

## Lecture 5:

## Lazy regime

1

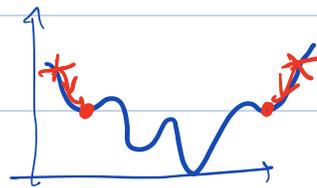
[Jacot, Gabriel, Hongler, 2018]

[Du et al. 2019] [Chizat, Bach, 2019] [Oymak et al. 2020] ...

To train NNs, need to solve a highly non-convex problem:

$$\min_{W_L \dots W_1} \frac{1}{n} \sum_{i=1}^n (y_i - W_L \circ \sigma \circ \dots \circ \sigma W_1 x_i)^2$$

Typically, expect to be hard



get trapped in local minima

→ to the point where most of the ML community abandoned NNs in the 2000s (until AlexNet 2012)

## Empirical observations:

1) as # parameters  $\uparrow$ , easier to optimize:  
 landscape appears to simplify greatly, enabling local search algo (e.g. SGD) to find global optima reliably.

"tractability via overparametrization"

2) Overparametrization does not harm generalization  
 (Lecture 3 & 4)

## This lecture:

\* Example of such a mechanism of tractability via overparametrization

"Lazy" or "linear" regime

\* General mechanism that has nothing to do a priori with neural nets

\* In this regime: NN behave as a linear model

→ neural tangent model

→ Neural tangent kernel in the limit

# Lazy optimization regime

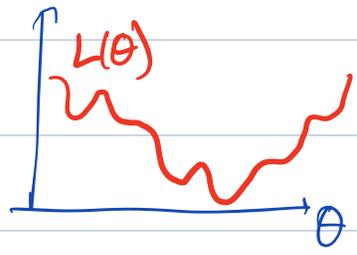
High-level intuition: very simple idea

1 data point  $(y, x)$   $y \in \mathbb{R}, x \in \mathbb{R}$

1-param model  $f(x; \theta)$   $\theta \in \mathbb{R}$

$$L(\theta) = (y - f(x; \theta))^2$$

→ might be very complicated



Intuition dynamics at  $\theta^0$

Idea: we can always rescale output of our model  $\alpha f(x; \theta)$   $\alpha > 0$  and  $\alpha \rightarrow \infty$  to linearizing the model and make the problem quadratic on a small neighborhood of  $\theta_0$ .

→ strongly convex → optimization is easy

2<sup>nd</sup> order Taylor expansion:

→ assume  $\neq 0$

$$f(x; \theta) = f(x; \theta_0) + f'(x; \theta_0)(\theta - \theta_0) + \frac{f''(\theta_0)}{2} (\theta - \theta_0)^2$$

$$L_\alpha(\theta) = (y - \alpha f(x; \theta))^2$$

4

For simplicity:  $v_\alpha := \alpha f(x; \theta_0) - y$

$$b_0' := f'(x; \theta_0)$$

$$b'' := \frac{1}{2} f''(x; \theta_0)$$

$$u := \theta - \theta_0$$

$$L_\alpha(\theta) = (v_\alpha + \alpha b_0' u + \alpha b'' u^2)^2$$

$$= v_\alpha^2 + 2v_\alpha b_0' (\alpha u) + (b_0')^2 (\alpha u)^2$$

$$+ 2v_\alpha b'' (\alpha u^2) + 2b_0' b'' (\alpha^2 u^3)$$

$$+ (b'')^2 (\alpha^2 u^4)$$

Assume  $f(x; \theta_0) \ll y$ , in fact  $f(x; \theta_0) = 0$

so that  $v_\alpha = v_0$

Consider a small neighborhood  $\theta - \theta_0 \lesssim \frac{1}{\alpha}$

$$\tilde{u} = \frac{\theta - \theta_0}{\alpha}$$

$$L_\alpha(\theta_0 + \alpha \tilde{u}) = v_0^2 + 2v_0 b_0' \tilde{u} + (b_0')^2 \tilde{u}^2$$

$$+ \frac{2v_0 b''}{\alpha} \tilde{u}^2 + \frac{2b_0' b''}{\alpha} \tilde{u}^3 + \frac{(b'')^2}{\alpha^2} \tilde{u}^4$$

$$\alpha \rightarrow \infty \quad L_\alpha(\theta_0 + \alpha \tilde{u}) \approx (v_0 + b_0' \tilde{u})^2$$

