# Understanding interpolation in machine learning

Stéphane Chrétien

University of Lyon 2 ERIC Laboratory

August 3, 2020

# Introduction

Modern machine learning and statistics deal with the problem of learning from data:

given a training dataset  $(y_i, x_i)$   $i \in I$  where

$$x_i \in \mathbb{R}^d$$
 is the input  $y_i \in \mathbb{R}$  is the output,

one seeks a function  $f: \mathbb{R}^d \mapsto \mathbb{R}$  from a certain function class  $\mathcal{F}$  that has good prediction performance on test data  $(y_t, x_t)$ ,  $t \in \mathcal{T}$ , i.e. which has small testing error

$$\sum_{t \in T} \ell(y_t, f(x_t)) \tag{1}$$

This problem is of fundamental significance and finds applications in numerous scenarios.

For instance, in image recognition,

the input *x* corresponds to the raw image the output *y* is the image category

and the goal is to find a mapping f that can classify new images with acceptable accuracy.

Decades of research efforts in statistical machine learning have been devoted to developing methods to find f efficiently with provable guarantees.

- Prominent examples include
  - linear classifiers (e.g., linear / logistic regression, linear discriminant analysis),
  - kernel methods (e.g., support vector machines),
  - tree-based methods (e.g., decision trees, random forests),
  - nonparametric regression (e.g., nearest neighbors, local kernel smoothing), etc.
- Roughly speaking, each aforementioned method corresponds to a different function class  $\mathcal{F}$  from which the final classifier f is chosen.

- Deep learning, in its simplest form, consists in looking for functions of the form

$$\mathcal{F} = \left\{ f(x,\theta) = W_L(\sigma_L(W_{L-1}(\sigma_{L-1}(\cdots \sigma_2(W_1(x)))))) \right\}.$$

where  $\sigma_l$  is a non-linear function which applies componentwise and  $W_l$  is an affine operator,  $l=1,\ldots,L$ .

- Evolution of the performances over the last 7 years . . .

Model	Year	# Layers	# Params	Top-5 error
Shallow	< 2012		_	> 25%
AlexNet	2012	8	61M	16.4%
VGG19	2014	19	144M	7.3%
GoogleNet	2014	22	7M	6.7%
ResNet-152	2015	152	60M	3.6%

- It is widely acknowledged that two indispensable factors contribute to the success of deep learning, namely
  - huge datasets that often contain millions of samples and
  - immense computing power resulting from clusters of graphics processing units (GPUs).
- Admittedly, these resources are only recently available.

- However, these two alone are not sufficient to explain the mystery of deep learning:
  - over-parametrization: the number of parameters in state-of-the-art models is very often much larger than the sample size,
    - ← which might make them prone to overfitting,



- nonconvexity does not seem to be a problem: even with the help of GPUs, training deep learning models is still NP-hard in the worst case due to the highly nonconvex loss function to minimize.

Nevertheless, standard incremental algorithms (Stochastic Gradient Descent, etc) often reach good minimisers of the Empirical Risk

- A lot remains to be understood! ...

- Deep learning is able to approximate complicated nonlinear maps through composing many simple nonlinear functions.
- The motivation for the multilayer architecture is that there are different levels of features and the layers might be able to properly account for these different levels independently.
- Here, we sample and visualize weights from a pre-trained AlexNet model.



This can be used to generate new images using for instance, Generative Adversarial Networks



# Why overparametrise?

- It is often observed that depth helps efficiently extract features from the dataset, whereas
- (recent studies found that increasing both depth and width in a shallow model leads to very nice continuous limits, where PDE tools can be put to work...)

# What consequences does overparametrisation have on learning?

- In deep neural networks, over-parametrization usually entails existence of many local minimisers with potentially different statistical performance.
- Common practice advises to runs stochastic gradient descent with random initialization and converges to parameters with very good practical prediction accuracy.

Why is this simple approach actually often working?

The goal of current research is to resolve these paradoxes!

- expressivity
- generalisation bounds (PAC + compression)
- optimisation algorithms
- interpolation

 ${\sf Expressivity}$ 

- Recent works have been devoted to the approximation accuracy of deep neural networks for various measures of the error (expressivity)
- Some notable works include
  - \* the approximation results of Yarotsky
  - \* the nonlinear approximation analysis of RELU networks by Daubechies, Devore Foucart, Hanin and Petrova
  - \* the approximation of analytic maps by Weinan E and Wang.
  - \* approximation of functions in Sobolev spaces by Guhring, Kutyniok and Petersen
  - \* the definition of new approximation spaces (Barron spaces) by Weinan E. et al. (which play the same role as Besov spaces for nonlinear approximation with wavelet bases)
  - \* etc



We will use the following theorem from Guhring, Kutyniok and Petersen.

