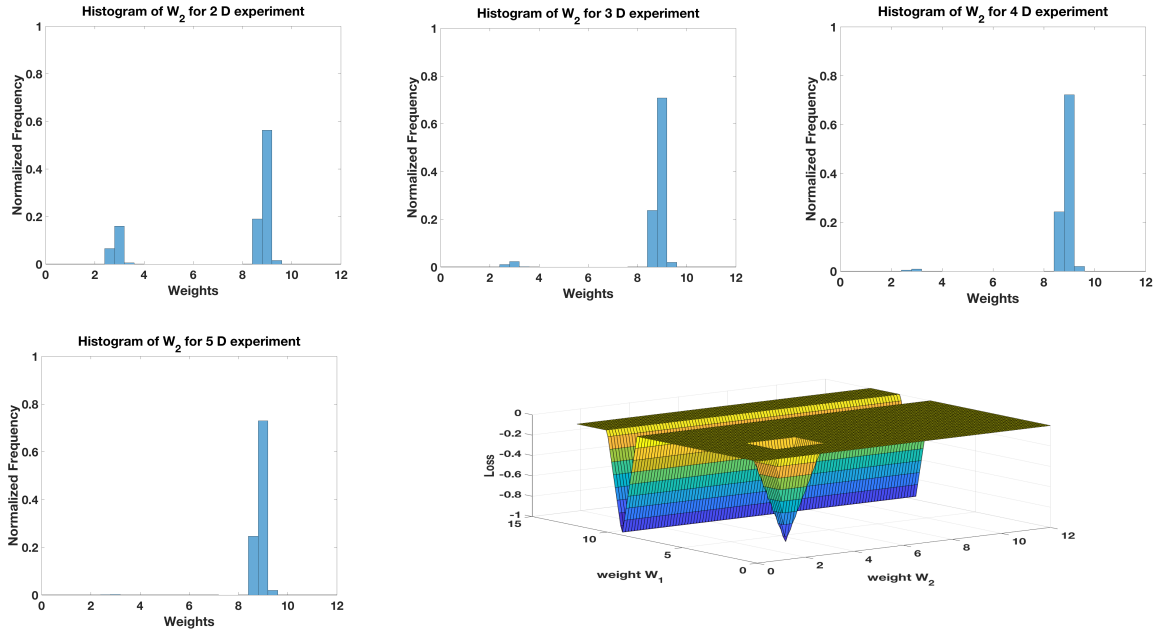


Figure 9: The figure shows the histogram of a one-dimensional slice of the asymptotic distribution obtained by running Langevin Gradient Descent (GDL) on the potential surface on the right. As expected from the form of the Boltzman distribution, the Langevin dynamics prefers degenerate minima to non-degenerate minima of the same depth. Furthermore, as dimensions increase the distribution concentrates strongly around the degenerate minima. This can be appreciated from the figure because the histogram density at  $W_1 = 2$  (the degenerate minimum) decreases in density rapidly as dimensions increase.

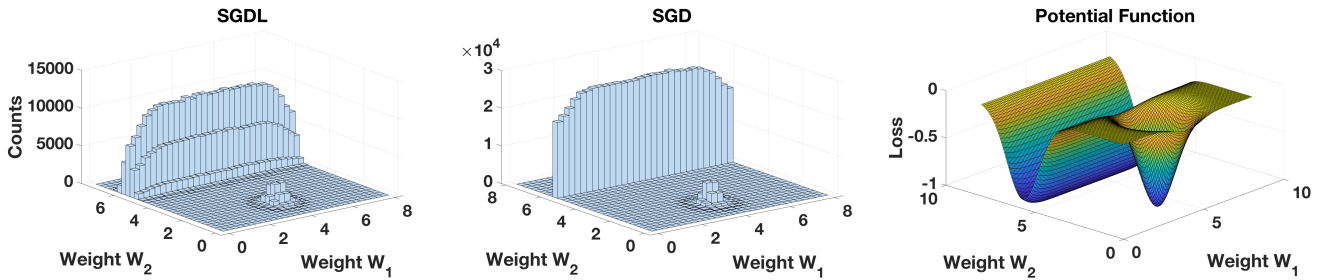


Figure 10: Stochastic Gradient Descent and Langevin Stochastic Gradient Descent (SGDL) on the 2D potential function shown above leads to an asymptotic distribution with the histograms shown on the left. As expected from the form of the Boltzman distribution, both dynamics prefers degenerate minima to non-degenerate minima of the same depth.

Another approach to show that SGD provides large margin is based on the equivalence between SGD and forms of robust optimization and is presented in Appendix 9.3. The proofs are rigorous only in the case of linear networks and approximate otherwise. The previous argument only holds within the validity of a linear approximation.

Thus SGD should select minimizers corresponding to maximum geometrical margin.

#### 5.4 “Equivalence” between perturbations in weights and in input

The following example presents a difficulty in any argument that establishes some sort of relation to robustness to perturbations in weights and in the input data.

Let  $x \in \mathbb{R}$ ,  $w \in \mathbb{R}^{2d}$ ,  $f(w, x) = x \cdot \sum_{j=1}^d w_{1,d} w_{2,d}$ .  $f$  may be viewed as a very simple 1-hidden layer neural network with linear activation functions, with  $d$  hidden units in the first hidden layer. Note that  $f$  is in fact fully connected.

Note that  $\frac{\partial f}{\partial x} = \sum_{j=1}^d w_{1,d} w_{2,d}$ , for  $i \in [d]$ ,

$$\frac{\partial f}{\partial w_{1,i}} = x w_{2,i}, \quad \frac{\partial f}{\partial w_{2,i}} = x w_{1,i}.$$

Now suppose that  $w_{j,i} = 1$  for all  $i \in [d], j \in [2]$  (this may correspond to the solution found for some training set by gradient descent). Note that for all  $|x| \leq b$ , we have that  $\left| \frac{\partial f}{\partial w_{j,i}} \right| \leq b$ , but that  $\frac{\partial f}{\partial x} = d \gg b$  for small constants  $b$ .

In light of this example, any result we prove cannot hold for layers of arbitrary sizes, and any sort of relation will have a factor of  $d$ , where  $d$  grows at least as quickly as the size of the first hidden layer.

## 6 SGD puzzles

Two puzzling properties of deep networks that we can now account for are

- the no generalization behavior tested with training data with random labels, see<sup>[4]</sup>, while expected error is much better than chance for natural labels;
- the behavior shown in Figure 13 and Figure 15: the expected error, which decreases with increasing training set size, is not affected by increasing the number of parameters beyond the training data size.

### 6.1 Random Labels

In the first case, Theory II predicts that it is in fact possible to interpolate the data on the training set, that is to achieve zero empirical error (because of overparametrization) and that this is in fact easy – because of the very high number of zeros of the polynomial approximation of the network– assuming that the target function is in the space of functions realized by the network. For  $n$  going to infinity we expect that the empirical error will converge to the expected (chance), as shown in Figures 1 and 2. For finite  $n$  when  $n < W$ , the fact that the empirical error (which is zero) is so different from the expected is puzzling, as observed by<sup>[4]</sup>, especially because the algorithm is capable of low expected error with the same  $n$  for natural labels (and also when trained with random labels but tested with perturbations of the training images, see Figures in Appendix 9.5).

The puzzle is explained in terms of the distribution-dependent geometrical margin in the linearized case. Our theory suggests that SGD maximizes margin since minimum norm corresponds to maximum margin in the linear case. A larger margin is in fact found for natural labels than for random labels as shown in Table 1 and in Figure 11 and Figure 12. Figure 11 shows “three-point interpolation” plots to illustrate the flatness of the landscape around global minima of the empirical loss found by SGD, on CIFAR-10, with natural labels and random labels, respectively. Specifically, let  $w_1, w_2, w_3$  be three minimizers for the empirical loss found by SGD. For  $\lambda = (\lambda_1, \lambda_2, \lambda_3)$  on the simplex  $\Delta_3$ , let

$$w_\lambda = \lambda_1 w_1 + \lambda_2 w_2 + \lambda_3 w_3 \quad (15)$$

We then evaluate the training accuracy for the model defined by each interpolated weights  $w_\lambda$  and make a surface plot by embedding  $\Delta_3$  in the 2D X-Y plane. As we can see, the natural label case depict a larger flatness region around each of the three minima than the random label case. Section 7.1 shows a direct relation between the range of flatness and the norm  $\lambda$  of the perturbations. Note that  $\lambda$  is also the regularization parameter (see section 9.4: *larger  $\lambda$  means greater regularization, which implies better stability ( $\beta$ -stability of regularization has  $\beta \approx \frac{K}{\lambda n}$ ) and better generalization bounds.*

The same phenomenon could be observed more clearly on the MNIST dataset, where the images of the same category are already quite similar to each other in the pixel space, making it more difficult to fit when random labels are used. Therefore, the difference in the characteristics of the landscapes is amplified. As shown in Figure 12, big flat regions could be observed in the natural label case, while the landscape for the random label experiment resembles sharp wells.

It is difficult to visualize the flatness of the landscape when the weights are typically in the scale of one million dimensions. To assess the isotropic flatness, we employ the following procedure around a minimum found by SGD: choose a random isotropic direction  $\delta w$  with  $\|\delta w\| = 1$ , perform a line search to find the “flatness radius” in that direction:

$$r(w, \delta w, \varepsilon) = \sup\{r : |\hat{I}(w) - \hat{I}(w + r\delta w)| \leq \varepsilon\} \quad (16)$$

The procedure is repeated  $T$  times and the average radius is calculated. The overall procedure is also repeated multiple times to test the average flatness at different minima. The results are shown in Table 1. For both CIFAR-10 and MNIST, we observe a difference between the natural label and random label.

### 6.2 No Overfitting with Increasing Network Size

In the second case, what is important for low expectation error is not parameter number but size of the training set. The intuition comes from Theory II: minimization by SGD finds global minima that are similar to minimum norm solutions – which are not changed by increasing the number of equations beyond the number of data points. Local minima instead are determined by the number of parameters since they determine the number of equations for the critical points of the gradient: this implies that increasing the number of parameters much beyond the number of data in the training set, should not affect generalization but may make it easier to find global minima relative to local minima. In summary overparametrization in Figure 4 for  $W > n$  does not affect the test error because SGD will disregard degenerate directions.

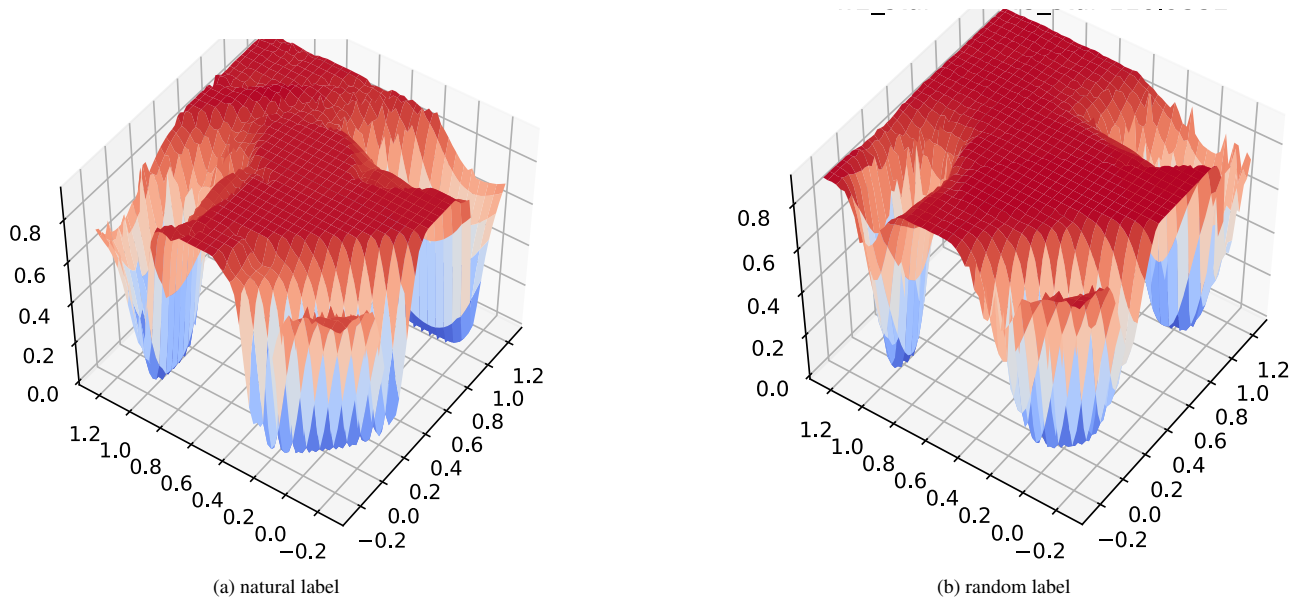


Figure 11: Illustration of the landscape of the empirical loss on CIFAR-10.

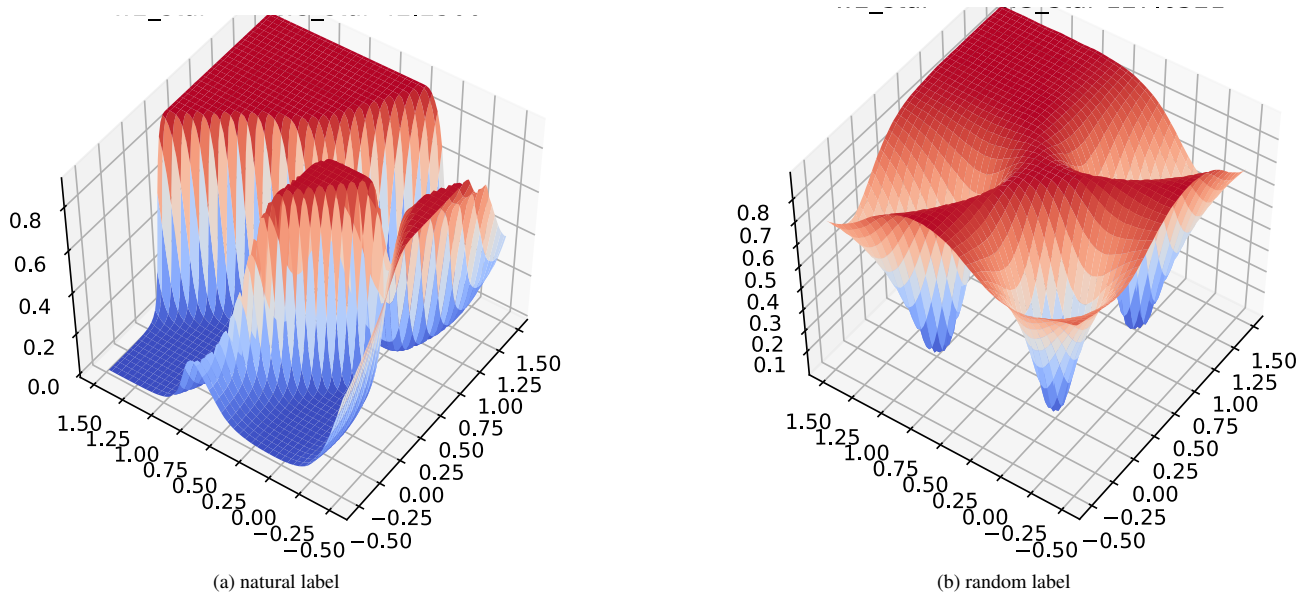


Figure 12: Illustration of the landscape of the empirical loss on MNIST.

	MNIST	CIFAR-10
all params	$45.4 \pm 2.7$	$17.0 \pm 2.4$
all params (random label)	$6.9 \pm 1.0$	$5.7 \pm 1.0$
top layer	$15.0 \pm 1.7$	$19.5 \pm 4.0$
top layer (random label)	$3.0 \pm 0.1$	$12.1 \pm 2.6$

Table 1: The “flatness test”: at the minimizer, we move the weights around in a random direction, and measure the furthest distance until the objective function is increased by  $\epsilon$  (0.05), and then measure the average distance.

Architecture	Number of Parameters	Training Accuracy	Test Accuracy
MLP 1x512	1,209,866	100.0	50.51
Alexnet	1,387,786	100.0	76.07
Inception	1,649,402	100.0	85.75
Wide Resnet	8,949,210	100.0	88.21

Table 2: A number of different network architectures are trained on CIFAR-10. We turn off all the regularizers in order to avoid implicitly constraining the hypothesis space size. We found that despite the network size continuously increases, the test performance does not drop, but even improves.