→ quadratic model:  $\Theta_t = \text{gradient flow}$

$$1) \quad L_\alpha(\Theta_t) \approx L_\alpha(\Theta_0) e^{-\frac{(b_\alpha)^2}{2} t}$$

↳ we started from a very complicated optimization landscape and proved convergence to global optima  $L_\alpha(\hat{\Theta}) = 0$  exponentially fast

2) Throughout the dynamics:

$$\alpha f(x; \Theta^t) \approx \alpha f(x; \Theta^0) + \alpha f'(x; \Theta^0) (\Theta^t - \Theta_0)$$

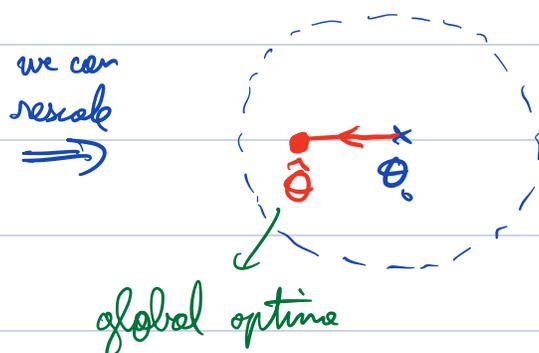
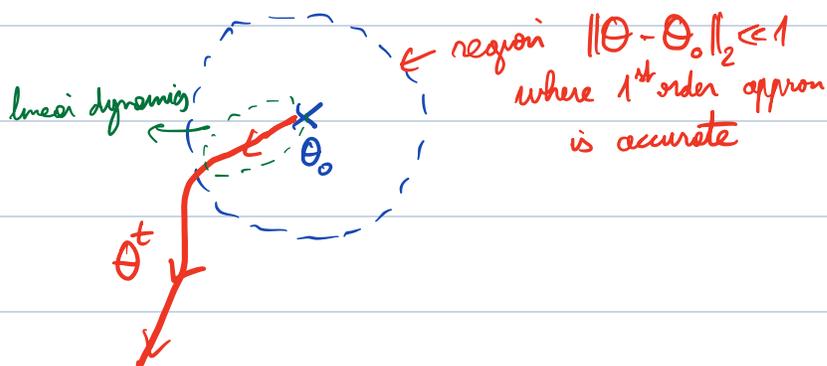
↳ behave as linear model (linear in parameter  $\Theta$ )

→ we can replace our non-linear model by a linear model during training & at test time

→ behave effectively as a linear model

Remark: Another way of thinking about it:

$$f(x; \Theta^t) = f(x; \Theta_0) + f'(x; \Theta_0) (\Theta^t - \Theta_0) + O(|\Theta^t - \Theta_0|^2)$$



3) A consequence is that  $\|\theta_t - \theta_0\|_2 \ll \theta_0$   
 $\rightarrow$  weights barely move

[Chizat, Bach, 2019] "Lazy regime"

Several papers showed global CV before (2018)  
 but 2) & 3) (weights barely moving + linearization of the NN)  
 were kind of hidden in the proof  
 $\Rightarrow$  there have important consequences on learning in this  
 regime  $\rightarrow$  essentially linear regression  
 which temper the achievement of showing global CV  
 and whether this is a good model for what happens  
 with NNs

**Next** precise conditions for global convergence

A meta theorem for global CV in lazy regime

$(y_i, x_i)_{i \leq m}$  data

param model  $f(x_i; \theta) \quad \theta \in \mathbb{R}^p$

$$\hat{R}(\theta) = \frac{1}{2m} \sum_{i=1}^m (y_i - f(x_i; \theta))^2 = \frac{1}{2m} \|y - f_m(\theta)\|_2^2$$

$$f_m: \mathbb{R}^p \rightarrow \mathbb{R}^m \quad f_m(\theta) = \begin{pmatrix} f(x_{1i}; \theta) \\ \vdots \\ f(x_{mi}; \theta) \end{pmatrix} \in \mathbb{R}^m$$

Denote  $\Phi_m(\theta) := D f_m(\theta) \in \mathbb{R}^{m \times p}$  the Jacobian

$$L_{m,0} := \sup_{\theta \neq \theta_0} \frac{\|\Phi(\theta) - \Phi(\theta_0)\|_{\text{op}}}{\|\theta - \theta_0\|_2}$$

Let's do a warm-up:

Condition so that there exists an interpolating solution in a small neighborhood of  $\theta_0$ .