## Theorem

Let  $k \in \mathbb{N}_{\geq 2}$ , 1 , <math>B > 0, and 0 < s < 1. Then, there exists a constant c = c(d, p, k, B, s) > 0 with the following properties: for any  $\varepsilon \in (0,1/2)$ , for any  $f:(0,1)^d \mapsto \mathbb{R}$  in the ball of radius B in W<sup>k,p</sup>, there exists a vector of weights W and an associated neural network fw such that

$$||f_W - f||_{W^{s,p}((0,1)^d)} \le \varepsilon$$

### Theorem

and

(i) the number L of layers is bounded by

$$L \le c \log_2 \left( \varepsilon^{-k/(k-s)} \right)$$

(ii) the number  $d + \sum_{l=1}^{L} N_l$  of neurons is bounded by

$$d + \sum_{l=1}^{L} N_l \le c \ \varepsilon^{-d/(k-s)} \cdot \log_2 \left( \varepsilon^{-k/(k-s)} \right).$$

Generalisation and compression

Introduction

# Really interesting work by Bartlett, Barron and others ...

## Theorem (B., Foster, Telgarsky, 2017)

With high probability over n training examples

$$(X_1,Y_1),\ldots,(X_n,Y_n)\in\mathcal{X} imes\{\pm 1\}$$
, every  $f_W$  with  $R_W\leq r$  has

$$\Pr(\operatorname{sign}(f(X)) \neq Y) \leq \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}[Y_i f(X_i) \leq \gamma] + \tilde{O}\left(\frac{rL}{\gamma \sqrt{n}}\right).$$

Here,  $f_W$  is computed in a network with L layers and parameters  $W_1, \ldots, W_\ell$ :

$$f_W(x) := \sigma_L(W_L \sigma_{L-1}(W_{L-1} \cdots \sigma_1(W_1 x) \cdots)),$$

where the  $\sigma_i$  are 1-Lipschitz, and we measure the scale of  $f_W$  using a product of norms of the matrices  $W_i$ ,

for example, 
$$r := \prod_{i=1}^{L} \|W_i\|_* \left(\sum_{i=1}^{L} \frac{\|W_i\|_{2,i}^{2/3}}{\|W_i\|_*^{2/3}}\right)^{3/2}$$
.

# New trends involve compression showing that intrinsic dimension of the deep networks is not as large as we think!

Stronger generalization bounds for deep nets via a compression approach

Sanjeev Arora\* Rong Ge<sup>†</sup> Behnam Neyshabur<sup>‡</sup> Yi Zhang<sup>§</sup>

#### Abstract

Deep nets generalize well despite having more parameters than the number of training samples. Recent works try to give an explanation using PAC-Bayes and Margin-based analyses, but do not as yet result in sample complexity bounds better than naive parameter counting. The current paper shows generalization bounds that re orders of magnitude better in practice. These rely upon new succinct reparametrizations of the trained net — a compression that is explicit and efficient. These yield generalization bounds via a simple compression-based framework introduced here. Our results also provide some theoretical justification for widespread empirical success in compressing deep nets.

Analysis of correctness of our compression relies upon some newly identified "noise stability" properties of trained deep nets, which are also experimentally verified. The study of these properties and resulting generalization bounds are also extended to convolutional nets, which had eluded earlier attempts on proving generalization.

# even the theory of coresets was used in order to compress!

Published as a conference paper at ICLR 2019

# DATA-DEPENDENT CORESETS FOR COMPRESSING NEURAL NETWORKS WITH APPLICATIONS TO GENERALIZATION BOUNDS

Cenk Baykal†, Lucas Liebenwein†, Igor Gilitschenski†, Dan Feldman‡, Daniela Rus†

#### ABSTRACT

We present an efficient coresets-based neural network compression algorithm that sparsifies the parameters of a trained fully-connected neural network in a manner that provably approximates the network's output. Our approach is based on an importance sampling scheme that judiciously defines a sampling distribution over the neural network parameters, and as a result, retains parameters of high importance while discarding redundant ones. We leverage a novel, empirical notion of sensitivity and extend traditional coreset constructions to the application of compressing parameters. Our theoretical analysis establishes guarantees on the size and accuracy of the resulting compressed network and gives rise to generalization bounds that may provide new insights into the generalization properties of neural networks. We demonstrate the practical effectiveness of our algorithm on a variety of neural network configurations and real-world data sets.

#### 1 Introduction



One of the main issues using the traditional approaches was to obtain bounds which do not blow up as the number of layers increases

⇒ this would go against empirical findings

Noah Golowich, Alexander Rakhlin, Ohad Shamir recently solved the problem using Rademacher complexity in a careful way

### Size-Independent Sample Complexity of Neural Networks

Noah Golowich Harvard University Alexander Rakhlin MIT Ohad Shamir Weizmann Institute of Science and Microsoft Research

#### Abstract

We study the sample complexity of learning neural networks, by providing new bounds on their Rademacher complexity assuming norm constraints on the parameter matrix of each layer. Compared to previous work, these complexity bounds have improved dependence on the network depth, and under some additional assumptions, are fully independent of the network size (both depth and width). These results are derived using some novel techniques, which may be of independent interest.

