

## Lecture 12 : Implicit Regularization & Noisy Dynamics

Recap from last lecture:

- \* Role of symmetries A geometric priors: while helpful to improve sample complexity, not sufficient in itself to break the CoD.
- \* Scale separation: → Breaks the curse, at the expense of reducing approximation properties.

↳ Inductive biases coming from architecture.

Today: Another form of "inductive bias": coming from optimisation algorithm.

## Implicit Regularization of GD

Next week: last lecture (Diffusion)

in two weeks: (May 5th) : project presentations.

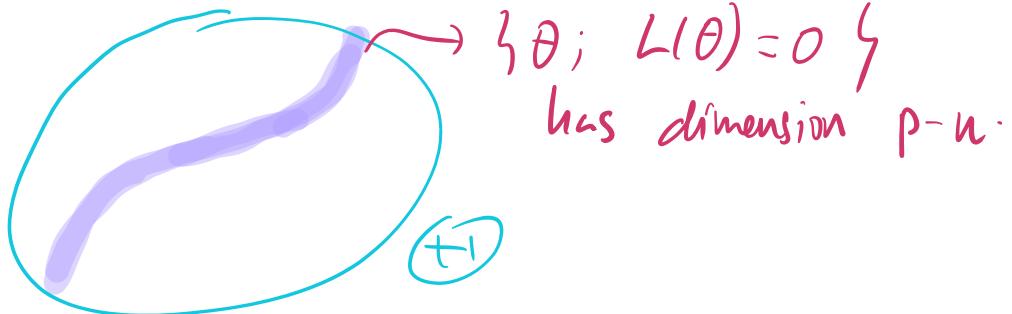
## Learning in the overparametrised Regime

→ Let  $\phi(\theta) : X \rightarrow \mathbb{R}$  be a parametrized (+ differentiable) hypothesis class,  $\theta \in \Theta \subseteq \mathbb{R}^P$ .

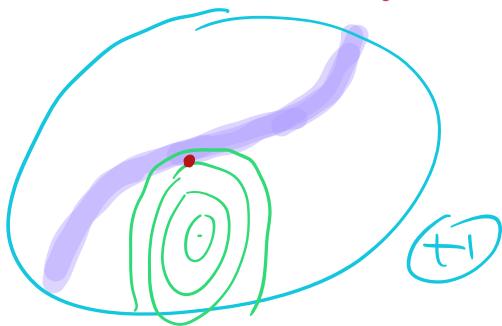
$$\rightarrow \text{Supervised Learning : } L(\theta) = \hat{\mathbb{E}}_D [ l(\phi(\theta)(x), y) ] = \\ = \frac{1}{n} \sum_{i=1}^n l(\phi(\theta)(x_i), y_i)$$

→ Overparametrized regime:  $p \gg n$ . For models that are non-degenerate  $(D\phi|\theta)(x_i)$  is a Jacobian matrix of size

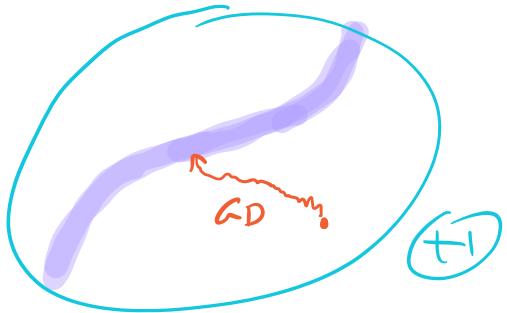
$p \times n$ , has rank  $n$  generically), finding the parameters  $\theta$  that fit the data becomes an underdetermined problem.



→ We saw in earlier lectures that adding regularization enables non-vacuous generalisation bounds:



→ Yet, in many experimental setups, we seem to rely on another form of capacity control / regularization, given by the algorithm instead.



- Natural Questions:
- Understand which solutions are picked by specific choices of algorithm.
  - Can we use this to establish generalisation bounds?

Part I : Implicit Bias of Gradient Descent.

• (Simplest) setup: linear regression:

$$(*) \quad \min_{x \in \mathbb{R}^p} E(x) = \frac{1}{2} \|Ax - y\|^2, \quad A \in \mathbb{R}^{n \times p} \text{ with } p > n.$$

• let's consider solving (\*) with AD:

$$x^{(k+1)} = x^{(k)} - \eta \cdot \nabla E(x^{(k)}) \quad ; \quad x^{(0)} = 0.$$

$$\nabla E(x) = A^T(Ax - y)$$

$$x^{(k+1)} = x^{(k)} - \eta \cdot A^T(Ax^{(k)} - y)$$

$$= [I - \eta A^T A] \cdot x^{(k)} + \eta A^T y.$$

$$= [I - \eta A^T A] \left( [I - \eta A^T A] x^{(k-1)} + \eta A^T y \right) + \eta A^T y$$

= ...