We want to find  $f_m(\theta) = y$

By Taylor's expansion:

$$\begin{aligned}
 y - f_m(\theta_0) &= f_m(\theta) - f_m(\theta_0) \\
 (*) \quad &= \Phi(\theta_0) (\theta - \theta_0) + \underbrace{\int_0^1 (\Phi(\theta_t) - \Phi(\theta_0)) (\theta - \theta_0) dt}_{e(\theta)}
 \end{aligned}$$

$\theta_t = (1-t)\theta_0 + t\theta$

Note that:

$$\|e(\theta)\|_2 \leq L_{m,0} \|\theta - \theta_0\|_2^2$$

Rearranging the terms in (\*)

$$\theta = \theta_0 + \Phi(\theta_0)^+ (y - f_m(\theta_0) + \delta)$$

where  $\delta = - \underbrace{e(\theta_0 + \Phi(\theta_0)^+ (y - f_m(\theta_0) + \delta))}_{=: F(\delta)}$

(fixed pt equation)

9

$S = F(S)$  When does such a solution exist?

Use Brouwer's fixed point theorem

$$\|F(S)\|_2 \leq L_{m,0} (\|\Phi_0^+ \tilde{y}\|_2 + \|\Phi_0^+\|_{op} \|S\|_2)^2$$

F maps the ball of radius  $r$  into the ball of radius  $L_{m,0} (\|\Phi_0^+ \tilde{y}\|_2 + \|\Phi_0^+\|_{op} r)^2$

Taking  $r := \frac{\|\Phi_0^+ \tilde{y}\|_2}{\|\Phi_0^+\|_{op}}$ ,

then if  $L_{m,0} \leq \frac{1}{4 \|\Phi_0^+ \tilde{y}\|_2 \|\Phi_0^+\|_{op}}$

F maps ball of radius  $r$  to ball of radius  $r$

and there exists a fixed pt  $S_*$

i.e.  $\exists$  interpolating solution  $\Theta_* = \Theta_0 + \Phi_0^+ (y - f_m(\Theta_0) + S_*)$

$\hookrightarrow$  approximately solut<sup>o</sup> of  $y = f_m(\Theta_0) + \Phi_0(\Theta - \Theta_0)$

$\leftarrow$  order  $\leq \|y - f_m(\Theta_0)\|_2$

# Global CV of gradient flow

$$\Phi(\theta_t) = \mathbb{D}f_m(\theta_t)$$

Gradient flow:

$$\begin{cases} \dot{\theta}_t = -\nabla_{\theta} \hat{R}(\theta) = \frac{1}{n} \Phi_t (y - f_m(\theta_t)) \\ \text{initialization at } \theta_0 \end{cases}$$

easy to  
extend  
analysis to  
discrete dynamics

We will compare this dynamic to the dynamics on the linearized model

1st order Taylor expansion

$$f_{lin}(x; \theta) = f(x; \theta_0) + \langle \theta - \theta_0, \nabla_{\theta} f(x; \theta_0) \rangle$$

$$\hat{R}_{lin}(\theta) = \frac{1}{2n} \|y - f_m(\theta_0) - \Phi_0(\theta - \theta_0)\|_2^2$$

$$\frac{d}{dt} \bar{\theta}_t = -\nabla_{\theta} \hat{R}_{lin}(\bar{\theta}_t) = \frac{1}{n} \Phi_0^T (y - f_m(\theta_0) - \Phi_0(\bar{\theta}_t - \theta_0))$$

Will assume  $p \geq n$   $\text{rank}(\Phi_0) = n$

↳ overparametrized

$$\text{ERM}_0^{\text{lin}} = \left\{ \theta : \mathbb{E}_0(\theta \cdot \theta_0) = y - f_m(\theta_0) \right\}$$

$$\bar{\theta}_t \rightarrow \bar{\theta}_\infty = \operatorname{argmin} \left\{ \|\bar{\theta} - \theta_0\|_2 : \bar{\theta} \in \text{ERM}_0^{\text{lin}} \right\}$$

