Statistical theory for deep neural networks

Lecture series



Johannes Schmidt-Hieber

The problem

general belief that no or little theory can be developed for modern deep network architectures

- $\bullet\,$ complex data structures $\rightsquigarrow\,$ no available statistical models
- combination of intricate network architectures with various regularization methods
- fitting a network is a non-linear problem in the network parameters
- non-convex function class
- . . .

Why theory?

What is the use of theoretical results in a field that is (successfully) driven by trial and error?

- understand why deep learning works
- deep learning is a chaotic field (thousands of publications)

 mathematical theory can be useful to extract key concepts
- comparison with other methods
- selection of tuning parameters
- detecting limitations of deep learning
- improvements
- hybrid methods

organization of the course

Lectures:

- Theory for shallow networks
- Advantages of additional layers
- Statistical theory for deep ReLU networks
- Overparametrization

Perceptron





updates

<ロ > < 回 > < 回 > < 三 > < 三 > < 三 > 三 の Q (~ 5/101

Shallow networks



• shallow neural network with one output is a function $f: \mathbb{R}^d \to \mathbb{R}$ of the form

$$f(\mathbf{x}) = \sum_{j=1}^{m} c_j \sigma(\mathbf{w}_j^\top \mathbf{x} + v_j), \quad \mathbf{w}_j \in \mathbb{R}^d, \ v_j, c_j \in \mathbb{R}.$$

• activation function $\sigma:\mathbb{R}\to\mathbb{R}$

Feedforward neural networks

• for $\mathbf{v} = (v_1, \dots, v_r)^\top$, $\mathbf{y} = (y_1, \dots, y_r)^\top \in \mathbb{R}^r$, define the shifted activation function $\sigma_{\mathbf{v}} : \mathbb{R}^r \to \mathbb{R}^r$ as

$$\sigma_{\mathbf{v}} = (\sigma(y_1 - v_1), \ldots, \sigma(y_r - v_r))^{\top}.$$

- network architecture (L, p)
 - positive integer L called number of hidden layers/depth
 - width vector $\mathbf{p} = (p_0, \dots, p_{L+1}) \in \mathbb{N}^{L+2}$

Neural network with architecture (L, p) is

$$f(\mathbf{x}) = W_L \sigma_{\mathbf{v}_L} W_{L-1} \sigma_{\mathbf{v}_{L-1}} \cdots W_1 \sigma_{\mathbf{v}_1} W_0 \mathbf{x},$$

- W_i is a $p_i \times p_{i+1}$ weight matrix
- $\mathbf{v}_i \in \mathbb{R}^{p_i}$ is a shift vector

Feedforward neural networks

Neural network:

$$f(\mathbf{x}) = W_L \sigma_{\mathbf{v}_L} W_{L-1} \sigma_{\mathbf{v}_{L-1}} \cdots W_1 \sigma_{\mathbf{v}_1} W_0 \mathbf{x},$$

Comments:

- \bullet feedforward \rightsquigarrow information is passed in one direction through the network
- $\bullet\,$ network functions are build by alternating matrix-vector multiplications with the action of the non-linear activation function $\sigma\,$
- network architecture is given
- parameters generating the underlying function class are the matrices W₀,..., W_L and the shift vectors v₁,..., v_L

Graph representation



- in CS, neural networks are introduced via graph representation
- nodes in the graph (also called units) are arranged in layers
- input layer is the first layer and the output layer the last layer
- layers that lie in between are called hidden layers
- number of hidden layers corresponds to L and the number of units in each layer generates the width vector p
- Each node/unit in the graph representation stands for operation σ(a^t · +b)

Special types

Neural network:

$$f(\mathbf{x}) = W_L \sigma_{\mathbf{v}_L} W_{L-1} \sigma_{\mathbf{v}_{L-1}} \cdots W_1 \sigma_{\mathbf{v}_1} W_0 \mathbf{x},$$

Comments:

- network is called sparse if W_i are sparse matrices
- *i*-th layer is fully connected $\rightsquigarrow W_i$ is dense
- for L = 1 network coincides with shallow networks
- if L > 1, network is called deep

Depth



Source: Kaiming He, Deep Residual Networks

Networks are deep

- version of ResNet with 152 hidden layers
- networks become deeper

High-dimensionality



Source: arxiv.org/pdf/1605.07678.pdf

- Number of network parameters is larger than sample size
 - AlexNet uses 60 million parameters for 1.2 million training samples

deep learning

deep learning denotes gradient based methods to fit neural networks to data

Neural network:

$$f(\mathbf{x}) = W_L \sigma_{\mathbf{v}_L} W_{L-1} \sigma_{\mathbf{v}_{L-1}} \cdots W_1 \sigma_{\mathbf{v}_1} W_0 \mathbf{x},$$

- cross-entropy loss (= -log-likelihood)
- e.g. for regression problems, we observe n i.i.d. pairs $(\mathbf{X}_i, Y_i) \in \mathbb{R}^d \times \mathbb{R}$ and consider least squares loss

$$\sum_{i=1}^{n} \left(Y_i - f(\mathbf{X}_i) \right)^2$$

- (stochastic) gradient descent based method
- neural network class is non-convex and deep learning will in general not find the global minimum

initialization

$f(\mathbf{x}) = W_L \sigma_{\mathbf{v}_L} W_{L-1} \sigma_{\mathbf{v}_{L-1}} \cdots W_1 \sigma_{\mathbf{v}_1} W_0 \mathbf{x},$

