

## Lecture 9: Learning Multi-index functions

In this lecture, we consider learning a specific class of functions: multi-index fcts, that is fcts that only depend on the input data through a low-dimensional projection.

The goal is to get more insights on

- When and how can NNs overcome the curse of dimensionality?
- How do they adapt dynamically to low-dimensional support?
- How do they compare to other methods? Optimal algorithms?

Data distribution: "multi-index functions"

$$x \sim \text{Unif}(S^{d-1}(\Gamma d))$$

$$y = f_*(x) + \varepsilon$$

mean 0  
indep noise  
e.g.  $N(0, \sigma^2)$

$$f_*(x) = \underline{g(U_*^T x)}$$

$$g: \mathbb{R}^d \rightarrow \mathbb{R}$$

$$U_* \in \mathbb{R}^{d \times d}$$

$$U_*^T U_* = I_d$$

$$= g(\langle u_1, x \rangle, \dots, \langle u_d, x \rangle)$$

$\Delta$  "indices"

(The support  $U_*$  is unknown)

We are interested in the high-dimensional regime where  $g$  is fixed (in particular  $s$ ) while  $d \rightarrow \infty$ .

Learn this data with a 2-layer NN.

We saw "three paradigms" of learning:

① Classical ERM with explicit regularization

→ assume algo output approximate minimizer of ERM

② Kernel methods

→ NNs trained in lazy regime

③ Feature learning

→ NNs trained non-linearly (e.g. MF regime)

↳ adapt "representation" to data

In these three cases: sample complexity? computational complexity? } → as  $d \rightarrow \infty$

① ERM with infinite-width 2 layer NN

"Breaking the curse of dimensionality with convex NNs"  
- Francis Bach, 2017.

Model:  $f(x) = \frac{1}{M} \sum_{j=1}^M a_j \sigma(\langle w_j, x \rangle + b_j)$

$\sigma(t) = \max(t, 0)$       RELU

Take  $z = (x, \sqrt{d}) \in \mathbb{R}^{d+1}$        $v = (w, \frac{b}{\sqrt{d}}) \in \mathbb{R}^{d+1}$

$\sigma(\langle w_j, x \rangle + b_j) = \sigma(\langle v_j, z \rangle)$

$\mu(dv) = \text{Unif}(S^d)$        $V := S^d$

Lecture 2: set of  $\infty$ -width 2-layer NN

$F_1(\mathbb{R}) = \left\{ f(x) = \int_V \sigma(v^T z) a(v) \mu(dv) : \int_V |a(v)| \mu(dv) \leq R \right\}$

In fact, closure of  $\mathcal{F}_1(R)$  can be written as

$$\overline{\mathcal{F}_1(R)} = \left\{ f(x) = \int_V \sigma(v^T z) \nu(dv) : \|\nu\|_{TV} \leq R \right\}$$

where  $\|\nu\|_{TV} = \nu_+(V) + \nu_-(V)$

signed measure

We saw that the Rademacher complexity

$$\text{Rad}_m(\overline{\mathcal{F}_1(R)}) \leq CR \sqrt{\frac{d}{m}}$$

We can define a functional norm:

$$\|f\|_{\mathcal{F}_1} = \inf \left\{ \|\nu\|_{TV} : f(x) = \int_V \sigma(v^T z) \nu(dv) \right\}$$

**ERM**

$$\hat{f} = \underset{\|f\|_{\mathcal{F}_1} \leq R}{\text{argmin}}$$

$$\hat{R}_m(f) = \frac{1}{m} \sum_{i=1}^m (y_i - f(x_i))^2$$

Test error:

(5)

$$R(\hat{f}) = \mathbb{E} [ (\hat{f}(x) - f_*(x))^2 ]$$

$$\leq \underbrace{\inf_{\|f\|_{\mathcal{F}_1} \leq R} \|f - f_*\|_{L^2}^2}_{\text{Approx error}} + \underbrace{2 \sup_{\|f\|_{\mathcal{F}_1} \leq R} |\hat{R}_n(f) - R(f)|}_{\text{generalized error}}$$

Approximation error:

Prop [Bach, 2017] For  $D \geq 1$  and  $R \geq C_d$  and any  
fct  $\varphi: \mathbb{R}^D \rightarrow \mathbb{R}$  s.t. for all  $\|x\|_2, \|y\|_2 \leq B$   
 $|\varphi(x)| \leq L$       $|\varphi(x) - \varphi(y)| \leq \frac{L}{B} \|x - y\|_2$

There exists  $f$  with  $\|f\|_{\mathcal{F}_1} \leq R$  and

$$\sup_{\|x\|_2 \leq B} |\varphi(x) - f(x)| \leq C_d L \left(\frac{L}{R}\right)^{\frac{2}{D+1}} \log\left(\frac{R}{L}\right)$$

In our case:  $D = d$  & we can use  $L = \mathcal{O}(1)$   
(using tail bounds)

For generalization error: we saw how to bound it from  $Rd_m(F)$  using Talagrand contraction inequality for Lipschitz loss

Here: squared loss.

→ assume  $\|y\|_\infty \leq C$

$$|\beta(x)| = \left| \int \langle v, z \rangle_+ v(dw) \right| \leq C\sqrt{d} \|v\|_{TV}$$

Hence  $|\ell(y, \beta) - \ell(y, \beta')| \leq C(1 + R\sqrt{d}) |\beta - \beta'|$

$$\begin{aligned} \mathbb{E} \left[ \sup_{\|\beta\|_{F_1} \leq R} |\hat{R}_m(\beta) - R(\beta)| \right] &\leq \frac{C}{\sqrt{m}} + C(1 + R\sqrt{d}) Rd_m(F_1(R)) \\ &\leq \frac{C}{\sqrt{m}} + C(1 + R\sqrt{d}) R \sqrt{\frac{d}{m}} \end{aligned}$$

Putting these two together, we obtain:

$$\mathbb{E}[R(\hat{\beta})] \lesssim R^{-\frac{2}{d+1}} \log R + R^2 \frac{d}{\sqrt{m}}$$

