10 - · Linearization - a.k.a. Neural -tongert Kernel - Last lecture, we computed risk curves associated with overperoversized linear models under implicit bias (e.g. min-norm solution.). \* Gradiant Flass on IRd - We want to train the model using a gradient-based algorithm. - Gradient flows are easier to analyze, but not practical. - Gradient descent is practical: (GD)(re-arrange)  $\iff \frac{\Theta_{t}}{\gamma} = \nabla R(\Theta_t)$  $\left(\begin{array}{c} lin\\ 100\end{array}\right) \implies \left(\begin{array}{c} \frac{d\theta_{t}}{dt} = \dot{\theta}_{t} = -\nabla \mathcal{R}(\theta_{t}) \right) \quad (6F)$ - GD is the Euler discretization of GF. For shall n, GD behaves like GF. Proposition:  $\int_{f} \frac{d}{dt} \Delta_{t} \leq -\mu \Delta_{t}$ , then  $\Delta_{t} \leq e^{-\mu t} \Delta_{0}$ . Proof:  $\frac{d}{dt}(e^{\mu t}\Delta_{t}) = \mu e^{\mu t}\Delta_{t} + e^{\mu t}\frac{d}{dt}\Delta_{t}$   $\leq \mu e^{\mu t}\Delta_{t} - \mu e^{\mu t}\Delta_{t}$ ( Ly (c)) Integrating both sides yields the result.

Assumption (Polyak - Lejancuites (PL) Inequality): For 
$$\mu > 0$$
,  
 $\forall \Theta$ ,  $R(\Theta) - R(\Theta_{k}) \leq \frac{1}{\mu} \|\nabla R(\Theta)\|^{2}$ .  
Proposition:  $PL \Rightarrow C \quad f = C \quad \Delta_{L} = R(\Theta_{k}) - R(\Theta_{k})$ .  
 $Prof: \frac{1}{dt} \left\{ R(\Theta_{k}) - R(\Theta_{k}) \right\} = \langle \nabla R(\Theta_{k}), -\nabla R(\Theta_{k}) \rangle$   
 $= -\|\nabla R(\Theta_{k})\|^{2}$ .  
 $\leq -\mu \left\{ R(\Theta_{k}) - R(\Theta_{k}) \right\}$   
 $= -\|\nabla R(\Theta_{k})\|^{2}$ .  
 $\leq -\mu \left\{ R(\Theta_{k}) - R(\Theta_{k}) \right\}$   
 $= -\|\nabla R(\Theta_{k})\|^{2}$ .  
 $\Delta_{L}$   
 $= -\mu \left\{ R(\Theta_{k}) - R(\Theta_{k}) \right\}$   
 $A = \frac{1}{2} \|\Theta_{L} - \Theta_{k}\|^{2}$ .  
 $\Theta_{L} \quad \langle \Theta - \Theta_{k} / \nabla R(\Theta) \rangle \geq \frac{\mu}{2} \|\Theta - \Theta_{k}\|^{2}$ .  
 $\Theta_{L} \quad \langle \Theta - \Theta_{k} / \nabla R(\Theta) \rangle \geq \frac{\mu}{2} \|\Theta_{L} - \Theta_{k}\|^{2}$ .  
 $Properties : SC \Rightarrow C \quad for \quad \Delta_{L} = \frac{1}{2} \|\Theta_{L} - \Theta_{k}\|^{2}$ .  
 $Properties : SC \Rightarrow C \quad for \quad \Delta_{L} = \frac{1}{2} \|\Theta_{L} - \Theta_{k}\|^{2}$ .  
 $\Phi_{L} \quad = \left\langle \Theta_{L} - \Theta_{k}, -\nabla R(\Theta_{L}) \right\rangle$   
 $\leq -\mu \cdot \frac{1}{2} \|\Theta_{L} - \Theta_{k}\|^{2}$ .  
 $E_{L} \quad E_{L}$ .  
 $Rouser L : - [n \quad fort, SC \Rightarrow PL \quad so \quad it is a stronger condition.
 $- \ln \quad hertin \quad Coscs, GF \quad converget \quad expressively \quad for f^{1}$ .$ 