- good initialization schemes of the network parameters is crucial for success of DL
- draw weight parameters independently such that the weight matrices are nearly orthogonal (in case they are square matrices)
- extra factors to compensate effect of the activation function

random initialization for wide shallow networks

A shallow network with m units and without biases can be written as

$$f_m(\mathbf{x}) = \sum_{i=1}^m c_i \sigma(w_i^\top \mathbf{x}), \quad i = 1, \dots, m$$

with parameters $c_i \in \mathbb{R}$ and $w_i \in \mathbb{R}^d$

random initialization: c_i i.i.d. and w_i i.i.d., such that $E[c_i] = 0$, $Var(c_i) = 1/m$

$$\operatorname{Cov}\left(f_{m}(\mathbf{x}), f_{m}(\mathbf{x}')\right) = \sum_{i=1}^{m} \operatorname{Cov}\left(c_{i}\sigma(w_{i}^{\top}\mathbf{x}), c_{i}\sigma(w_{i}^{\top}\mathbf{x}')\right)$$
$$= \sum_{i=1}^{m} E\left[c_{i}^{2}\sigma(w_{i}^{\top}\mathbf{x})\sigma(w_{i}^{\top}\mathbf{x}')\right]$$
$$= E\left[\sigma(w_{1}^{\top}\mathbf{x})\sigma(w_{1}^{\top}\mathbf{x}')\right]$$

15/101

random initialization for wide shallow networks (ctd.)

A wide, randomly initialized shallow neural network is approximately a Gaussian process with covariance

$$K(\mathbf{x}, \mathbf{x}') = E\left[\sigma(w^{\top}\mathbf{x})\sigma(w^{\top}\mathbf{x}')\right]$$

- expectation with respect to distribution of w
- covariance depends on activation function
- For ReLU with d = 1, the limiting Gaussian process on [0,∞) is

 $x \mapsto \xi x$

for $\xi \sim \mathcal{N}(0, E[(w)^2_+])$

random initialization for deep neural networks

- what is the limit for a wide deep neural network?
- depends on how the limit is taken
- if the number of units in all hidden layers is m and $m \to \infty$, then, proper initialization converges again to a Gaussian process (Theorem 4 in Matthews et al. 2018)
- convergence of finite dimensional distributions can be proved using CLT for exchangeable sequences
- covariance structure of limiting process can be expressed via recursion formula on the layers
- edge of chaos phenomenon: if depth increases, correlation of $f(\mathbf{x})$ and $f(\mathbf{x}')$ does almost not depend on $(\mathbf{x}, \mathbf{x}')$ anymore

random initialization for deep neural networks (ctd.)



- also non-Gaussian limits are possible for deep networks
- if the width of some of the hidden layers is fixed and the number of units in all the other hidden layers tends to infinity, then, the limit will be a composition of Gaussian processes (deep Gaussian process)

a lot of probability theory missing to understand properties of the process $\mathbf{x} \mapsto f(\mathbf{x})$

functions generated by shallow networks

Consider function class

$$\mathcal{F}_{m,\sigma} := \Big\{ f = \sum_{j=1}^m c_j \sigma \big(\mathbf{w}_j^\top \cdot + v_j \big) : \mathbf{w}_j \in \mathbb{R}^d, v_j, c_j \in \mathbb{R} \Big\}.$$

problems:

- how large is this class?
- how well can we approximate functions of a specific smoothness?
- or the function $f(x_1, x_2) = x_1 x_2$?

universal approximation

$$\mathcal{F}_{m,\sigma} := \Big\{ f = \sum_{j=1}^m c_j \sigma \big(\mathbf{w}_j^\top \cdot + v_j \big) : \mathbf{w}_j \in \mathbb{R}^d, v_j, c_j \in \mathbb{R} \Big\}.$$

• functions in the class $\mathcal{F}_{m,\sigma}$ have m(d+2) real parameters

• nested spaces, e.g. $\mathcal{F}_{m,\sigma} \subseteq \mathcal{F}_{m',\sigma}$ whenever $m' \geq m$.

Universal approximation property: Shallow networks with activation function σ have the universal approximation property if for any $\varepsilon > 0$ and any continuous function f on $[0,1]^d$, there exists an integer $m = m(f, \varepsilon)$, such that

$$\inf_{g\in\mathcal{F}_{m,\sigma}}\|f-g\|_{L^{\infty}([0,1]^d)}\leq\varepsilon.$$

reduction to ridge functions

- many proofs first show universal approximation in dimension one
- univariate functions {σ(w · +v) : w, v ∈ ℝ} span the space of continuous functions
- statement does not involve scalar products anymore

afterwards, it is enough to show that the function space spanned by so called ridge functions

$$f = \sum_{j=1}^m g_j(\mathbf{w}_j^\top \cdot)$$

with g_j univariate and continuous has the universal approximation property

universal approximation for univariate functions

Theorem: Shallow networks with smooth activation function that is not a polynomial have universal approximation property for d = 1.

Proof:

•
$$\Delta_h^1 \sigma(t) := (\sigma(t + xh) - \sigma(t))/h$$

•
$$\Delta_h^k \sigma(t) := \Delta_h^1(\Delta_h^{k-1}\sigma)(t)$$

• definition of the k-th derivative \rightsquigarrow

$$\Big| rac{\Delta_h^k \sigma(t)}{x^k} - \sigma^{(k)}(t) \Big| o 0, \quad ext{as} \ h o 0$$

universal approximation for univariate functions (ctd.)

- σ not a polynomial \rightsquigarrow there exists for each k a real number t_k with $\sigma^{(k)}(t_k) \neq 0$
- multiplying with x^k and division by $\sigma^{(k)}(t_k)$ yields

$$\left. rac{\Delta_h^k \sigma(t_k)}{\sigma^{(k)}(t_k)} - x^k
ight| o 0, \quad ext{as } h o 0.$$