→ take  $R_* = \left(\frac{\sqrt{m}}{d}\right)^{\frac{d+1}{2d+4}}$

we deduce

$$\mathbb{E}[R(\hat{\beta})] \lesssim \frac{\sqrt{d}}{m^{\frac{1}{2D+4}}} \log m$$

(7)

As long as  $m \gtrsim d^{\Delta+2}$  test error is vanishing

[Note that this is not an optimized proof]

↳ e.g. with a little bit more care we can get

$$\text{bound} \lesssim \left(\frac{d^{\frac{1+\frac{2}{\Delta+1}}{\Delta+1}}}{m}\right)^{\frac{1}{\Delta+1}} \log^C m \quad (\text{can be probably improved})$$

This learning guarantee is for any  $O(1)$ -Lipschitz fct

$$f_{\star}(\alpha) = g(U_{\star}^{\top} \alpha) \quad U_{\star} \in \mathbb{R}^{d \times \Delta}$$

as long as we have  $\hat{f}$  a good enough approximate minimizer

Is there an efficient way of constructing approximate minimizers over  $\mathcal{F}_1$ ?

1<sup>st</sup> natural approach "random feature approximation"

$$\int \sigma(\langle v, z \rangle) a(v) \mu(dv) \rightarrow \frac{1}{M} \sum_{j=1}^M a_j \sigma(\langle v_j, z \rangle) =: f(x; a)$$

$$v_j \stackrel{iid}{\sim} \mu(dv)$$

$$\hat{f} = \underset{a \in \mathbb{R}^M}{\text{argmin}} \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i; a))^2 + \frac{\lambda}{M} \sum_{j=1}^M |a_j|$$

Approximation lower bound  $d^{k+s} \leq M \leq d^{k+1-s}$

$$R(\hat{f}) \geq \inf_{a \in \mathbb{R}^M} \mathbb{E}[(f_*(x) - f(x; a))^2] = \underbrace{\|P_{\leq k} b_*\|_{L^2}^2}$$

can only fit  $P_{\leq k} g(U_*^\top x)$   
best degree-k polynomial approx

When fitting  $f_*(x) = (\langle w_*, x \rangle + b_*)_+$

$$\|w_*\|_2 = d^2$$
$$|b_*| \leq 6d^4 + 1$$

[Yehudai, Shamir, '19]

Then test error must be  $\geq \frac{1}{50}$  unless  $M = e^{\Omega(d)}$   
or  $\|a\|_1 = e^{\Omega(d)}$

This is because we are trying to fit  $f(\langle w_0, x \rangle)$  with  $\sigma(\langle w_j, x \rangle)$  and in high-dim  $|\langle w_j, w_0 \rangle| \leq \frac{1}{\sqrt{d}}$  for all  $j \in [M]$  w.h.p unless  $M$  superpolynomial in  $d$ .

**2<sup>nd</sup> natural approach** adding one neuron at a time

Rule: [ We can choose  $\hat{f}$  supported on  $(n+1)$  neurons ]

(a type of "representer" theorem)

$$\inf_{\|f\|_{F_n} \leq R} \frac{1}{m} \sum_{i=1}^m (y_i - f(x_i))^2 = \inf \left\{ \frac{1}{m} \sum_{i=1}^m (y_i - \hat{y}_i)^2 : \hat{y} \in \mathbb{R}^m, \gamma_1(\hat{y}) \leq R \right\}$$

where

$$\gamma_1(\hat{y}) = \inf \left\{ \|v\|_{TV} : \int \sigma(\langle v, z_i \rangle) v(dv) = \hat{y}_i \right\}_{i=1, \dots, m}$$

By Carathéodory theorem\*, there optimal function can be decomposed into at most  $(n+1)$   $\sigma(\langle v_j, z \rangle)$  i.e. we can write

[\* if  $x \in \text{Conv}(P) \subset \mathbb{R}^m$  then  $x$  can be written as combinat<sup>n</sup> of  $n+1$  elements of  $P$ ]

$$\hat{f}(x) = \sum_{j=1}^{m+1} a_j \sigma(\langle v_j, x \rangle) \quad \text{with } \sum_{j=1}^{m+1} |a_j| \leq R$$

10

Of course, we do not know these  $\{v_j\}_{j=1}^{m+1}$  in advance

Instead of choosing  $v_j$ 's at random as above, maybe we can choose them one at a time, that makes biggest progress in decreasing the train error.

Algo to minimizing a functional  $\Psi(f)$  over  $f \in L^2(\rho)$   
where  $\rho$  is some probe measure on  $\mathcal{X}$

Conditional gradient algo (Frank-Wolfe algorithm):

Assume  $\Psi$  is convex and  $L$ -smooth and denote  $S\Psi$  its functional gradient

$$0 \leq \Psi(f) - \Psi(h) - \langle f - h, S\Psi \rangle_{L^2(\rho)} \leq \frac{L}{2} \|f - h\|_{L^2(\rho)}^2$$

FW algo: iterative algo

• initialization  $f_0 \in \overline{F}_1(R)$

•  $\overline{f}_t \in \operatorname{argmin}_{f \in \overline{F}_1(R)} \langle f, \delta^2 \Psi(f_t) \rangle_{L^2(\mathcal{C})}$

$$f_{t+1} := (1 - c_t) f_t + c_t \overline{f}_t$$

Note that  $f_t \in \overline{F}_1(R)$  for all  $t$

If we choose  $c_t = \frac{2}{t+1}$ , it is known that we have CV rate:

$$\Psi(f_t) - \inf_{f \in \overline{F}_1(R)} \Psi(f) \leq \frac{2L}{t+1} \sup_{f, h \in \overline{F}_1(R)} \|f - h\|_{L^2(\mathcal{C})}^2$$

Using that  $\|f\|_\infty \leq C \sqrt{d} \|f\|_{F_1}$ , we get

$$\Psi(f_t) - \inf_{f \in \overline{F}_1(R)} \Psi(f) \leq \frac{C d R^2 L}{t+1}$$

So if we can solve the update problem, this will CV to optimum

(12)

$$\langle f, g \rangle_{L^2(\mathcal{E})} = \int_V \left( \int_{\mathcal{X}} \sigma(\langle z, v \rangle) g(x) p(dx) \right) v(dv)$$

$$\leq \|v\|_{TV} \cdot \max_{v \in V} \left| \int_{\mathcal{X}} \sigma(\langle z, v \rangle) g(x) p(dx) \right|$$

Hence  $\max_{\|f\|_{F_1} \leq R} \langle f, g \rangle_{L^2(\mathcal{E})} = R \max_{v \in V} \left| \int_{\mathcal{X}} \sigma(\langle z, v \rangle) g(x) p(dx) \right|$

which is achieved by taking  $\pm R \cdot \sigma(\langle \cdot, v \rangle)$  where  $v$  is a maximizer of RMS.