#### 1 Introduction

One of the major challenges involving neural networks is explaining their ability to generalize well, even if they are very large and have the potential to overfit the training data [Nevshabur et al., 2014, Zhang et al.,



Where do gradient methods end up landing in this wild lanscape ?

# Overparameterized Nonlinear Learning: Gradient Descent Takes the Shortest Path?

Samet Oymak\* and Mahdi Soltanolkotabi<sup>†</sup> December 2018

#### Abstract

Many modern learning tasks involve fitting nonlinear models to data which are trained in an overparameterized regime where the parameters of the model exceed the size of the training dataset. Due to this overparameterization, the training loss may have infinitely many global minima and it is critical to understand the properties of the solutions found by first-order optimization schemes such as (stochastic) gradient descent starting from different initializations. In this paper we demonstrate that when the loss has certain properties over a minimally small neighborhood of the initial point, first order methods such as (stochastic) gradient descent have a few intriguing properties: (1) the iterates converge at a geometric rate to a global optima even when the loss is nonconvex, (2) among all global optima of the loss the iterates converge to one with a near minimal distance to the initial point, (3) the iterates take a near direct route from the initial point to this global optima. As part of our proof technique, we introduce a new potential function which captures the precise tradeoff between the loss function and the distance to the initial point as the iterations progress. For Stochastic Gradient Descent (SGD), we develop novel martingale techniques that guarantee SGD never leaves a small neighborhood of the initialization, even with rather large learning rates. We demonstrate the utility of our general theory for a variety of problem domains spanning low-rank matrix recovery to neural network training.

# Stochastic Gradient Langevin Dynamics (SGLD) avoids spurious local minimisers!

Proved in the work of Charikar et al. for the Langevin approximation

## A Hitting Time Analysis of Stochastic Gradient Langevin Dynamics

Yuchen Zhang\* Percy Liang† Moses Charikar‡

April 10, 2018

#### Abstract

We study the Stochastic Gradient Langevin Dynamics (SGLD) algorithm for non-convex optimization. The algorithm performs stochastic gradient descent, where in each step it injects appropriately scaled Gaussian noise to the update. We analyze the algorithm's hitting time to an arbitrary subset of the parameter space. Two results follow from our general theory: First, we prove that for empirical risk minimization, if the empirical risk is pointwise close to the (smooth) population risk, then the algorithm finds an approximate local minimum of the population risk in polynomial time, escaping suboptimal local minima that only exist in the empirical risk. Second, we show that SGLD improves on one of the best known learnability results for learning linear classifiers under the zero-one loss.

#### 1 Introduction

A central challenge of non-convex optimization is avoiding sub-optimal local minima. Although escaping all local minima is NP-hard in general [e.g. 7], one might expect that it should be possible to escape "appropriately shallow" local minima, whose basins of attraction have relatively low barriers. As an illustrative



L| 1 May 2019

#### Wide Neural Networks of Any Depth Evolve as Linear Models Under Gradient Descent

Jachoon Lee<sup>\*12</sup> Lechao Xiao<sup>\*12</sup> Samuel S. Schoenholz<sup>1</sup> Yasaman Bahri<sup>1</sup> Jascha Sohl-Dickstein<sup>1</sup> Jeffrey Pennington<sup>1</sup>

#### Abstract

A longstanding goal in deep learning research has been to precisely characterize training and generalization. However, the often complex loss landscapes of neural networks have made a theory of learning dynamics clusive. In this work, we show that for wide neural networks the learning dynamics simplify considerably and that, in the infinite width limit, they are governed by a linear model obtained from the first-order Taylor expanded

systems can often shed light on these hard problems. For neural networks, one such limit is that of infinite width, which refers either to the number of hidden units in a fully-connected layer or to the number of channels in a convolutional layer. Under this limit, the output of the network at initialization is a draw from a Gaussian process (GP); moreover, the network output remains governed by a GP after exact Bayesian training using squared loss (Neal, 1994; Lee et al., 2018; Matthews et al., 2018; Novak et al., 2019; Garriga-Alonso et al., 2019a). Aside from its theoretical simplicity, the infinite-width limit is also of practical inter-

# But recent work of Chizat, Oyallon and Bach showed that deep neural networks then reach in a tangent kernel regime . . .

# On Lazy Training in Differentiable Programming

Lénaïc Chizat CNRS, Université Paris-Sud Orsay, France lenaic.chizat@u-psud.fr Edouard Oyallon CentraleSupelec, INRIA Gif-sur-Yvette, France edouard.oyallon@centralesupelec.fr

Francis Bach
INRIA, ENS, PSL Research University
Paris, France
francis.bach@inria.fr

#### Abstract

In a series of recent theoretical works, it was shown that strongly overparameterized neural networks trained with gradient-based methods could converge exponentially fast to zero training loss, with their parameters hardly varying. In this work, we show that this "lazy training" phenomenon is not specific to overparameterized neural networks, and is due to a choice of scaling, often implicit, that makes the model behave as its linearization around the initialization, thus yielding a model equivalent to learning with positive-definite kernels. Through a theoretical analysis, we exhibit various situations where this phenomenon arises in non-convex optimization and we provide bounds on the distance between the lazy and linearized optimization paths. Our numerical experiments bring a critical note, as we observe that the performance of commonly used non-linear deep con-