- for any h > 0, (σ^(k)(t_k))⁻¹Δ^k_hσ(t_k) can be realized by a shallow network with k + 1 units
- → build networks approximating the function x → x^k arbitrarily well in sup-norm
- apply Weierstrass approximation theorem

some comments on the proof

- proof provides explicit construction of networks that closely resemble polynomials
- construction requires that some parameters are extremely small and others are very large
- uses only one point of the activation function to generate a specific power
- ~> small perturbations of the activation function can lead to completely different properties
- networks can "zoom in" at local features of the activation function
- the universal approximation theorem can be extended to continuous activation functions using local smoothing
- \bullet universal approximation theorem does not hold if σ is a polynomial

universal approximation via Fourier transform

- Fourier transform $\mathcal{F}f(\boldsymbol{\xi}) = \int e^{-i\boldsymbol{\xi}^{\top} \mathbf{x}} f(\mathbf{x}) \, d\mathbf{x}$
- inverse Fourier transform $\mathcal{F}^{-1}f(\mathbf{x}) = (2\pi)^{-d} \int e^{i\mathbf{x}^{\top}\boldsymbol{\xi}} f(\boldsymbol{\xi}) d\boldsymbol{\xi}$

•
$$f = \mathcal{F}^{-1}\mathcal{F}f$$

- for any complex number $z, z = |z|e^{i\phi}$ for some real number $\phi = \phi(z)$
- \rightsquigarrow there exists a real valued function $\phi(\mathbf{w})$ such that $\mathcal{F}f(\mathbf{w}) = e^{i\phi(\mathbf{w})}|\mathcal{F}f(\mathbf{w})|$
- Fourier inversion ~→

$$f(\mathbf{x}) = \frac{1}{(2\pi)^d} \operatorname{Re} \int e^{i\mathbf{w}^\top \mathbf{x}} e^{i\phi(\mathbf{w})} |\mathcal{F}f(\mathbf{w})| d\mathbf{w}$$
$$= \frac{1}{(2\pi)^d} \int \cos\left(\mathbf{w}^\top \mathbf{x} + \phi(\mathbf{w})\right) |\mathcal{F}f(\mathbf{w})| d\mathbf{w}$$

- discretization of the integral on the right hand side gives the structure of a shallow network with activation function cos()
- ~ will be used later for approximation rates

Approximation rates for shallow networks

How well can we approximate a function in dependence on smoothness etc. ?

- smooth activation functions
- approximation rates using multivariate polynomials
- Barron's class

approximation rates for smooth activation function

- Mhaskar '96
- smooth activation function
- β -smooth function (in L^2 -Sobolev sense)
- rate of approximation over all shallow networks with m units is $m^{-\beta/d}$ with d the dimension
- proof first approximates polynomials of ridge functions and then continues with polynomial approximation

approximation rates for arbitrary activation function

- Petrushev '99
- good approximation rates can be obtained for functions that are smoother than the activation function

Theorem: if activation function is *s*-smooth (Sobolev), optimal approximation rates are obtained for s + (d - 1)/2-smooth functions

- $\bullet \ \rightsquigarrow$ effect becomes better as input dimension increases
- **proof:** reduce to ridge functions + approximation of Radon inversion + polynomial eigenbasis
- proof is constructive \rightsquigarrow several interesting conclusions

Barron's approximation theorem

- for any sigmoidal activation function
- any $m \geq 1$,
- any function f
- define $C_f := \int |\mathbf{w}|_1 \mathcal{F}(f)(\mathbf{w}) d\mathbf{w}$
- there exist shallow network such that

$$\left\|f(\cdot)-f(\mathbf{0})-\sum_{j=1}^m c_j\sigma(\mathbf{w}_j^\top\cdot+\mathbf{v}_j)\right\|\leq \frac{2C_f}{(2\pi)^d\sqrt{m}},$$

Remarks:

- proof is based on Maurey's theorem
- rate $m^{-1/2}$ does not depend on the dimension d
- do neural networks avoid curse of dimensionality?

On the rate

• Recall: $C_f = \int |\mathbf{w}|_1 |\mathcal{F}f(\mathbf{w})| d\mathbf{w}$

- indeed there is nothing special about neural networks here
- Candes '02 shows that truncated Fourier series achieves faster approximation rate

 $m^{-1/2-1/d}$

for the same function class $\{f: C_f < \infty\}$

• gain is related to loss in Maurey's theorem

Up to now, no approximation problem has been found where shallow networks outperform Fourier series or polynomial approximation

statistical model

- combine approximation theory with statistical analysis
- given an i.i.d. sample (X_i, Y_i) ∈ ℝ^d × ℝ, i = 1,..., n with bounded responses |Y_i| ≤ 1,
- want to recover the regression function

$$f(\mathbf{x}) = E[Y_i | \mathbf{X}_i = \mathbf{x}]$$

• covers binary classification $\rightsquigarrow Y_i \in \{0, 1\}$ and $f(\mathbf{x}) = P(Y_i = 1 | \mathbf{X}_i = x)$

oracle inequality

• \widehat{f} be the empirical risk minimizer

$$\widehat{f} \in \operatorname{argmin}_{f_{\theta}: \theta \in \Theta} \sum_{i=1}^{n} (Y_i - f_{\theta}(\mathbf{X}_i))^2.$$

standard exponential inequalities
 → if Θ is a discrete set with cardinality |Θ|, then

$$E_f \left[\|\widehat{f} - f\|_2^2 \right] \le C \inf_{\theta \in \Theta} \|f - f_\theta\|_2^2 + C \frac{\log |\Theta|}{n}$$

<ロト < 回ト < 巨ト < 巨ト < 巨ト < 巨ト < 三 へのの 32/101

statistical bounds for shallow networks

- Barron '94
- discretizes network parameters
- study empirical risk minimizer
- number of parameters is m(d+2)
- $\log |\Theta| \lesssim m(d+2) \log n$
- oracle inequality + approximation theory \rightsquigarrow

$$E_f\big[\|\widehat{f}-f\|_2^2\big] \lesssim m^{-1} + \frac{m\log n}{n}$$

if $C_f = \int |\mathbf{w}|_1 |\mathcal{F}f(\mathbf{w})| d\mathbf{w} < \infty$.