In our case,  $\mathcal{X} = \{x_1, \dots, x_m\}$  and  $p$  is uniform measure

$$\Psi(f) = \hat{R}_m(f)$$

$$\nabla \Psi(f) = \frac{2}{m} \begin{pmatrix} f(x_1) - y_1 \\ \vdots \\ f(x_m) - y_m \end{pmatrix}$$

Therefore in our case, it reduces to solving the problem

argmax  $\sum_{i=1}^m a_i (v^T z_i)_+$   
 $\|v\|_2 \leq 1$

This is a NP-hard problem even to solve approximately

Remark: This is shown by doing a reduction to the problem of Maximum Agreement for Halfspaces problem

For  $(x_1^{(1)}, \dots, x_{m_+}^{(1)}) (x_1^{(2)}, \dots, x_{m_-}^{(2)}) \in \{\pm 1\}^d$

$$M(w, a) = \sum_{i=1}^{m_+} \mathbb{1}[\langle w, x_i^{(1)} \rangle - a > 0] + \sum_{i=1}^{m_-} \mathbb{1}[\langle w, x_i^{(2)} \rangle - a < 0]$$

Then it is NP-hard to distinguish between the two cases:

- (1)  $\exists w, a$  s.t.  $M(w, a) \geq (1 - \epsilon)(m_+ + m_-)$
- (2) For any  $w, a$ :  $M(w, a) \leq (\frac{1}{2} + \epsilon)(m_+ + m_-)$

classified at random

Remark: Does that mean FW algo will not work?

(1) We do not need to solve exactly the update problem possible there is a convex relaxation

[Open problem]

(2) NP-hard  $\rightarrow$  but might still be tractable on random data [Open]

Two natural approaches failed. Is  $F_1$ -minimization intrinsically hard? (i.e., we can't hope to have a polynomial time algo that is guaranteed to minimize ERM problem with  $F_1$ )

See next lecture!

\* Optimization hardness: getting  $O(n^c)$  approximate minimizer is NP-hard

[Celentano, M., Montanari, 2021]

\* Learning hardness: if we could solve  $F_1$ -ERM then we could learn intersection of half spaces

(under strong random CSP assumption, this is impossible)

[Neyshabur, Tomioka, Srebro, 2015]

(More next lecture!)

Summary

$F_1$ -ERM breaks the curse of dim when learning any Lipschitz multi-index fct in sample complexity. However, no polynomial time algo to construct  $\hat{f}$

# 2 Kernel methods

In the lazy regime, 2-layer NNs behave as kernel method with inner-product kernel.

(For simplicity, we take all biases  $b_j = 0$ )

We saw in lecture 6 that if  $d^{k+s} \leq n \leq d^{k+1-s}$

$$R(\hat{f}) = \| P_{>k} f \|_{L^2}^2 \quad \left. \vphantom{\| P_{>k} f \|_{L^2}^2} \right\} \text{only fit } d^k \text{ polynomial approx}$$

More generally, using the dimension lower bound, when learning multi-index model, no kernel method can do better with <sup>1</sup> data uniform on the sphere

Rank: [Learning with  $F_1$  vs  $F_2$  regularization]

$$F_1(R) = \left\{ f(x, a) = \int \sigma(\langle w, x \rangle) a(w) \mu(dw) : \int |a(w)| \mu(dw) \leq R \right\}$$

$$F_2(R) = \left\{ f(x, a) = \int \sigma(\langle w, x \rangle) a(w) \mu(dw) : \left( \int |a(w)|^2 \mu(dw) \right)^{\frac{1}{2}} \leq R \right\}$$

*uniform on the sphere*

Similarly to  $F_1$ , we can define a functional norm

$$\|f\|_{F_2} = \inf \left\{ \|a\|_{L^2(\mu)} : f = \int \sigma(\langle \omega, \alpha \rangle) a(\omega) \mu(d\omega) \right\}$$

Then 
$$\min_f \hat{R}_n(f) + \lambda \|f\|_{F_2}$$

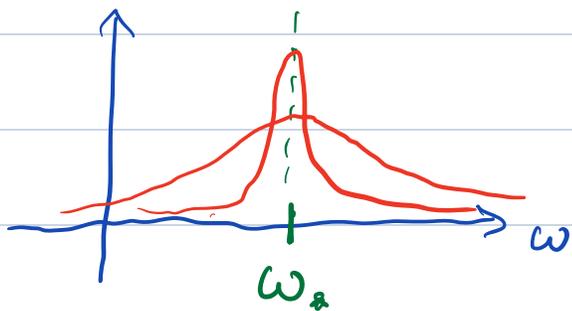
is a kernel method associated to the kernel

$$K(\alpha_1, \alpha_2) = \int \sigma(\langle \omega, \alpha_1 \rangle) \sigma(\langle \omega, \alpha_2 \rangle) \mu(d\omega)$$

and  $F_2$  is an RKHS with RKHS-norm  $\|f\|_{F_2}$

So: when  $+ \lambda \|f\|_{F_1} \rightarrow \text{learn}$   
 $+ \lambda \|f\|_{F_2} \rightarrow \text{fail to learn}$  } multi-index fcts

Some heuristic intuition: if we want to learn  $\sigma(\langle \omega, \alpha \rangle)$   
 then to get a good approx using  $a(\omega) \mu(d\omega) \rightarrow \delta_{\omega_*}$



sequence of  $\overset{\text{positive}}{a_k(\omega)} \mu(d\omega) \rightarrow \delta_{\omega_*}$

s.t.  $\int a_k(\omega) \mu(d\omega) = 1$

but  $\int |a_k(\omega)|^2 \mu(d\omega) \nearrow \infty$

We can approximate  $S_{w_a}$  with  $\|a_k\|_{L^1} = 1$

but we must have  $\|a_k\|_{L^2} \rightarrow \infty$

worse & worse generalization error

### 3 Feature learning

Consider:  $f(x; \Theta) = \sum_{j=1}^M a_j \sigma(\langle w_j, x \rangle + b_j)$  ← ReLU

and consider training  $(a_j, w_j, b_j)_{j=1}^M$  using GD

on squared loss  $\hat{R}(\Theta) = \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i; \Theta))^2$