# The importance of being flat!

# **Understanding Generalization through Visualizations**

W. Ronny Huang

University of Maryland wrhuang@umd.edu

Micah Goldblum

University of Maryland goldblum@math.umd.edu

Liam Fowl

University of Maryland lfowl@math.umd.edu

Justin K. Terry
University of Maryland
justinkterry@gmail.com

**Furong Huang** 

University of Maryland furongh@cs.umd.edu Tom Goldstein

Zevad Emam

University of Maryland

zevad@math.umd.edu

University of Maryland tomg@cs.umd.edu

#### Abstract

The power of neural networks lies in their ability to generalize to unseen data, yet the underlying reasons for this phenomenon remain clusive. Numerous rigorous attempts have been made to explain generalization, but available bounds are still quite loose, and analysis does not always lead to true understanding. The goal of

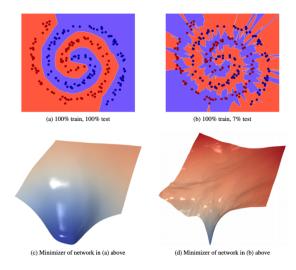


Figure: Flat minimisers are easier to reach and have better generalisation properties (empirical)

# Yet another fascinating direction . . .

Recent work of Lavaei and co-authors show that incremental methods can avoid spurious minimisers

# Absence of Spurious Local Trajectories in Time-Varying Optimization: A Control-Theoretic Perspective

Salar Fattahi, Cedric Josz, Reza Mohammadi, Javad Lavaei, and Somayeh Sojoudi

Abstract-In this paper, we study the landscape of an optimization problem whose input data vary over time. This time-varying problem consists of infinitely-many individual optimization problems, whose solution is a trajectory over time rather than a single point. To understand when it is possible to find a global solution of a time-varying non-convex optimization problem, we introduce the notion of spurious (i.e., non-global) local trajectory as a generalization to the notion of spurious local solution in nonconvex (time-invariant) optimization. We develop an ordinary differential equation (ODE) which, at limit, characterizes the spurious local solutions of the time-varying optimization problem. By building upon this connection, we prove that the absence of spurious local trajectory is closely related to the transient behavior of the proposed ODE. In particular, we show that: (1) if the problem is time-invariant, the spurious local trajectories are ubiquitous since any strict local minimum is a locally stable equilibrium point of the ODE, and (2) if the ODE is time-varying. the data variation may force all ODE trajectories initialized at arbitrary local minima at the initial time to gradually converge to the global solution trajectory. This implies that the natural data variation in the problem may automatically trigger escaping local

minima over time.

optimization. This observation naturally gives rise to the following question:

Would fast local-search algorithms escape spurious local minima in online nonconvex optimization, similar to their timeinvariant counterparts?

In this paper, we attempt to address this question by developing a control-theoretic framework for analyzing the landscape of online and time-varying optimization. In particular, we demonstrate that even if a time-varying optimization may have undesired point-wise local minima at almost all times, the variation of its landscape over time would enable simple local-search algorithms to escape these spurious local minima. Inspired by this observation, this paper provides a new machinery to analyze the global landscape of online decision-making problems by drawing tools from optimization and control theory.

We consider a class of nonconvex and online optimization problems where the input data varies over time. First, we introduce the notion of spurious local trajectory as a generalization



Interpolation

- Overparametrisation works despite contradicting the intuition that "overfitting makes no sense"



# Reconciling modern machine-learning practice and the classical bias-variance trade-off

Mikhail Belkina,b,1, Daniel Hsuc, Siyuan Maa, and Soumik Mandala

\*Department of Computer Science and Engineering, The Ohio State University, Columbus, OH 43210; \*Department of Statistics, The Ohio State University, Columbus, OH 43210; \*Department of Statistics, The Ohio State University, Columbus, OH 43210; \*Department of Statistics, The Ohio State University, Columbus, OH 43210; \*Department of Statistics, The Ohio State University, Columbus, OH 43210; \*Department of Statistics, The Ohio State University, Columbus, OH 43210; \*Department of Statistics, The Ohio State University, Columbus, OH 43210; \*Department of Statistics, The Ohio State University, Columbus, OH 43210; \*Department of Statistics, The Ohio State University, Columbus, OH 43210; \*Department of Statistics, The Ohio State University, Columbus, OH 43210; \*Department of Statistics, The Ohio State University, Columbus, OH 43210; \*Department of Statistics, The Ohio State University, Columbus, OH 43210; \*Department of Statistics, The Ohio State University, Columbus, OH 43210; \*Department of Statistics, The Ohio State University, Columbus, OH 43210; \*Department of Statistics, The Ohio State University, OHIO STATE OF The Ohio State Uni

Edited by Peter J. Bickel, University of California, Berkeley, CA, and approved July 2, 2019 (received for review February 21, 2019)