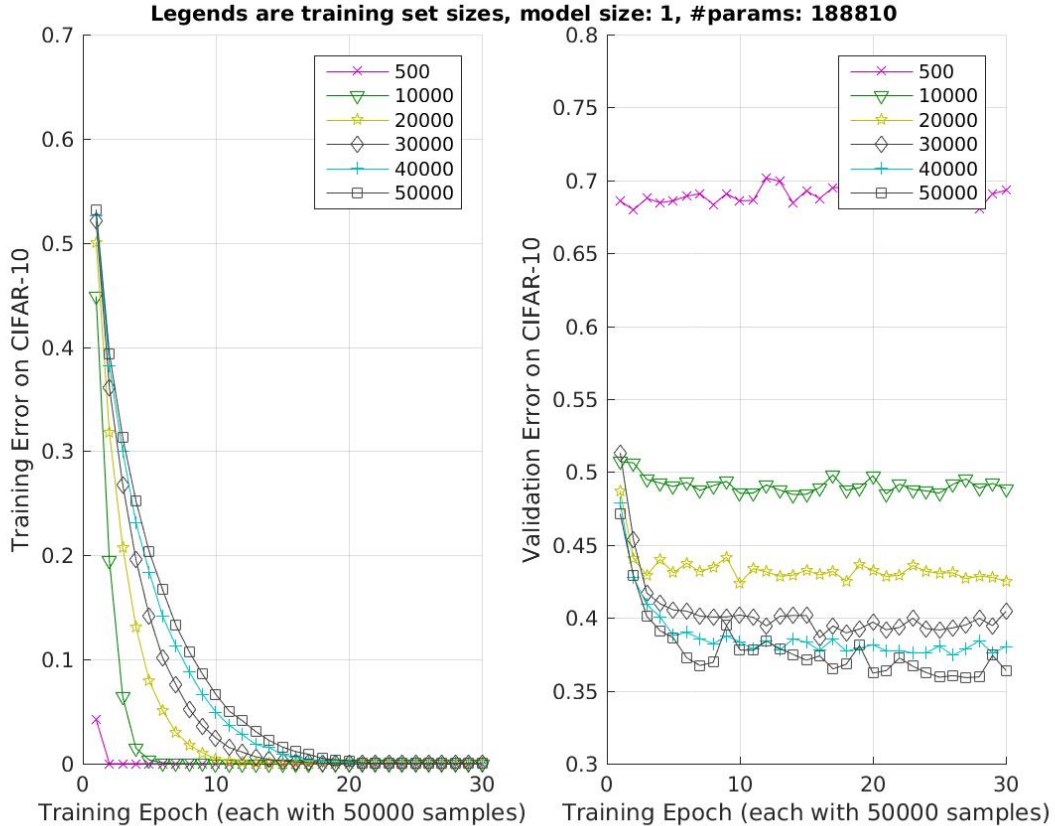


Figure 13: Training with different training set sizes (from 500 to 50,000, one per curve) on CIFAR-10. It shows what happens to empirical and testing error if one severely reduces the number of training examples. There is overfitting and little generalization with small training sets.

## 7 Discussion

### 7.1 Comments on Robustness wrt Weights and Robustness wrt Data

A natural intuition, that we discuss further in Appendix 9.7, is that several forms of stability are closely related. In particular, *perturbations of the weights follow from perturbations of the data points in the training set* because the function  $f(x)$  resulting from the training is parametrized by the weights that depend on the data points  $S_n$ :  $f(x) = f(x; w(z_1, \dots, z_n))$ , with  $z = (x, y)$ <sup>4</sup>. Thus *isotropic flat regions in weight space of the global minimum of the empirical loss indicate stability with respect to the weights; the latter in turn indicates stability with respect to training examples*.

<sup>4</sup>In a differentiable situation, one would write  $df = \frac{df}{dw} \frac{dw}{dS} dS$  and  $dw = \frac{dw}{dS} dS$ . The latter equation would show that perturbations in the weights depend on perturbations of the training set  $S$ .



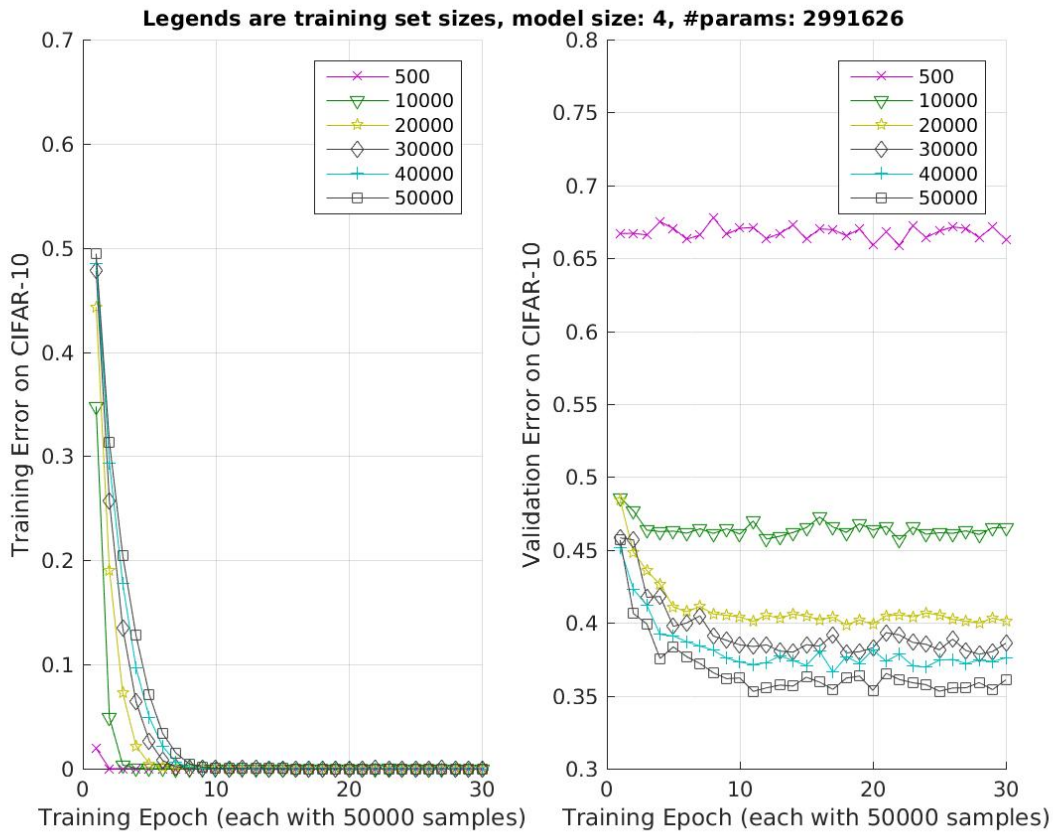


Figure 14: Training with different training set sizes (from 500 to 50,000, one per curve) on CIFAR-10. This is a larger model with more parameters.

## 7.2 Summary: why and when SGD has low expected error and why and when it generalizes

We conjecture that low expected error in the regime of  $N < W$  and of generalization in the regime  $N > W$  exhibited by deep convolutional networks in multi-class tasks such as CIFAR10 and Imagenet are due to three main factors:

- the implicit regularizing effect of SGD
- the task is compositional
- the task is multiclass.

Speed of convergence of SGD is not the deep reason for good generalization. We expect the speed of convergence to be correlated with good generalization because convergence will depend on the relative size of the basins of attraction of the minima of the empirical risk, which in turn depend on the ratio between the effective dimensionality of the minima and the ambient dimensionality. Theory III in this paper, together with Theory II, answers the question about the unusual properties of Stochastic Gradient Descent used for training overparametrized deep convolutional networks. SGDL and SGD select with high probability solutions with zero or small empirical error (Theory II) – because they are flatter. It is interesting that among these solutions GD and SGD select the best in achieving low expected error – because they are the minimum norm solutions and thus have a large margin, are the most stable wrt parameters, most stable wrt training data and most stable with respect to the  $x$  value.

## 7.3 Open questions

The main task is to quantify the qualitative characterization of generalization and consistency of convolutional deep networks given here. In particular, which non-vacuous bounds may relate margin to expected error? Can we spell conditions on data distributions that yield large margin and low expected error? Can we improve SGD or find optimal parameters such as mini-batch size?

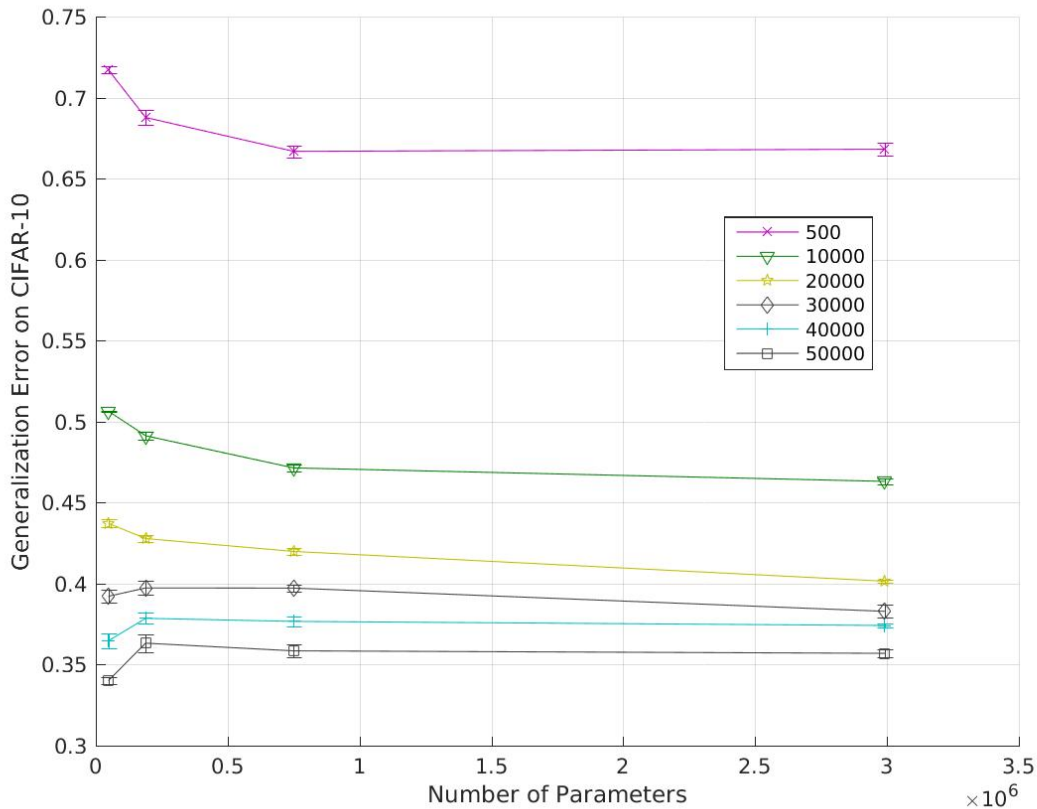


Figure 15: Generalization error as a function of number of parameters in the network. Generalization error is defined as the difference between training error and validation error. We train a 4-hidden-layer convolutional network with different training set sizes on CIFAR-10. Different curves correspond to different training set sizes (from 500 to 50,000). All models are trained 30 epochs. Increasing the number of parameters does not hurt generalization.

## Acknowledgment

This work was supported by the Center for Brains, Minds and Machines (CBMM), funded by NSF STC award CCF – 1231216. We gratefully acknowledge the support of NVIDIA Corporation with the donation of the DGX-1 used for this research.

## References

- [1] T. Poggio, H. Mhaskar, L. Rosasco, B. Miranda, and Q. Liao, “Why and when can deep - but not shallow - networks avoid the curse of dimensionality: a review,” tech. rep., MIT Center for Brains, Minds and Machines, 2016.
- [2] T. Poggio and Q. Liao, “Theory ii: Landscape of the empirical risk in deep learning,” *arXiv:1703.09833, CBMM Memo No. 066*, 2017.
- [3] S. Shalev-Shwartz and S. Ben-David, *Understanding Machine Learning: From Theory to Algorithms*. Cambridge eBooks, 2014.
- [4] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, “Understanding deep learning requires rethinking generalization,” in *International Conference on Learning Representations (ICLR)*, 2017.
- [5] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *Nature*, pp. 436–444, 2015.
- [6] K. Fukushima, “Neocognitron: A self-organizing neural network for a mechanism of pattern recognition unaffected by shift in position,” *Biological Cybernetics*, vol. 36, no. 4, pp. 193–202, 1980.
- [7] D. Hubel and T. Wiesel, “Receptive fields, binocular interaction and functional architecture in the cat’s visual cortex,” *The Journal of Physiology*, vol. 160, no. 1, p. 106, 1962.
- [8] M. Riesenhuber and T. Poggio, “Hierarchical models of object recognition in cortex,” *Nature Neuroscience*, vol. 2, pp. 1019–1025, Nov. 1999.
- [9] P. Chaudhari, A. Choromanska, S. Soatto, Y. LeCun, C. Baldassi, C. Borgs, J. Chayes, L. Sagun, and R. Zecchina, “Entropy-SGD: Biasing Gradient Descent Into Wide Valleys,” *arXiv:1611.01838 [cs]*, Nov. 2016. arXiv: 1611.01838.

- [10] N. S. Keskar, D. Mudigere, J. Nocedal, M. Smelyanskiy, and P. T. P. Tang, "On Large-Batch Training for Deep Learning: Generalization Gap and Sharp Minima," *arXiv:1609.04836 [cs, math]*, Sept. 2016. arXiv: 1609.04836.
- [11] L. Dinh, R. Pascanu, S. Bengio, and Y. Bengio, "Sharp minima can generalize for deep nets," *arXiv preprint arXiv:1703.04933*, 2017.
- [12] G. Karolina Dziugaite and D. M. Roy, "Computing Nonvacuous Generalization Bounds for Deep (Stochastic) Neural Networks with Many More Parameters than Training Data," *ArXiv e-prints*, Mar. 2017.
- [13] B. Neyshabur, R. Tomioka, and N. Srebro, "Geometry of Optimization and Implicit Regularization in Deep Learning," *arXiv preprint arXiv:1705.03071*, 2017.
- [14] B. Neyshabur, R. Tomioka, and N. Srebro, "Norm-Based Capacity Control in Neural Networks," *arXiv:1503.00036 [cs, stat]*, Feb. 2015. arXiv: 1503.00036.
- [15] V. Koltchinskii and D. Panchenko, "Empirical margin distributions and bounding the generalization error of combined classifiers," *Annals of Statistics*, pp. 1–50, 2002.
- [16] P. L. Bartlett, "The sample complexity of pattern classification with neural networks: the size of the weights is more important than the size of the network," *IEEE transactions on Information Theory*, vol. 44, no. 2, pp. 525–536, 1998.
- [17] P. Bartlett, D. J. Foster, and M. Telgarsky, "Spectrally-normalized margin bounds for neural networks," *ArXiv e-prints*, June 2017.
- [18] Y. Freund and R. E. Schapire, "A decision-theoretic generalization of on-line learning and an application to boosting," *J. Comput. Syst. Sci.*, vol. 55, pp. 119–139, Aug. 1997.
- [19] S. Ma and M. Belkin, "Diving into the shallows: a computational perspective on large-scale shallow learning," *ArXiv e-prints*, Mar. 2017.
- [20] S. Gelfand and S. Mitter, "Recursive stochastic algorithms for global optimization in  $R^d$ ," *Siam J. Control and Optimization*, vol. 29, pp. 999–1018, September 1991.
- [21] L. Bottou, "Online algorithms and stochastic approximations," in *Online Learning and Neural Networks* (D. Saad, ed.), Cambridge, UK: Cambridge University Press, 1998. revised, oct 2012.
- [22] D. Bertsekas and J. Tsitsiklis, "Gradient Convergence in Gradient methods with Errors," *SIAM J. Optim.*, vol. 10, pp. 627–642, Jan. 2000.
- [23] D. P. Bertsekas and J. N. Tsitsiklis, *Neuro-dynamic Programming*. Athena Scientific, Belmont, MA, 1996.
- [24] B. Gidas, "Global optimization via the Langevin equation," *Proceedings of the 24th IEEE Conference on Decision and Control*, pp. 774–778, 1985.
- [25] L. Rosasco and S. Villa, "Learning with incremental iterative regularization," in *Advances in Neural Information Processing Systems*, pp. 1630–1638, 2015.
- [26] M. Anthony and P. Bartlett, *Neural Network Learning - Theoretical Foundations*. Cambridge University Press, 2002.
- [27] S. Bubeck *et al.*, "Convex optimization: Algorithms and complexity," *Foundations and Trends® in Machine Learning*, vol. 8, no. 3-4, pp. 231–357, 2015.
- [28] H. Xu and S. Mannor, "Robustness and generalization," *CoRR*, vol. abs/1005.2243, 2010.
- [29] N. Srebro, K. Sridharan, and A. Tewari, "Smoothness, low noise and fast rates.," in *NIPS* (J. D. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R. S. Zemel, and A. Culotta, eds.), pp. 2199–2207, Curran Associates, Inc., 2010.
- [30] S. Gelfand and S. Mitter, "Recursive Stochastic Algorithms for Global Optimization in  $\mathbb{R}^d$ ," *SIAM J. Control Optim.*, vol. 29, pp. 999–1018, Sept. 1991.
- [31] S. B. Gelfand and S. K. Mitter, "Metropolis-Type Annealing Algorithms for Global Optimization in  $R^d$ ," *SIAM Journal on Control and Optimization*, vol. 31, no. 1, pp. 111–131, 1993.
- [32] S. Rajasekaran, "On the Convergence Time of Simulated Annealing," *Technical Reports (CIS)*, Nov. 1990.
- [33] M. Raginsky, A. Rakhlin, and M. Telgarsky, "Non-convex learning via stochastic gradient langevin dynamics: A nonasymptotic analysis," *arXiv:180.3251 [cs, math]*, 2017.
- [34] S. Mandt, M. D. Hoffman, and D. M. Blei, "A Variational Analysis of Stochastic Gradient Algorithms," *arXiv:1602.02666 [cs, stat]*, Feb. 2016. arXiv: 1602.02666.
- [35] F. Anselmi, L. Rosasco, and T. Tan, C. and Poggio, "Deep Convolutional Networks are Hierarchical Kernel Machines," *Center for Brains, Minds and Machines (CBMM) Memo No. 35*, also in *arXiv*, 2015.
- [36] A. Ben-Tal, L. El Ghaoui, and A. Nemirovski, *Robust Optimization*. Princeton Series in Applied Mathematics, Princeton University Press, October 2009.
- [37] H. Xu, C. Caramanis, and S. Mannor, "Robustness and regularization of support vector machines," *J. Mach. Learn. Res.*, vol. 10, pp. 1485–1510, Dec. 2009.
- [38] M. Hardt, B. Recht, and Y. Singer, "Train faster, generalize better: Stability of stochastic gradient descent," *CoRR*, vol. abs/1509.01240, 2015.
- [39] S. Mukherjee, P. Niyogi, T. Poggio, and R. Rifkin, "Learning theory: stability is sufficient for generalization and necessary and sufficient for consistency of empirical risk minimization," *Advances in Computational Mathematics*, vol. 25, no. 1, pp. 161–193, 2006.
- [40] P. L. Bartlett and S. Mendelson, "Rademacher and gaussian complexities: Risk bounds and structural results," *Journal of Machine Learning Research*, vol. 3, no. Nov, pp. 463–482, 2002.
- [41] S. M. Kakade, K. Sridharan, and A. Tewari, "On the complexity of linear prediction: Risk bounds, margin bounds, and regularization," in *Advances in neural information processing systems*, pp. 793–800, 2009.
- [42] B. M. Lake, R. Salakhutdinov, and J. B. Tenenbaum, "Human-level concept learning through probabilistic program induction," *Science*, vol. 350, no. 6266, pp. 1332–1338, 2015.
- [43] A. Maurer, "Bounds for Linear Multi-Task Learning," *Journal of Machine Learning Research*, 2015.
- [44] C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, "Understanding deep learning requires rethinking generalization," *CoRR*, vol. abs/1611.03530, 2016.

