



UNDERSTANDING DEEP LEARNING REQUIRES RETHINKING GENERALIZATION

Aldo Lamarre ¹ Matthew C. Scicluna ²

Feb 21 2018

¹Département d'Informatique et de Recherche Opérationnelle
Université de Montréal

²Montréal Institute of Learning Algorithms
Université de Montréal

Table of contents

1. Introduction
2. Background
3. Results
4. Technical dive
5. Discussion

Introduction

Main Question

What distinguishes Neural Networks that generalize well from those that don't?

- Capacity ?
- Regularization ?
- How we train the model?

Traditional View

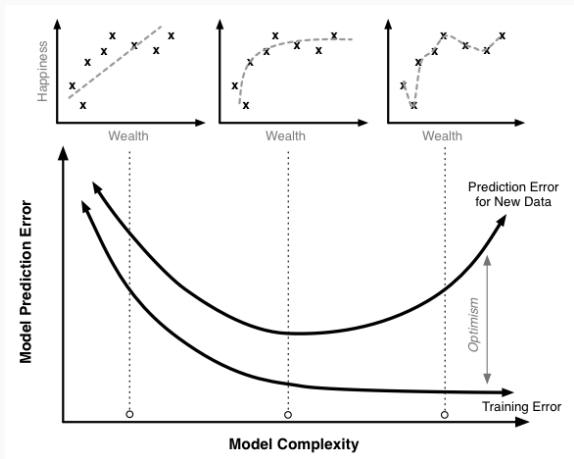


Figure 1: Traditional view of generalization. Image taken from [1]

Why do we care about the problem?

- Make neural networks more interpretable
- May lead to more principled and reliable model architecture design

Background

We can bound the Generalization Error using measures of complexity such as:

- VC Dimension
- Rademacher Complexity
- Uniform Stability

Additionally, regularization can help (including Early Stopping)

In 2016 Hardt et al. gives an Upper bound on Generalization error on model using SGD using uniform stability [2]

BUT

Uniform stability is a property of a learning algorithm and is not affected by the labelling of the training data.

Main Message

Classic results (e.g. PAC bounds) are insufficient in that they cannot distinguish between neural networks with dramatically different generalization performance.

This is demonstrated in the paper [3]. The central finding:

Deep neural networks easily fit random labels

Results

Setup: trained several standard architectures on the data with various modifications:

1. True labels → No modifications
2. Random labels → randomly changed some labels
3. shuffled pixels → apply some fixed permutation of pixels to all images
4. Random pixels → apply some random permutation of pixels to all images
5. Gaussian → Generate pixels for all images from a Gaussian

Main Results

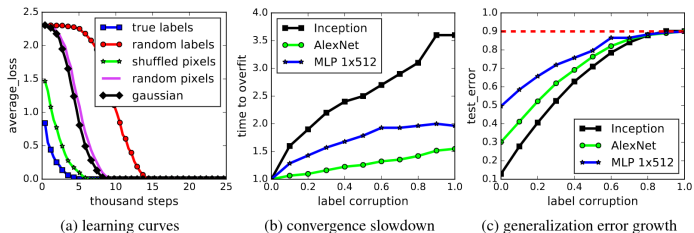


Figure 2: Fitting random labels and random pixels on CIFAR10.

In most cases, the training error went to zero while test error was high

Notice:

the model capacity, hyperparameters, and the optimizer remained the same!

Explicit regularization may improve generalization performance, but is neither necessary nor by itself sufficient for controlling generalization error

Table 4: Results on fitting random labels on the CIFAR10 dataset with weight decay and data augmentation.

Model	Regularizer	Training Accuracy
Inception	Weight decay	100%
Alexnet		Failed to converge
MLP 3x512		100%
MLP 1x512		99.21%
Inception	Random Cropping ¹	99.93%
	Augmentation ²	99.28%

Technical dive

The empirical observations are complemented with a theoretical construction showing that generically large neural networks can express any labelling of the training data.