Breakthroughs in machine learning are rapidly changing science and society, yet our fundamental understanding of this technology has lagged far behind. Indeed, one of the central tenets of the field, the bias-variance trade-off, appears to be at odds with the observed behavior of methods used in modern machine-learning practice. The bias-variance trade-off implies that a model should balance underfitting and overfitting: Rich enough to express underlying structure in data and simple enough to avoid fitting spurious patterns. However, in modern practice, very rich models such as neural networks are trained to exactly fit (i.e., interpolate) the data. Classically, such models would be considered overfitted, and vet they often obtain high accuracy on test data. This apparent contradiction has raised questions about the mathematical foundations of machine learning and their relevance to practitioners. In this paper, we reconcile the classical understanding and the modern practice within a unified performance curve. This "double-descent" curve subsumes the textbook U-shaped bias-variance trade-off curve by showing how increasing model capacity beyond the point of interpolation results in improved performance. We provide evidence for the existence and ubiquity of double descent for a wide spectrum of models and datasets, and we posit a mechanism for its emergence. This connection between the performance and the structure of machine-learning models delineates the limits of classical analyses and has implications for both the theory and the practice of

machine learning | bias-variance trade-off | neural networks

ing data (i.e., have large empirical risk) and hence predict poorly on new data. 2) If H is too large, the empirical risk minimizer may overfit spurious patterns in the training data, resulting noor accuracy on new examples (small empirical risk but large true risk).

The classical thinking is concerned with finding the "sweet" between underfitting and overfitting. The control of the function class capacity may be explicit, via the choice of  $\mathcal{H}$  (e.g. see that the control of the function class capacity may be explicit, via the choice of  $\mathcal{H}$  (e.g. see that the control of the co

However, practitioners coulinely use modern machine learning methods, such a large neural networks and other nonlinear predictors that have very low or zero training risk. Despite the high function class capacity and near-perfect fit to training data, these predictors often give very accurate predictions on new data. Indeed, this behavior has guided a best practice in deep learning for choosing neural network architectures, specifically that the network should be large enough to permit effortless that the relevoit should be large enough to permit effortless that the network should be large sought to permit effortless that the network should be large sought to permit effortless that the network should be large to the bias-avainance tradeoff rhislooder, revert empirical evidence indicates that releval



machine learning.

- Belkin et al. introduced the "double descent curve"



# Montanari et al. resolved this paradox ... for the linear model! (Uses a lot of random matrix theory in the asymptotic regime)

## Surprises in High-Dimensional Ridgeless Least Squares Interpolation

Trevor Hastie Andrea Montanari Saharon Rosset Ryan J. Tibshirani

#### Abstract

Interpolators—estimators that achieve zero training error—have attracted growing altention in machine learning, mainly because state-of-the art neural networks appear to be models of this type. In this paper, we study minimum  $\ell_2$  norm ("ridgeless") interpolation in high-dimensional least squares regression. We consider two different models for the feature distribution: a linear model, where the feature vectors  $x_i \in \mathbb{R}^p$  are obtained by applying a linear transform to a vector of i.i.d. entries,  $x_i = \Sigma^{1/2} z_i$  (with  $z_i \in \mathbb{R}^p$ ); and a nonlinear model, where the feature vectors are obtained by passing the input through a random one-layer neural network,  $z_i = \varphi(Wz_i)$  (with  $z_i \in \mathbb{R}^d$ ,  $W \in \mathbb{R}^{p \times d}$  a matrix of i.i.d. entries, and  $\varphi$  an activation function acting componentwise on  $Wz_i$ ). We recover—in a precise quantitative way—several phenomena that have been observed in large-scale neural networks and kernel machines, including the "double descent" behavior of the prediction risk, and the potential benefits of overparametrization.

#### 1 Introduction

Introduction

Modem deep learning models involve a huge number of parameters. In nearly all applications of these models, current practice suggests that we should design the network to be sufficiently complex so that the model (as training, by pigally, by gradient descent) interpolates the data, i.e., achieves zero training error. Indeed, in a thought-provoking experiment, Zhang et al. (2016) showed that state-of-the-art deep neural network architectures can be trained to interpolate the data even when the actual labels are replaced by entirely random ones.

Despite their enormous complexity, deep neural networks are frequently seen to generalize well, in meaningful practical problems. At first sight, this seems to defy conventional statistical wisdom; interpolation (vanishing training

# Other models in non-parametric regression have been addressed by Belkin, Rakhlin and Tsybakov

Does data interpolation contradict statistical optimality?

Mikhail Belkin The Ohio State University Alexander Rakhlin MIT Alexandre B. Tsybakov CREST, ENSAE

#### Abstract

We show that learning methods interpolating the training data can achieve optimal rates for the problems of nonparametric regression and prediction with square loss.