We follow [Dannan, Lee, Sotkinolkotabi, 2022]

(paper is for Gaussian data, but can be modified for spherical data)

Assumption [Data]  $x_i \sim N(0, Id)$   $y_i = f_*(x_i) + \epsilon_i$   
 $\epsilon_i \sim \text{Unif}(\pm 3)$

- $f_*(x) = g(\langle u_1, x \rangle, \dots, \langle u_s, x \rangle)$ 
  - $\rightarrow d^0 p$  polynomial
  - $\rightarrow$  orthogonal to  $d^1$  polynomial in  $x$
  - can be learned by pretraining
- $M = \mathbb{E}_x [\nabla^2 f_*]$  has rank  $s$   $\kappa := \frac{1}{\sqrt{\Delta} \sigma_{\min}(M)}$

Algo: 1) Initialization  $a_j^0 = -a_{M-j}^0 \sim \text{Unif}(\pm 1)$   $\omega_j^0 = \omega_{M-j}^0 \sim N(0, \frac{Id}{\lambda})$   
 $b_j^0 = 0$   
 $\rightarrow$  symmetric initialization s.t.  $f(x, \Theta^0) = 0$

2) One gradient step on  $W$   $l_2$  regularization  

$$\begin{cases} W^{(1)} = W^{(0)} - \eta_1 [\nabla_W \hat{R}(\Theta^0) + \lambda_1 W^{(0)}] \\ a^{(1)} = a^{(0)} \end{cases}$$

3) Re-initializing  $b_j \sim N(0, 1)$   
 Train 2nd layer weights for  $t=2, \dots, T$   

$$\begin{cases} a^{(t)} = a^{(t-1)} - \eta_2 [\nabla_a \hat{R}(\Theta^{(t)}) + \lambda_2 a^{(t-1)}] \\ W^{(t)} = W^{(t-1)} \end{cases}$$

Hyperparameters:  $\eta_1 = \tilde{O}(\sqrt{d})$      $\lambda_1 = \frac{1}{\eta_1}$

Thm: Assume  $n = \tilde{\Omega}(d^2 k^2 s)$  and  $d = \tilde{\Omega}(k s^{\frac{3}{2}})$

Then there exists  $\eta_2$  and  $\lambda_2$  sufficiently small,

$$T = \tilde{\Omega}\left(\frac{1}{\eta_2 d_2}\right) \text{ s.t. } w.p. \geq 0.99$$

$$\mathbb{E}[|f(x; \hat{\omega}^{(T)}) - y|] - 3 \lesssim \sqrt{\frac{d s^p k^{2p}}{n}} + \sqrt{\frac{s^p k^{2p}}{M}} + \frac{1}{n^{\frac{1}{4}}}$$

Proof: Learn with  $n = \tilde{\Omega}_d(d^2)$  and  $M = \tilde{\Omega}_d(1)$

Proof idea

$$\nabla_{\omega_j} \mathbb{E}[(y - f(x; \omega^0))^2] = -2 \mathbb{E}[f_*(x) \nabla_{\omega_j} f(x; \omega^0)]$$

$$= -2 a_j \mathbb{E}[f_*(x) \mathbb{1}_{\omega_j, x \geq 0}]$$

$$\approx -\frac{2}{\omega_j} a_j M \omega_j + \tilde{O}\left(\frac{\sqrt{s}}{d}\right)$$

For  $n = \tilde{\Omega}(d^2)$  population  $\approx$  empirical gradient

$$W^{(1)} = W^{(0)} - \eta_1 (\nabla_W \hat{R}(\Theta^{(0)})) + \frac{1}{\eta_1} W^{(0)}$$

$$= -\eta_1 \nabla_W \hat{R}(\Theta^{(0)})$$

$$\approx -\eta_1 \underbrace{M W^{(0)}}_{\text{the } \omega_j^{(1)} \text{ are in the span of } U \in \mathbb{R}^{d \times s}} \text{diag}(a_j)$$

the  $\omega_j^{(1)}$  are in the span of  $U \in \mathbb{R}^{d \times s}$

$M W^{(0)}$  has singular value  $\approx \frac{1}{\sqrt{d}}$

and we take  $\eta_1 = \mathcal{O}(\sqrt{d})$

Hence after this first step, we fix  $W^{(1)}$  and get

$$f(x; a) = \sum_j a_j \sigma(\underbrace{\langle \omega_j, x \rangle}_{s\text{-dimensional span}} + b_j)$$

s-dimensional span

↳ then optimizing over  $a_j$ : kernel methods with dim-s input

↳ can use approx + generalization type analysis

(but for SGD)

does not depend on d

With this specific choice of GD:  $W^{(H)}$  align with the support  $U$  after one step. Learn a good "feature represent<sup>o</sup>" with  $\sigma(\langle w_j, x \rangle)$   $w_j \in \text{span}(U)$

→ once learned the represent<sup>o</sup>, problem becomes low-dimensional and learning is easy, e.g., linear regression on 2<sup>nd</sup> layer weights

Shortcomings of this analysis:

•  $\mathbb{E}[\nabla^2 g]$  full rank : is it necessary?  
(e.g. not verified by  $\text{He}_3(\langle u_1, x \rangle)$ )

•  $n = \Omega(d^2)$  is it necessary?

↳ e.g.  $g(\langle u_1, x \rangle) = \langle u_1, x \rangle$  only need  $n = \Theta(d)$

Can we get a more precise picture of how GD/SGD learn the low-dimensional support?