- bias variance trade-off $\rightsquigarrow m = \sqrt{n/\log n}$
- yields the rate

$$\sqrt{\frac{\log n}{n}}$$

summary

shallow networks:

- universal approximation
- approximation rates
- estimation risk bounds

no gain in terms of rates with respect to series estimators

advantages of additional layers

- Iocalization
- approximation of polynomials with deep networks
- Kolmogorov-Arnold representation theorem
- advantages of deep ReLU networks

localization with Heaviside activation function

- no localization for shallow networks in dimension d > 1 (?)
- for commonly used activation functions, taking two hidden layers allows us to localize in arbitrary dimensions
- Heaviside activation function $\sigma_0 = \mathbf{1}(\cdot \geq \mathbf{0})$,

$$\mathbf{1}(\mathbf{x} \in [-1,1]^d) = \sigma_0 \Big(\sum_{i=1}^d \sigma_0(x_i+1) + \sigma_0(-x_i+1) - 2d + \frac{1}{2} \Big)$$

- ~> outer neuron only gets activated iff all the inner neurons output one
- this is the case iff $-1 \le x_i \le 1$ for all $i = 1, \ldots, d$

<ロト < 回 ト < 巨 ト < 巨 ト < 巨 ト 三 の Q (~ 36 / 101
localization by other activation functions

• for sigmoidal activation function

 $\sigma(\alpha x) \approx \sigma_0(x)$, for large α .

• for the ReLU $\sigma(x) = (x)_+$,

$$\sigma(\alpha x) - \sigma(\alpha x - 1) \approx \sigma_0(x)$$
, for large α .

- ${\, \bullet \, }$ approximation quality depends on α
- \rightsquigarrow results using neural networks with sigmoidal activation often have conditions on the speed at which $|\sigma(x)| \to 0$ for $x \to -\infty$, and $|1 \sigma(x)| \to 0$ for $x \to +\infty$
- localization might be a useful property for approximation, being non-local might be helpful for the (stochastic) gradient descent

approximation of x^{2^k} with deep networks

- stacking layers on top of each other, this can be reduced to O(k) units in k layers
- rescaled finite second order differences

$$\frac{\sigma(t+2xh)-2\sigma(t+xh)+\sigma(t)}{\sigma''(t)h^2}\approx x^2.$$

a graphical proof





<ロ > < 回 > < 回 > < 三 > < 三 > < 三 > 三 の Q (~ 39/101

advantages of additional layers

- localization
- approximation of polynomials with deep networks
- Kolmogorov-Arnold representation theorem
- advantages of deep ReLU networks

KA representation theorem

- for any continuous function $f:[0,1]^d \to \mathbb{R},$
- there exist univariate continuous functions g_q , $\psi_{p,q}$ such that

$$f(x_1,\ldots,x_d)=\sum_{q=0}^{2d}g_q\Big(\sum_{p=1}^d\psi_{p,q}(x_p)\Big).$$

- used to prove Hilbert's 13th problem
- very different from other representation/approximation schemes, e.g. wavelets, splines, ...

modern KA representation theorem

Theorem 2.14 in Braun '09: There are real numbers a, b_p, c_q and a continuous and monotone function $\psi : \mathbb{R} \to \mathbb{R}$, such that for any continuous function $f : [0, 1]^d \to \mathbb{R}$, there exists a continuous function $g : \mathbb{R} \to \mathbb{R}$ with

$$f(x_1,\ldots,x_d) = \sum_{q=0}^{2d} g\left(\sum_{p=1}^d b_p \psi(x_p + qa) + c_q\right).$$

comparison with two hidden layer neural networks

A two hidden layer neural network is a function of the form

$$\mathsf{x} \mapsto \sum_{\ell=1}^{m_2} e_\ell \, \sigma \bigg(\sum_{j=1}^{m_1} c_{\ell j} \sigma(\mathbf{a}_j^\top \mathsf{x} + b_j) + d_\ell \bigg), \, \, \mathsf{a}_j \in \mathbb{R}^d, b_j, c_{\ell,j}, d_\ell, e_\ell \in \mathbb{R}$$

this has a similar structure as the KA representation

$$f(x_1,\ldots,x_d)=\sum_{\ell=0}^{2d}g\Big(\sum_{j=1}^d c_j\psi(x_j+\ell b)+d_\ell\Big).$$

• interior function is independent of $f \rightsquigarrow$ pre-training

Representation Properties of Networks: Kolmogorov's Theorem Is Irrelevant

Federico Girosi Tomaso Poggio

Massachusetts Institute of Technology, Artificial Intelligence Laboratory, Cambridge, MA 02142 USA and Center for Biological Information Processing, Whitaker College, Cambridge, MA 02142 USA

Many neural networks can be regarded as attempting to approximate a multivariate function in terms of one-input one-output units. This note considers the problem of an exact representation of nonlinear mappings in terms of simpler functions of fewer variables. We review Kolmogorov's theorem on the representation of functions of several variables in terms of functions of one variable and show that it is irrelevant in the context of networks for learning.

NOTE

Communicated by Halbert White

Kolmogorov's Theorem Is Relevant

Věra Kurková

Institute of Computer Science, Czechoslovak Academy of Sciences, P. O. Box 5, 182 07 Prague 8, Czechoslovakia

We show that Kolmogorov's theorem on representations of continuous functions of *n*-variables by sums and superpositions of continuous functions of one variable is relevant in the context of neural networks. We give a version of this theorem with all of the one-variable functions approximated arbitrarily well by linear combinations of compositions of affine functions with some given sigmoidal function. We derive an upper estimate of the number of hidden units.

approximation theory based on KA representation

$$f(x_1,\ldots,x_d)=\sum_{\ell=0}^{2d}g\Big(\sum_{j=1}^d c_j\psi(x_j+\ell b)+d_\ell\Big).$$