#### 1 Introduction

In this paper, we exhibit estimators that interpolate the data, yet achieve optimal rates of convergence for the problems of nonparametric regression and prediction with square loss. This curious observation goes against the usual (or, folklore?) intuition that a good statistical procedure should forego the exact fit to data in favor of a more smooth representation. The family of estimators we consider do exhibit a bias-variance trade-off with a tuning parameter, yet this "regularization" co-exists in harmony with data interpolation.

Motivation for this work is the recent focus within the machine learning community on the out-of-sample performance of neural networks. These flexible models are typically trained to fit the data exactly (either in their sign or in the actual value), yet they predict well on unseen data. The conundrum has served both as a source of excitement about the "magical" properties of neural networks, as well as a call for the development of novel statistical techniques to resolve it.

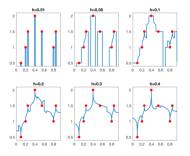
The aim of this short note is to show that not only can interpolation be a good statistical procedure, but it can even be optimal in a minimax sense. To the best of our knowledge, such outimality has not been exhibited before.

Let (X,Y) be a random pair on  $\mathbb{R}^d \times \mathbb{R}$  with distribution  $P_{XY}$ , and let  $f(x) = \mathbb{E}[Y|X=x]$ 



The Nadaraya-Watson estimator for a singular kernel K is defined as

$$f_n(x) = \left\{ \begin{array}{ll} Y_i & \text{if } x = X_i \text{ for some } i = 1, \dots, n, \\ 0 & \text{if } \sum_{i=1}^n K\left(\frac{x - X_i}{h}\right) = 0, \\ \frac{\sum_{i=1}^n Y_i K\left(\frac{x - X_i}{h}\right)}{\sum_{i=1}^n K\left(\frac{x - X_i}{h}\right)} & \text{otherwise.} \end{array} \right.$$



Interpolation with  $K(u) = ||u||^{-a} \mathbf{I}\{||u|| \le 1\}$ , a = 0.49, and various values of h.

## Figure: Singular Kernel estimators that interpolate!

 $\label{eq:Asymptotic} A \ simple \ analysis \ of \ interpolation$ 

## Mathematical Model

Let  $Z_i = (X_i, Y_i)$  in  $\mathbb{R}^{d+1} \times \mathbb{R}$ ,  $i = 1, \ldots, n$  be observations drawn from the following model

$$Y_i = f^*(X_i) + \varepsilon_i \tag{2}$$

 $i = 1, \ldots, n$ , where we assume that

the vectors  $X_i$  are random and i.i.d., taking values in  $\mathbb{R}^d$ and the noise vector  $\varepsilon = [\varepsilon_1, \dots, \varepsilon_n]^t$  is sub-Gaussian, with sub-Gaussian constant denoted by  $K_{\epsilon}$ .

The goal is to estimate  $f^*$  based on the observation  $Z_1, \ldots, Z_n$ .

The estimation of f\* will be based on restricting the search to a subset  $\mathcal{F}$  of functions of a Banach space  $\mathcal{B}$ . 4□ > 4□ > 4□ > 4□ > □
900

In order to generalise, the estimator should be chosen in the set of stationary points of the empirical version of the risk  $R: \mathcal{F} \to \mathbb{R}$ defined by

$$R(f) = \mathbb{E}\left[\ell(Y, f(X))\right],$$

where  $\ell: \mathbb{R} \times \mathbb{R} \to \mathbb{R}$  satisfies

Recent trends

$$\ell(y,y)=0$$
 for all  $y\in\mathbb{R}$  and

 $\ell(y,\cdot)$ :  $\mathbb{R} \mapsto \mathbb{R}$  is a strictly convex twice continuously differentiable nonnegative function.

Let  $\hat{R}_n(f)$  denote the empirical risk defined by

$$\hat{R}_n(f) = \frac{1}{n} \sum_{i=1}^n \ell(Y_i, f(X_i)).$$
 (3)

Then, the Empirical Risk Minimizer  $\hat{f}^{ERM}$  will be a solution to

$$\hat{f}^{ERM} \in \operatorname{argmin}_{f \in \mathcal{F}} \hat{R}_n(f).$$
 (4)

## Assumption

The sample satisfies the following separation

$$\min_{i,i'=1}^{n} \|X_i - X_{i'}\|_2 \ge cn^{-1/\nu}$$
 (5)

with probability larger than or equal to  $1-\delta$ , for some positive constants c,  $\nu$  and for  $\delta \in (0,1)$ .

The Holder exponent  $\nu$  is usually interpreted as a surrogate for the intrinsic dimension of the data manifold. E.g., this intrinsic dimension was estimated to be less than 20 for the MNIST dataset

Intrinsic Dimensionality Estimation of Submanifolds in  $\mathbb{R}^d$ 

Matthias Hein

Max Planck Institute for Biological Cybernetics, Tübingen, Germany

Jean-Yves Audibert

AUDIBERT@CERTIS.ENPC.FR

MH@TUEBINGEN.MPG.DE



An handy result from Neuberger

# Theorem (Neuberger's theorem)

Recent trends

Suppose that  $\mathcal{B}$ ,  $\mathcal{J}$ , and  $\mathcal{K}$  are three Banach spaces and that  $\mathcal{B}$  is compactly embedded in  $\mathcal{J}$ .