# The staircase mechanism

[Albe, Boix, Advers, M., '22, '23]

For simplicity, consider  $x \sim \text{Unif}(\{\pm 1\}^d)$

Remark: [Functions on the hypercube]

The Fourier-Walsh basis is an orthogonal basis of  $L^2(\{\pm 1\}^d)$

$$\chi_S(x) = \prod_{i \in S} x_i \quad S \subseteq [d] \quad \text{with } \chi_\emptyset(x) = 1$$

Indeed, we have

"parity fct on subset  $S$ "

$$\mathbb{E}[\chi_S(x) \chi_{S'}(x)] = \mathbb{E}[\chi_{S \Delta S'}(x)] = \begin{cases} 0 & \text{if } S \Delta S' \neq \emptyset \\ 1 & \text{if } S \Delta S' = \emptyset \end{cases}$$

So any function  $f: \{\pm 1\}^d \rightarrow \mathbb{R}$  can be decomposed as

$$f(x) = \sum_{S \subseteq [d]} \hat{f}(S) \chi_S(x) \quad \hat{f}(S) = \mathbb{E}_x[f(x) \chi_S(x)]$$

We consider learning fcts with sparse support:

$$y_i = f_\theta(x_i) + \varepsilon_i$$

where  $f_{\star}(x) = g(z_1, z_2, \dots, z_n)$

with  $(z_1, \dots, z_n) = (\underbrace{x_{i_1}, x_{i_2}, \dots, x_{i_s}}_{\text{unknown subset of } s \text{ coordinates}})$

unknown subset of  $s$  coordinates

From above:

$$g(z) = \sum_{S \subseteq [n]} \hat{g}(S) \prod_{i \in S} z_i$$

Motivating examples:

$$\begin{cases} g_1(z) = z_1 + z_1 z_2 + z_1 z_2 z_3 \\ g_2(z) = z_1 z_2 z_3 \end{cases}$$

Q: Which one is easier to learn with GD-trained NNs?

Rank:

- Both are simple fcts that only depend on 3 coordinates and can be approximated easily by a 2-layer NN
- $\text{rank}(\mathbb{E}[\nabla^2 g_1(z)]) = 2 < 3$   
 $\text{rank}(\mathbb{E}[\nabla^2 g_2(z)]) = 0 < 3$

} both don't satisfy  
Denman et al condition  
(let's pretend  $z_i$  are iid  $N(0,1)$   
→ theory apply to this case too)

To build intuition, we consider population gradients:

$$R(\Theta) = \mathbb{E}_x [(y - f(x; \Theta))^2]$$

$$f(x; \Theta) = \sum_{j=1}^M a_j \sigma(\langle \omega_j, x \rangle)$$

From earlier computation, for NNs to learn the target function, the weights  $W$  need to align with the support of  $f_*$ .

Let's freeze  $a_j$  and consider a few gradient step on  $\omega_j$ :

$$\omega_j^{(t+1)} = \omega_j^{(t)} - \eta \nabla_{\omega_j} R(\Theta^{(t)})$$

$$= \omega_j^{(t)} + \eta a_j \mathbb{E} [ f_*(x) \sigma'(\langle \omega_j, x \rangle) x ]$$

$$- \eta a_j \underbrace{\mathbb{E} [ f(x; \Theta^{(t)}) \sigma'(\langle \omega_j, x \rangle) x ]}$$

NN self interaction term that do not contain signal,  $f(x; \Theta^0) \approx 0$

→ at the beginning of the dynamics  $\approx 0$

and we will neglect it

Let's set  $\eta a_j = 1$  for simplicity

At initialization, consider  $\omega_j \sim \text{Unif}(\xi \pm \frac{1}{\sqrt{d}})^d$

(could be Gaussian too)

Lemma: At initialization

$=: \mu_{|S|+1}(\sigma)$  "Hermite coeff"

$$\mathbb{E}[\sigma(\langle \omega_j^0, x \rangle) \prod_{i \in S} x_i] \approx \mathbb{E}[\sigma^{(|S|+1)}(G)] \prod_{i \in S} \langle \omega_j^0 \rangle_i$$

Proof sketch:

$$\mathbb{E}_{\alpha_1}[\sigma(\langle \omega, x \rangle) \alpha_1] = \frac{1}{2} [\sigma(\omega_1 + \langle \omega_{-1}, \alpha_{-1} \rangle) - \sigma(-\omega_1 + \langle \omega_{-1}, \alpha_{-1} \rangle)]$$

$$\approx \omega_1 \sigma''(\langle \omega_{-1}, \alpha_{-1} \rangle)$$

$|\omega_1| = \frac{1}{\sqrt{d}}$  is small

$$\frac{1}{\sqrt{d}} \sum_{j=1}^d u_j \Rightarrow N(0,1) \quad \square$$

Let's use this lemma to get a heuristic picture of what happens at the beginning of the dynamics when learning either  $g_1$  or  $g_2$

Assume that  $z = (x_1, x_2, x_3)$

① Learning  $g_1(z) = x_1 + x_1 x_2 + x_1 x_2 x_3$

1<sup>st</sup> GD step:  $\omega^{(1)} = \omega^{(0)} + \mathbb{E} \left[ g_1(z) \sigma'(\langle \omega^{(0)}, x \rangle) \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_1 \\ \vdots \end{pmatrix} \right]$

$$\omega_1^{(1)} = \omega_1^{(0)} + \mathbb{E} \left[ \sigma'(\langle \omega^{(0)}, x \rangle) g_1(z) x_1 \right]$$

$$= \omega_1^{(0)} + \mathbb{E} \left[ \sigma'(\langle \omega^{(0)}, x \rangle) (1 + x_2 + x_2 x_3) \right]$$

$$= \omega_1^{(0)} + \underbrace{\mu_1}_{\approx \frac{1}{\sqrt{d}}} + \underbrace{\mu_2 \omega_2^{(0)}}_{\approx \frac{1}{\sqrt{d}}} + \underbrace{\mu_3 \omega_2^{(0)} \omega_3^{(0)}}_{\approx \frac{1}{d}}$$

$$\boxed{\omega_1^{(1)} \approx 1}$$

similarly  $\omega_2^{(1)} = \omega_2^{(0)} + \mathbb{E} \left[ \sigma'(\langle \omega^{(0)}, x \rangle) (x_2 + x_1 + x_1 x_3) \right]$

$$\boxed{\omega_2^{(1)} \approx \frac{1}{\sqrt{d}}}$$

and  $\boxed{\omega_3^{(1)} \approx \frac{1}{\sqrt{d}}}$   $\boxed{\omega_j^{(1)} \approx \frac{1}{\sqrt{d}}$  for  $j \geq 4$