- [45] H. N. Mhaskar, "Approximation properties of a multilayered feedforward artificial neural network," *Advances in Computational Mathematics*, vol. 1, pp. 61–80, 1993.
- [46] T. Poggio, "On optimal nonlinear associative recall," *Biological Cybernetics*, vol. 19, pp. 201–209, 1975.
- [47] J. Lin and L. Rosasco, "Optimal learning for multi-pass stochastic gradient methods," in *Advances in Neural Information Processing Systems*, pp. 4556–4564, 2016.
- [48] M. Hardt, B. Recht, and Y. Singer, "Train faster, generalize better: Stability of stochastic gradient descent," *Proceedings of the International Conference on Machine Learning*, 2016.
- [49] M. Hardt and T. Ma, "Identity matters in deep learning," *CoRR*, vol. abs/1611.04231, 2016.
- [50] B. Neyshabur, S. Bhojanapalli, D. McAllester, and N. Srebro, "Exploring generalization in deep learning," *arXiv:1706.08947*, 2017.
- [51] C. Zhang, Q. Liao, A. Rakhlin, K. Sridharan, B. Miranda, N. Golowich, and T. Poggio, "Theory of deep learning iii: Generalization properties of sgd," *CBMM Memo No. 067*, 2017.
- [52] M. Espig, W. Hackbush, and A. Khachatryan, "On the convergence of alternating least squares optimisation in tensor format representations," *arXiv:1506.00062*, 2015.
- [53] T. Poggio, "On optimal nonlinear associative recall," *Biological Cybernetics*, vol. 19, no. 4, pp. 201–209, 1975.
- [54] Y. Cho and L. K. Saul, "Kernel methods for deep learning," in *Advances in Neural Information Processing Systems 22* (Y. Bengio, D. Schuurmans, J. Lafferty, C. Williams, and A. Culotta, eds.), pp. 342–350, Curran Associates, Inc., 2009.

## 8 First Appendix: classical generalization bounds for $N > W$

### 8.1 Classical generalization bounds

The simplest and oldest bound is provided by theorem 16.2 in <sup>[26]</sup> which provides the following sample bound for a generalization error  $\epsilon_G$  with probability at least  $1 - \delta$  in a network in which the  $W$  parameters (weights and biases) that minimize the empirical error (the theorem is stated in the standard ERM setup) are expressed in terms of  $k$  bits:

$$M(\epsilon_G, \delta) \leq \frac{2}{\epsilon_G^2} \left( kW \log 2 + \log \left( \frac{2}{\delta} \right) \right) \quad (17)$$

An interesting comparison between shallow and deep compositional networks based on this bound is described in Appendix 9.1.

### 8.2 Covering number of the space explored by SGD

**Definition 2** (Covering Number). *The covering number  $N_\epsilon(K)$  for a set  $K$  is the minimum cardinality of a set  $C$ , such that*

$$K \subset \bigcup_{x \in C} B(x; \epsilon)$$

where  $B(x; \epsilon)$  denotes a metric ball of radius  $\epsilon$  centered at  $x$ . The metric should be clear from the context unless otherwise specified.

The motivation is that when we run an actual learning algorithm (e.g. SGD) on a hypothesis space  $\mathcal{H}$ , the algorithm that runs in finite time might only be able to explore a subset  $K$ . Therefore the effective hypothesis space is  $K$  instead of  $\mathcal{H}$ . We would like to characterize  $N_\epsilon(K)$ , which might be much smaller than  $N_\epsilon(\mathcal{H})$ .

Consider the SGD algorithm as defined in (4). If there is an explicit projection step onto the set  $K$ , then we could compute  $N_\epsilon(K)$  directly. If there is no explicit projection, then for an algorithm that runs for  $T$  steps:

$$\begin{aligned} f_T &= f_{T-1} - \gamma_{T-1} \nabla V(f_{T-1}, z_{T-1}) \\ &= f_0 - \sum_{t=0}^{T-1} \gamma_t \nabla V(f_t, z_t) \end{aligned}$$

The effectively explored space can be controlled if we have a bound on the size of the accumulated gradients,

$$\left\| \sum_{t=0}^{T-1} \gamma_t \nabla V(f_t, z_t) \right\| \leq G_T \quad (18)$$

where  $G_T$  should be a universal bound that is independent of the actual sampled training data, and of the randomness of the algorithm. Our formulation should apply to both the case of pure SGDs and multi-epoch SGDs. In this case, it is clear that

$$f_T \in K = B(f_0; G_T)$$

**Lemma 2.** In a  $d$  dimensional Euclidean space, for  $\varepsilon \in (0, r)$ , the covering number of a Euclidean ball is bounded as

$$N_\varepsilon(B(\cdot; r)) \leq \left(\frac{3r}{\varepsilon}\right)^d \quad (19)$$

**Remarks** Applying the lemma directly with the bound  $G_T$  on SGD movements, we get

$$N_\varepsilon(K) \leq \left(\frac{3G_T}{\varepsilon}\right)^d$$

The uniform bounds on generalization are roughly of order

$$O\left(\sqrt{\frac{\log N_\varepsilon(K)}{n}}\right) \leq O\left(\sqrt{\frac{d \log(3G_T/\varepsilon)}{n}}\right)$$

As long as the bound  $G_T$  grows slower than  $\exp(n)$ , then generalization should be possible. Here  $T > n$  in the multiple epoch case.

**Convex learning with Lipschitz Loss** Consider the case of convex learning problem with an  $L$ -Lipschitz loss function. With a constant step size  $\gamma = R/(L\sqrt{T})$ , where  $T$  is the total number of (full) gradient steps, and  $\|f_0 - f^*\| \leq R$ , then an optimization convergence rate on the loss of  $O\left(\frac{RL}{\sqrt{T}}\right)$  could be obtained [27]. In this case,

$$\begin{aligned} \left\| \sum_{t=0}^{T-1} \gamma \nabla V(f_t, z_t) \right\| &\leq \gamma \sum_{t=0}^{T-1} \|\nabla V(f_t, z_t)\| \\ &\leq \frac{R}{L\sqrt{T}} \sum_{t=0}^{T-1} L \\ &\leq R\sqrt{T} \end{aligned}$$

By plugging  $G_T \leq R\sqrt{T}$ , we get the following generalization bound

$$O\left(\sqrt{\frac{d \log(T/\varepsilon)}{n}}\right)$$

As a result, as long as the number of optimizing steps grows less than exponentially in the number of training samples, generalization could be guaranteed.

**Convex learning with smooth loss** For the Lipschitz loss, the scale of gradients does not necessarily decay even when approaching the optimal solution (e.g. the absolute value function). For smooth loss, the gradients becomes smaller when we get closer to the optimal. Therefore, the bound on the movements by the gradient methods could be improved.

Furthermore, when we have an explicitly (and fast enough) decay rate of the size of the gradients, we might be able to get finer-grain control by realizing that after some fixed number of steps, the gradients are so small that they will remain within an  $\varepsilon$ -ball.

### 8.3 Robustness-based Bounds for Generalization

A way to theoretically connect ‘‘flatness’’ in weight space to existing generalization results, is to invoke the notion of *weak robustness* and to show that it is equivalent to flatness of a minimizer. Theorem 18 in [28] proves that *weak robustness* is necessary and sufficient for generalization. For simplicity, we focus on robustness (instead of weak robustness) since robustness implies generalization which is our focus here. We use the notation of [28].

**Definition 3.** Algorithm  $A$  is  $(K, \epsilon(\cdot))$ -robust, for  $K \in \mathcal{N}$  and  $\epsilon(\cdot) : \mathcal{Z}^n \mapsto \mathcal{R}$ , if  $\mathcal{Z}$  can be partitioned into  $K$  disjoint sets, denoted by  $\mathcal{C} = (C_i)_{i=1}^K$ , such that the following holds: for any dataset  $S \in \mathcal{Z}^n$ ,  $\forall s \in S$  and  $\forall z \in \mathcal{Z}$ ,  $\forall i = 1, \dots, K$ : if  $s, z \in C_i$ , then

$$\|V(\mathcal{A}_S, s) - V(\mathcal{A}_S, z)\| \leq \epsilon(S) \quad (20)$$

**Theorem 1.** If a learning algorithm  $A$  is  $(K, \epsilon(\cdot))$ -robust, and the training sample set  $S$  is generated by  $N$  IID draws from  $\mu$ , then for any  $\delta > 0$ , with probability at least  $1 - \delta$  we have,

$$|\mathcal{L}(\mathcal{A}_s) - \ell_{emp}(\mathcal{A}_S)| \leq \epsilon(S) + M \sqrt{\frac{2K \log 2 + 2 \log(\frac{1}{\delta})}{n}}, \quad (21)$$

with  $M$  being an upper bound on the loss function  $V$ .

The key step in the argument about flatness, stability and generalization is to prove that flatness implies robustness as defined above.

We now extend theorem 1 to a data dependent version using a standard union bound. Suppose algorithm  $A$  is robust for a countable collection  $\mathcal{K}$  of pairs  $(K_C, \epsilon_C(\cdot))$ , indexed by  $\mathcal{C}$ , where  $\mathcal{C}$  is the partition of size  $K_C$  in the definition of robustness. We have in mind a situation where a coarse partition can be set (in a data-independent way) with a potentially larger function  $\epsilon_C$  but smaller  $K_C$ , while a more fine partition would yield a smaller value of  $\epsilon_C$  but a larger  $K_C$ . Since  $\epsilon_C(S)$  is a data-dependent quantity, the optimal choice of the partition can only be performed after seeing the data.

**Theorem 2.** *If a learning algorithm  $A$  is  $(K_C, \epsilon_C(\cdot))$ -robust for a collection  $\mathcal{K}$  of pairs,*

$$\mathbb{P} \left\{ |\mathcal{L}(\mathcal{A}_s) - \ell_{temp}(\mathcal{A}_s)| \leq \inf_{(K_C, \epsilon_C(\cdot)) \in \mathcal{K}} \left[ \epsilon_C(S) + M \sqrt{\frac{2K_C \log 2 + 2 \log((\pi(\mathcal{C})\delta)^{-1})}{n}} \right] \right\} \geq 1 - \delta \quad (22)$$

for any prior distribution  $\pi(\mathcal{C})$  on the collection  $\mathcal{K}$ .

As an example, consider a collection of axis-aligned grid partitions with granularity  $\gamma$ . If the input space is  $[0, 1]^d$ , the size of the partition is  $(1/\gamma)^d$ , and the corresponding function  $\epsilon_C$ , of course, depends on the algorithm. One can then assign a prior probability proportional to  $1/k^2$  for discretization at level  $\gamma = 2^{-k}$ , in which case the penalty  $\pi(\mathcal{C})$  is  $\log k$  (as, for instance, in <sup>[15]</sup>).

In our CIFAR case the key question is how large  $K$  is and whether it is smaller than  $N$ .

### Remarks

The authors in <sup>[28]</sup> assume that  $f$  is a  $L$ -layer neural network of the binary tree type with  $P$  ReLU units per node and  $D = 2^L$  inputs, trained on  $\mathcal{Z}_n$  by SGD algorithm  $\mathcal{A}_{S_n}$  with loss  $V(\mathcal{A}_S, z) = |y - \mathcal{A}_S(x)| = |y - f_S(x)|$ , with  $z = (y, x)$  and weights taking values on the torus. They define the network as follows. The input to the first layer is

$$x^0 := x; \quad (23)$$

the output of layer  $v - 1$  is  $x^v$  with components

$$\forall v = 1, \dots, L - 1: x_i^v := \left[ \sum_{j=1}^{2P} w_{i,j}^{v-1} x_j^{v-1} \right]_+; i = 1, \dots, d_v; \quad (24)$$

where  $d^v = P \times 2^{L-v}$

$$f_S(x) = \left[ \sum_{j=1}^{d_{L-1}} w_j^{L-1} x_j^{L-1} \right]_+. \quad (25)$$

The ReLU nonlinearity is Lipschitz continuous since  $[a]_+ - [b]_+ \leq \beta|a - b|$  with  $\beta = 1$ . If the weights are bounded then it follows that

$\sum_{j=1}^{d_v} w_{i,j}^v \leq \alpha \forall v, i$  (thus the largest  $\alpha$  is  $\alpha \leq 2P\pi$ ). As an aside, note that, in the case of compositional networks, the use of the path norm should give better bounds than <sup>[28]</sup>. Lemma 22 in <sup>[28]</sup> proves then that under these conditions such a multilayer network is robust in the infinity norm and therefore generalizes. As noted by <sup>[28]</sup> the total number of hidden units does not play a role in the bound. Notice also that the local connectivity of deep local networks (such as convolutional network and more in general networks matched to compositional functions) avoids the appearance of the dimensionality of the inputs  $x_i$ .

## 8.4 Distribution Dependent Generalization because of functional margin

In repeat SGD there is a training set and an empirical risk that we write in the form

$$I_{S_n}(f) = \frac{1}{n} \sum_i^n V(f, z_i) \quad (26)$$

It is clear from the previous sections that repeat SGD provably generalizes under relatively weak assumptions – similar to the usual conditions of convergence of SGD. The problem is of course to provide bounds on the expected error that are non-vacuous for the regime of  $W > N$ . Several of the classical bounds are meaningless in this case (see section 8.1), but other proofs are possible (see 8.5).

The observations relevant here are:

- We assume that the space of functions exposed by SGD is compact. One way is to assume Lipschitz and convexity. Another is to set as an hypothesis “If  $f_i$  remains bounded then...”. In practice, boundness of some kind must be assumed. It is in principle guaranteed by numerical limits of computer simulations. A simple way to ensure compactness – as a theoretical trick which will not make any difference in the usual simulations – is to assume that the coordinates are periodic – that is that the domain is a torus. Then standard covering numbers arguments guarantee generalization for  $N > W$ .
- The following arguments are needed to ensure bounds that are independent of the explicit ratio  $\frac{W}{N}$ . We remark that SGDL should asymptotically select among the zero-minimizers the ones which have minimum norm (and thus large margin). We provide empirical support for large empirical margin in the case of CIFAR in the Figures 17. We introduce the zero-one loss and the  $\gamma$ -margin empirical zero-one loss which is defined as  $I_{S,\gamma}(h) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}_{y_i h(x_i) < \gamma}$ . We then use Theorem 5 in [29]

**Theorem 3.** *For any hypothesis class  $(F)$ , with  $|f| \leq b$ , and any  $\delta > 0$ , with probability at least  $1 - \delta$ , simultaneously for all margins  $\gamma > 0$  and all  $f \in \mathcal{F}$*

$$I(f) \leq 1.01 I_{S,\gamma}(f) + K \left( \frac{2 \log^3 n}{\gamma^2} \mathcal{R}_n^2 + \frac{2 \log \left( \frac{\log \left( \frac{4b}{\gamma} \right)}{\delta} \right)}{n} \right) \quad (27)$$

where  $K$  is an appropriate numeric constant,  $\mathcal{R}$  is the Radamacher complexity of the space of functions  $\mathcal{F}$ .