- for approximation theory it is sufficient to approximate outer function g
- unfortunately, g does not have good smoothness properties
- best bound so far derived in Kurkova '92 requires m^{4+d} network parameters to approximate β -Hölder function up to an error $m^{-\beta}$
- the 4 in the exponent seems to be suboptimal

a simple KA representation

There exists a monotone function $\psi : \mathbb{R} \to \mathbb{R}$ such that for any function $f : [0,1]^d \to \mathbb{R}$, we can find a function $g : \mathbb{R} \to \mathbb{R}$ with

$$f(x_1,\ldots,x_d)=g\Big(\sum_{p=1}^d 2^{-p}\psi(x_p)\Big).$$

proof sketch for d = 2:

- dyadic expansion $x_1 = [0.a_1a_2a_3...]_2$ and $x_2 = [0.b_1b_2b_3...]_2$
- set $\Psi(x_1, x_2) = [0.a_1b_1a_2b_2a_3b_3...]_2$
- $\Psi(x_1, x_2) = \sum_{p=1}^2 2^{-p} \psi(x_p)$ for suitable ψ

$$f(x_1, x_2) = \underbrace{f \circ \Psi^{-1}}_{=:g} \circ \Psi(x_1, x_2) = g\left(\sum_{p=1}^2 2^{-p} \psi(x_p)\right)$$

Morton order



Source: Bader, Space-Filling Curves, Springer 2013

- \bullet the inverse $\Psi^{-1}:[0,1]\to [0,1]^2$ is a space-filling curve, called the Morton order
- the outer function g is much rougher than f

generalization



Source: wikipedia

what we need:

- $\Psi^{-1}:[0,1]\to [0,1]^d$ is a space-filling curve
- $\Psi: [0,1]^d
 ightarrow [0,1]$ must be an additive function
- we can also replace [0,1] by a subset

the Lebesgue curve:

•
$$x_1 = [0.a_1a_2a_3...]_2$$
 and $x_2 = [0.b_1b_2b_3...]_2$

$$\Psi(x_1, x_2) = [0.(2a_1)(2b_1)(2a_2)(2b_2)\dots]_3$$

• $\Psi : [0,1]^2 \to \mathcal{C}$ (Cantor set), Ψ is invertible and additive

smoothness of the outer function

For
$$x = [0.a_1a_2...]_2$$
, set

$$\phi(x) := \sum_{j=1}^{\infty} \frac{2a_j}{3^{d(j-1)}}$$

Lemma: for any function $f : [0,1]^d \to \mathbb{R}$, we can find a function $g : \mathcal{C} \to \mathbb{R}$ such that

$$f(x_1,\ldots,x_d)=g\Big(\sum_{p=1}^d 3^{-p}\phi(x_p)\Big);$$

smoothness of the outer function

$$f(x_1,\ldots,x_d)=g\Big(\sum_{p=1}^d 3^{-p}\phi(x_p)\Big);$$

Lemma: If for $\beta \leq 1$, $|f(\mathbf{x}) - f(\mathbf{y})| \leq Q|\mathbf{x} - \mathbf{y}|_{\infty}^{\beta}$, for all $\mathbf{x}, \mathbf{y} \in [0, 1]^d$, then,

$$|g(x) - g(y)| \le 2^{\beta}Q|x - y|^{\frac{\beta \log 2}{d \log 3}}, \text{ for all } x, y \in \mathcal{C};$$

If f is piecewise constant on dyadic intervals, then, g is Lipschitz

 $\log 2/\log 3$ is the Hausdorff dimension of the Cantor set C

approximation theory

- based on the previous result we can build an approximation theory
- truncate interior function: for $x = [0.a_1a_2...]_2$ set

$$\phi_{\mathbf{K}}(\mathbf{x}) := \sum_{j=1}^{\mathbf{K}} \frac{2a_j}{3^{d(j-1)}}$$

• because of the smoothness of the outer function we can control the induced error

$$\left|f(\mathbf{x}) - g\left(\sum_{p=1}^{d} 3^{-p} \phi_{\mathcal{K}}(x_p)\right)\right| \leq Q 2^{-\beta(\mathcal{K}-4)}$$

• the result can be used to build a neural network approximating *f* with the optimal rate

implementation of the interior function

the interior function

$$\phi_{\boldsymbol{K}}(\boldsymbol{x}) := \sum_{j=1}^{\boldsymbol{K}} \frac{2a_j}{3^{d(j-1)}}$$

can be represented by a deep network with

- K hidden layers
- linear $\sigma(x) = x$ and threshold activation function $\sigma(x) = \mathbf{1}(x \ge 1/2)$
- 3 units in each hidden layer doing bit extraction

this can be well approximated by a deep ReLU network with the same architecture

deep ReLU approximation



add one hidden layer to construct a deep ReLU network computing approximately

$$g\left(\sum_{q=1}^{d} 3^{-q} \phi_{\kappa}(x_q)\right) \approx f(\mathbf{x})$$

<ロト < 団ト < 巨ト < 巨ト < 巨ト 三 の Q @ 54/101

approximation by deep ReLU networks

if $\beta \leq 1,$ there exists a deep ReLU network \widetilde{f} with

- K + 3 hidden layers
- $O(2^{Kd})$ network parameters (fully connected)

such that the approximation error is bounded by

 $\lesssim 2^{-\beta K}$.

differently speaking we need m^d parameters to achieve approximation error $m^{-\beta}$

improved representation theorems

- Kolmogorov-Arnold approximation theorem shows that every continuous function can be represented by a specific two-layer network
- very different structure if compared with the universal approximation theorem for shallow networks
- indicates that additional layers can lead to new features of network functions

Theorem (Braun '09): Fix $d \ge 2$. There are real numbers a, b_p, c_q and a continuous and monotone function $\psi : \mathbb{R} \to \mathbb{R}$, such that for any continuous function $f : [0, 1]^d \to \mathbb{R}$, there exists a continuous function $g : \mathbb{R} \to \mathbb{R}$ with

$$f(x_1,\ldots,x_d) = \sum_{q=0}^{2d} g\Big(\sum_{p=1}^d b_p \psi(x_p+qa) + c_q\Big).$$