↑ Lecture 3

Notations:

$$L_m := \sup_{\theta_1 \neq \theta_2} \frac{\|\mathbb{D}f_m(\theta_1) - \mathbb{D}f_m(\theta_2)\|_{\text{op}}}{\|\theta_1 - \theta_2\|_2}$$

$$\sigma_{\min} := \sigma_{\min}(\mathbb{D}f_m(\theta_0))$$

$$\sigma_{\max} := \sigma_{\max}(\mathbb{D}f_m(\theta_0))$$

$$\|\mathbb{D}f(\theta)\|_{\text{op}} = \sup_{a \in \mathbb{R}^p} \frac{\|a^\top \nabla f(\cdot, \theta)\|_{L^2(\mathcal{H})}}{\|a\|_2}$$

$$\text{Lip}(\mathbb{D}f) = \sup_{\theta_1 \neq \theta_2} \frac{\|\nabla f(\cdot, \theta_1) - \nabla f(\cdot, \theta_2)\|_{\text{op}}}{\|\theta_1 - \theta_2\|_2}$$

Thm: [Berkhett, Monkenari, Rokhlin, '21]

Assume

$$L_m \|y - f_m(\theta_0)\|_2 < \frac{1}{4} \sigma_{\min}^2 \quad (*)$$

Then  $\forall t > 0$

$$(1) \hat{R}(\theta_t) \leq \hat{R}(\theta_0) e^{-\lambda_0 t} \quad \lambda_0 := \frac{\sigma_{\min}^2}{2m}$$

$$(2) \|\theta_t - \theta_0\|_2 \leq \frac{2}{\sigma_{\min}} \|y - f_m(\theta_0)\|_2$$

$$\|\theta_t - \bar{\theta}_t\|_2 \leq C L_m \frac{\sigma_{\max}^2}{\sigma_{\min}^2} \|y - f_m(\theta_0)\|_2^2$$

$$(3) \|f(\theta_t) - f_{\min}(\bar{\theta}_t)\|_{L^2(P)} = \mathbb{E}_x \left[ (f(x; \theta_t) - f_{\min}(x; \bar{\theta}_t))^2 \right]^{1/2} \\ \leq C \left\{ \frac{\text{Lip}(Df)}{\sigma_{\min}^2} + \|Df(\theta_0)\|_{\text{op}} \frac{L_m \sigma_{\min}^2}{\sigma_{\max}^5} \right\} \|y - f_m(\theta_0)\|_2^2$$

Remark 1: Condition (\*)  $\Rightarrow$  condition  $\exists$  interpolating solut<sup>o</sup> earlier

$$L_{m,0} \leq \frac{1}{4 \| \Phi_0^+(y - f_m(\theta_0)) \|_2 \| \Phi_0^+ \|_{\text{op}}}$$

$$\|\Phi_0^+(y - f_m(\Theta_0))\|_2 \|\Phi_0^+\|_{\text{op}} \leq \|\Phi_0^+\|_{\text{op}}^2 \|y - f_m(\Theta_0)\|_2$$

$$\hookrightarrow = \frac{1}{\sigma_{\min}(\Phi_0)^2}$$

$$L_{m,0} \leq L_m < \frac{\sigma_{\min}^2}{4 \|y - f_m(\Theta_0)\|_2} \leq \frac{1}{4 \|\Phi_0^+(y - f_m(\Theta_0))\|_2 \|\Phi_0^+\|_{\text{op}}}$$

Remark 2: (1) shows global cv to 0!! (success!)

(2) shows  $\Theta_t$  stays close to  $\Theta_0$

and  $\Theta_t \approx \bar{\Theta}_t$  dynamics where we replace  $f \rightarrow f_{\text{lin}}$

(3) From a statistical perspective, this is the most important: shows that  $f(\cdot; \Theta^t)$  and  $f_{\text{lin}}(\cdot; \bar{\Theta}^t)$  behave the same on test data

$\Rightarrow$  show that learning  $\equiv$  to learning with linear model  $f_{\text{lin}}(\cdot; \Theta)$

