

GOING BEYOND THE LINEAR REGIME

1) Trackability via overparametrization

2) Double descent

3) Benign overfitting

(self-induced regularization)

→ (A) Separation between linearized NNs vs NNs

→ (B) Fixed features vs feature learning:

breaking the curse of dimensionality using feature learning

→ (C) An example: learning parities

→ (D) Match approaches beyond linear regime

① Separation between linearized NNs vs NNs

LINEAR REGIME: training regime where network can be approximated by a linear model during the whole training dynamics

$$\Rightarrow f(x, \theta) \hookrightarrow f^{\text{lin}}(x, \theta) = f(x, \theta_0) + \langle \theta - \theta_0, \nabla_{\theta} f(x, \theta_0) \rangle$$

$$\theta^0 = \bar{\theta}^0 = \theta_0$$

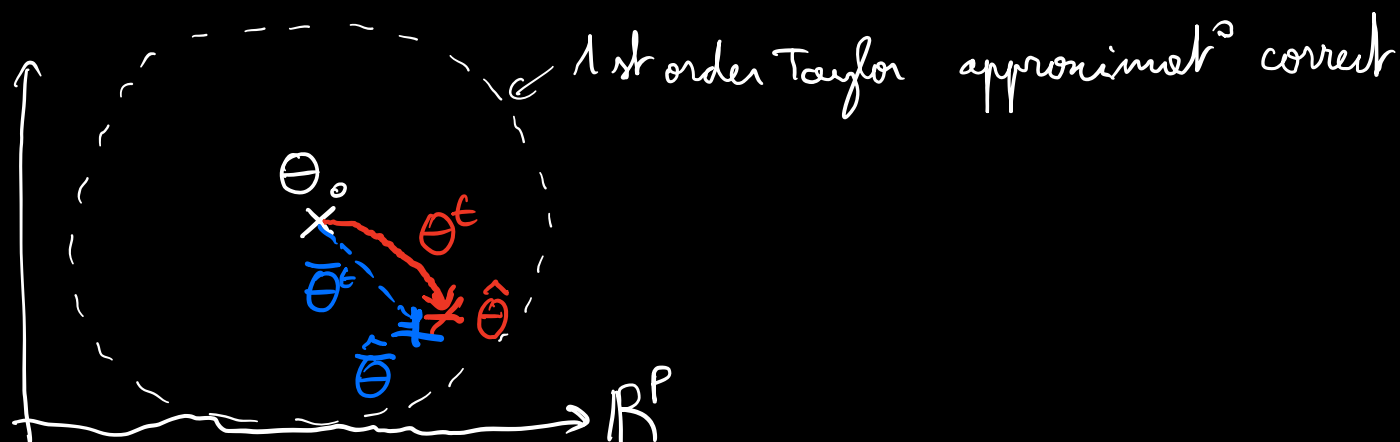
$$\textcircled{1} \dot{\theta}^t = -\nabla \hat{R}_n(f(x, \theta^t))$$

$$\textcircled{2} \dot{\bar{\theta}}^t = -\nabla \hat{R}_n(f^{\text{lin}}(x, \bar{\theta}^t))$$

$$\frac{\|\theta^t - \bar{\theta}^t\|_2}{\|\bar{\theta}^t - \bar{\theta}^0\|_2} \ll 1$$

$$\Rightarrow f(x, \theta^t) \approx f^{\text{lin}}(x, \bar{\theta}^t)$$

⇒ in this regime: NNEDs can effectively be replaced by linear models



1) Are NNEDs in practice trained in the linear regime?

→ Sometimes, mostly not.

2) Does linear theory capture what can be achieved by NNEDs?

→ No.

3) Do we have a better theory? → understand both optimization and generalization.

→ Not yet.

Linear regime: → explains why GD/SGD can find a global optima of a highly non-convex problem

Successfully illustrated: **TRACTABILITY VIA OVERPARAMETRIZATION**

→ problem becomes more tractable as # of parameters ↑

\Rightarrow since then: lots of work to show the limitation of linear
require theory to explain good generalization of NNETs

Most of the (theoretical) work: show in specific examples that
NNETs outperform linearized NNETs.

This results are called "separation results".