× Linearization - Reall: We want to train a function s.t.  $\gamma \approx f(x; \theta)$  e.g. f is NN. output input parometer Data:  $(Y_i, x_i)$  for i = 1, ..., n, let  $x_i \in \mathbb{R}^d$  and  $\theta \in \mathbb{R}^p$ .  $Y_{n} = \begin{bmatrix} Y_{i} \\ Y_{n} \end{bmatrix} + IR^{n} \qquad f_{n} (\Theta) = \begin{bmatrix} f(x_{i}, \Theta) \\ f(x_{n}; \Theta) \end{bmatrix} + IR^{n}$ Empirical risk:  $\widehat{R}(\Theta) = \frac{1}{2n} \| \gamma_n - f_n(\Theta) \|^2$ . Minimize R(0) with GFI  $\dot{\Theta}_{t} \triangleq \frac{d\Theta_{t}}{dt} = \frac{1}{n} \nabla_{f_{n}}(\Theta_{t}) (\gamma_{n} - f_{n}(\Theta_{t}))$ · ∈ IR<sup>p×n</sup> (Jacobien<sup>T</sup> of fn: IR<sup>P</sup>→IR<sup>n</sup>). - JG.H18 argues that, in highly overperametrined regime, Of changes only slightly w.r.t. initializertion Ob. Therefore we compare of with its linearization at to:  $\overline{\Theta}_{t} \stackrel{\text{d}}{=} \frac{d \theta_{t}}{dt} = \frac{1}{n} \nabla_{f_{n}} (\theta_{o}) \left( \gamma_{n} - \frac{1}{f_{n}} (\theta_{o}) - \nabla_{f_{n}} (\theta_{o})^{\mathsf{T}} (\overline{\Theta}_{t} - \Theta_{o}) \right)$  $= \int_{n}^{\ln} \left(\overline{\Theta}_{t}\right) \quad \text{first-order Taylor} \\ \text{approx around } \Theta_{0}. \\ \text{which minimizes } \quad \widehat{R}^{\ln}(\overline{\Theta}) = \frac{1}{2n} \left\| \gamma - \int_{n} (\Theta_{0}) - \nabla_{fn} (\Theta_{0})^{T} (\overline{\Theta} - \Theta_{0}) \right\|^{2}.$ \* Why expect generalization?.  $- \int \nabla_{f_n}(\theta_0)$  is full-rank,  $\mathcal{E}_{\mathfrak{D}} = \left\{ \overline{\Theta} : \widehat{R}^{h_n}(\overline{\Theta}) = 0 \right\}$  is an affine subspace

- 
$$\overline{\Theta}_{00} = uggin \left\{ 1\overline{\Theta} - \Theta_{0} \right\} : \nabla f_{n}(\Theta_{0})^{2} (\overline{\Theta} - \Theta_{0}) = Y - f_{n}(\Theta_{0})^{2}$$
  
 $\in \mathcal{G}$   $\Rightarrow$  restriction solution  $\Rightarrow$  implicit, regularization  $\left( \text{ cer previous letters} \right) (AS)$   
- But how realistic is this approximation:  
 $1 \cdot \Theta_{1} \approx \overline{\Theta}_{1}$   $\forall t ?$   $2 \cdot f_{n}(\Theta_{1}) \approx f_{n}^{\text{the}}(\overline{\Theta}_{1}) ?$   
This regular is called the linear regular.  
 $= \overline{\Theta}_{\infty}$  generalizes  $+ \Theta_{E} \approx \overline{\Theta}_{1}$   $\forall t$ .  
 $\Rightarrow QF$  trained MW generalizes!  
 $\Rightarrow QF$  trained MW generalizes!  
 $= \frac{1}{2} QF$  trained MW generalizes!  
 $= \frac{1}{2} QF$  trained  $MW$  generalizes!  
 $= \frac{1}{2} QF$  trained  $\frac{1}{2} \frac{\nabla f_{n}(\nabla f_{n}(\Theta_{0}))}{\frac{1}{2} \frac{1}{2}}$   
 $= \frac{1}{2} QF$  trained  $\frac{1}{2} \frac{\nabla f_{n}(\nabla f_{n}(\Theta_{0}))}{\frac{1}{2} \frac{1}{2}}$   
 $= \frac{1}{2} QF$  trained  $\frac{1}{2} \frac{\nabla f_{n}(\nabla f_{n}(\Theta_{0})}{\frac{1}{2} \frac{1}{2}}$   
 $= \frac{1}{2} \frac{1}$ 

Proof: 
$$L$$
,  $Lat$ ,  $f_{t} \triangleq f_{n}(\theta_{t})$ , so  $\hat{R}(\theta_{t}) = \frac{1}{2} \|f_{t} - Y_{t}\|^{L}$ .  
 $f_{t} = \nabla f_{n}(\theta_{t}) \hat{\Theta}_{t} = -\frac{1}{n} \nabla f_{n}(\theta_{t}) \nabla f_{n}(\theta_{t})^{T} [f_{t} - Y_{n}]$   
 $= -\frac{1}{n} k_{t} (f_{t} - Y_{n})$   $\approx k_{0}$  : Mound touch book  
 $= -\frac{1}{n} k_{t} (f_{t} - Y_{n})$   $\approx k_{0}$  : Mound touch book  
 $Aliss, \frac{d}{dt} \hat{R}(\theta_{t}) = \frac{d}{dt} \frac{1}{2} \|f_{t} - Y_{n}\|^{2} = \cdot \langle f_{t} - Y_{n}, f_{t} \rangle$   
 $= -\frac{1}{n} \langle f_{t} - Y_{n}, K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} - Y_{n}, K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} - Y_{n}, K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} - Y_{n}, K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} - Y_{n}, K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} - Y_{n}, K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} - Y_{n}, K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} - Y_{n}, K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} - Y_{n} | \xi - \frac{1}{2n} \langle f_{t} - Y_{n}, K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} - Y_{n} | K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} - Y_{n} | K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} - Y_{n} | K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} - Y_{n} | K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} - Y_{n} | K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} - Y_{n} | K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} - Y_{n} | K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} - Y_{n} | K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} - Y_{n} | K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} (f_{t} - Y_{n} | K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} (f_{t} - Y_{n} | K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} (f_{t} - Y_{n} | K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} (f_{t} - Y_{n} | K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} (f_{t} - Y_{n} | K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} (f_{t} - Y_{n} | K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} (f_{t} - Y_{n} | K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} (f_{t} - Y_{n} | K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} (f_{t} - Y_{n} | K_{t} (f_{t} - Y_{n}) \rangle$   
 $= -\frac{1}{n} \langle f_{t} (f_{t} - Y_$ 