Under the same hypotheses that guarantee compactness for the space of solution of SGD,  $\mathcal{R}^2$  goes to zero with increasing  $N$  because the covering numbers provides an upper bound on Radamacher averages. Additionally, by checking Equation 27, we can observe that zero-minimizers have a better generalization bound if they have larger margin.

Theory I results (see 8.1) imply that  $\mathcal{N}(F)$  for  $d$ -dimensional functions of a Sobolev space with smoothness  $s$  grow as  $\log \mathcal{N}(\mathcal{F}) \sim \left(\frac{1}{\epsilon}\right)^{\frac{d}{s}}$ . On the other hand for the subset of hierarchically local compositional functions ( $\mathcal{F}_c$  – and for the set of functions generated by the corresponding networks of which convolutional networks are a subset – the covering numbers are much smaller, only growing as  $\log \mathcal{N}(\mathcal{F}_c) \sim d \left(\frac{1}{\epsilon}\right)^{\frac{h}{s}}$  as a function of  $d$ .

Since

$$\mathcal{R}_n \leq C \cdot \inf_{\epsilon > 0} \left\{ \epsilon + \frac{1}{\sqrt{n}} \int_{\epsilon}^{\|\mathcal{F}\|_{\infty}} \sqrt{\log \mathcal{N}_{\delta}(\mathcal{F})} d\delta \right\} \quad (28)$$

we conclude that the term  $\mathcal{R}_n^2$  is smaller for convolutional networks than for equivalent shallow networks. This may be the reason for the good prediction properties of deep networks when the underlying function to be learned is compositional. Furthermore, the covering numbers of the subset of functions in a ball around the global minimizer should be even smaller, because the effective  $d$  is dictated by the number of data and not by the much larger number of weights. Apart from this non-rigorous argument, generalization is also controlled by the margin  $\gamma$ . We now collect together the results from Theory I and II and from the previous sections in this paper to claim the following

**Proposition 1.** • *Good expected error by deep networks relative to shallow networks in tasks that correspond to learn a compositional function (or a linear combination of a “small” number of them) follows from the smaller Radamacher complexity.*

- *Consistency is further controlled by the empirical margin of minimizers that depends not only on the distributions of inputs  $x_i$  but also on the distribution of  $y_i$ .*
- *Later we will show that SGDL prefers among the minima the degenerate zero-minimizers – if they exist; among the degenerate zero-minimizers it prefers the ones with flattest minima, that, as we will show later, correspond to the largest margin.*

*In summary, SGDL (and likely also SGD) selects with high probability the minimizers that do the best job in terms of minimizing expected error. This suggests that the large-margin argument to explain consistency in the regime  $W > N$  also explains generalization in the regime  $W < N$ .*

## 8.5 Global minimization and stability under regularization-like assumptions

The most interesting proofs for our goal are available for the case of Langevin equations and annealing techniques where results about global optimization and stability of the algorithm have become available. Gelfand and Mitter [30, 31] prove convergence of Langevin equations to global minima. Related results on simulated annealing for finite Markov chains [32] prove polynomial convergence time (meaning probability of visiting the global optimum at least once). A very recent paper by Raginsky et al. [33] considers a version of SGD – called SGDL – in which isotropic Gaussian noise is added to an estimate of the gradient. They also prove global convergence but in addition they prove generalization with rates independent of parameter numbers or size and valid in the case of general loss functions. In order to apply their results to standard SGD one should first prove that SGD is a good approximation of a Langevin equation (and that their results still hold approximatively). We skip this step because there are empirical reasons to believe that SGDL may be the better way to train neural networks. In particular our numerical analysis suggests that the behavior of SGDL is very similar to SGD but SGDL consistently shows a slightly better generalization performance. In the next section, however, we provide theoretical and numerical arguments for the intuition that SGD is similar to SGDL and in fact to a Langevin equation (GDL). This claim is not new: there are previous analyses of it. As an example see [34].

## 9 Second Appendix: Various Topics

### 9.1 Generalization Bounds: comparison with compositional networks

The following is from section 8 in Theory I<sup>[1]</sup>, provided here for completeness.

Assume a network size that ensure the same approximation error  $\epsilon$ . Then in order to achieve the same generalization error  $\epsilon_G$ , the sample size  $M_{shallow}$  of the shallow network must be much larger than the sample size  $M_{deep}$  of the deep network:

$$\frac{M_{deep}}{M_{shallow}} \approx \epsilon^n. \quad (29)$$

This implies that for largish  $N$  there is a (large) range of training set sizes between  $M_{deep}$  and  $M_{shallow}$  for which deep networks will not overfit (corresponding to small  $\epsilon_G$ ) but shallow networks will (for dimensionality  $N \approx 10^4$  and  $\epsilon \approx 0.1$  Equation 29 yields  $m_{shallow} \approx 10^{10^4} m_{deep}$ ).

A similar comparison is derived if one considers the best possible expected error obtained by a deep and a shallow network. Such an error is obtained finding the architecture with the best trade-off between the approximation and the estimation error. The latter is essentially of the same order as the generalization bound implied by inequality (17), and is essentially the same for deep and shallow networks, that is

$$\frac{rn}{\sqrt{M}}, \quad (30)$$

where we denoted by  $M$  the number of samples. For shallow networks, the number of parameters corresponds to  $r$  units of  $N$  dimensional vectors (plus off-sets), whereas for deep compositional networks the number of parameters corresponds to  $r$  units of 2 dimensional vectors (plus off-sets) in each of the  $N - 1$  units. Using our previous results on degree of approximation, the number of units giving the best approximation/estimation trade-off is

$$r_{shallow} \approx \left( \frac{\sqrt{M}}{n} \right)^{\frac{n}{m+n}} \quad \text{and} \quad r_{deep} \approx \left( \sqrt{M} \right)^{\frac{2}{m+2}} \quad (31)$$

for shallow and deep networks, respectively. The corresponding (excess) expected errors  $\mathcal{E}$  are

$$\mathcal{E}_{shallow} \approx \left( \frac{n}{\sqrt{M}} \right)^{\frac{m}{m+n}} \quad (32)$$

for shallow networks and

$$\mathcal{E}_{deep} \approx \left( \frac{1}{\sqrt{M}} \right)^{\frac{m}{m+2}} \quad (33)$$

for deep networks. For the expected error, as for the generalization error, deep networks appear to achieve an exponential gain. The above observations hold under the assumption that the optimization process during training finds the optimum parameters values for both deep and shallow networks. Taking into account optimization, e.g. by stochastic gradient descent, requires considering a further error term, but we expect that the overall conclusions about generalization properties for deep vs. shallow networks should still hold true.

Notice that independently of considerations of generalization, deep compositional networks are expected to be *very efficient memories* – in the spirit of hierarchical vector quantization – for associative memories reflecting compositional rules (see Appendix and<sup>[35]</sup>). Notice that the advantage with respect to shallow networks from the point of view of memory capacity can be exponential (as in the example after Equation 29 showing  $m_{shallow} \approx 10^{10^4} m_{deep}$ ).

### 9.2 Empirical results on margin

### 9.3 SGD is equivalent to a form of Robust Optimization and corresponds to large margin

The claim is that SGD while minimizing the empirical error effectively maximizes robustness to perturbations in the weights (and similarly in the training data). In particular, section 9.4 shows that this is similar to regularization with a regularizer consisting of a path norm in the weights of the deep network. In the subsequent section 8.3 we outline results – old and new – showing that robustness yield generalization bounds that are independent of number of parameters. Unfortunately these bounds are very loose.

In the following sections we argue that flatness in the minimum provided by SGD yields better generalization. We do not have a practical bound but we present strong arguments based on a close connection between flatness, robust optimization and regularization.



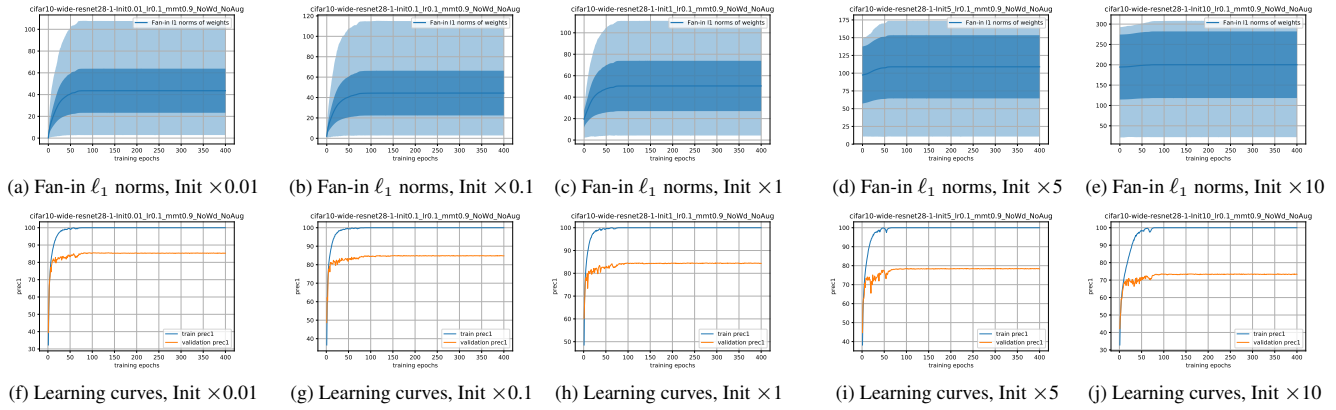


Figure 16: Fan-in  $\ell_1$  norms of the weights for training on CIFAR-10.

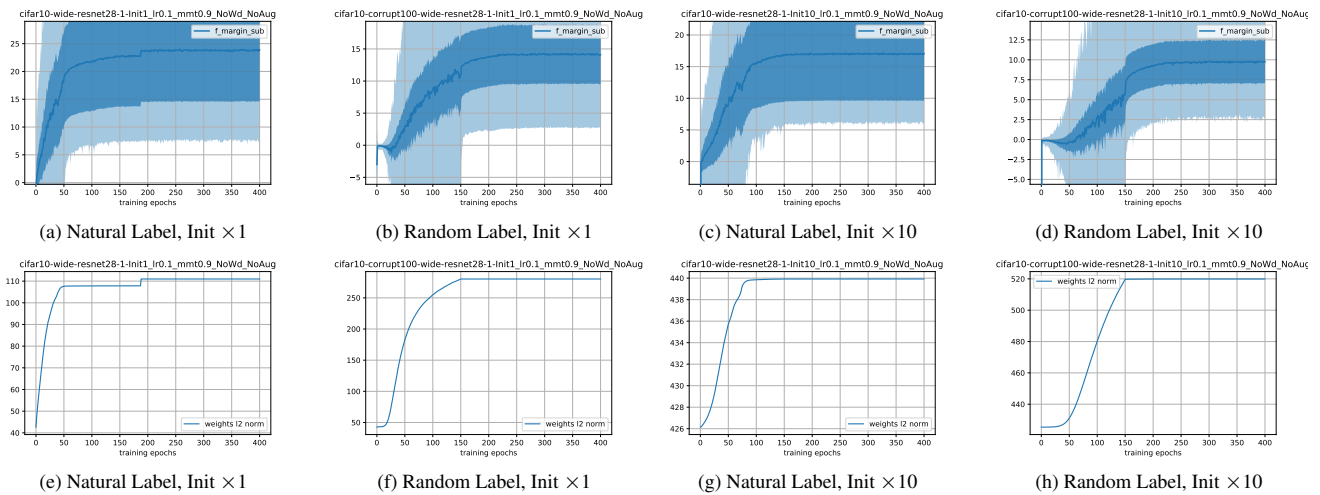


Figure 17: Functional margins for training on CIFAR-10. Patterns: random labels has smaller margin and larger l2 norm than natural labels. Larger init ( $\times 10$ ) leads to smaller margin and larger l2 norm than normal init ( $\times 1$ ).

Let us first introduce a notation to describe the deep networks that we will use in the following sections.

### Networks and weights

We assume that  $f$  is a  $L$ -layer neural network with  $d$  inputs, trained on  $z_1, \dots, z_d$  with  $z = (x, y)$ . We assume one of the components of the effective input vector  $x$  is a dummy input equal to 1 (this allows to drop the parameter  $b$  in describing the RELU activity). For convenience we repeat the definitions in (23), (24), and (25). The input to the first layer is

$$x^0 := x;$$

the output of layer  $v - 1$  is  $x^v$  with components

$$\forall v = 1, \dots, L - 1: x_i^v := \left[ \sum_{j=1}^{N_{v-1}} w_{i,j}^{v-1} x_j^{v-1} \right]_+ \quad i = 1, \dots, d_v$$

and

$$f(x) = \left[ \sum_{j=1}^{N_{L-1}} w_j^{L-1} x_j^{L-1} \right]_+.$$

Notice that if we neglect the RELUs (or equivalently assume that all the RELUs are active), the  $L$  layer network with  $d$  inputs and one output, is equivalent to a linear network with a  $d$ -dimensional weight vector  $\tilde{\mathbf{w}} = \sum_{h,t,j \dots} w_{1,h}^{L-1} \dots w_{t,k}^2 w_{k,j}^1 w_{j,i}^0$  with  $i = 1, \dots, d$ . Notice that  $\|\tilde{\mathbf{w}}\|$  is the path norm defined elsewhere<sup>[13]</sup>. The vector  $\tilde{\mathbf{w}}$  involved in the path norm has components which are each a sum of monomials in the weights, each monomial corresponding to a path. Suppose there are a total of  $D$  paths. We define  $\hat{\mathbf{w}}$  as the  $D$ -dimensional vector in which each component consist of the monomial corresponding to a path and as  $\hat{\mathbf{x}}$  the corresponding ‘‘augmented’’  $x$  vector (for instance in Figure 18  $D = d$  if we assume one RELU per node; with 2 RELUs per node  $D = 4d$ ). The corresponding ‘‘Dpath norm’’ is  $\|\hat{\mathbf{w}}\|$ . Thus

$$f(x) = \Lambda_{w,x} \hat{\mathbf{w}} \hat{\mathbf{x}} \quad (34)$$

where  $\Lambda$  is a diagonal matrix with 0, 1 components indication which paths are switched on. Notice that  $\frac{df}{dw} \approx 0$  around a flat minimum  $w^*$  of  $f$ .

## 9.4 SGDL, RO and generalization

In the next subsection we extend previous results<sup>[36]</sup> to relate minimizers that corresponds to flatness in the minimum to robustness: thus SGD performs robust optimization of a certain type. In particular, we will show that robustness to perturbations in multilayer deep networks is related to regularization (see<sup>[36, 37]</sup>). We first describe the simplest separable and linear case.

### 9.4.1 For linear networks and separable data, SGD converges to a large-margin classifier

In linear classification we try to separate the two classes of a binary classification problem by a hyperplane  $\mathcal{H} = \{x : w^\top x + b = 0\}$ , where  $w \in \mathbb{R}^d$ . There is a corresponding decision rule of the form  $y = \text{sign}(w^\top x + b)$ . Let us consider the separable condition in which the decision rule makes no error on the data set. This corresponds to the following set of inequalities

$$y_i(w^\top x_i + b) > 0, \quad i = 1, \dots, n \quad (35)$$

We assume that the inequalities are feasible, that is the data are separable. Assume now to impose robustness of the classifier wrt to perturbations  $\delta w$  in the weights which we assume here to be such that  $\|\delta w\|_2 \leq \rho \|w\|_2$  with  $\rho \geq 0$ . In other words, we require that  $\forall \delta w$  such that  $\|\delta w\|_2 \leq \rho \|w\|_2$ :

$$y_i((w + \delta w)^\top x_i + b) \geq 0, \quad i = 1, \dots, n \quad (36)$$

Further assume that  $\|x_i\|_2 \approx 1$ , then for  $i = 1, \dots, n$ , if we let  $\delta w = -\rho y_i x_i \|w\|_2 / \|x_i\|_2$ ,

$$y_i(w^\top x_i + b) \geq -y_i \delta w^\top x_i = \rho \|w\|_2 \|x_i\|_2 \approx \rho \|w\|_2$$

As a result, an approximate *robust counterpart* of Equation (35) is

$$y_i(w^\top x_i + b) \geq \rho \|w\|_2, \quad i = 1, \dots, n \quad (37)$$

Maximizing  $\rho$  – the margin – subject to the constraints Equation (37) leads to minimizing  $w$ , because we can always enforce  $\rho \|w\|_2 = 1$  and thus to

$$\min_{w,b} \{ \|w\|_2 : y_i(w^\top x_i + b) \geq 1 \quad i = 1, \dots, n \} \quad (38)$$

Notice that the same Equation 38 follows if the maximization is with respect to spherical perturbations of radius  $\rho$  around each data point  $x_i$ . In either case the resulting optimization problems is equivalent to hard margin SVMs<sup>5</sup>. Notice that the classification problem is similar to using the loss function  $V(y, f(x)) = \log(1 + e^{-yf(x)})$  which penalizes errors but is otherwise very small in the zero classification error case. Section 5.2 implies that SGD maximizes flatness at the minimum that is SGD maximizes  $\delta w$ .