"Obvious" separation in approximation power:

2-layer NNET: $f_{\text{NN}}(\alpha, a, W) = \sum_{i=1}^N \underline{a_i} \underline{\sigma}(\underline{w_i}, \alpha)$
 $a_i \in \mathbb{R}, w_i \in \mathbb{R}^d$

2-layer linearized NNET: fix $W^0 = (w_1^0, \dots, w_N^0)$
(RF) $f_{\text{RF}}(\alpha, a) = \sum_{i=1}^N \underline{a_i} \underline{\sigma}(\underline{w_i^0}, \alpha)$ \rightarrow also apply full linearizat.

* $F_{\text{NN}}(\mathcal{B}) = \{ f_{\text{NN}}(\alpha, a, W) : \|\alpha\|_2, \|W\|_F \leq \mathcal{B} \}$

* $F_{\text{RF}}(W^0) = \{ f_{\text{RF}}(\alpha, a) : a \in \mathbb{R}^N \}$ \rightarrow test error $n \rightarrow \infty$

Approximation error:

$R_{\text{App}}(f_*, F) = \inf_{f \in F} \|f_* - f\|_{L^2}^2$

↳ best you can hope for any number of samples

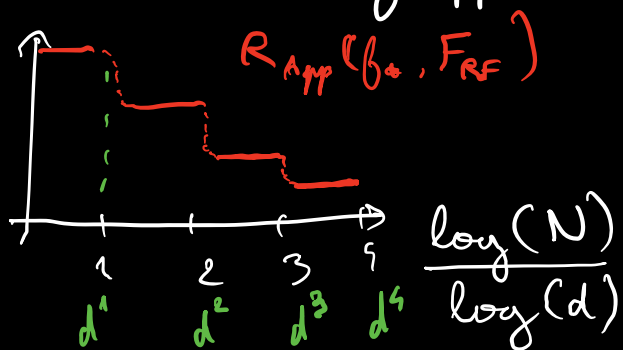
Take: $x \sim \text{Unif}(\mathbb{S}^{d-1})$
 $w_i^0 \sim \text{Unif}(\mathbb{S}^{d-1})$ fixed

Thm: [Misiakiewicz et al., 2019]

For any $f_* \in L^2(\mathbb{S}^{d-1})$. If $d^l \ll N \ll d^{l+1}$, then

$$R_{\text{app}}(f_*, F_{\text{RF}}(W^0)) = \underbrace{\|P_{>l} f_*\|_{L^2}^2}_{\text{best degree-}l \text{ poly. approx.}} + o_d(1).$$

→ can only approximate degree l polynomials



Staircase phenomena.

⇒ Simple example: one neuron $f_*(x) = \sigma(\langle w_*, x \rangle)$

App. with RF: $R_{\text{app}}(f_*, F_{\text{RF}}) \approx \underbrace{\|P_{>l} \sigma\|_{L^2}^2}_{\text{best degree-}l \text{ poly. approx.}}$ if $N \ll d^l$

App. with NNets: $R_{\text{app}}(f_*, F_{\text{NN}}) = 0$ for $N \geq 1$.

↳ simply take $\alpha_i = 1, w_i = w_*$
+ rest set to 0

⇒ F_{NN} much richer class of fcts

Intuition: in high-dim, $\sup_{i \in [N]} |\langle \omega_i^\circ, \omega_* \rangle| \approx \frac{1}{\sqrt{d}}$
 $\sigma(\langle \omega_i, \cdot \rangle)$ using $\sigma(\langle \omega_i, \cdot \rangle)$

→ no 'good' feature $\sigma(\langle \omega_i, \cdot \rangle)$ to approximate $\sigma(\langle \omega_*, \cdot \rangle)$

→ when 1st layer is not fixed, can select 'good' features
(i.e. ω_i high correlat^o with ω_*)

Kind of obvious and not interesting: the fact that you can approximate does not mean that you can efficiently find these good networks.

Want a separation between linearized NNETs and NNETs that can be "constructed" in practice, e.g., using GD.