56 / 101

remarks

$$f(x_1,\ldots,x_d)=\sum_{q=0}^{2d}g\Big(\sum_{p=1}^d b_p\psi(x_p+qa)+c_q\Big).$$

- ullet one inner function ψ and one outer function g
- inner function is independent of f
- *q*-dependence in the first layer comes through the shifts *qa*.
- right hand side can be realized by a network with two hidden layers, architecture $\mathbf{p} = (d, d, 2d + 1, 1)$, and ψ being the activation function in the first layer.

link to pre-training

- inner function in the Kolmogorov-Arnold representation theorem is independent of the represented function *f*
- in deep learning it has been observed that the first layers build function systems which can be used for other classification problems
- exploited in pre-training where a trained deep network from a possibly completely different classification problem is taken and only the last layer is learned by the new dataset
- fact that pre-training works shows that deep networks build generic function systems in the first layers.

we discuss several advantages of deep ReLU networks

- representation of identity
- growth of number of linear pieces
- approximation by ReLU networks with small parameters

deep ReLU networks can learn skip connections

 $\sigma(x) = \max(x, 0)$

projection property

 $\sigma \circ \sigma = \sigma$

- $\bullet \rightsquigarrow$ pass a signal without change through several layers in the network
- $\bullet \, \rightsquigarrow$ network synchronization by adding hidden layers
- related to skip connections and ResNets
- for other activation functions it is much harder to approximate the identity

number of linear pieces of deep ReLU networks

- deep ReLU networks are piecewise linear functions of the input
- adding layers \rightsquigarrow highly oscillating functions with few parameters
- consider ReLU network with two hidden layers and width vector (1, m, 1, 1) of the form

$$\Big(\sum_{j=1}^m c_j(w_jx+v_j)_+\Big)_+$$

- ~> number of added pieces by outer ReLU is proportional to number of zero crossings of inner function
- any ReLU network with width vector (1, p₁,..., p_L, 1) has at most

$$\left(\frac{3}{2}\right)^{L}\prod_{j=1}^{L}(p_{j}+1)$$

pieces

an open problem

Problem: Show that the function

$$(x_1,\ldots,x_{2^k})\mapsto \max(0,x_1,\ldots,x_{2^k})$$

can be exactly represented by a ReLU network with k + 1 hidden layers.

Conjecture: The same function cannot be represented by any ReLU network with less than k + 1 hidden layers.

• Except for k = 1, no prove yet.

example of a highly oscillating function

Functions:

- let $T : [0,1] \to [0,1],$ $T(x) := (2x) \land (2-2x) = (2x)_+ - (4x-2)_+$
- can be realized by shallow network with two units • $R^k: [0,1] \to [0,1],$

$$R^k := \underbrace{T \circ T \circ \dots T}_{k \text{ times}}$$

network representation



<ロト < 回 ト < 三 ト < 三 ト ミ の < () 64/101

multiplication

how can we (approximately) multiply two inputs with a network?

- crucial problem for approximation theory
- for deep networks this can be reduced to approximation of square function x → x² via

$$xy = \left(\frac{x+y}{2}\right)^2 - \left(\frac{x-y}{2}\right)^2$$

• has a surprising answer for ReLU networks

network approximation of the function $x \mapsto x^2$ is very important!

• for twice differentiable activation function, we used

$$rac{\sigma(t+2xh)-2\sigma(t+xh)+\sigma(xh)}{h^2\sigma''(t)}
ightarrow x^2 ext{ for } h
ightarrow 0$$

- ~ network parameters become large
- for deep ReLU networks we use a different construction

ReLU approximation of the square function

Functions:
• let
$$T^k : [0, 2^{2-2k}] \to [0, 2^{-2k}],$$

 $T^k(x) := (x/2) \land (2^{1-2k} - x/2) = (x/2)_+ - (x - 2^{1-2k})_+$
• $R^k : [0, 1] \to [0, 2^{-2k}],$
 $R^k := T^k \circ T^{k-1} \circ \dots T^1.$

Lemma (Telgarsky '16, Yarotski '18, SH '17):

$$\left|x(1-x)-\sum_{k=1}^{m}R^{k}(x)\right|\leq 2^{-m}.$$

<ロト < 回 ト < 巨 ト < 巨 ト < 巨 ト 三 の へ () 67/101

ReLU approximation of the square function



- (Left plot) The functions R^1 (red), R^2 (orange), R^3 (yellow).
- (Right plot) Approximation of x(1 x) (blue) by R¹ (red) and R¹ + R² (orange).

rewriting approximation as network

$$\left|x(1-x)-\sum_{k=1}^{m}R^{k}(x)\right|\leq 2^{-m}.$$

deep ReLU approximation:

- m hidden layers
- O(m) network parameters
- bounded parameters
- approximation 2^{-m}

shallow ReLU network

• for x(1-x) a shallow ReLU network needs at least $O(2^{m/2})$ parameters to achieve approximation error 2^{-m}

multiplication with deep ReLU networks

Lemma: There exists a network $Mult_m$ with m + 4 hidden layers, width vector $(2, 6, 6, \ldots, 6, 1)$ and all network parameters bounded by one, such that

$$\big|\operatorname{\mathsf{Mult}}_m(x,y)-xy\big|\leq 2^{-m},\quad ext{for all }x,y\in[0,1].$$

Proof:

• use polarization identity

$$xy = \left(\frac{x+y}{2}\right)^2 - \left(\frac{x-y}{2}\right)^2$$

- separation of positive and negative part
- compute (x + y)/2 and (x y)/2 in first layer (non-negativity!)
- square network has to be incorporated twice (inefficient)

a step in the proof



localization and approximation

- we have seen that with two hidden layers we can localize
- how can this be done for ReLU networks?
- goes back to Yarotsky '18
- define $\mathbf{D}(M)$ as all grid points on the grid

$$\left\{(\ell_j/M)_{j=1,\ldots,r}: \boldsymbol{\ell}=(\ell_1,\ldots,\ell_r)\in\{0,1,\ldots,M\}^r\right\}$$

partition of unity on unit cube

$$\sum_{\mathbf{x}_{\ell}\in\mathbf{D}(M)} \underbrace{\prod_{j=1}^{r} (1-M|x_j-x_j^{\ell}|)_+}_{\text{localized functions}} = \prod_{j=1}^{r} \sum_{\ell=0}^{M} (1-M|x_j-\ell/M|)_+ = 1,$$
local Taylor approximation