In the non-separable case, the hinge loss<sup>6</sup> leads to the following robust minimization problem

$$\min_{w,b} \frac{1}{n} \sum_{i=1}^n [1 - y_i(w^\top x_i + b) + \rho \|w\|_2]_+ \quad (40)$$

Note that the robust version of this worst-case loss minimization is not the same as in classical SVM because the regularization term is inside the hinge loss. Note also that standard regularization is an upper bound for the robust minimum since

$$\min_{w,b} \frac{1}{n} \sum_{i=1}^n [1 - y_i(w^\top x_i + b) + \rho \|w\|_2]_+ \leq \min_{w,b} \frac{1}{n} \sum_{i=1}^n [1 - y_i(w^\top x_i + b)]_+ + \rho \|w\|_2. \quad (41)$$

In the case of the square loss, robust optimization gives with  $\|\delta w\|_2 \leq \rho \|w\|_2$

$$\min_{w,b} \frac{1}{n} \sum_{i=1}^n [(y_i - w^\top x_i)^2 + \rho^2 \|w\|_2]_+ \quad (42)$$

The summary here is that depending on the loss function and on the uncertainty set allowed for the perturbations  $\delta w$  one obtains a variety of robust optimization problems. In general they are not identical to standard regularization but usually they contain a regularization term. Notice that a variety of regularization norms (for instance  $\ell_1$  and  $\ell_2$ ) guarantee  $CV_{i\infty}$  stability and therefore generalization. The conclusion is that *in the case of linear networks and linearly separable data, SGD provides a solution closely related to hinge-loss SVM.*

#### 9.4.2 Deep networks, SGD, RO and regularization

In the case of deep (convolutional) networks we can *only approximate* the robust optimization problems described above *at a minimum* by using the pseudo-linear representation of a deep networks provided by Equation 34, where  $\Lambda$  depends on the weights and on the input  $x$ . For a deep network the problem around the minimum for the square loss under the assumptions of the previous section reads

$$\min_w \max_{(\delta w)} \frac{1}{n} \sum_{i=1}^n (y_i - f_w(\mathbf{x}_i))^2 \quad (43)$$

which, since  $f(x) = \Lambda_{w,x} \hat{\mathbf{w}} \hat{\mathbf{x}}$ , leads to

$$\min_{\hat{\mathbf{w}}} \frac{1}{n} \sum_{i=1}^n [(y_i - \hat{\mathbf{w}}^\top \Lambda_{w,x} \hat{\mathbf{x}}_i)^2 + \frac{1}{n} \rho^2 \hat{\mathbf{w}}^\top \Lambda_{w,x} \hat{\mathbf{w}}] \quad (44)$$

<sup>5</sup>If we start from the less natural assumption that  $\|\delta w\|_\infty \leq \rho \|w\|_2$  with  $\rho \geq 0$  and assume that each of the  $d$  components  $|(x_i)_j| \approx 1$ ,  $\forall i, d$ , then the *robust counterpart* of Equation 35 is, since  $(\delta w)^\top x \leq \|(\delta w)^\top\|_{\ell_\infty} \|x\|_{\ell_1}$ ,

$$y_i(w^\top x_i + b) \geq (\delta w)^\top \sup \delta w = \rho \|w\|_2, \quad i = 1, \dots, n \quad (39)$$

<sup>6</sup>The hinge loss  $V(y, f(x)) = [1 - yf(x)]_+$ , with  $y$  binary (used by [37]) is similar to the cross-entropy loss (for classification)  $V(y, f(x)) = \log(1 + e^{-yf(x)})$ .

Here the diagonal matrix  $\Lambda$  weighs each path by the times it is active over the training set  $S_n$ . Of course the perturbation will switch on or off several of the paths, switching on or off associated RELUs, depending on the training sample  $y_i, (x_i)$ . We consider the simple case of Figure 18, where the number of paths is equal to the dimensionality of the input vector. Let us assume that there are several zero minima in the square loss, all close to zero loss. Then maximization over  $\delta$  selects the flattest minima – the ones where perturbation can be largest without increasing the errors over the training set. Let us assume as before that all the explored minima are at close to zero loss, that is  $y_i - f(x_i) \approx 0$ . In this case the robust optimum is close to the minimum of a weighted regularization problem.

In fact the penalization term above is an upper bound so that

$$\min_w \max_{(\delta w)} \frac{1}{n} \sum_{i=1}^n (y_i - f_{w+\delta w}(\mathbf{x}_i))^2 \leq \min_{\hat{\mathbf{w}}} \frac{1}{n} \sum_{i=1}^n [(y_i - \hat{\mathbf{w}}^\top \Lambda_{w,x} \hat{\mathbf{x}}_i)^2 + \frac{1}{n} \rho^2 \hat{\mathbf{w}}^\top \Lambda_{w,x} \hat{\mathbf{w}}]. \quad (45)$$

More in general

**Proposition 2.** For a compositional network with several units per node the following holds

$$\min_{\hat{\mathbf{w}}} \frac{1}{n} \sum_{i=1}^n [(y_i - \hat{\mathbf{w}}^\top \Lambda_{w,x} \hat{\mathbf{x}}_i)^2 + \frac{1}{n} \rho^2 \hat{\mathbf{w}}^\top \hat{\mathbf{w}}] \leq \min_w \max_{\delta \in \mathcal{N}} \frac{1}{n} \sum_{i=1}^n V(y_i, f(\mathbf{x}_i - \delta)) \leq \min_{\hat{\mathbf{w}}} \frac{1}{n} \sum_{i=1}^n V(y_i, f(\mathbf{x}_i)) + c|\hat{\mathbf{w}}| \leq \min_{\hat{\mathbf{w}}} \sum_{i=1}^n [(y_i - \hat{\mathbf{w}}^\top \Lambda_{w,x} \hat{\mathbf{x}}_i)^2 + \rho^2 \hat{\mathbf{w}}^\top \Lambda_{w,x} \hat{\mathbf{w}}] \quad (46)$$

where the lower and upper bounds depend on the assumption that each path is active at least once over the training set.

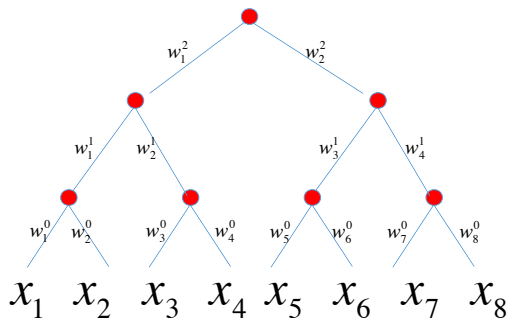


Figure 18: Assume that the binary network in the figure has one relu unit per node. The path norm is  $|\hat{\mathbf{w}}|$  and  $\hat{\mathbf{w}} = w_1^0 w_1^1 w_1^2, w_2^0 w_1^1 w_1^2, w_3^0 w_2^1 w_1^2, \dots$ . If one of the unit is switched off then the two weights in the input to the unit can be put to zero and one or more of the monomials representing the components of  $\hat{\mathbf{w}}$  will be zero. In the case of several units per node, the components of  $\hat{\mathbf{w}}$  will be sum of monomials, some of which can be put to zero if corresponding RELU units are switched off.

It seems likely that we can bound  $CV_{loo}$  stability (for any  $N$ ) and guarantee generalization. This may require to make  $\rho$  decrease with  $N$  by controlling the power of a Gaussian noise added to SGD (we assume to be in the SGDL case). In summary we can establish a close connection between the flatness maximized by SGD, maximizing robustness and regularization *but only* approximately and only in some special cases.

### Remarks

- The proposition suggests that a deep network optimized by SGD has bounded weights, uniform stability and generalization. The associate bounds are *independent* of the number of parameters.

## 9.5 Training with random labels, testing with perturbed images

In this section, we empirically compare the robustness of the classifier learned on natural labels and random labels. Specifically, since for both natural labels and random labels, the models could perfectly fit the labels, we are interested in comparing how the fitted model respond to small perturbations on the training examples. We define the following notion of perturbation robustness:

$$\tilde{E}(f, S_n, \gamma) = \frac{1}{n} \sum_{i=1}^n \mathbb{I}[f(x_i + \gamma \varepsilon_i) \neq y_i] \quad (47)$$

where  $S_n = \{(x_i, y_i)\}_{i=1}^n$  is the training set,  $f$  is a model that was trained to fit the training set.  $\varepsilon_1, \dots, \varepsilon_n$  are i.i.d. standard Gaussian vectors. And  $\gamma$  is the size of the perturbation. When the model fits the training set perfectly,  $\tilde{E}(f, S_n, 0) = 0$ .

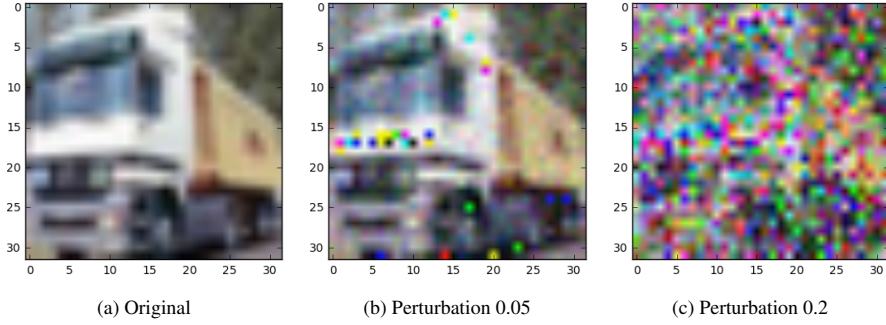


Figure 19: Example of an image on the CIFAR-10 training set and its perturbed versions with different perturbation size.

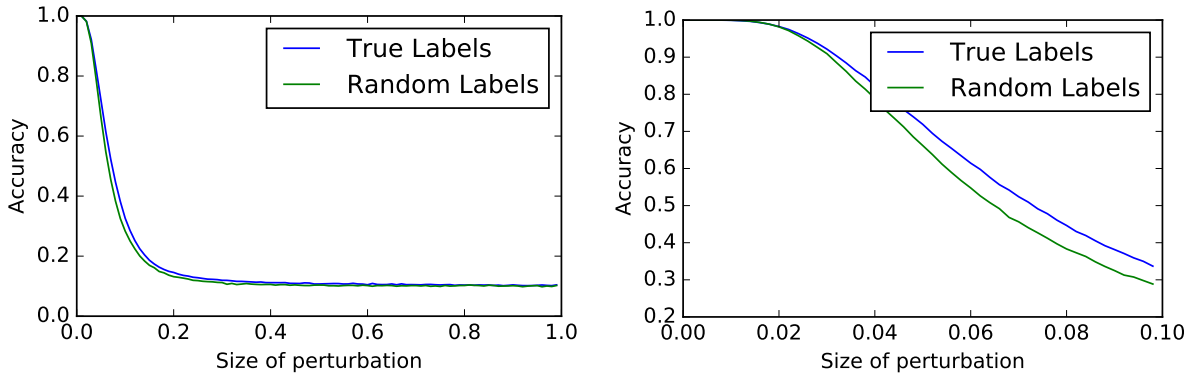


Figure 20: Perturbation robustness of models trained by SGD on CiFAR with correct and random labels and then tested on perturbed images. The right hand side figure is a zoomed-in view of the initial part of the left hand side figure.

Figure 19 shows an example of the training image and its perturbed version with different perturbation sizes. Figure 20 shows the perturbation robustness with models trained on natural labels and random labels. As we can see, the model trained on random labels has slightly worse robustness than the model trained on natural labels.

## 9.6 GD and SGD converges to Minimum Norm Solution for Underdetermined Linear System

Consider the following setup:  $X = (x_1, \dots, x_n)^\top \in \mathbb{R}^{n \times d}$  are the data points, with  $d > n$ . We further assume that the data matrix is of full row rank:  $\text{rank}(X) = n$ . Let  $y \in \mathbb{R}^n$  be the labels, and consider the following linear system:

$$Xw = y \quad (48)$$

where  $w \in \mathbb{R}^d$  is the weights to find. This linear system has infinite many solutions because  $X$  is of full row rank and we have more parameters than the number of equations. Now suppose we solve the linear system via a least square formulation

$$L(w) = \frac{1}{2n} \|Xw - y\|^2 \quad (49)$$

by using gradient descent (GD) or stochastic gradient descent (SGD). In this section, we will show that both GD and SGD converges to the minimum norm solution, for the  $\ell_2$  norm. A similar analysis (with fewer details) was presented in [4].

**Lemma 3.** *The following formula defines a solution to (48)*

$$w_{\dagger} \triangleq X^\top (XX^\top)^{-1} y \quad (50)$$

and it is the minimum norm solution.

*Proof* Note since  $X$  is of full row rank, so  $XX^\top$  is invertible. By definition,

$$Xw_{\dagger} = XX^\top (XX^\top)^{-1} y = y$$

Therefore,  $w_{\dagger}$  is a solution. Now assume  $\hat{w}$  is another solution to (48), we show that  $\|\hat{w}\| \geq \|w_{\dagger}\|$ . Consider the inner product

$$\begin{aligned}\langle w_{\dagger}, \hat{w} - w_{\dagger} \rangle &= \langle X^{\top} (X X^{\top})^{-1} y, \hat{w} - w_{\dagger} \rangle \\ &= \langle (X X^{\top})^{-1} y, X \hat{w} - X w_{\dagger} \rangle \\ &= \langle (X X^{\top})^{-1} y, y - y \rangle \\ &= 0\end{aligned}$$

Therefore,  $w_{\dagger}$  is orthogonal to  $\hat{w} - w_{\dagger}$ . As a result, by Pythagorean theorem,

$$\|\hat{w}\|^2 = \|(\hat{w} - w_{\dagger}) + w_{\dagger}\|^2 = \|\hat{w} - w_{\dagger}\|^2 + \|w_{\dagger}\|^2 \geq \|w_{\dagger}\|^2$$

**Lemma 4.** *When initializing at zero, the solutions found by both GD and SGD for problem (49) live in the span of rows of  $X$ . In other words, the solutions are of the following parametric form*

$$w = X^{\top} \alpha \quad (51)$$

for some  $\alpha \in \mathbb{R}^n$ .

*Proof*

The gradient for (49) is

$$\nabla_w L(w) = \frac{1}{n} X^{\top} (X w - y) = X^{\top} e$$

where we define  $e = (1/n)(X w - y)$  to be the error vector. GD use the following update rule:

$$w_{t+1} = w_t - \eta_t \nabla_w L(w_t) = w_t - \eta_t X^{\top} e_t$$

Expanding recursively, and assume  $w_0 = 0$ . we get

$$w_t = \sum_{\tau=0}^{t-1} -\eta_{\tau} X^{\top} e_{\tau} = X^{\top} \left( -\sum_{\tau=0}^{t-1} \eta_{\tau} e_{\tau} \right)$$

The same conclusion holds for SGD, where the update rule could be explicitly written as

$$w_{t+1} = w_t - \eta_t (x_{i_t}^{\top} w - y_{i_t}) x_{i_t}$$

where  $(x_{i_t}, y_{i_t})$  is the pair of sample chosen at the  $t$ -th iteration. The same conclusion follows with  $w_0 = 0$ .

Q.E.D.

**Theorem 4.** *Let  $w_t$  be the solution of GD after  $t$ -th iteration, and  $w_t$  be the (random) solution of SGD after  $t$ -th iteration. Then  $\forall \varepsilon > 0$ ,  $\exists T > 0$  such that*

$$L(w_t) \leq \varepsilon, \quad \mathbb{E}L(w_t) \leq \varepsilon$$

where the expectation is with respect to the randomness in SGD.