Suppose that  $F: \mathcal{B} \to \mathcal{K}$  is a function that is continuous with respect to the topologies of  $\mathcal{J}$  and  $\mathcal{K}$ . For r > 0 and u in  $\mathcal{B}$ ,  $B_r(u)$ and  $\bar{B}_r(u)$  will denote the open and closed balls in  $\mathcal{B}_r$ , respectively, with center u and radius r. Suppose that  $f \in \mathcal{B}$ , that r > 0, and that for each g in  $B_r(f)$  there is an h in  $\bar{B}_r(0)$  such that

$$\lim_{t \to 0^+} \frac{1}{t} (F(g + th) - F(g)) = -F(f).$$

Then there is  $\hat{f}$  in  $\bar{B}_r(f)$  such that  $F(\hat{f}) = 0$ .

Recent trends

We recall that  $f \in \mathcal{F}$ , and  $d' \in \mathcal{B}$  such that  $\mathcal{F} \subset \mathcal{B}$ . Let us compute the directional derivative of  $\hat{R}_n$ 

$$D\hat{R}_{n}(f) \cdot h' = \lim_{t \to 0} \frac{\hat{R}_{n}(f + th') - \hat{R}_{n}(f)}{t}$$

$$= \lim_{t \to 0} \frac{\frac{1}{n} \sum_{i=1}^{n} \ell(Y_{i}, f(X_{i}) + th'(X_{i})) - \ell(Y_{i}, f(X_{i}))}{t}$$

$$= \lim_{t \to 0} \frac{\frac{1}{n} \sum_{i=1}^{n} \partial_{2} \ell(Y_{i}, f(X_{i})) th'(X_{i}) + c \partial_{2}^{2} \ell(Y_{i}, f(X_{i}) t^{2}h'^{2}(X_{i})}{t}$$

with  $c \in [0,1]$ , and thus

$$D\hat{R}_n(f)\cdot h'=\frac{1}{n}\sum_{i=1}^n \partial_2\ell(Y_i,f(X_i))h'(X_i).$$

In the same spirit, we get

$$D^{2}\hat{R}_{n}(f)\cdot(h',h)=\frac{1}{n}\sum_{i=1}^{n}\partial_{2}^{2}\ell(Y_{i},f(X_{i}))h'(X_{i})h(X_{i}).$$

Based on these computations, Neuberger's theorem resorts to obtaining a bound on the norm of an appropriate solution d' to the following linear system

$$\frac{1}{n}\sum_{i=1}^{n} \partial_{2}^{2} \ell(Y_{i}, f(X_{i})) h'(X_{i})h(X_{i}) = -\frac{1}{n}\sum_{i=1}^{n} \partial_{2}\ell(Y_{i}, f^{*}(X_{i}))h'(X_{i})$$

for all  $f \in B_r(f^*)$  and for all  $h' \in \mathcal{B}$ .

Let  $\psi$  denote the bump function

$$\psi(x) = \begin{cases} \exp\left(1 - \frac{1}{1 - \|x\|_2^2}\right) & \text{if } \|x\|_2^2 \le 1, \\ 0 & \text{otherwise} \end{cases}$$
 (6)

We will decouple the problem and first solve it in a Sobolev space, and then approximate the solution by a deep neural network using the Guhring, Kutyniok and Petersen theorem . . .

### Theorem

Set  $\ell$  to be the  $\ell_2^2$  loss, i.e.  $\ell(y,z) = \frac{1}{2}(y-z)^2$  for all y, z in  $\mathbb{R}$ . Let Assumption 1 hold. and let  $\psi_{\sigma} = \psi(\cdot/\sigma)$ . Take any  $\sigma \leq c n^{-1/\nu}$  such that the ball in  $\mathcal{B}$  centered at  $f^*$  with radius  $6K_{\epsilon} \ n \|\psi_{\sigma}\|_{\mathcal{B}} \subset \mathcal{F}$ . Then, there exists with high probability a mapping  $\hat{f}^{ERM}$ :  $\mathbb{R}^d \mapsto \mathbb{R}$  which is a stationary point of the empirical risk minimisation problem and which lies at a distance at most

$$6K_{\epsilon} n \|\psi_{\sigma}\|_{\mathcal{B}}$$

from the neural network f\*.