Remark 3: Rescaling  $f_\alpha(\cdot; \theta) = \alpha f(\cdot; \theta)$

assume  $f(\alpha; \theta_0) = 0$

$$\text{Lip}(\mathbb{D}f_{\alpha, m}) = \alpha L_m$$

$$\|y - f_\alpha(\theta_0)\|_2 = \|y\|_2$$

$$\sigma_{\min}(\mathbb{D}f_{\alpha, m}) \approx \alpha$$

Condition become

$$\alpha \ll \alpha^2$$

i.e.  $\alpha \gg 1$

Application: 2-layer neural networks

$$f(x; \theta) = \frac{\alpha}{\sqrt{M}} \sum_{j=1}^M a_j \sigma(\langle w_j, x \rangle) \quad x \stackrel{\text{iid}}{\sim} \text{Unif}(S^{d-1})$$

To simplify, fix  $a_j \in \{\pm 1\}$

$$\theta = (w_1, \dots, w_M) \in \mathbb{R}^{Md} \quad p = Md$$

$$w_1^0, \dots, w_{\frac{M}{2}}^0 \stackrel{\text{iid}}{\sim} \mathcal{N}(0, \text{Id}_d)$$

$$w_{\frac{M}{2}+1}^0 = w_1^0$$

$$a_1 = \dots = a_{M/2} = 1$$

$$a_{M/2+1} = \dots = a_M = -1$$

so that  $f(x; \theta^0) = 0$ .

$$\left[ D f_m(\theta^0) \right]_{i, (j,k)} = \frac{\alpha}{\sqrt{M}} a_j \sigma'(\langle w_j, x_i \rangle) (x_i)_k$$

Assume  $y_i$  are  $O(1)$ -sub-Gaussian.

Lemma: Assume  $Md \geq C n \log n$ . Then w.h.p

$$\|Y - f_m(\theta_0)\| \lesssim \sqrt{m} \quad \sigma_{\max} \lesssim \alpha (\sqrt{m} + \sqrt{d})$$

$$\sigma_{\min} \gtrsim \alpha \sqrt{d} \quad L_m \lesssim \alpha \sqrt{\frac{d}{M}} (\sqrt{m} + \sqrt{d})$$

Condition  $L_m \cdot \|Y - f_m(\theta_0)\|_2 < \frac{1}{4} \sigma_{\min}^2$

$$\alpha \sqrt{\frac{d}{M}} (\sqrt{m} + \sqrt{d}) \sqrt{m} \lesssim \alpha^2 d \quad \stackrel{m \geq d}{\implies} \quad \alpha \gtrsim \sqrt{\frac{m^2}{Md}}$$

Corollary: If  $n \geq d$   $Md \gtrsim n \log n$   $\alpha \gtrsim \sqrt{\frac{m^2}{Md}}$

Then (1)  $\hat{R}_m(\theta_t) \lesssim \sqrt{m} e^{-\alpha^2 \frac{d}{m} t}$

(2)  $\|\theta_t - \theta_0\|_2 \lesssim \frac{1}{\alpha} \sqrt{\frac{m}{d}} = \frac{1}{\alpha} \sqrt{\frac{m}{Md^2}} \|\theta_0\|_2$

(3)  $\|f(\theta_t) - f_{\text{lin}}(\bar{\theta}_t)\|_{L^2(P)} \lesssim \frac{1}{\alpha^2} \sqrt{\frac{m^5}{Md^4}}$

RMK: Hence linear theory is accurate when

(1) Fixed  $\alpha$ :  $p = Md \rightarrow \infty$   
(scaling  $\frac{\alpha}{\sqrt{M}}$ )

(2) Fixed  $M$ :  $Md \gtrsim m \log m$        $\alpha \rightarrow \infty$

RMK: For  $\alpha = 1$ , need  $p = Md \gtrsim m^2$   
Necessary  $p = Md \geq m$       bound might be suboptimal.