• on each localized bit $(\mathbf{a} \in \mathbf{D}(M))$ do a Taylor approximation

$$f(\mathbf{x}) \approx P_{\mathbf{a}}^{\beta} f(\mathbf{x}) := \sum_{0 \le |\boldsymbol{\alpha}| < \beta} (\partial^{\boldsymbol{\alpha}} f)(\mathbf{a}) \frac{(\mathbf{x} - \mathbf{a})^{\boldsymbol{\alpha}}}{\boldsymbol{\alpha}!} = \sum_{0 \le |\boldsymbol{\gamma}| < \beta} \mathbf{x}^{\boldsymbol{\gamma}} c_{\boldsymbol{\gamma}}$$

- ${\scriptstyle \bullet} \,$ this can be approximately realized by a deep ReLU network
- many technicalities occur (see the article SH '17)

approximation rate

Theorem: For any β -smooth function $f : [0,1]^r \to \mathbb{R}$ and any integers $m, N \ge 1$, there exists a ReLU network with

- depth $L \asymp m$
- ullet width in each layer bounded by $\lesssim N$
- ${\scriptstyle \bullet}$ number of non-zero network parameters $s \lesssim {\it Nm}$ such that



remarks



- for deep networks first term is of smaller order
- second term becomes suboptimal for large depth
- trade-off
- sparse networks

risk bounds for deep ReLU networks

Framework:

- we now study a statistical problem
- requires that we first need to specify a statistical model
- we study nonparametric regression

mathematical problem

 $X = \{images\}$ $f: X \longrightarrow Y$ $Y = \{"cat", "dog"\}$

The data are used to fit a network, i.e. estimate the network parameters.

How fast does the estimated network converge to the truth f as sample size increases?

statistical analysis

• we observe *n* i.i.d. copies $(\mathbf{X}_1, Y_1), \ldots, (\mathbf{X}_n, Y_n),$

$$Y_i = f(\mathbf{X}_i) + \varepsilon_i, \quad \varepsilon_i \sim \mathcal{N}(0, 1)$$

• $\mathbf{X}_i \in \mathbb{R}^d, \ Y_i \in \mathbb{R},$

 ${\scriptstyle \bullet }$ goal is to reconstruct the function $f:\mathbb{R}^{d}\rightarrow \mathbb{R}$

 has been studied extensively (kernel smoothing, wavelets, splines, ...)

the estimator

- choose network architecture (*L*, **p**) and sparsity *s*
- denote by $\mathcal{F}(L,\mathbf{p},s)$ the class of all networks with
 - architecture (L, p)
 - number of active (e.g. non-zero) parameters is s
- our theory applies to any estimator \hat{f}_n taking values in $\mathcal{F}(L,\mathbf{p},s)$
- prediction error

$$R(\widehat{f}_n, f) := E_f \big[\big(\widehat{f}_n(\mathbf{X}) - f(\mathbf{X}) \big)^2 \big],$$

with $\mathbf{X} \stackrel{\mathcal{D}}{=} \mathbf{X}_1$ being independent of the sample • study the dependence of n on $R(\hat{f}_n, f)$

- $\bullet\,$ classical idea: assume that regression function is $\beta\text{-smooth}\,$
- optimal nonparametric estimation rate is $n^{-2\beta/(2\beta+d)}$
- suffers from curse of dimensionality
- to understand deep learning this setting is therefore useless
- \rightsquigarrow make a good structural assumption on f

hierarchical structure



- Important: Only few objects are combined on deeper abstraction level
 - few letters in one word
 - few words in one sentence

function class

We assume that

$$f = g_q \circ \ldots \circ g_0$$

with

- $g_i : \mathbb{R}^{d_i} \to \mathbb{R}^{d_{i+1}}$.
- each of the d_{i+1} components of g_i is β_i-smooth and depends only on t_i variables
- t_i can be much smaller than d_i
- effective smoothness

$$eta_i^* := eta_i \prod_{\ell=i+1}^q (eta_\ell \wedge 1).$$

we show that the rate depends on the pairs

$$(t_i, \beta_i^*), \quad i=0,\ldots,q.$$

 similar conditions have been proposed by Horowitz & Mammen (2007), Kohler & Kryzak (2017), Bauer & Kohler (2017), Kohler & Langer (2018) example

$$f_0(x_1, x_2, x_3) = g_{11}(g_{01}(x_3), g_{02}(x_2))$$

•
$$f_0 = g_1 \circ g_0$$

• $d_0 = 3, t_0 = 1, d_1 = t_1 = 2, d_2 = 1$

main result

Theorem: If

(i) depth $\asymp \log n$

(ii) width \geq network sparsity $\asymp \max_{i=0,...,q} n^{\frac{t_i}{2\beta_i^*+t_i}} \log n$ Then, for any network reconstruction method \hat{f}_n ,

prediction error $\approx \phi_n + \Delta_n$

(up to log *n*-factors) with

$$\Delta_n := E\Big[\frac{1}{n}\sum_{i=1}^n (Y_i - \widehat{f}_n(\mathbf{X}_i))^2 - \inf_{f \in \mathcal{F}(L,\mathbf{p},s)} \frac{1}{n}\sum_{i=1}^n (Y_i - f(\mathbf{X}_i))^2\Big]$$

and

$$\phi_n := \max_{i=0,...,q} n^{-\frac{2\beta_i^*}{2\beta_i^*+t_i}}.$$

consequences

- empirical risk minimizer is optimal in this class
- fitted neural networks are high-dimensional (no upper bound on the width)
- network sparsity induces regularization
- the assumption that depth $\asymp \log n$ appears naturally
- in particular the depth scales with the sample size

important for statistical performance is not the size of the network but the amount of regularization

consequences (ctd.)