**Corollary 1.** *When initialize with zero, both GD and SGD converges to the minimum-norm solution.*

*Proof* Combining Lemma 4 and Theorem 4, GD is converging  $w_t \rightarrow w_{\star} = X^{\top} \alpha_{\star}$  as  $t \rightarrow \infty$  for some optimal  $\alpha_{\star}$ . Since  $w_{\star}$  is a solution to (48), we get

$$y = X w_{\star} = X X^{\top} \alpha_{\star}$$

Since  $X X^{\top}$  is invertible, we can solve for  $\alpha_{\star}$  and get

$$w_{\star} = X^{\top} \alpha_{\star} = X^{\top} (X X^{\top})^{-1} y = w_{\dagger}$$

Similar argument can be made for SGD with expectation with respect to the randomness of the algorithm.

## 9.6.1 Weight Perturbations

We could define in different ways the basic perturbations of the data that we consider<sup>7</sup> but here we assume independent perturbation of each component of a vector of perturbations  $\delta_i = \delta$  which is the same for each example  $x_i$ .

7

**Definition 4.** *A set  $\mathcal{N}_0 \in \mathcal{R}^n$  is called an Atomic Uncertainty Set if*

1.  $\mathbf{0} \in \mathcal{N}_0$
2. For any  $\mathbf{w}_0 \in \mathcal{R}^n$   $\sup_{\delta \in \mathcal{N}_0} [\mathbf{w}_0^{\top} \delta] = \sup_{\delta' \in \mathcal{N}_0} [-\mathbf{w}_0 \delta'] < \infty$

Based on this<sup>[37]</sup> then define more complex models of perturbations in which the behavior across multiple samples  $i = 1, \dots, n$  is controlled.

Suppose that there is zero-minimizer  $f^*$  of the empirical loss with parameters  $w^*$ , yielding  $I_{S_n}(w^*) = 0$ . Suppose further that the minimizer is  $\Delta w$ -flat, e.g.  $I_{S_n}(w^* + \Delta w) = 0$ , where  $\Delta w$  is a vector of perturbations of the weights  $w^*$  and  $\min |\Delta w_i| \geq C$ .

For an illustration of the basic idea, consider just the first layer with  $P = 1$ . Then calling  $x_j^0 = x_j$  we obtain with  $\Delta$  being a percentage perturbation of  $w$

$$x_i^1 = \sum_{j=1}^2 (w_{i,j}^0 \pm \Delta w_j) x_j = \sum_{j=1}^2 (w_{i,j}^0 (x_j \pm \Delta x_j)) \quad (52)$$

implying that a delta percentage of weight perturbation corresponds to the same percentage of x-data perturbation, that is to  $\Delta$ -robustness.

Consider now the case of a multilayer network with a graph corresponding to a binary tree and  $P = 1$ . Then a “flatness”  $\pm \Delta$  in all weights – we now consider a fraction  $\Delta$  of the maximum range of the weights (which are between  $\pi$  and  $-\pi$ ) – corresponds to the following maximum perturbation in the contribution of component  $x_j$  to the output (when it is not zero because of ReLU)

$$\left[ \prod_{d=1}^L (1 \pm (\Delta) w_j^d) \right] x_j. \quad (53)$$

Equation 53 has to be compared with the effect of perturbed data

$$\left[ \prod_{d=1}^L w_j^d \right] (1 \pm \Delta) x_j. \quad (54)$$

This suggests the following

**Conjecture 3.** *Invariance of the empirical error to perturbations of all weights are equivalent to larger, identical perturbations of all training examples  $x_i$  with  $i = 1, \dots, n$ .*

As an example, supposed there are  $d = 8$  layers and that  $\Delta = \frac{1}{10}$ . Then this tolerance in all the weights corresponds to a tolerance in the data perturbation of a factor 2 (1.1 vs 2.14). In addition, note also that there are usually many more weights than data. This relatively small  $\delta$  “flatness” of the zero-minimizer in weight space yields large robustness in data space (in this picture, to the same perturbations for all data points). The product of the weights that appears in the equations above is related to the path norm that we will discuss later (see also Figure 18).

## 9.7 Stability of SGD

In this section we discuss the equivalence of the following properties (some have appeared in the literature):

- flat minima in the weights of the minimizer
- stability of the minimizer wrt perturbations of weights
- stability wrt perturbations of the data
- fewer bits needed for parameters
- “small” local complexity (Ramacher averages) of functions that are flat
- stability of SGD in the sense of<sup>[38, 39]</sup>

Stability plays an important role in several different situations such as in inverse problems, in dynamical systems – under the form of “structural stability” – and in learning – as stability under perturbations of the training set. The intuition is that these are very similar forms of stability but formally the different cases requires different definitions and their mathematical equivalence remains in general elusive.

We discuss here at the intuitive level the relation between different types of perturbations in the specific case of SGD, that we define as before as

$$f_{t+1} = f_t - \gamma_n \nabla V(f_t, z_t), \quad (55)$$

- *SGD can be written as Gradient Descent with additive noise*

We rewrite

$$\nabla V(f_t, z_t) = \nabla I_{S_n}(f_t) + \xi_t \quad (56)$$

As we discussed earlier, the SGD update step is then rewritten as the following Langevin equation

$$f_{t+1} = f_t - \gamma_n(\nabla I_{S_n}(f_t) + \xi_t) \quad (57)$$

- *Perturbing the training examples in SGD – which we described in the context of perturbation stability or robustness – is equivalent to a pointwise change in the noise process in the equivalent noisy GD* Suppose that  $z_j$  is replaced by  $z'_j$  for a given  $j$  in  $S_n$ . Then  $U(f_t)$  is replaced by  $U'(f_t) = U(f_t) - V(f_t, z_j) + V(f_t, z'_j) = U(f_t) + \xi'_j$  where  $\xi'_j$  is a “noise” term.

- *Perturbing the hypothesis  $f_t$ , that is the weights, is equivalent to adding an additional noise term to noisy GD*

Suppose the parameters  $w$  are by replacing them component-wise with  $w + \tau w$ . Then  $\nabla_w U(f_f) \rightarrow \nabla_w U(f_f)(I + \tau) = \nabla_w U(f_f) + \xi$ .

The three cases end up being equivalent to noisy GD with slightly different noise terms. Notice that GD and SGD converge to a minimum norm solution in the linear degenerate case.

We conjecture that under some assumptions each of the cases above can be transformed in another one. This means that GD with replacement can be seen as noisy GD and this in turn may be seen under appropriate conditions as GD with replacement. The relevance of this observation is that stability wrt data is thus ensured by the presence of some form of noise in GD such as provided by the SGD update.

Notice that convergence to the same minima or to a class of equivalence is a form of *structural stability of the dynamical system* associated with SGD. Our intuition is that “repeat SGD” generalizes because the “noisy” update of the gradient makes it equivalent to “pure SGD” for  $N \rightarrow \infty$  in  $S_n$ .

- The key to generalization by “repeat SGD” is the intrinsic noise associated with SGD. This is different wrt GD. Because of properties of the Boltzman distribution in high dimensions, SGD finds with high probability *structurally stable* flat zeros (that is robust, wrt to perturbations of the weights). These via intrinsic regularization of the gradient descent steps coincide with stable solutions that generalize.
- The relative amount of noise can be controlled by changing the size of the minibatches.
- SGD solutions generalize because each example  $z_n$  with its repetitions (multiple passes) provides effectively different data points because of the noise term. This is formally correct if we could prove that perturbing the sample  $z_n$  is equivalent to a change in the noise term which does not affect the properties of SGD and associated theorems. Then the pure SGD theorem applies and convergence to the expected risk is guaranteed.
- Another intuition is that the solutions found by SGD are not just zeros of the polynomial set of equations of Theory II but “stable zeros” because of the noise and the fact that SGD is a Langevin equation converging in probability to the minima of the cost functions and preferring the stable ones.

## 9.8 Compositionality Improves Generalization Bounds

It has been long recognized that capacity of the set of neural networks can be controlled not only through the number of units (and, hence, the VC dimension) but also through the size of the weights <sup>[16]</sup>. In particular, the argument of <sup>[40]</sup> leads to a bound on the generalization error in terms of  $\ell_1$ -norms of the weights into each neuron. If each weight in the network is bounded, the result yields a non-vacuous generalization guarantee for any network with number of incoming connections into each neuron being  $o(n)$ , the sample size. For completeness, we provide the statement of <sup>[40]</sup> here, with minor modifications.

Define

$$\mathcal{F}_1 = \{x \mapsto \langle w^1, x \rangle : \|w^1\|_1 \leq B_1\} \quad (58)$$

and

$$\mathcal{F}_i = \left\{ x \mapsto \sum_j w_j^i \sigma(f_j(x)) : f_j \in \mathcal{F}_{i-1}, \|w^i\|_1 \leq B_i \right\} \quad (59)$$

for  $i \geq 2$ , and suppose  $x \in \mathcal{X} \subset [-1, 1]^d$ .



**Lemma 5.** Consider a  $k$ -layer neural network defined recursively as in (59). Let  $V : \mathbb{R} \times \{\pm 1\} \rightarrow \mathbb{R}$  be 1-Lipschitz loss function, and write  $V(f, z) = V(f(x), y)$  for  $z = (x, y) \in \mathcal{X} \times \{\pm 1\}$ . Let  $\sigma : \mathbb{R} \rightarrow \mathbb{R}$  be 1-Lipschitz and  $\sigma(0) = 0$ .<sup>8</sup> Then

$$\mathbb{E} \sup_{f \in \mathcal{F}_k} \left\{ \mathbb{E} V(f, Z) - \frac{1}{n} \sum_{t=1}^n V(f, Z_t) \right\} \leq 2\mathbb{E} \widehat{\mathcal{R}}(\mathcal{F}_k) \leq \frac{\sqrt{2 \log d} \cdot \prod_{i=1}^k (2B_i)}{\sqrt{n}}, \quad (60)$$

where  $\widehat{\mathcal{R}}(\mathcal{F}_k)$  are empirical Rademacher averages of  $\mathcal{F}_k$ .

*Proof.* Using the standard symmetrization argument followed by contraction, the uniform deviations in (60)

$$2\mathbb{E} \sup_{f \in \mathcal{F}_k} \frac{1}{n} \sum_{t=1}^n \epsilon_t V(f, Z_t) \leq 2\mathbb{E} \sup_{f \in \mathcal{F}_k} \frac{1}{n} \sum_{t=1}^n \epsilon_t f(X_t). \quad (61)$$

Now, for any  $i \geq 2$ ,

$$\widehat{\mathcal{R}}(\mathcal{F}_i) = \mathbb{E} \sup_{f_j \in \mathcal{F}_{i-1}, \|w^i\| \leq B_i} \frac{1}{n} \sum_{t=1}^n \epsilon_t \sum_j w_j^i \sigma(f_j(X_t)) \leq 2B_i \widehat{\mathcal{R}}(\mathcal{F}_{i-1}) \quad (62)$$

where the last step uses the Cauchy-Schwartz inequality. Finally,

$$\widehat{\mathcal{R}}(\mathcal{F}_1) = B_1 \left\| \frac{1}{n} \sum_{t=1}^n \epsilon_t X_t \right\|_{\infty} \leq B_1 \sqrt{\frac{2 \log d}{n}}$$

using the usual estimates.  $\square$

Lemma 5, together with the empirical observation that neural networks can be trained to achieve zero error on the data, implies that the expected risk of the minimizer is converging to zero. Lemma 5 assumes that  $V$  is Lipschitz. However, the standard argument from the theory of large margin classifiers yields a guarantee on the expected zero-one loss as well, in terms of the Rademacher averages of the class and the empirical margin. This result, which can be found in [15, Theorem 2] (see also [41] for the present version), is stated below:

**Theorem 5.** For all  $\delta > 0$ , with probability at least  $1 - \delta$ , for all  $f \in \mathcal{F}_k$  and  $\gamma > 0$ ,

$$P(Yf(X) \leq 0) \leq \frac{\text{card}(i : Y_i f(X_i) < \gamma)}{n} + \frac{4}{\gamma} \mathbb{E} \widehat{\mathcal{R}}(\mathcal{F}_k) + \sqrt{\frac{\log \log(4C/\gamma)}{n}} + \sqrt{\frac{\log \delta^{-1}}{2n}}, \quad (63)$$

where  $C \geq \sup_{f,x} f(x)$ .

The bound of (5) on Rademacher averages can be used in Theorem 5 to obtain an upper bound on the expected zero-one loss of a solution that has a large margin  $\gamma$  on the data. Notice that  $\gamma$  optimizing the above bound is based on the empirical margin for the particular draw of the data.

We observe that the covering numbers are much better for compositional networks than for densely connected networks consistently with section 9.1; the covering number are an upper bound to the Rademacher averages.

## 9.9 Compositionality and Multiclass Tasks

Most of the successful neural networks exploit compositionality not only to avoid the curse of dimensionality but also for better generalization in an important way (see [42]). Suppose that the mappings to be learned in a family of classification tasks (for instance classification of different object classes in Imagenet) may be approximated by compositional functions such as  $f(x_1, \dots, x_n) = h_l \cdots (h_{21}(h_{11}(x_1, x_2), h_{12}(x_3, x_4)), h_{22}(h_{13}(x_5, x_6), h_{14}(x_7, x_8) \cdots)) \cdots$ , where  $h_i$  depends on the task (for instance to which object class) but all the other constituent functions  $h$  are common across the tasks. Under such an assumption, multi-task learning, that is training simultaneously for the different tasks, forces the deep network to “find” *common constituent functions*. Multi-task learning has theoretical advantages that depends on compositionality: the sample complexity of the problem can be significantly lower (see [43]). The Maurer’s approach is in fact to consider the overall function as the composition of a preprocessor function common to all task followed by a task-specific function. As a consequence, the *generalization error*, defined as the difference between expected and empirical error, averaged across the  $T$  tasks, is bounded with probability at least  $1 - \delta$  (in the case of finite hypothesis spaces) by

$$\frac{1}{\sqrt{2M}} \sqrt{\ln |\mathcal{H}| + \frac{\ln |\mathcal{G}| + \ln(\frac{1}{\delta})}{T}}, \quad (64)$$

<sup>8</sup>This assumption can be easily removed.

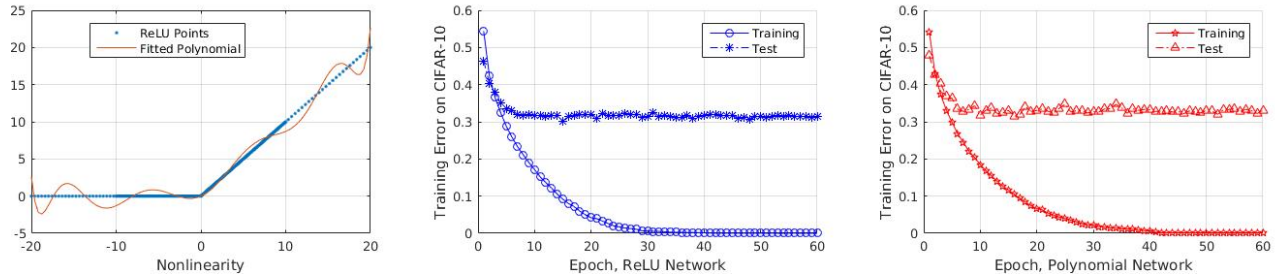


Figure 21: A standard convolutional deep network is converted into a polynomial function by replacing each of the REL units with a univariate polynomial approximation of the RELU function. As long as the nonlinearity approximates ReLU well (especially near 0), the “polynomial network” performs similarly to a ReLU net. The polynomial shown in the inset on the left is of degree 10.

where  $M$  is the size of the training set,  $\mathcal{H}$  is the hypothesis space of the common classifier and  $\mathcal{G}$  is the hypothesis space of the system of constituent functions, common across tasks.

The improvement in generalization error because of the multitask structure can be in the order of the square root of the number of tasks (in the case of Imagenet with its 1000 object classes the generalization error may therefore decrease by a factor  $\approx 30$ ). It is important to emphasize the dual advantage here of compositionality, which a) reduces generalization error by decreasing the complexity of the hypothesis space  $\mathcal{G}$  of compositional functions relative the space of non-compositional functions and b) exploits the multi task structure, that replaces  $\ln|\mathcal{G}|$  with  $\frac{\ln|\mathcal{G}|}{T}$ .

## 9.10 Polynomial approximation and deep polynomial networks

Nonlinear activations in a deep networks can be approximated up to an arbitrary accuracy over a bounded interval by a univariate polynomial of appropriate degree (see <sup>[2]</sup>). Since so much is known about polynomials several interesting results follow from these approximation properties such as a characterization of when deep convolutional networks can be exponential better than shallow networks in terms of representational power (<sup>[1]</sup>). However, not all properties of the polynomial approximants hold for the target function – the deep network. An example is the exact number of zeros of function, which does not follow directly from approximation arguments.