2.  

$$\|\hat{\Theta}_{t}\| = \frac{1}{n} \|\nabla f_{n}(\Theta_{t})^{T}(f_{t} \cdot \gamma_{n})\|$$

$$= \frac{1}{dt} \|f_{t} \cdot \gamma_{n}\| = \frac{1}{(f_{t} \cdot \gamma_{n})} \frac{\|\nabla f_{n}(\Theta_{t})^{T}(f_{t} - \gamma_{n})\|}{\|f_{t} - \gamma_{n}\|}$$

$$= -\frac{1}{2n} \|\nabla f_{n}(\Theta_{t})^{T}(f_{t} - \gamma_{n})\|$$

$$= -\frac{1}{2n} \|\Theta_{t} - \Theta_{t}\|^{2} \|\Theta_{t} - \Theta_{t}\|^{2}$$

$$= -\frac{1}{2n} \left(\frac{1}{2} \|\Theta_{t} - \Theta_{t}\| + \frac{1}{2n} \|\Theta_{t}\| + \frac{1}{2n} \|\Theta_{t} - \Theta_{t}\| + \frac{1}{2n} \|\Theta_{t} - \Theta_{t}\| + \frac{1}{2n} \|\Theta_{t} - \Theta_{t}\| + \frac{1}{2n} \|\Theta_{t}\| + \frac{1}{2n}$$

Technical defail

$$\begin{array}{l} \hline \label{eq:product} \hline \end{tabular} \\ \hline$$

$$= 2e^{-pt} \|y_{n} - f_{0}\| \left\{ \frac{1}{n} \frac{1}{\nabla_{max}} \|y_{n} + f_{0}\| + \frac{\nabla_{max}}{n} \right\}.$$
In leg relay this, we get (note  $\theta_{0} = \overline{\theta}_{0}$ )
$$\|\theta_{k} - \overline{\theta}_{k}\| \leq \frac{2}{pk} \|y_{n} - f_{0}\| \left\{ \frac{1}{n} \frac{1}{\nabla_{max}} \|y_{n} + f_{0}\| + \frac{\nabla_{max}}{n} \right\}.$$

$$\# Two - layer Neural Networks$$

$$- Consider  $f(x; \theta) = \frac{1}{16} \sum_{j=1}^{m} q_{j} \sigma((\omega_{j}, z)), \quad \theta = (\omega_{1, \dots, j}, \omega_{m}).$ 

$$p = md.$$

$$wight \qquad wight \qquad wight \qquad unput \\ \text{wight} \qquad \text{to grave } \frac{1}{p} \int_{j=1}^{m} q_{j} \sigma((\omega_{j}, z)), \quad \theta = (\omega_{1, \dots, j}, \omega_{m}).$$

$$p = md.$$

$$= for suplicity close half of  $q_{j}$ 's  $+l$  and other  $log - l.$ 

$$= lnitialise \quad \Theta = \omega_{j} \sim llag(S^{l-1})$$

$$= Train is a \quad \Theta F:$$

$$= we can compute: [\nabla f(x; \theta)]_{i_{1}(j_{1}k)} = \frac{1}{m} q_{j} \sigma^{-1}(\langle \omega_{j}, z_{2} \rangle) z_{1k}$$

$$= leg(x_{1}, \gamma_{k}) \in lni \times lid)$$

$$p = nd$$

$$leguna (BMR21): Under contains conditions, whe p.$$

$$l. \|y_{n} - f_{1}(\theta_{0})\| \lesssim in$$

$$= log - \overline{\theta}_{k} \| \rightarrow \theta_{n}(l).$$
Remarks is  $1 - Since = \overline{\theta}_{0}$  is the min-norm solution  $\Rightarrow$  unput these isotres.  

$$2 - T_{log} = result con be charpored. See [SH221].$$$$$$

Proof														
	the	theo	nerl,	رب و •	t		° C	L .		•	° )			
•		∥θ <b>⊦</b> -	- <mark>⊕<sub>t</sub> ∥</mark>	<u>د</u> - 6		/n-fn(6	&)∥{	Jun 1	Y-fn(Q	י מייי ייי	ūax J			
	(using in	the e	stivete: una	S.) ∠ ~	е-	Ju Z	J d l	[[]+]	<u>n)</u> 1	n +	In +1	J P		
				- ~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~	•	° N    +								
						Ŵ	, d					, B		