Separation between linearized NNETs and gradient trained NNETs:

Inner-Prod kernel: "infinite-width" linearized NNET

$$\frac{1}{N} \sum_{i=1}^N \sigma(\langle x, \omega_i^\circ \rangle) \sigma(\langle z, \omega_i^\circ \rangle) \rightarrow \mathbb{E}_{\omega^\circ} [\sigma(\langle x, \omega^\circ \rangle) \sigma(\langle z, \omega^\circ \rangle)]$$

$N \rightarrow \infty \quad =: \quad \underline{h(\langle x, z \rangle)}$

Test error of KRR (more details later about Kernel methods) ✓

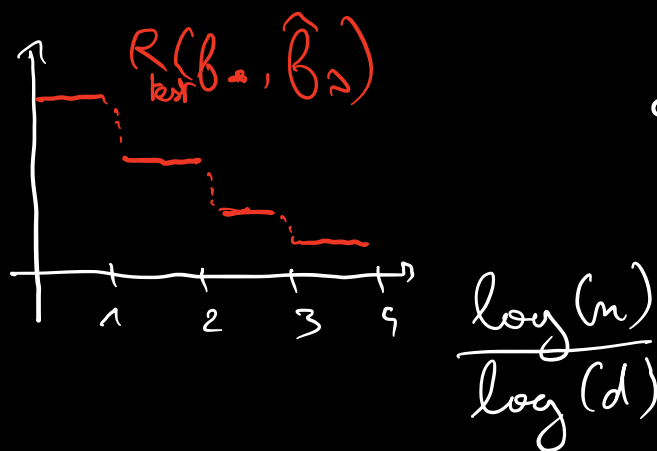
$$R_{\text{test}}(f_*, \hat{f}_d) = \mathbb{E}_n [(f_*(n) - \hat{f}_d(n))^2] \leftarrow$$

Take: $n \sim \text{Unif}(\mathbb{S}^{d-1})$ $f_* \in L^2(\mathbb{S}^{d-1})$

If $d^k \ll \underline{n} \ll d^{k+1}$ (now # training samples)

Then [Minchewicz et al., 2019]

$$R_{\text{test}}(f_*, \hat{f}_d) = \underline{\|P_{>d} f_*\|_{L^2}^2} + o_d(1)$$



Staircase descent

Again: $f_*(n) = \underline{\sigma(\langle w_*, a \rangle)}$

$$P_{>d} \sigma = 0 \quad n \gg d$$

then $R_{\text{test}}(f_*, \hat{f}_d) \approx \underline{\|P_{>d} \sigma\|_{L^2}^2}$ $n \propto d^k$

Take: $\hat{f}_d(n) = \underline{\sigma(\langle w_n, a \rangle)}$ 1 hidden unit

\rightarrow learn w_n using GD then if $\underline{n \geq d \log d}$

$$\underline{R_{\text{test}}(f_*, \hat{f}_d) \approx 0} \quad [\text{Montanari et al., 2017}]$$

\Rightarrow GD in underparametrized regime ($\underline{N=1} \ll n$)

↳ very simple case but already very technical proof.

⇒ In general, studying NNETs trained by GD is currently out of reach except in the linear regime.

⇒ Can we understand the benefit of training more abstractly?

③ Fixed features vs feature learning

$$\rightarrow f_{RF}(x, a) = \sum_{i=1}^N a_i \sigma(\langle w_i^0, x \rangle)$$

$$\rightarrow \text{Kernel} \quad \frac{1}{N} \sum_{i=1}^N \sigma(\langle w_i^0, x \rangle) \sigma(\langle w_i^0, y \rangle) = \underline{H_0(x, y)}$$

$$f_{NN}(x, \underline{\theta^t}) = \sum_{i=1}^N a_i^t \sigma(\langle w_i^t, x \rangle)$$

$$\rightarrow \text{"Kernel"}: \quad \frac{1}{N} \sum_{i=1}^N \sigma(\langle w_i^t, x \rangle) \sigma(\langle w_i^t, y \rangle) = \underline{H_t(x, y)}$$

⇒ GD is a way to "training" a good kernel

i.e., learning "good" features adapted to the data.

H_0
↓
 H_t

$$\underline{\text{e.g.:}} \quad f_* = \sigma(\langle w_*, \cdot \rangle) \rightarrow \underline{H_*(x, y) = \sigma(\langle w_*, x \rangle) \sigma(\langle w_*, y \rangle)}$$

- Linear regime: "kernel regime", "lazy regime"
- Outside linear regime: "feature learning regime"
"rich regimes"

Vastly different behavior between fixed feature (fixed kernel)
and methods that allow 'feature learning'

\Rightarrow they are "adaptive" and can vastly outperform
fixed feature models.