We show in this section that the puzzling generalization properties of deep convolutional networks hold for networks with the same architecture in which all the RELU’s have been replaced by a univariate polynomial approximation of degree 10. The resulting networks are obviously polynomial networks. Figure 21 demonstrates that polynomial and RELU networks with the same architecture have very similar performance in terms of training and testing error. Figures 21 and 23 show the same qualitative behavior of testing and training error on natural labeled data and on random data. The same puzzling lack of overfitting in the overparametrized case can be seen in Figure 24. Together Figures 22,23, 24 include all the puzzles discussed in <sup>[44]</sup>. The computational overhead posed by calculation of a polynomial rather than RELU is negligible, because computation of the activations is a very small fraction of the total computational budget involved in training.

The rest of the paper is dedicated to explaining this puzzling behavior of polynomial networks. Obviously the explanation for polynomial networks is likely to apply to standard deep networks.

### 9.10.1 Deep Puzzles

Classical learning theory characterizes generalization behavior of a learning system as a function of the number of training examples  $n$ . From this point of view deep learning networks behave as expected: the more training data the better the test error. But other aspects of their learning curves are puzzling:

- the test error decreases for increasing  $n$  even when the training error is zero.
- the same network which shows nontrivial generalization for normally labeled training examples has zero training error for randomly labeled data
- there is no apparent overfitting for overparametrization of the network – that is when the number of weights is much larger than the number of data points.

The figures in Appendix 1 (Figures 1, 2, 3 ) show these phenomena for a standard convolutional network with RELUs units, trained on CIFAR10.

### 9.10.2 Polynomial approximation of RELUs and deep polynomial networks

The mapping from input to output of a deep polynomial network is

Model #params: 9370

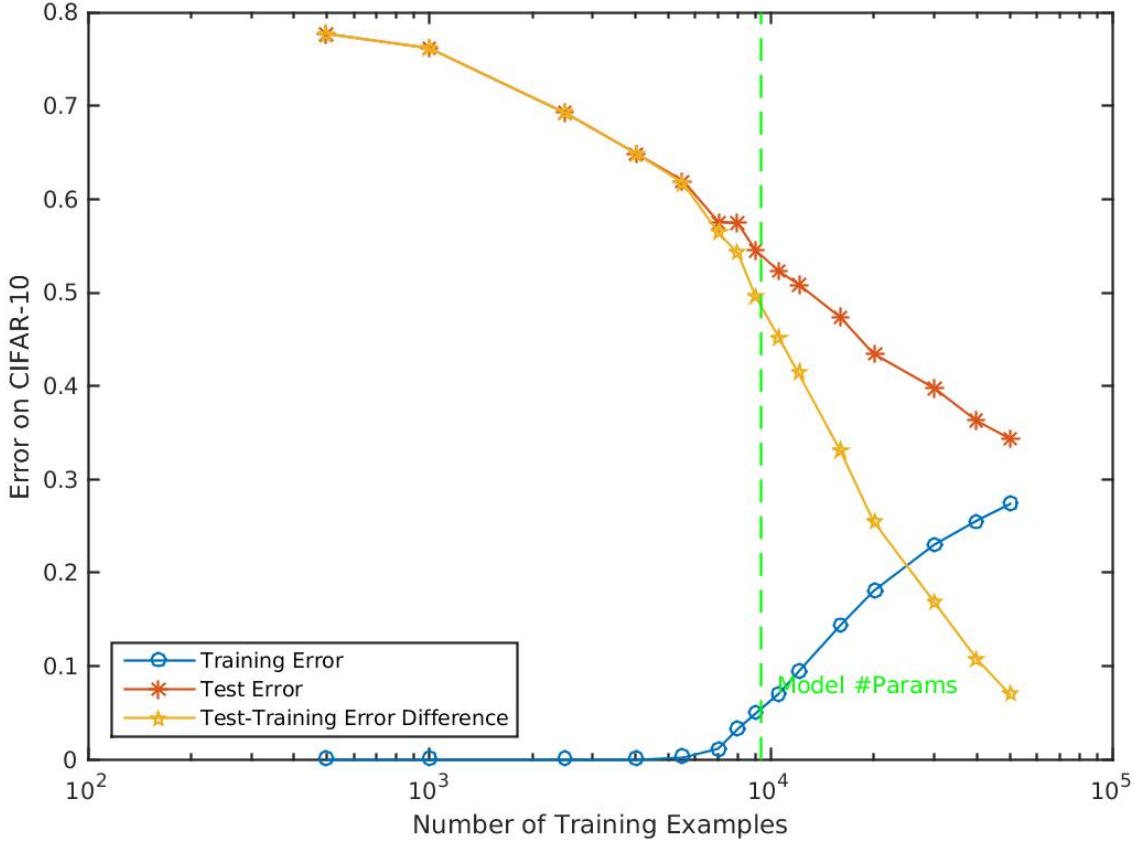


Figure 22: The figure shows the behavior of a polynomial deep network trained on subsets of the CIFAR database. The network is a 5-layer all convolutional network (i.e., no pooling) with 16 channels per hidden layer, resulting in only  $W \approx 10000$  weights instead of the typical 300,000. Neither data augmentation nor regularization is performed.

$$y = P_k(x) \quad y \in \mathcal{R}^m, x \in \mathcal{R}^d, \quad (65)$$

where  $P_k$  is a polynomial in  $x$  of degree  $k = 10 \times l$  – since 10 is the degree chosen for the univariate approximation of the RELU of Figure 21) and  $l$  is the number of layers.

A generic polynomial  $P_k(x)$  of degree  $k$  in  $d$  variables is in the linear space  $\mathcal{P}_k \cup_{s=k}^k H_s$  composed by the union of homogeneous polynomials  $\mathcal{H}_k$  of degree  $k$ .  $\mathcal{H}_k$  is of dimensionality  $r = \dim \mathcal{H}_k = \binom{d-1+k}{k}$ : the latter is the number of monomials and thus the number of coefficients of the polynomials in  $\mathcal{H}_k$ . A generic polynomial in  $\mathcal{P}_k$  can always be written (see Proposition 3.5 in [45]) as

$$P_k(x) = \sum_{i=1}^r p_i(\langle w_i, x \rangle). \quad (66)$$

where  $p_i$  is a univariate polynomials of degree at most  $k$  and  $\langle w_i, x \rangle$  is a scalar product. In fact

$$P(x) = \sum_j a_j x^j = \sum_{k=1}^N c_k ((w_k, x) + b_k)^n = \sum_{k=1}^N c_k \sum_{|j|=n} \binom{n}{j} u_k^j y^j \quad (67)$$

with  $y = (x, 1)$  and  $u_k = (w_k, b_k)$  so that  $a_j = \binom{n}{j} \sum_{k=1}^N c_k u_k^j$ .

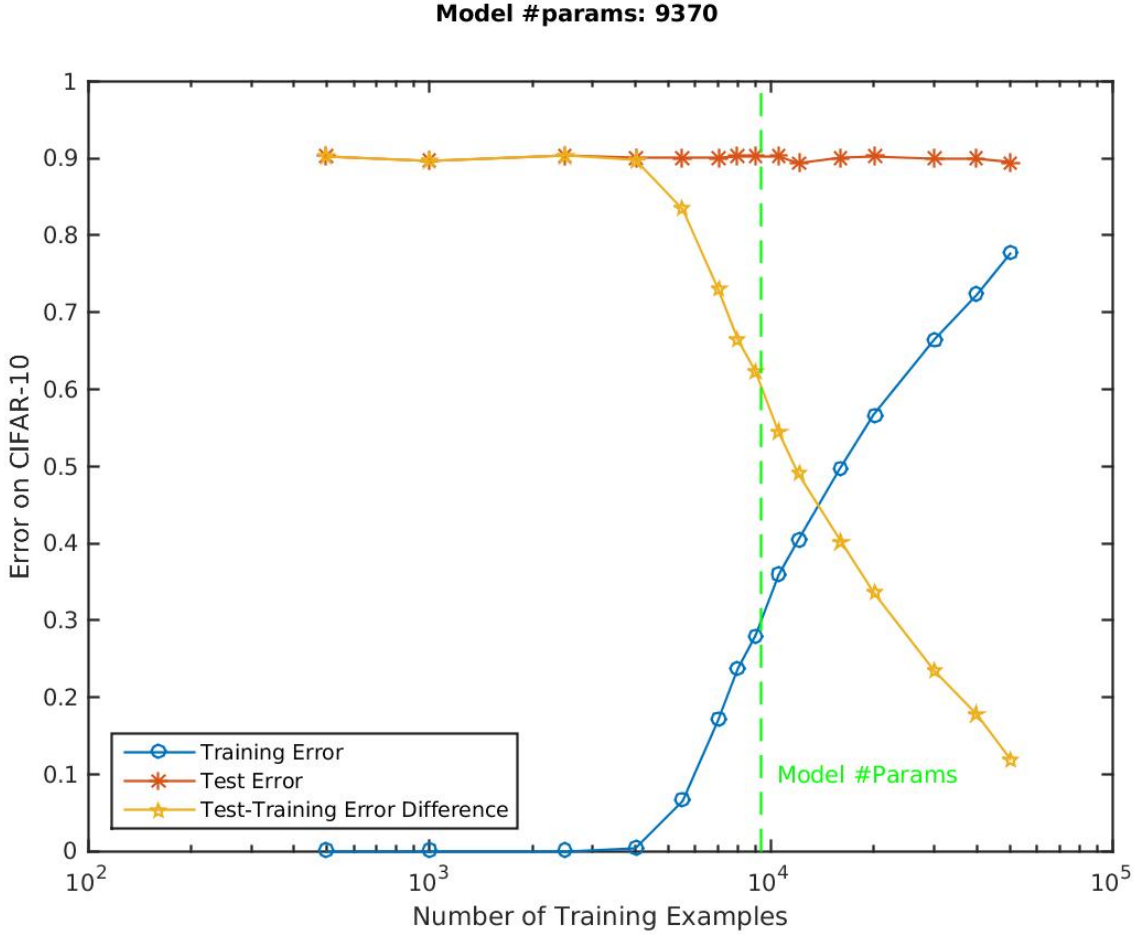


Figure 23: The figure shows the behavior of a polynomial deep network trained on subsets of the CIFAR database in which the labels have been randomly scrambled. The network is a 5-layer all convolutional network (i.e., no pooling) with 16 channels per hidden layer, resulting in only  $W \approx 10000$  weights instead of the typical 300,000. Neither data augmentation nor regularization is performed.

Two remarks about Figure 25 are in order:

- in the case of hierarchically local<sup>[1]</sup> functions and networks – convolutional networks are a special case – the associated polynomials are very sparse with a number of monomials that grows linearly with  $d$  instead of exponentially ( $\binom{d-1+k}{k} \approx k^d$ ).
- Notice that the effective  $r$ , is closely related to the *separation rank* of a tensor (see <sup>[1]</sup> which provides the mathematical foundations for the connection between deep nets and the HT tensor decomposition).

Using the notation of <sup>[46]</sup>, we can regard Equation 65 as a polynomial operator on the vector space  $\mathcal{R}^d$ . For instance, consider in Figure 25 the node in the second layer that receives  $x_1, x_2, x_3, x_4$  inputs. The outputs of the node are of the form

$$\sum_{j=1}^{r'} p_j [W_1 \sum_{i=1}^r p_i (w_{1,i}x_1 + w_{2,i}x_2) + W_2 \sum_{i=1}^r p_i (w_{3,i}x_3 + w_{4,i}x_4)]. \quad (68)$$

with  $p_h(x) = c_h p(x)$  and  $p(x)$  is the polynomial activation function.

The key point in our argument is that a polynomial operator of degree  $k$  can be regarded as a linear mapping  $\mathcal{L}$  from the (feature) space of all tensor products of  $x$  up to degree  $k$  to the space  $\mathcal{Y}$ .

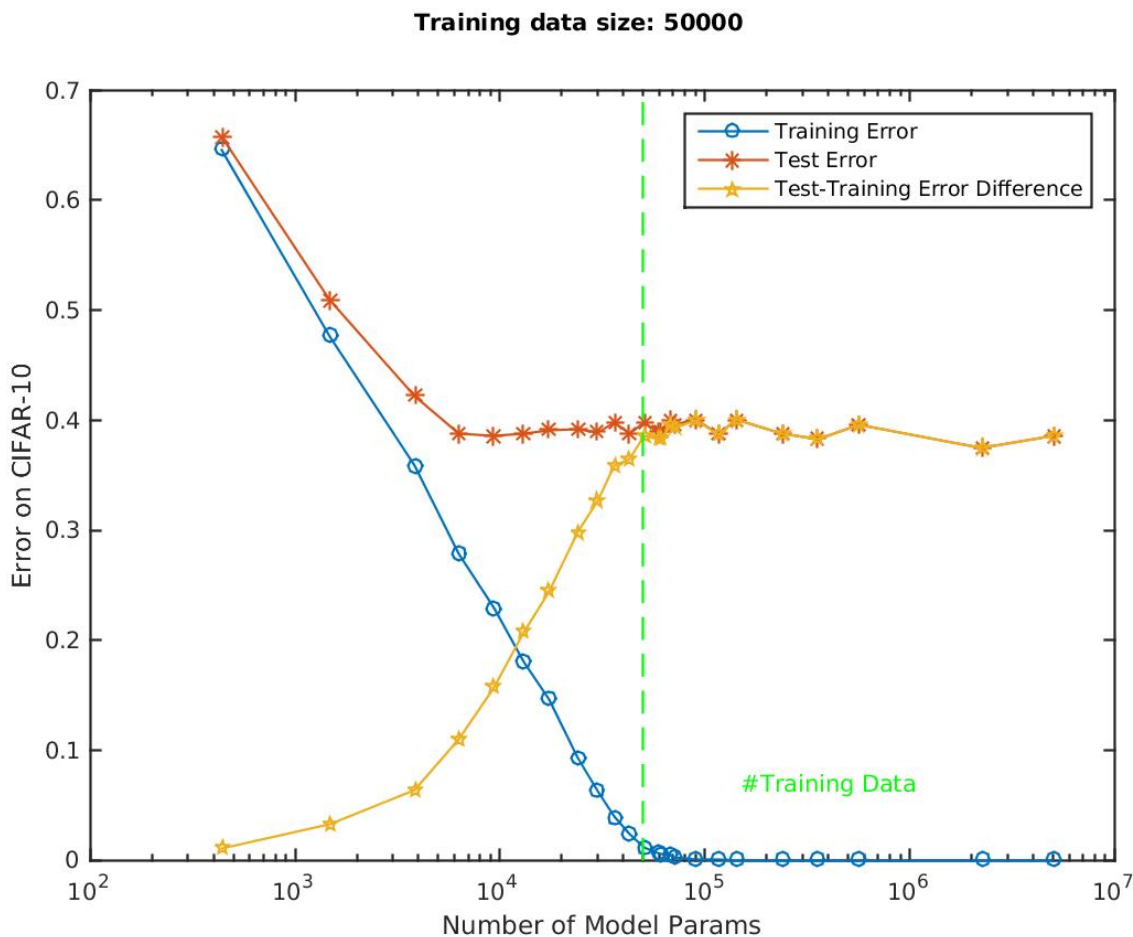


Figure 24: The previous figures show dependence on  $N$  – number of training examples – for a fixed architecture with  $W$  parameters. This figure shows dependence on  $W$  for a fixed training set with  $N$  examples. The network is again a 5-layer all convolutional polynomial network. All hidden layers have the same number of channels. Neither data augmentation nor regularization is performed. The classical theory explains the generalization behavior on the left; the challenge is to explain the lack of overfitting for  $W > n$ . As shown here, there is zero error for  $W \geq n$ .