The empirical observations are complemented with a theoretical construction showing that generically large neural networks can express any labelling of the training data.

Definition

Finite-sample expressivity : is the expressive power of neural networks on a finite sample of size n .

NB : It is possible to transfer population level results to finite sample results using uniform convergence theorems.

Theorem

There exists¹ a two-layer neural network with ReLU activations and $2n + d$ weights that can represent any function on a sample of size n in d dimensions.

¹NOT all networks satisfy this

Lemma 1

For any two interleaving sequences of n real numbers

$b_1 < x_1 < b_2 < x_2 \cdots < b_n < x_n$, the $n \times n$ matrix

$A = [\max\{x_i - b_j, 0\}]_{ij}$ has full rank. Its smallest eigenvalue is

$\min_i \{x_i - b_i\}$

For weight vectors $w, b \in R^n$ and $a \in R^d$, consider the function $c : R^n \rightarrow R$,

$$c(x) = \sum_{j=1} w_j \max\{a^T x - b_j, 0\}$$

For weight vectors $w, b \in R^n$ and $a \in R^d$, consider the function $c : R^n \rightarrow R$,

$$c(x) = \sum_{j=1} w_j \max\{a^T x - b_j, 0\}$$

- This can be done trivially with a depth 2 neural network with relu.

For weight vectors $w, b \in R^n$ and $a \in R^d$, consider the function $c : R^n \rightarrow R$,

$$c(x) = \sum_{j=1} w_j \max\{a^T x - b_j, 0\}$$

- Now, fixing a sample $S = z_1, \dots, z_n$ of size n and a target vector $y \in R_n$. We need to find weights a, b, w so that $y_i = c(z_i)$ for all $i \in \{1, \dots, n\}$

For weight vectors $w, b \in R^n$ and $a \in R^d$, consider the function $c : R^n \rightarrow R$,

$$c(x) = \sum_{j=1} w_j \max\{a^T x - b_j, 0\}$$

- First, choose a and b such that with $\xi_i = a^T z_i$ we have the interleaving property $b_1 < x_1 < b_2 < \dots < b_n < x_n$. Next, consider the set of n equations in the n unknowns w ,

$$y_i = c(z_i), i \in \{1, \dots, n\}$$

We have $c(z_i) = Aw$, where $A = [\max\{\xi_i - b_j, 0\}]_{ij}$ is the matrix of Lemma 1.

For weight vectors $w, b \in R^n$ and $a \in R^d$, consider the function $c : R^n \rightarrow R$,

$$c(x) = \sum_{j=1} w_j \max\{a^T x - b_j, 0\}$$

- Now, fixing a sample $S = z_1, \dots, z_n$ of size n and a target vector $y \in R_n$. We need to find weights a, b, w so that $y_i = c(z_i)$ for all $i \in \{1, \dots, n\}$
- We chose a and b so that the lemma applies and hence A has full rank. We can now solve the linear system $y = Aw$ to find suitable weights w .

Discussion

What We Were Talking About...




To recap:

1. We just showed that generically large neural networks can express any labelling of the training data
2. And so it is not surprising to see the networks learn the training data perfectly
3. but it is surprising that we can't explain this well!

Some Thoughts...

1. Our favourite papers are the ones that shed light on truths that are taken for granted.
2. Its obvious that randomizing the labels would eliminate generalizability, but finding precise mathematical statements about this is not!
3. Models used in practice have the capability of memorizing the training data. Is it somehow easier not to?
4. The interplay between generalization and ease of optimization seems like an interesting thing to explore...

Any Questions?

-  D. Sowinski, “What is generalization in machine learning?.” Post.
-  M. Hardt, B. Recht, and Y. Singer, “Train faster, generalize better: Stability of stochastic gradient descent,” *CoRR*, vol. abs/1509.01240, 2015.
-  C. Zhang, S. Bengio, M. Hardt, B. Recht, and O. Vinyals, “Understanding deep learning requires rethinking generalization,” *CoRR*, vol. abs/1611.03530, 2016.