"Breaking the curse of dimensionality using conven NNets"
- Francis BACH (2017)

Background on KRR / RKHS:

$$\{(y_i, x_i)\}_{i=1}^n \quad y_i \in \mathbb{R}, \quad x_i \in \mathbb{R}^d = \mathcal{X} \quad \leftarrow$$

$(\mathcal{X}, \mathbb{P})$: proba space of the data covariates \leftarrow

(Ω, μ) : proba space of the features weights \leftarrow weights

Featurization map: $\phi: \mathcal{X} \times \Omega \rightarrow \mathbb{R}$

$$(x, \omega) \mapsto \underline{\phi(x, \omega)}$$

$$(\phi \in L^2(\mathcal{X} \times \Omega))$$

Model: $f(x, a) = \int_{\Omega} a(\omega) \sigma(\langle x, \omega \rangle) \mu(d\omega)$

→ infinitely-wide 2 layers NNet with
 $a: \Omega \rightarrow \mathbb{R}$

e.g. $a(\omega) = \sum_{i=1}^N a_i \delta_{\omega=\omega_i}$ (discrete) ←

then $f(x, a) = \sum_{i=1}^m a_i \sigma(\langle x, \omega_i \rangle)$

Define norm: $\|f(\cdot, a)\|_{F_2} = \left(\int_{\Omega} |a(\omega)|^2 \mu(d\omega) \right)^{\frac{1}{2}}$
 $= \|a\|_{L^2}$ F_2 -norm

→ $F_2 = \{ f(\cdot, a) \text{ such that } \|f(\cdot, a)\|_{F_2} < \infty \}$

→ Reproducing Kernel Hilbert Space (RKHS) ←

with kernel: $K(x, z) = \int_{\Omega} \sigma(\langle x, \omega \rangle) \sigma(\langle z, \omega \rangle) \mu(d\omega)$

$\|f\|_H = \|a\|_{L^2}$

Kernel Ridge Regression:

$\hat{a} = \underset{a: \Omega \rightarrow \mathbb{R}}{\operatorname{argmin}} \left\{ \sum_{i=1}^m (y_i - f(x_i, a))^2 + \lambda \|a\|_{L^2}^2 \right\}$
 $(a \in L^2(\Omega))$

→ convex problem in $a \in L^2(\Omega) \Rightarrow$ but on ∞ dimensional space

→ Tractable: celebrated representer theorem

usually not tractable

\Rightarrow the solution $\hat{a} \in \text{Span}\{\sigma(\langle \alpha_i, \cdot \rangle) : i \leq n\}$ $\in L^2(\Omega)$
 n -dim linear subspace

proof: $a \in L^2(\Omega)$, consider subspace $V = \text{span}\{\sigma(\langle \alpha_i, \cdot \rangle) : i \leq n\}$
 $\subset L^2(\Omega)$

let $a = \underbrace{a_V}_{\in V} + \underbrace{a_\perp}_{\in V^\perp}$ space orthogonal to V

We have $f(\alpha_i, a) = \int \sigma(\langle \alpha_i, \omega \rangle) a(\omega) \mu(d\omega)$
 $= \langle \sigma(\alpha_i, \cdot), a \rangle_{L^2(\mu)} = \langle \sigma(\alpha_i, \cdot), a_V \rangle_{L^2}$

and $\|a\|_{L^2}^2 = \|a_V\|_{L^2}^2 + \|a_\perp\|_{L^2}^2$

$\Rightarrow \hat{a} = \underset{a = a_V + a_\perp}{\text{argmin}} \left\{ \sum_{i=1}^n \underbrace{(y_i - \langle \sigma(\alpha_i, \cdot), a_V \rangle)_{a_V}}^2 + \lambda \underbrace{\|a_V\|_{L^2}^2}_{a_V} + \lambda \underbrace{\|a_\perp\|_{L^2}^2}_{a_\perp} \right\}$

$\Rightarrow \hat{a}_\perp = 0$ hence $\hat{a} \in \text{Span}\{\sigma(\langle \alpha_i, \cdot \rangle) : i \leq n\}$ \square
 $\hat{a} = \hat{a}_V$