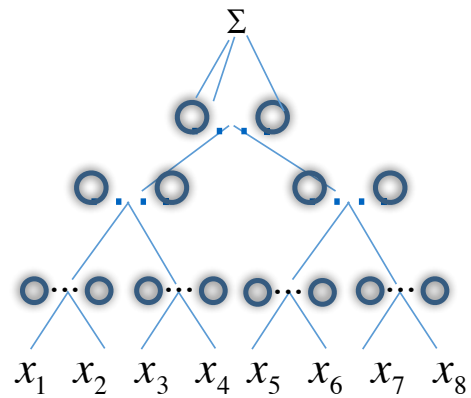
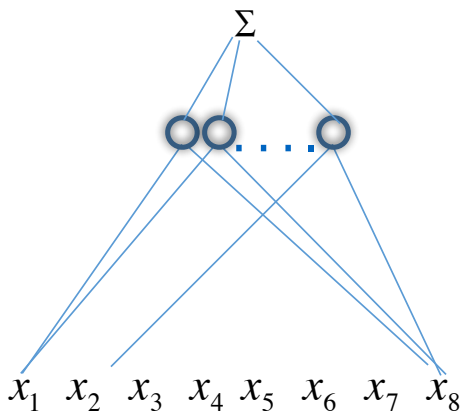
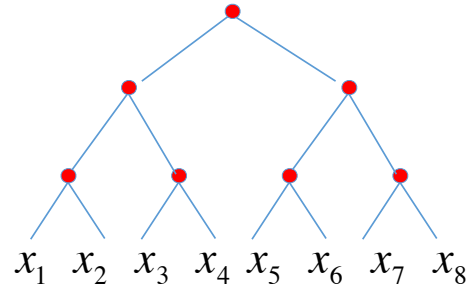
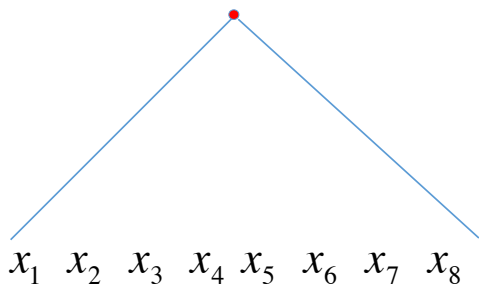
### 9.10.3 Solving the generalization puzzle for the standard representation of polynomial networks

Consider learning  $\mathcal{L}$  from a set of example pairs  $x_i, y_i$  represented as columns of the matrices  $X, Y$  with  $X$  a  $d, n$  matrix and  $Y$  a  $m, n$  matrix. In this case  $\mathcal{L}$  is a linear operator  $\mathcal{L} : V^{\otimes k} \mapsto Y$  where  $V^{\otimes k} = \cup_{i=1}^k V^i$  and  $V^i$  is the linear space of the homogeneous polynomials in  $x$  of degree  $i$ . Notice that the dimensionality of each of the spaces of homogeneous polynomials  $V^i$  is high:  $\dim V^i = D = \binom{d+i}{i} \approx i^d$ . Let us denote with  $L$  the  $m, D$  matrix corresponding to  $\mathcal{L}$ , with  $V$  the  $D, n$  matrix corresponding to example vectors  $x_i$  and with  $Y$  the  $m, n$  matrix of the  $y$  examples. The best solution in the  $L_2$  sense of equation  $LV = Y$  is then  $L = YV^\dagger$ . Since the matrix  $V$  is so large, only iterative techniques can be used to provide a solution. Two different constructions are sketched in the Appendix.

The arguments above apply to polynomials in the usual format. As an example consider the polynomial  $P(x_1, x_2)$  of degree 2 with  $P = a_0 + a_1x_1 + a_2x_2 + a_3x_1^2 + a_4x_2^2 + a_5x_1x_2$  with  $\binom{2+2}{2} = 6$  monomials. We will refer to this as the *standard parametrization* of a polynomial mapping where there is one coefficient  $a_j$  for each monomial in  $x$ .

### 9.10.4 Weight parametrizations of polynomial networks

Such a polynomial arises from the hierarchical network of Figure 25. For instance the output of the first node on the bottom left of part (b) of the Figure is the linear combination of 6 terms such as  $(w_1^i x_1 + w_2^i x_2 + b^i)^2$  (assuming a nonlinear activation to be a quadratic polynomial). The network is directly parametrized in terms of *sums of products of powers of the weights* for each monomial. We call this parametrization *the weights parametrization*. Each  $a_i$  can be expressed in terms of the  $w_i$  and viceversa. Notice (this observation is due to H. Mhaskar) that the representation of a polynomial in  $d$  variables with total degree  $\leq k$  (so dimension  $N$ ) with monomial or any other fixed basis involves exactly  $N$  parameters to optimize. In the representation that is produced by a single hidden layer in a network as  $\sum_k a_k (w_k \cdot x + b_k)^n$ ,



*a*

*b*

Figure 25: The top graphs are associated to functions; each of the bottom diagrams depicts a polynomial network approximating the function above. In a) a shallow polynomial in 8 variables and  $N$  units approximates a function of 8 variables  $f(x_1, \dots, x_8)$ . If the activation for each of the hidden units is a polynomial of degree 10 then  $N = \binom{8+10}{8}$  since a polynomial in  $d$  variables of degree  $k$  can be written as a linear combination of the same number  $N$  of univariate polynomials of degree up to  $k$  in  $\langle w, x \rangle$ . The bottom of inset b) shows a binary tree hierarchical network in  $n = 8$  variables, which approximates well functions of the form  $f(x_1, \dots, x_8) = h_3(h_{21}(h_{11}(x_1, x_2), h_{12}(x_3, x_4)), h_{22}(h_{13}(x_5, x_6), h_{14}(x_7, x_8)))$  as represented by the binary graph above. In the approximating network each of the  $n - 1$  nodes in the graph of the function corresponds to a set of polynomial units. In a binary tree with  $d$  inputs, there are  $\log_2 d$  levels and a total of  $d - 1$  nodes. The locality of the constituent functions – that is their low dimensionality – means that the overall polynomial computed by the network is sparse.

there are  $(d + 2)N$  parameters which are not independent of each other. In particular, there is no uniqueness for this representation, let alone stability and other such nice attributes. A simple way to make it unique is to fix the  $w_k$ 's and  $b_k$ 's, thereby giving a fixed basis again. An example of the connection of the two representations is  $a_5 = 2 \sum_{i=1}^6 (w_1^i w_2^i)$ . Because of the hierarchical (in particular, convolutional) architecture of the networks, it is enough to have at each layer enough units to express with the required accuracy the relevant constituent polynomial<sup>[2]</sup>.

### 9.10.5 Making sense of the puzzle (and a technical conjecture)

The global minima (in terms of classes of equivalence) for the standard vs weight parametrization of deep polynomial networks are equivalent (SGD finds global minima in both cases<sup>[33]</sup>). Assuming that the loss function is convex, adding a strictly convex (for instance quadratic) regularization term ensures that the the minimum generalizes in both cases (see for instance <sup>[25, 47, 48]</sup>). The minimizers will in general be different in the standard vs the weight parametrization case, since (in the  $L_2$  case) the constraint  $\min \|a\|^2$  is different from  $\min \|w\|^2$ . Notice that SGD and GD can find the  $\min \|a\|^2$  even without an explicit regularizer *but only if the initial conditions are appropriate*. Empirical results showing that a  $\min \|w\|^2$  constraint does not significantly improve the SGD solution suggest that something similar to the linear case is happening here: optimization by SGD and GD for initial conditions with small norm provides a small norm solution<sup>[49]</sup>. Thus it seems likely that a small norm solution provides good generalization. Empirical results on simple multivariate polynomials suggest that solutions found by SGD are really bad only when the initial conditions have very large norm. Notice that minimizing  $\|a\|^2$  corresponds to minimize a "path-like" norm<sup>[50]</sup> since each  $a_i$  corresponds to  $\sum_{j,k} w_1^p \dots w_l^q$  (see <sup>[50, 13]</sup>). It is important to notice that in all cases – in the linear, polynomial and weight parametrization cases – *the well-posed, robust way to obtain generalization is to add a regularization term*. It is intriguing – *though, in a sense, a theoretical footnote* – that regularization may not be strictly needed in the absence of noise and with good initial conditions.

Thus we have the following summary of our empirical and theoretical argument:

#### If we assume that

1. the loss function is the square loss;
2. there is overparametrization in terms of  $w_i$ , wrt to number of data points, ensuring that degenerate global minimizers with zero error exists;
3. additionally, each significant  $a_j$  may be expressed redundantly in terms of  $w_i$  (this assumption is not strictly necessary for the main results);

#### then

1. there is a very large number of global minimizers in the  $w_i$  with zero loss – because of Bezout theorem<sup>[2]</sup> – which are degenerate when the equations are not inconsistent (zero training error implies that the equations  $y_i = f(x_i)$  are consistent);
2. SGD selects large volume minima which tend to be flat wrt all or most  $w$  (experiments validate this prediction<sup>[51]</sup>, implying that SGD selects preferentially degenerate minima that are global, avoiding local minima and saddle points);
3. theory shows that a small amount of quadratic regularization in the  $w_i$  (such as weight decay) guarantees a minimum norm solution that generalizes (this is because if  $\min_w F(w, x)$ , with  $F$  convex, has degenerate solutions in  $w^*$  with  $F(w^*, x) = 0$ , then  $\min_w F(w, x) + \lambda \|w\|^2$  yield a unique  $w^*$ );

The theoretical arguments above are simple and classical. Empirically, a regularization term is used most of the times via weight decay. Of course a puzzle remains: empirical evidence, as in several figures of this paper, shows that generalization is obtained even without an explicit regularization term. How can this happen? For the standard parametrization of a polynomial, the arguments in this note show that the problem is linear. In the linear case it is well known, as shown in the SI, that generalization without regularization can be achieved with initial conditions of small norm of the coefficients  $a_i$ . Notice that, in a strict sense, such generalization is not robust (e.g. independent of initial conditions and of noise).

Though the above argument does not apply directly to optimization wrt  $w_i$ , it suggests the following scenario. In the linear case degenerate directions - determined by the data – have zero gradient and the corresponding coefficients are therefore not changed by gradient descent (and are not changed in expectation by stochastic gradient descent) as shown in SI section 9.6. This is the underlying reason for the minimum norm solution (with initial conditions of zero norm) for the degenerate directions of the polynomial in the standard parametrization case. Assume for simplicity that a degenerate direction corresponds to the  $k$ th monomial in the standard parametrization (more in general it will correspond to a linear combination of monomials determined by the SVD of the correlation matrix between the input matrix  $V^\dagger$  and the output  $Y$ ). Then the associated coefficient  $a_k$  will not be changed by GD. In particular, its norm will remain small if it was set to be small at the beginning of optimization. Consider now what happens in the  $w_i$  parametrization of the same monomial, where  $a_k$  is now represented as the sum of products of several  $w_i$  (we assume there is more than one weight per each product term; notice that the total number of  $w$  weights is expected to be orders of magnitude larger than the number of *effective*  $a_i$  – that is those  $a_i$ , or more in general those linear combinations of  $a_i$ , associated with significant monomials). The gradient wrt each of the weights involved in the specific (degenerate) monomial will not affect the corresponding component in the loss and therefore will not be updated in the process of GD. Of course, one of

those weights may be updated because it is also present in another term  $a_j$  that is not degenerate. However, the value of the degenerate product will remain small (if all weights were initialized to be small  $\approx \epsilon$ ) because the probability that they are all present in other non degenerate terms is very small, given their redundancy. In other words, we assume that for all degenerate coefficients  $a_i$  there will be at least one  $w_j$  that does not appear in any of the non-degenerate  $a_i$ . We expect then that only the products of  $w_i$  associated with a nondegenerate direction have a significant probability to be updated during GD. Notice that we know that there is degeneracy in the weights (because of theoretical reasons<sup>[2]</sup> and experimental measurements of flatness<sup>[51]</sup>). Thus, SGD and GD wrt to the weight parametrization would provide implicit regularization in a way which is very similar to the situation for the standard parametrization of the same polynomial.

This conjecture above would imply that each sums of products of powers of the  $w_j$ , corresponding to a degenerate direction associated with  $a_k$ , will maintain with high probability the same small norm it had at the beginning of the GD or SGD iterations. This control of the associated path-like norm would of course correspond to a regularization effect, exactly as in the linear case of the standard parametrization of the polynomial. In the SI we suggest how it could be formally proved. It should be noted that a somewhat similar situation arises in optimizing tensor representations, which are related to the networks described here, with the alternating least square (ALS) technique (see SI and for instance <sup>[52]</sup>).

A possible way to prove the conjecture is based on the following two observations on networks with one hidden layer. The first is for units with linear activation, the second for units with polynomial activation.

- Take the gradient wrt  $W_1$  of the loss  $L(w) = \|W_2 W_1 X - Y\|^2$ . We denote  $E = W_2 W_1 X - Y$ ,  $E \in \mathbb{R}^{d',n}$ ,  $W_2 \in \mathbb{R}^{d',p}$ ,  $W_1 \in \mathbb{R}^{p,d}$  and obtain

$$dW_1 = -\gamma \nabla_{W_1} L(w) = W_2^\top E X^\top = E' X^\top \quad (69)$$

where  $E' = W_2^\top E$ ,  $E' \in \mathbb{R}^{p,n}$ . The matrix  $E'$  is projected onto the span of the data  $X$  and thus there is no increment  $dW_1$  at each iteration that is in the space orthogonal to the space spanned by  $X$ . If the norm of  $W_1$  in the space orthogonal to the span of the data was zero at the start, it will remain zero.

- Consider the loss  $L(w) = \|P_m(X) - Y\|^2$  where

$$P_m(x) = (W_2)_{1,i} (\langle (W_1)_{i,j}, x_j \rangle)^{\textcircled{m}}. \quad (70)$$

where repeated indices denote summation and the power  $m$  is elementwise.

We obtain, denoting  $E = P_m(X) - Y$ , with  $E \in \mathbb{R}^{d',n}$ ,  $W_2 \in \mathbb{R}^{d',N}$ ,  $W_1 \in \mathbb{R}^{N,d}$ ,  $E' \in \mathbb{R}^{d',d}$

$$\nabla_{W_1} L(w) = (W_2^\top E (W_1 X)^{\textcircled{m}1}) X^\top = E' X^\top \quad (71)$$

As in the first case, this implies that  $\nabla L(w)$  is in the span of the rows of  $X^\top$  and thus is a linear combination of the columns of  $X$ , that is a linear combination of the data points. We now observe that each column  $x$  of  $X$  can be lifted to a vector comprising all monomials up to degree  $k$  spanning  $P_k(x)$ . The same lifting can be done with the rows  $w$  of  $W$ . Lemma 2.1 in Poggio <sup>[53]</sup> observes that  $(WX)^{\textcircled{K}} = \mathbb{W}\mathbb{X}$  (this is an early version of the “kernel trick” for polynomial kernels). For any  $x \in X^\perp$ ,  $Wx = 0$  and thus  $(Wx)^{\textcircled{K}} = 0$  implying that  $\mathbb{W}$  projects into the span of  $\mathbb{X}$ . Thus if  $\mathbb{W}$  is the solution that minimizes the loss, it is also the minimum norm solution in the space of the  $\mathbb{X}$ .

## Remarks

- It seems that one can take the limit for  $m \mapsto \infty$  In Equations 70 and 71. This would imply an extension of the result to activation functions such as the sigmoid and the ReLU function.
- On a similar vein, the limit argument may be replaced by an argument based on the kernel associated with the specific activation function (a kernel associated with the ReLU activation function is the order 1 arc-cosine kernel function<sup>[54]</sup>). The linear space of the monomial vectors then should be generalizable to the linear space of the features  $\phi(x)$  associated with the kernel. An explicit representation of such features is not needed to reach the main conclusions of the argument.

## 9.10.6 Discussion

The main conclusion of the paper is that robust, well-posed generalization in deep networks explaining most of the puzzles around their non-overfitting and predictive properties is ensured by a quadratic regularization term added to the loss function which is then optimized by SGD. We also suggest an answer to why, as suggested by empirical evidence, such a regularization term is not strictly needed and is implicitly provided by SGD optimization. It is well known for linear networks that, even without a regularization term, SGD and GD, when initialized with close-to-zero norm, yield a minimum norm solution which generalizes. We show that the same result applies to polynomial networks in the standard representation. It does not automatically follow, however, that the result still holds when the minimization is performed with respect to the weight representation in a hierarchical architecture. In this paper we provide evidence – both theoretical and empirical – in support of the conjecture that the implicit regularization provided by GD and SGD optimization wrt the weights  $w_i$  is



equivalent to the implicit regularization provided by optimization wrt the coefficients  $a_j$  that parametrize the same polynomial. It is satisfying that the implicit regularization by GD and SGD can be seen as completely equivalent to an explicit regularization term usually added during training. In fact, as shown in <sup>[25]</sup> GD and SGD perform incremental iterative regularization, effectively computing the pseudoinverse according to the following definition

$$A^\dagger = \lim_{\delta \rightarrow 0} A^\top (AA^\top + \delta^2 I)^{-1}$$

The conclusion is that there is nothing magic in deep learning that requires an explanation different from the classical linear one <sup>[25]</sup> with respect to generalization, intended as convergence of the empirical to the expected error, and with respect to the absence of overfitting in the case of overparametrization. Our result that deep learning networks optimized with gradient descent techniques do generalize is consistent with a number of previous results such as Bartlett's <sup>[26]</sup> and Recht <sup>[38]</sup> (without the need of additional technical assumptions such as a bound in the weights). The property of generalization, though important, is however of academic importance only, since deep networks are typically used in the overparametrized case, that is in a regime far from convergence. It is more important that our results provide a simple explanation – completely analogous to the pseudoinverse behavior for linear networks – for the absence of overfitting in deep polynomial networks. They also extend to nonlinear multi-layer networks the use linear techniques to analyze the convergence to the pseudoinverse solution and its dynamics.

The main difference between linear and deep networks is in terms of approximation power rather than generalization: deep convolutional network can avoid the curse of dimensionality <sup>[1]</sup> for the important class of compositional functions and are thus capable of using efficiently large training data set, thereby achieving good predictive performance. The results of this paper on robustness with respect to overparametrization imply that the architecture of the convolutional network does not need to be exactly matched to the underlying function graph: a relatively large amount of overparametrization does not affect the predictive performance of the network.