paradox:

- good rate for all smoothness indices
- existing piecewise linear methods only give good rates up to smoothness two
- Here the non-linearity of the function class helps

 \rightsquigarrow non-linearity is essential!!!

additive models

• functions are of the form

$$f(x_1,\ldots,x_d)=f_1(x_1)+\ldots+f_d(x_d)$$

•
$$f_i \text{ are } \beta\text{-smooth}$$

• $f = g_1 \circ g_0$ with
 $g_0(\mathbf{x}) = (f_1(x_1), \dots, f_d(x_d))^\top \text{ and } g_1(\mathbf{y}) = \sum_{j=1}^d y_j$
• $\rightsquigarrow d_0 = d, t_0 = 1, d_1 = t_1 = d, d_2 = 1$

rate achieved by a neural network

$$R(\widehat{f}_n, f_0) \lesssim n^{-rac{2eta}{2eta+1}} \log^3 n + \Delta(\widehat{f}_n, f_0).$$

on the proof

oracle inequality (roughly)

$$R(\widehat{f}, f) \lesssim \inf_{f^* \in \mathcal{F}(L,\mathbf{p},s)} \left\| f^* - f \right\|_{\infty}^2 + \frac{\log \mathcal{N}_n}{n}.$$

• $\log N_n$ denotes the covering entropy

shows the trade-off between approximation and model size

• for networks we obtain a bound of the type

 $\log N_n \lesssim sL \log(n)$

 $\bullet \ \rightsquigarrow$ trade-off between approximation and network sparsity

lower bounds on the network sparsity

the convergence theorem implies a deterministic lower bound on the network sparsity required to approximate β -smooth functions on $[0,1]^d$

Result:

 $\bullet~$ if for $\varepsilon>0,$ $s\lesssim \frac{\varepsilon^{-d/\beta}}{L\log(1/\varepsilon)}$

then

$$\sup_{f_0 \text{ is } \beta-\text{H\"older } f \text{ a } s-\text{sparse network}} \|f - f_0\|_{\infty} \geq \varepsilon.$$

has been proved via a different technique in Bölcskei et al. '17

Network sparsity is crucial in the proof but classical deep learning produces dense networks. Recently many new methods have been proposed generating sparsely connected networks.

- sparsifying as post-processing step \rightsquigarrow compression
- starting with sparse network topology
- evolutionary methods inspired by human brain

other statistical results

- piecewise smooth functions, Imaizumi and Fukumizu '18
- binary classification with hinge loss, Kim, Ohn, Kim '18
- causal models, Farell, Liang, Misra '18
- deep Q-learning, Fan et al. '20
- and others ...

suboptimality of wavelet estimators

•
$$f(\mathbf{x}) = h(x_1 + \ldots + x_d)$$

- for some α -smooth function h
- Rate for DNNs $\lesssim n^{-lpha/(2lpha+1)}$ (up to logarithmic factors)
- Rate for best wavelet thresholding estimator $\gtrsim n^{-lpha/(2lpha+d)}$
- Reason: Low-dimensional structure does not affect the decay of the wavelet coefficients

double descent and implicit regularization



overparametrization generalizes well \rightsquigarrow implicit regularization

overfitting

- training error = 0 implies that $\Delta_n = 0$
- Δ_n does not fully characterize the statistical properties anymore
- because of implicit regularization, SGD will pick interpolant with good statistical properties

can implicit regularization avoid sparsity?

we conjecture the answer is **no** for the regression problem!

interpolation properties

- consider continuous activation function that is not a polynomial
- given data $(\mathbf{X}_k, Y_k) \in \mathbb{R}^d imes \mathbb{R}$ with distinct design vectors \mathbf{X}_k
- shallow networks: one can perfectly interpolate *n* data points with *n* units in the hidden layer
- related to the universal approximation theorem (therefore same condition appears)

theory for vanishing training error

smooth activation function

- Du et al. '18 consider highly over-parametrized setting
- number of units in each layer has to be of some (unspecified
 ?) polynomial order in the sample size
- setup is regression with least-squares loss
- show that gradient descent with randomly initialization converges to zero training error

ReLU networks

- Allen-Zhu et al. '18 shows a similar result
- one assumptions is that the network width scales at least with the 30-th power of the sample size

does data interpolation contradict statistical optimality?



Source: Belkin, Rakhlin, Tsybakov, 2018

in principle it is possible to interpolate and to denoise simultaneously

more details



we can show that for a simplified model and properly chosen learning rate, SGD converges to natural cubic spline interpolant \rightsquigarrow inconsistent estimator

main idea

A shallow network (L = 1) can be written as

$$x\mapsto \sum_{j=1}^m a_j(b_jx-c_j)_+, \quad a_j, b_j, c_j\in \mathbb{R}.$$

Taylor expansion in one dimension

$$g(x) = g(0) + xg'(0) + \int g''(u)(x-u)_+ du$$

If, say g(0) = g'(0) = 0, we have that approximately

$$g(x) \approx \frac{1}{m} \sum_{j=1}^{m} g''\left(\frac{j}{m}\right) \left(x - \frac{j}{m}\right)_+.$$

<ロト < 団ト < 巨ト < 巨ト < 巨ト 三 の Q () 99/101

denoising vs. interpolation

- implicit regularization is not sufficient to do denoising
- it still works in practice because standard datasets have a lot of structure in common (classification with few misclassified data points)

All statements that start with "In deep learning" are wrong, what matters is the structure of the data. To describe for which data structures such claims are true is a major challenge for research in statistics.

outlook

deep networks are an exciting field with many open problems

- new phenomena, ...
- network types: CNNs, RNNs, autoencoders, ...
- generative adversarial networks (GANs)

Thank you for your attention!