Closed form solution: $\hat{a} = \sum_{i=1}^n \underbrace{c_i}_{\in V} \underbrace{\sigma(\langle \alpha_i, \cdot \rangle)}_{\in V}$

$f(\alpha, \hat{a}) = \sum_{i=1}^n c_i K(\alpha, \alpha_i)$ $(K(\alpha, \alpha_i) = \int \sigma(\langle \alpha, \omega \rangle) \sigma(\langle \alpha_i, \omega \rangle) \mu(d\omega))$

Denote $K_n = (K(\alpha_i, \alpha_j))_{i,j \leq n} \in \mathbb{R}^{n \times n}$ \leftarrow

$y = (y_1, \dots, y_n) \in \mathbb{R}^n$ \leftarrow

$\hat{c} = \underset{c \in \mathbb{R}^n}{\text{argmin}} \left\{ \|y - K_n c\|_2^2 + \lambda c^T K c \right\}$ \leftarrow

$\hat{c} = (K_n + \lambda \text{Id})^{-1} y$ $\in \mathbb{R}^n$

What is the performance of KRR?

→ consider $\mathcal{G} = \{f_* \text{ } L\text{-Lipschitz}\}$

$$\sup_{f_* \in \mathcal{G}} R_{\text{test}}(f_*, \hat{f}_\lambda) \asymp \underline{n^{-\frac{1}{d}}}$$

→ for the 'work case', to get error $\leq \varepsilon$
we need $\underline{n \geq \left(\frac{1}{\varepsilon}\right)^d}$

CURSE OF DIMENSIONALITY

KRR is adaptive to smoothness of the function

↳ smoother fcts will be easier to fit

e.g.: previous theorem: to fit degree- λ polynomial

→ need $\underline{n \geq d^\lambda}$ ← not here $\lambda \in [0, C]$
 $\lambda = 0 \quad \lambda \uparrow \quad \text{test error} \uparrow \quad \underline{\lambda \geq C}$

2) $\mathcal{G}_S = \{f_* \text{ with } S \text{ first derivatives bounded}\}$

good kernel → $\sup_{f_* \in \mathcal{G}_S} R_{\text{test}}(f_*, \hat{f}_\lambda) \asymp n^{-\frac{S}{S+d}}$ ← $\underline{\lambda^*}$
← Hölder space

⇒ still need smoothness S to grow with d (curse of dim.)

$$S \propto d$$

Can we hope to do better?

→ No: these classes of fcts are too big ('plague by the curse of dim')

Need to restrict to a smaller class of fct

Interesting class of fcts: $f_{\alpha}(Ux)$ $U \in \mathbb{R}^{d \times d}$ $s \ll d$

→ fcts that only depend on x on a low-dimensional projection

Why?

$$f(x, a) = \int \sigma(\langle x, w \rangle) a(w) \mu(dw) \quad \leftarrow$$

→ put all the weights $a(w)$ on $w \in \text{Im}(U^T)$

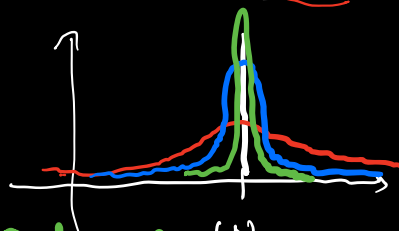
→ problem effectively s -dimensional $d \rightarrow s$
↳ hope to get $\mathcal{O}(n^{-\frac{1}{s}})$ rate.

However: Kernel methods are not adaptive to fcts that depend only on a low-dimensional projection of data

Recall: if $n \propto d^k$, $R_{\text{test}}(f_*, \hat{f}) = \|P_{>2} f_*\|_{L^2}^2 + o_d(1)$
↳ no matter the structure on f_* (e.g. $= \sigma(\langle w_*, x \rangle)$)

Intuition: in order to have $f(a, a) \rightarrow \sigma(\langle \omega_*, a \rangle)$
 we need $a(\omega) \rightarrow \delta_{\omega=\omega_*}$ ($\sigma \notin F_2$)

a density $\| \delta_{\omega=\omega_*} \|_{L^2} = \infty$
 $\sigma(\langle \omega_*, \cdot \rangle)$ not in RKHS, σ_* not in F_2
 we must have $\|a\|_{L^2} \rightarrow \infty$



HOWEVER: $\int_{\Omega} |a(\omega)| \mu(d\omega) =: \|a\|_{L^1}$ remains bounded
 $\hat{a}_n \rightarrow \delta_{\omega, \omega_*}$
 $\|\hat{a}_n\|_{L^1} \leq 1$

\Rightarrow so... why not replacing $\|a\|_{L^2}$ by $\|a\|_{L^1}$

i.e.