RMK: If  $\alpha \downarrow 0$  with  $M$ , we will get different regime

e.g.  $\alpha = \frac{\alpha_0}{\sqrt{M}}$

$$f(\alpha; \theta) = \frac{1}{M} \sum_j a_j \sigma(\langle w_j, \alpha \rangle)$$

$$= \int a \sigma(\langle w, \alpha \rangle) \rho(d\alpha dw)$$

$M \rightarrow \infty$

→ non-linear dynamics (Mean-Field regime)

# Learning in the Lazy regime

$$f_{\text{lin}}(\alpha; \theta) = \underbrace{f(\alpha; \theta_0)}_{\text{effect not learned}} + \langle \nabla f(\alpha; \theta_0), \theta - \theta_0 \rangle$$

NT model:  $b_{\text{NT}}(\alpha; b) = \langle \nabla f(\alpha; \theta_0), b \rangle$

$$f(\alpha; \theta) = \frac{1}{\sqrt{M}} \sum_{j \in [M]} a_j \sigma(\langle \omega_j, \alpha \rangle)$$

$$b_{\text{NT}} = \underbrace{\langle a, \Phi_{\text{RF}}(\alpha) \rangle}_{\text{lineareigt } 2^{\text{nd}} \text{ layer}} + \underbrace{\langle b, \Phi_{\text{NT}}(\alpha) \rangle}_{\text{lineareigt } 1^{\text{st}} \text{ layer}}$$

$$\Phi_{\text{RF}}(\alpha) = \frac{1}{\sqrt{M}} \begin{pmatrix} \sigma(\langle \omega_1, \alpha \rangle) \\ \vdots \\ \sigma(\langle \omega_M, \alpha \rangle) \end{pmatrix} \in \mathbb{R}^M$$

$$\Phi_{\text{NT}}(\alpha) = \frac{1}{\sqrt{M}} \begin{pmatrix} \sigma'(\langle \omega_1, \alpha \rangle) \alpha \\ \vdots \\ \sigma'(\langle \omega_M, \alpha \rangle) \alpha \end{pmatrix} \in \mathbb{R}^{Md}$$

Associated kernel:

$$K_M(\alpha_1, \alpha_2) = \langle \Phi_{\text{RF}}(\alpha_1), \Phi_{\text{RF}}(\alpha_2) \rangle + \langle \Phi_{\text{NT}}(\alpha_1), \Phi_{\text{NT}}(\alpha_2) \rangle$$

$$= \frac{1}{M} \sum_{j \in [M]} \sigma(\langle \omega_j, \alpha_1 \rangle) \sigma(\langle \omega_j, \alpha_2 \rangle) + \langle \alpha_1, \alpha_2 \rangle \sigma'(\langle \omega_j, \alpha_1 \rangle) \sigma'(\langle \omega_j, \alpha_2 \rangle)$$

$$\begin{matrix} \text{WLLN} \\ \xrightarrow{M \rightarrow \infty} \end{matrix} K(\alpha_1, \alpha_2)$$

$$= \mathbb{E}_{\omega} [\sigma(\langle \omega, \alpha_1 \rangle) \sigma(\langle \omega, \alpha_2 \rangle)] + \langle \alpha_1, \alpha_2 \rangle \mathbb{E}_{\omega} [\sigma'(\langle \omega, \alpha_1 \rangle) \sigma'(\langle \omega, \alpha_2 \rangle)]$$

## NEURAL TANGENT KERNEL (NTK)

Saor, Gabriel, Mongler, 2018.

If  $\|\alpha_1\|_2 = \|\alpha_2\|_2 = 1$        $\omega \sim N(0, Id_d)$

$$K(\alpha_1, \alpha_2) = h(\langle \alpha_1, \alpha_2 \rangle) \quad h(t) = \mathbb{E} [\sigma(G_1) \sigma(G_2) + t \sigma'(G_1) \sigma'(G_2)]$$

$$\begin{pmatrix} G_1 \\ G_2 \end{pmatrix} \sim N\left(0, \begin{pmatrix} 1 & t \\ t & 1 \end{pmatrix}\right)$$

$M \rightarrow \infty$  in this limit: learning with NN in this regime  
 $\equiv$  kernel method with NTK

L see next lecture!

Does the lazy regime explain NNs used in practice?