After 1 step,  $\omega^{(1)}$  align with 1<sup>st</sup> coordinate

2<sup>nd</sup> GD step:

$$\begin{aligned} \omega_2^{(2)} &= \omega_2^{(1)} + \mathbb{E}[\sigma'(k\omega^{(1)}, \alpha) (\alpha_2 + \alpha_1 + \alpha_1 \alpha_3)] \\ &= \underbrace{\omega_2^{(1)}}_{\approx \frac{1}{\sqrt{d}}} + \underbrace{\mu_2 \omega_2^{(1)}}_{\approx \frac{1}{\sqrt{d}}} + \underbrace{\mu_2 \omega_1^{(1)}}_{= 1} + \underbrace{\mu_3 \omega_1^{(1)} \omega_3^{(1)}}_{\approx \frac{1}{\sqrt{d}}} \end{aligned}$$

$$\boxed{\omega_2^{(2)} \approx 1}$$

and  $\boxed{\omega_3^{(2)} \approx \frac{1}{\sqrt{d}}}$   $\boxed{\omega_j^{(2)} \approx \frac{1}{\sqrt{d}}$   $j \geq 4$

3<sup>rd</sup> GD step:  $\omega_3^{(3)} = \omega_3^{(2)} + \mathbb{E}[\sigma'(k\omega^{(2)}, \alpha) \alpha_1 \alpha_2]$

$$= \underbrace{\omega_3^{(2)}}_{\approx \frac{1}{\sqrt{d}}} + \underbrace{\mu_3 \omega_1^{(2)} \omega_2^{(2)}}_{= 1}$$

$$\boxed{\omega_3^{(3)} \approx 1}$$

GD sequentially align to the support, by using intermediary monomials to learn one coordinate at a time

$n = \Theta(d)$  samples is enough so that we can replace population by empirical gradient and have similar computation

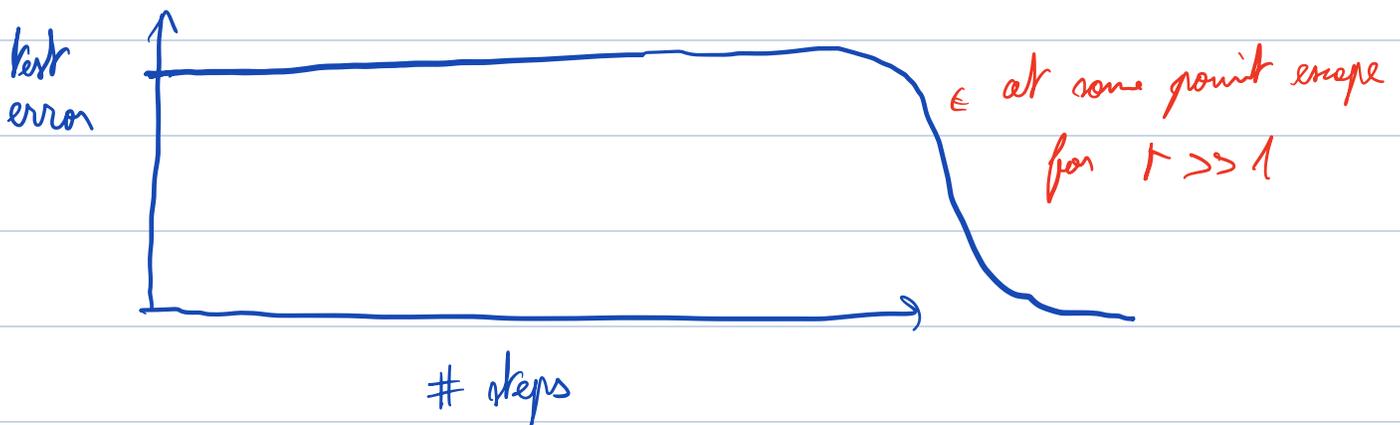
② Learning  $g_2(z) = \alpha_1 \alpha_2 \alpha_3$

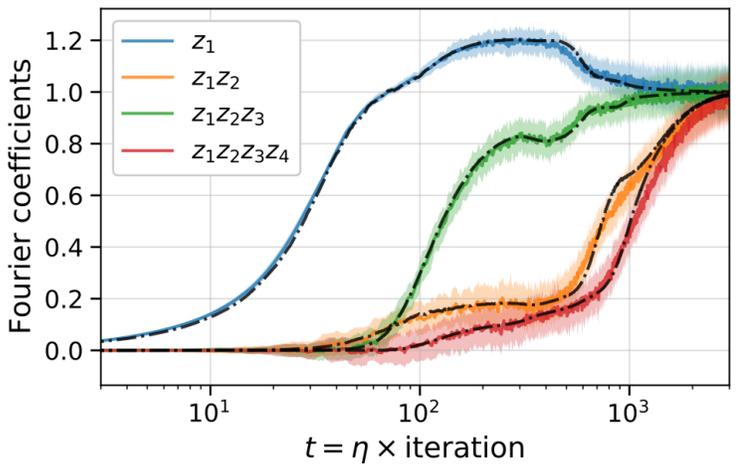
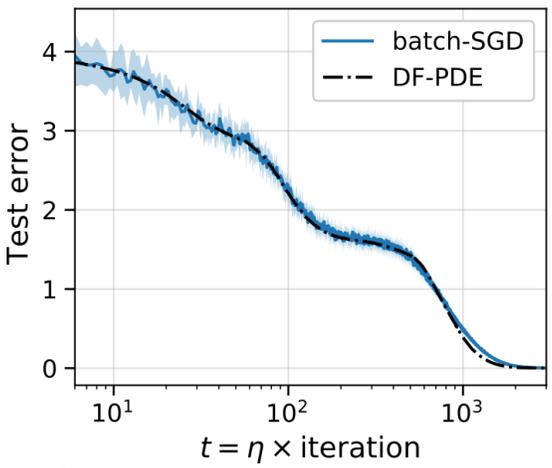
In that case, we don't have the intermediary "stairs" (monomials)

$$\omega_1^{(1)} = \omega_1^{(0)} + \underbrace{\mathbb{E}[\sigma'(\langle \omega^{(0)}, x \rangle) \alpha_2 \alpha_3]}_{\approx \frac{1}{d}}$$

same for  $\omega_2^{(1)}$  and  $\omega_3^{(1)} \approx \frac{1}{\sqrt{d}}$  and stay small for any constant # of steps