$\rightarrow (F_2 - P) \quad \hat{a} = \underset{a: \Omega \rightarrow \mathbb{R}}{\operatorname{argmin}} \left\{ \sum_{i=1}^n (y_i - f(x_i, a))^2 + \lambda \|a\|_{L^2}^2 \right\}$
 (Kernel method)

$\rightarrow (F_1 - P) \quad \hat{a} = \underset{a: \Omega \rightarrow \mathbb{R}}{\operatorname{argmin}} \left\{ \sum_{i=1}^n (y_i - f(x_i, a))^2 + \lambda \|a\|_{L^1} \right\}$

\rightarrow still convex
 (might not be tractable)
 \rightarrow no representer thm

CONVEX NNETS

$f(a, \hat{a}_{F_1})$

If $f_* = g(Ux)$,
 $U \in \mathbb{R}^{n \times d}$

$$R(f_*, \hat{f}_{F_1}) \leq \underline{\underline{\underline{n^{-\frac{1}{d}}}}}$$

⇒ Convn NNets break the curse of dimensionality on fcts that only depend on a low-dim projection of the data

⇒ adaptive to latent linear structure (U is unknown)

However F_1 is not tractable (hard problem)

→ can think about GD as approximately solving F_1

⇒ in general do not expect GD to solve F_1 -problems (not the right implicit bias)

↳ However, one situation where GD was proven to solve approximately F_1 problem

"Implicit bias of GD for wide 2-layers NNets"
— Chizat and Bach, 2020.

Linear regime → F_2 : problem

- ↗ curve of dim
- ↘ adaptive to smoothness
- ↘ not adaptive to low-dim project fcts

Sometimes

Non-linear dynamics → F_1 : problem

- ↗ adaptive to smoothness
- ↘ adaptive to low-dim project
- ⇒ break curse of dim. on these fct classes.

③ An example: learning parities

So far, we saw

- Limitation of kernel methods/linear models
 - Feature learning necessary to break the curse of dimensionality
 - One "classical regime" example: GD fitting a single neuron with another neuron.
- underparametrized* →

More realistic example where we can study feature learning with GD

$$\alpha \sim \text{Unif}(\{\pm 1\}^d)$$

learning class of k -parity fcts

$$\mathcal{C}_k = \left\{ f_A(x) = \prod_{i \in A} x_i : \forall A \subseteq [d], |A|=k \right\}$$

EA
parity of subset A .

Hardness result of learning parity fcts with kernel methods

Prop: [Allen-Zhu et al., 2020] If for any $f_A \in \mathcal{C}_k$

$$R_{\text{test}}(f^*, \hat{f}_\lambda) \leq \frac{1}{9}$$

then we must have $n \geq \frac{3}{4} \binom{d}{k} \quad (\propto d^k)$

Remark: 1) $f_A(x)$ is a degree- k polynomial, already implied by previous result if kernel is an inner-product kernel
then need $n \geq d^k$ samples

→ here very elementary proof for any kernel

2) $f_A(x)$ only depend on a low-dimensional projection of dimension k

→ expect F_1 problem to be able to efficiently learn \mathcal{C}_k

Proof: [If time, probably not: very nice proof using only elementary algebra]

"Learning parities with neural networks"

[Amit Daniely and Eran Melech, 2020]

With slightly different distribution

+ classification setting: $l(\hat{y}, y) = \max(1 - y\hat{y}, 0)$

+ 2 layers NNets with ReLU activations

Thm: for any linear model $\hat{f}(u) = \langle \psi(u), \hat{a} \rangle$ with $\psi(u) \in \mathbb{R}^N$ and $\|\hat{a}\|_2 \leq B$, then there exists $f_A \in \mathcal{C}_k$ such that

$$R_{\text{test}}(f_A, \hat{f}) \geq \frac{1}{2} - \frac{\sqrt{N} B}{2^k \sqrt{2}}$$