Community got very excited: scaling  $\frac{1}{\sqrt{M}}$  seemed to correspond to what is used in practice, but:

- Lazy regime:
- \*  $\|\theta^t - \theta^0\|_2 \ll \|\theta^0\|$  weights barely move
  - \*  $f(x; \theta^t) \approx f(x; \theta^0) + \langle \theta^t - \theta^0, \nabla_{\theta} f(x; \theta^0) \rangle$   
effectively behave as a linear model (kernel method)

Extensive literature checking numerically validity of this regime:

- weights  $\theta^t$  don't stay near initialized  $\theta^0$  (e.g. filters in CNNs)
- if we replace  $f(x; \theta)$  by  $f_{lin}(x; \theta) \Rightarrow$  drop in performance
- ⚠ Depends on architecture & data set: NTK sometimes match performance

Theory: linear models are much less powerful methods than non-linear NNs (in terms of approximation + generalization capabilities)

⇒ many "separation" results

Proof of Theorem Part (1):

$$y_t := f_m(\theta_t)$$

$$\dot{\theta}_t = -\frac{1}{n} \Phi_t^T (y_t - y_0) \quad [\text{param space}]$$

$$\dot{y}_t = \Phi_t \dot{\theta}_t = -\frac{1}{n} K_t (y_t - y) \quad [\text{feb space}]$$

"Kernel"  $K_t := \Phi_t \Phi_t^T \in \mathbb{R}^{m \times m}$

$$\frac{d}{dt} \underbrace{\|y_t - y\|_2^2}_{= 2m \hat{R}(\theta_t)} = -\frac{2}{n} \langle y_t - y, \underbrace{K_t (y_t - y)}_{\lambda_{\min}(K_t) = \sigma_{\min}^2 = \sigma_{\min}(\Phi_t)^2} \rangle$$

$$\sigma_{\min}(\Phi_t) \geq \sigma_{\min}(\Phi_0) - L_m \cdot \|\theta_t - \theta_0\|_2$$

If  $\|\theta_t - \theta_0\|_2 \leq r_* := \frac{\sigma_{\min}}{2L_m}$  then  $\lambda_{\min}(K_t) \geq \left(\frac{\sigma_{\min}}{2}\right)^2$

Let  $t_* := \inf \{t : \|\theta_t - \theta_0\|_2 > r_*\}$

$$t < t_* \implies \|y_t - y\|^2 \leq \|y_0 - y\|^2 e^{-\lambda_0 t}$$

$\hookrightarrow \lambda_0 = \frac{\sigma_{\min}^2}{2m}$

→ Need to show  $t_* = \infty$

$$\|\dot{\theta}_t\|_2 = \frac{1}{n} \|\Phi_t^T (y_t - y_0)\|_2$$

$$\frac{d}{dt} \|y_t - y\|_2 = \frac{1}{2 \|y_t - y\|_2} \frac{d}{dt} \|y_t - y\|_2^2$$

$$= -\frac{1}{n} \frac{\langle y_t - y, K_t (y_t - y) \rangle}{\|y_t - y\|_2}$$

$$= -\frac{1}{n} \frac{\|\Phi_t^T (y_t - y)\|_2^2}{\|y_t - y\|_2}$$

$t \leq t_*$   
 $\sigma_{\min}(\Phi_t)$

$$\leq -\frac{\sigma_{\min}}{2n} \|\Phi_t^T (y_t - y)\|_2$$

$\geq \frac{\sigma_{\min}}{2}$

$$\|y_t - y_0\|_2 \leq \frac{2}{\sigma_{\min}} \|\Phi_t^T (y_t - y)\|_2$$

$$\frac{d}{dt} \|y_t - y\|_2 + \frac{\sigma_{\min}}{2} \|\dot{\theta}_t\|_2 \leq 0$$

$$\frac{d}{dt} \left[ \|y_t - y\|_2 + \frac{\sigma_{\min}}{2} \|\theta_t - \theta_0\|_2 \right] \leq 0$$

$$\Rightarrow \|y_t - y\|_2 + \frac{\sigma_{\min}}{2} \|\theta_t - \theta_0\|_2 \leq \|y_0 - y\|_2$$

$$\Rightarrow \|\theta_t - \theta_0\|_2 \leq \frac{2}{\sigma_{\min}} \|y_0 - y\|_2$$

$L_n \|y - y_0\|_2 < \frac{1}{4} \sigma_{\min}^2$   
↑

$$\text{If } t_* < \infty, \|\theta_{t_*} - \theta_0\|_2 \leq \frac{2}{\sigma_{\min}} \|y - y_0\|_2 < \frac{\sigma_{\min}}{2L_n} = r_*$$

$\Rightarrow t_* = \infty$  This proves (1) □ ↙