Recent trends

Our main result is the following

Recent trends

#### $\mathsf{Theorem}$

Set  $\ell$  to be the  $\ell_2^2$  loss, i.e.  $\ell(y,z) = \frac{1}{2}(y-z)^2$  for all y, z in  $\mathbb{R}$ . Let Assumption 1 hold. Assume that  $||f^*||_{\mathcal{W}^{k,p}} < B$  for some  $k \in \mathbb{N}$  and  $p \in [1, +\infty]$ . Assume that d, p,  $\nu$  and n are such that  $6K_{\varepsilon}n^{1/2-d/(\nu p)}\|\psi\|_{\mathcal{M}^{k,p}} \leq B$ . Then, for any  $s \in [0,1]$ , there exists with high probability a deep neural network  $f_{\hat{\mathcal{M}}} : \mathbb{R}^d \mapsto \mathbb{R}$  with

$$\|f_{\hat{W}} - f^*\|_{L^p(\mathcal{D})} \le CK_{\epsilon} n^{1-d/(\nu p)} \|\psi\|_{\mathcal{W}^{k,p}(\mathcal{D})}$$

for some positive constant C = C(d, p, k, B, s)

#### Theorem

and with

(i) a number L of layers upper bounded by

$$L \leq c \log_2 \left( \left( CK_{\epsilon} n^{1-d/(\nu p)} \|\psi\|_{\mathcal{W}^{k,p}(\mathcal{D})} \right)^{-k/(k-s)} \right)$$

(ii) a number  $d + \sum_{l=1}^{L} N_l$  of neurons upper bounded by

$$d + \sum_{l=1}^{L} N_{l} \leq c \left( CK_{\epsilon} n^{1-d/(\nu p)} \|\psi\|_{\mathcal{W}^{k,p}(\mathcal{D})} \right)^{-d/(k-s)}$$
$$\cdot \log_{2} \left( \left( CK_{\epsilon} n^{1-d/(\nu p)} \|\psi\|_{\mathcal{W}^{k,p}(\mathcal{D})} \right)^{-k/(k-s)} \right)$$

which approximately solves the empirical risk minimisation problem (4) over the Sobolev class  $\mathcal{W}^{k,p}(\mathcal{D})$  with  $\mathcal{W}^{s,p}(\mathcal{D})$ -distance at most  $K_{\epsilon}n^{1-d(\nu p)}\|\psi\|_{\mathcal{W}^{k,p}(\mathcal{D})}$  to the solution set.

 ${\sf Sketch} \,\, {\sf of} \,\, {\sf the} \,\, {\sf proof} \,\,$ 

Notice that for all  $f \in B_s(f_{W^*})$ , we have

$$\frac{1}{n}\sum_{i=1}^{n}\frac{\partial \ell}{\partial 2}(Y_{i},f(X_{i}))\ h'(X_{i})=-\frac{1}{n}\sum_{i=1}^{n}\left(Y_{i}-f(X_{i})\right)\ h'(X_{i}),$$

and that

Introduction

$$\frac{1}{n}\sum_{i=1}^{n}\frac{\partial^{2}\ell}{\partial_{2}^{2}}(Y_{i},f(X_{i}))\ h'(X_{i})h(X_{i})=\frac{1}{n}\sum_{i=1}^{n}h'(X_{i})h(X_{i}).$$

Then, using the fact that  $\ell$  is the  $\ell_2^2$  loss, Neuberger's condition reads

$$\frac{1}{n}\sum_{i=1}^n h'(X_i)h(X_i) = \frac{1}{n}\sum_{i=1}^n h'(X_i)(Y_i - f_{W^*}(X_i)).$$

One possible solution can be obtained by setting

$$h(X_i) = Y_i - f_{W^*}(X_i) = \varepsilon_i$$

 $i=1,\ldots,n$ .

One simple option is to take

$$h(x) = \sum_{i=1}^{n} \epsilon_{i} \psi \left( \frac{x - X_{i}}{\sigma} \right)$$

where  $\psi: \mathbb{R}^p \to \mathbb{R}$  is a kernel function and  $\sigma > 0$  is a bandwidth.

Let

$$\psi_{\sigma} = \psi\left(\cdot/\sigma\right).$$

Now, observe that, based on Assumption 1, the functions  $\psi((x-X_i)/\sigma)$ , and their successive derivatives up to k,  $i=1,\ldots,n$ , have disjoint supports for with probability larger than or equal to  $1-\delta$  as long as  $\sigma \leq c n^{-1/\nu}$ . We thus obtain that

$$||h||_{\mathcal{B}} = ||\epsilon||_1 ||\psi_{\sigma}||_{\mathcal{B}}$$

Moreover, as is well known for subGaussian vectors, the norm is controlled by

$$\|\epsilon\|_2 \leq 6K_{\epsilon}n.$$

with probability at least  $1 - \exp(-n)$ , combining the conclusion of Theorem 4 follows from Neuberger's Theorem 3.

The proof for the deep neural network case is completed by using the approximation result of Guhring, Kutyniok and Petersen. The number of layers may have to increase logarithmically with the number of samples

The total number of parameters blows up polynomially in the number of samples and exponentially in the dimension of the problem This simple exercice in using quantitative zero finding theorems such as Neubergers shows that we can easily prove results that do not blow up with the number of layers with interpolating networks

We can easily study local minimisers as well using the same technique

We would need to explore approximation theory in unusual/non standard directions:

improve the Guhring, Kutyniok and Petersen theorem by introducing the constraint that the network be a flat minimiser

This would explain that Stochastic Gradient methods can find the correct approximation with large probability (?)