⇒ The dynamics is initialized near a saddle pt ( $\nabla_{\omega} R(\omega^0) \approx 0$ ) and takes many more steps to escape



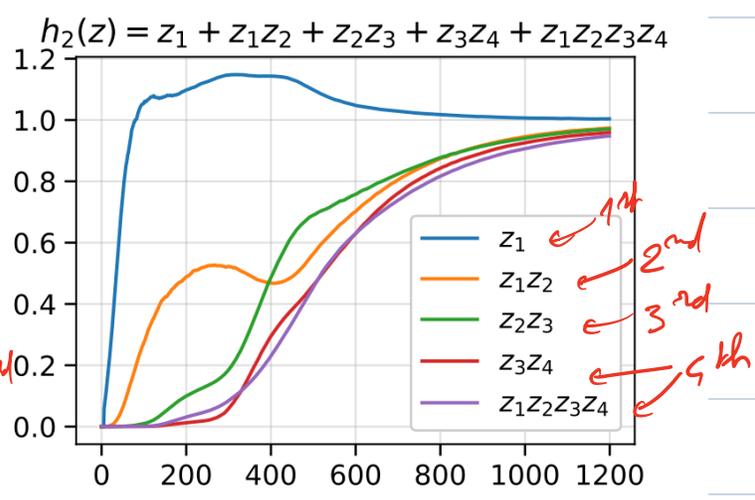
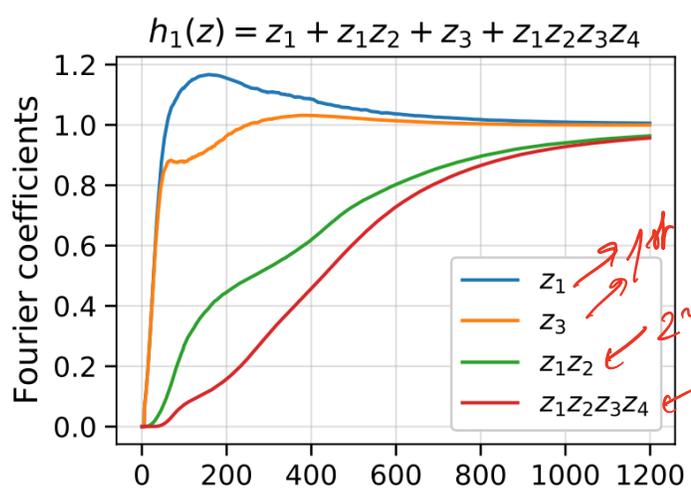


this is online SGD

these are the Fourier coefficients of the NN throughout the dynamics, i.e.

$$\mathbb{E} \left[ \beta_{NN}(\alpha; \Theta^t) \prod_{i \in S} \alpha_i \right]$$

Other examples:



We can show that  $T = \mathcal{O}_d(d)$  online SGD iterations

are enough to learn these functions, hence  $n = \mathcal{O}_d(d)$  samples

Def: [Merged - Staircase fcts]

$$g(z) = \sum_{S \subseteq [n]} \hat{g}(S) \prod_{i \in S} z_i$$

Denote  $\Sigma = \{S \subseteq [n] : \hat{g}(S) \neq 0\}$ . We say that  $\Sigma$  has MS property, if there exists an ordering

$$\Sigma = \{S_1, S_2, S_3, \dots, S_n\}$$

such that

$$|S_i \setminus \bigcup_{j=1}^{i-1} S_j| \leq 1 \quad i \in [1, \dots, n]$$

that is we can order the non zero monomials such that we add at most one coordinate at a time

Following, the heuristic argument from above, we expect these fct to be learnable in  $\tilde{O}(d)$  samples

Thm: MSP fcts are learnable with  $\tilde{O}(d)$  online SGD steps, except for a Lebesgue measure-0 set of them

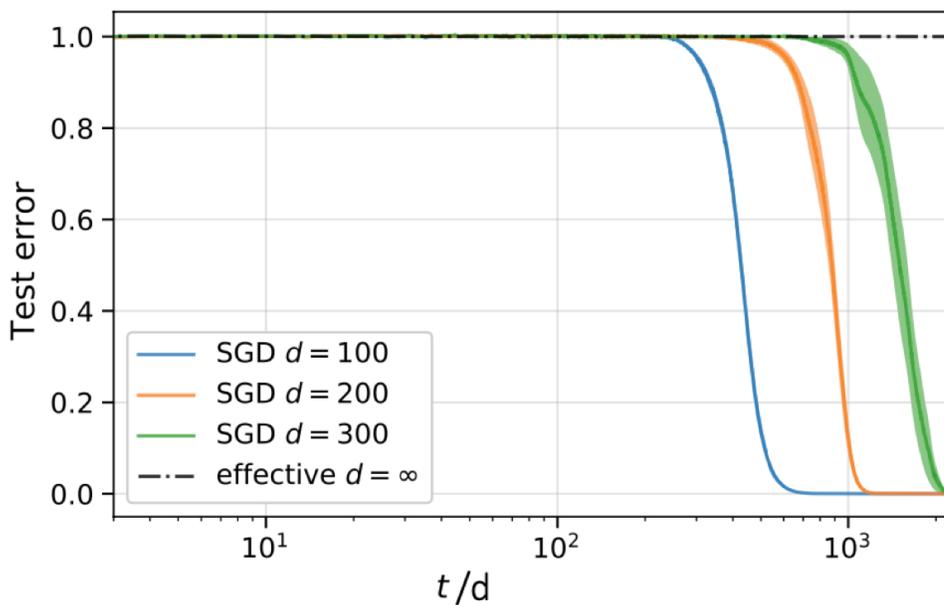
Conversely: if  $g$  is not MSP, then NNs trained in the mean-field regime will not learn  $g$  with  $T = O(1)$  continuous time