Thm: (informal) GD with some initialization and learning steps on population loss, for T steps with high probability, for any $f_A \in \mathcal{C}_k$

$$R_{\text{test}}(f_A, \hat{f}^{(T)}) \lesssim \frac{k^8}{\sqrt{N}} + \frac{Nk}{\sqrt{d}} + \frac{k^2 \sqrt{N}}{T}$$

Rule: 1) $k \propto d^{\frac{1}{64}}$, $N \propto k^{16} \propto d^{\frac{1}{4}}$, $T \propto d^{\frac{3}{4}}$

Then $\exists f_A \in \mathcal{C}_k$ such that

$$R_{\text{test}}(f_A, \hat{f}^{(\text{lin})}) \geq \frac{1}{4} \checkmark$$

$$R_{\text{test}}(f_A, \hat{f}^{(\text{GD})}) \lesssim 0 \checkmark$$

2) Still unsatisfactory: here $m = \infty$ (or very large) + artificial GD learning steps schedule

Proof idea: Initialization

$$(0) \quad f(x, \theta^{(0)}) = \sum_{j=1}^N a_j^{(0)} \sigma(\langle \underline{w_j^{(0)}}, x \rangle) \leftarrow$$

1 large GD step

$$(1) \quad f(x, \theta^{(1)}) = \sum_{j=1}^N a_j^{(1)} \sigma(\langle \underline{w_j^{(1)}}, x \rangle) \leftarrow$$

* learn good weights $w_j^{(1)}$ with large correlation with A

* show that if we fix $w_j^{(1)}$ and only know $a_j^{(1)}$, can fit f_A

(2 \rightarrow T) Following learning steps sufficiently small such that we are in the linear regime

Summary:

(0) Initializing at $a_i^{(0)}, w_i^{(0)}$

GD steps

(1) One gradient step learn good $w_i^{(1)}$

(2-T) Fit second layer $a_i^{(1)}$ while $w_i^{(1)}$ almost fixed

□

①

Goin beyond the Linear regime:
mathematical approaches

1) Higher-order Taylor expansion around initialization



(Dan Roberts et al. 2021)

$$GD \rightarrow K_t \rightarrow K_t \rightarrow K_t^{(2)}$$

2) Can see GD dynamics as kernel dynamics with time varying kernel $K_t(x, y)$



→ can write ODE for K_t

→ hierarchy of ODEs with higher order kernels $K_t^{(k)}$

→ can truncate at some level $K_t^{(k)} = K_0^{(k)}$ is fixed

$$\begin{matrix} K_t \\ \downarrow \\ K_t^{(2)} \\ \downarrow \\ \dots \\ K_t^{(k)} \end{matrix}$$

completely non linear

3) Mean-Field dynamics: $\Theta^{(t)}$ weights after t SGD steps

$$f_N(x, \Theta^{(t)}) = \frac{1}{N} \sum_{i=1}^N a_i^{(t)} \sigma(\langle w_i^{(t)}, x \rangle) = \int a(w) \sigma(\langle w, x \rangle) \hat{\rho}_t(dw)$$

$$\left(\text{where } \hat{\rho}_t(dw) = \frac{1}{N} \sum_{i=1}^N \delta_{w=w_i^{(t)}} \right)$$

$$\longrightarrow f(x, \rho_t) = \int a(w) \sigma(\langle w, x \rangle) \rho_t(dw)$$

PDE on ρ_t : evolution in the space of measures

→ send paper: on Multi-layer MF.

→ Mario Mondelli

generalizat^o error of

$$\underline{f(x, \theta^t)} \rightarrow \theta^t$$

NTK: FC multilayer NN
weights $\sim N(0, Id)$

→ NTK $h(\langle x, y \rangle)$ inner prod

1) learn only low d^o poly on the sphere

2) F_2 -P → not adaptive

→ NTK: KRR with NT kernel

$$\textcircled{1} \quad x \in \mathbb{R}^d \quad x = U z + U_{\perp} z^{\perp}$$

$$x = U z \quad U \in \mathbb{R}^{d \times d} \quad \uparrow \text{low variance}$$

$$h(\langle x_1, x_2 \rangle) = h(\langle z_1, U^T U z_2 \rangle)$$

$$= h(\langle z_1, z_2 \rangle)$$

\rightarrow denoising