with corresponds to  $\mathcal{O}(d)$  discrete SGD steps

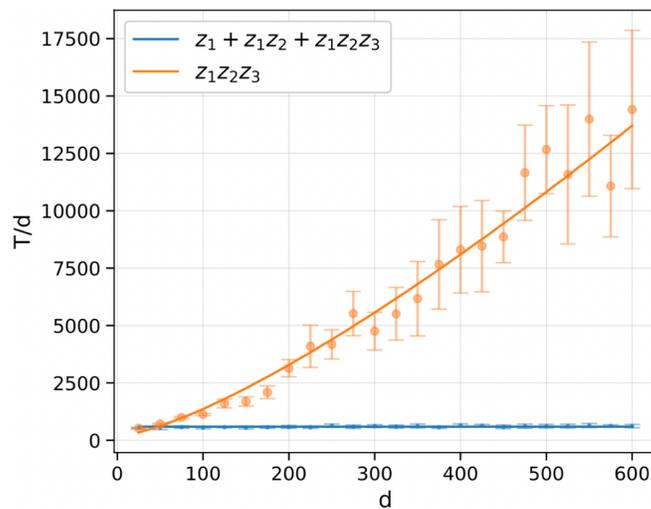
31

What happens for non-MSP functions?

$$h_{*,2}(z) = z_1 z_2 z_3 .$$



$T$  = number of SGD steps to reach 0.05 test error.



$$\underbrace{h_{*,1}(z) = z_1 + z_1 z_2 + z_1 z_2 z_3,}_{T=n=\mathcal{O}(d) \text{ SGD steps to learn}}$$

$$\underbrace{h_{*,2}(z) = z_1 z_2 z_3}_{\text{needs } T = n = \tilde{\mathcal{O}}(d^2) \text{ steps}}$$

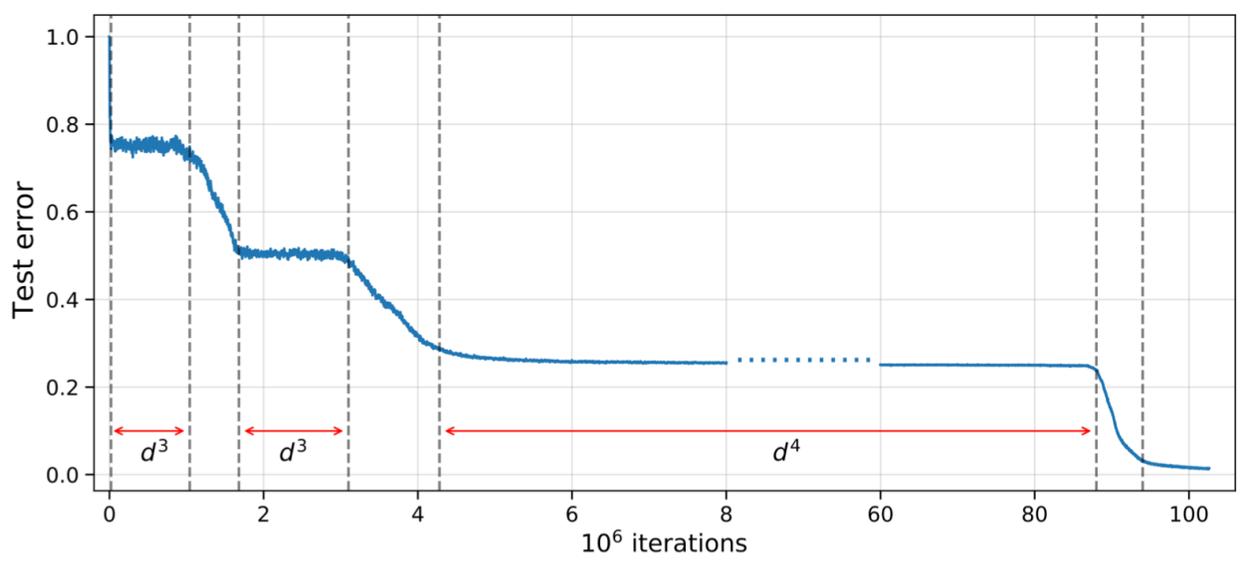
E.g.  $g_2(z) = z_1 z_2 z_3 \dots z_k$



SGD requires  $\tilde{O}(d^{k-1})$  step size to escape the saddle at initialization and align  $w$ 's with the  $k$  coordinates in the support  
 → as soon  $w$ 's aligned with the support, NNs can quickly fit this set using 2<sup>nd</sup> layer weights

More generally:

$h_*(z) = z_1 + z_1 z_2 \dots z_5 + z_1 z_2 \dots z_9 + z_1 z_2 \dots z_{14} .$



SGD sequentially aligns the weights with the sparse support with saddle - to saddle dynamics, and require

$\tilde{\Theta}(d^{k-1})$  SGD steps to align with  $k \geq 2$  new coordinates at once

def: [Leap-complexity]  $g$  is a leap- $k$  fct if it is the smallest integer such that we can order the non-zero monomials s.t

$$|S_i \setminus \bigcup_{j=1}^{i-1} S_j| \leq k$$

(MSP  $\equiv$  leap-1 fcts)

with squared loss

Conjectural picture: online SGD learns  $g$  with # steps

{	$\Theta(d)$	iff	$g$ leap-1
	$\tilde{\Theta}(d \log^c d)$	iff	$g$ leap-2
	$\tilde{\Theta}(d^{k-1} \log^c d)$	iff	$g$ leap- $k$

Summary:

\* Functions of low-dimensional projection can be represented efficiently by 2-layer NNs

\* If we could solve ERM efficiently, we could learn these fcts with  $\text{poly}(d)$  sample complexity (regardless of leap complexity of the fct)

↳ However, evidence that there cannot be an algo that provably solves  $F_1$ -ERM in  $\text{poly}(d)$  time

\* Kernel methods cannot exploit this low-dimensional structure

\* In contrast, NNs trained non-linearly:  $W$  align with the support, learn "good representation", and dynamics become low dimensional.

↳ learning is sequential, following a staircase

⇒ if function has high leap complexity, GD fail to learn it