$$= \sum_{j=0}^k [I - \eta A^T A]^j \eta A^T y = \eta \left( \sum_{j=0}^k [I - \eta A^T A]^j \right) A^T y.$$

→  $A^T A$  is symmetric + psd.

$$\nabla^2 E$$

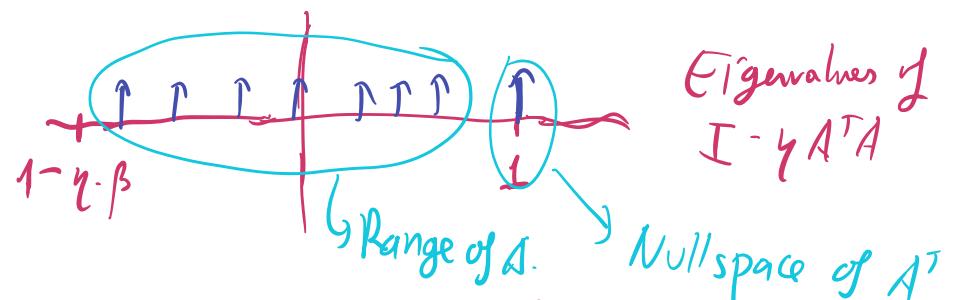
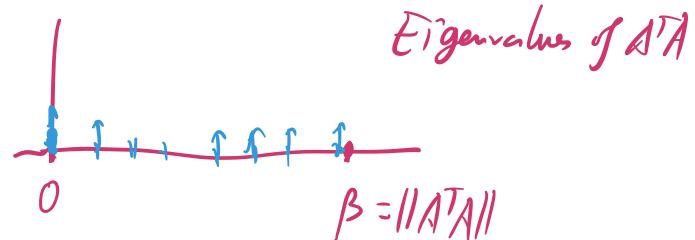
→ IF  $1 - \eta \beta > -1$

$$\Leftrightarrow \eta < \frac{2}{\beta}, \text{ then}$$

$I - \eta A^T A$  has eigenvalues in  $(-1, 1]$ .

⇒ In the range of  $A$ ,  $\sum_{j=0}^k [I - \eta A^T A]^j \rightarrow (\eta A^T A)^{-1} = \eta^{-1}(A^T A)^{-1}$   $k \rightarrow \infty$ .

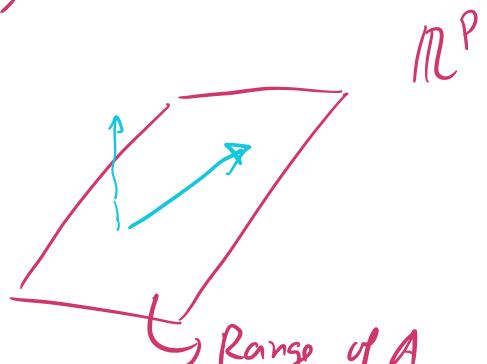
→ so  $x^{(k)} = (A^T A)^{-1} A^T y$  ... it.



$$\xrightarrow[k \rightarrow \infty]{} (\underbrace{AA^\top}_{\text{pseudo-inverse of } A} A^\top y = A^\top y$$

→ variational characterisation:  $A^\top y$  is the minimum  $L_2$  norm solution of  $\min_x \|Ax - y\|^2$ ; in other words,

$$\min_x \|Ax - y\|^2 \text{ st } Ax = y.$$



↳ So, for this particular setting, GD finds the minimum  $L_2$  solution that fits the data.

→ Note: GD in this overparametrised regime does not forget initial conditions → GD initialised at  $x^{(0)} = z$  converges to

$$\min_x \|x - z\|^2 \text{ st } Ax = y.$$

Q: How general is this phenomena?

- Different loss function
- Different parameterization
- Model
- Different flavor of GD?

→ Extension to logistic regression (binary classification)

[Soudry et al '17].

Using linear models. Consider a dataset  $\{(x_i, y_i)\}_{i=1..n}$  with  $y_i = \pm 1$ , and a loss function

$$L(\theta) = \sum_{i=1}^n l(y_i \langle x_i, \theta \rangle) = \sum_{i=1}^n l(\langle \theta, x_i, y_i \rangle)$$

$$l(u) = \log(1 + e^{-|u|}) \text{ logistic loss.}$$

↳ Assume linearly separable data:  $\exists \theta^* \in \mathbb{R}^d$  such that

$$\langle \theta, x_i y_i \rangle > 0 \text{ if } y_i = +1$$

Theorem [Soudry et al.] GD on any

initial condition satisfies

$$\lim_{k \rightarrow \infty} \frac{\theta^{(k)}}{\|\theta^{(k)}\|} = \hat{\theta} \quad \text{when } \hat{\theta} \text{ is the L2 max-margin}$$

Solution:  $\hat{\theta} = \underset{\theta \in \mathbb{R}^d}{\operatorname{argmin}} \|\theta\|^2 \text{ st } \langle \theta, y_i x_i \rangle \geq 1 \text{ for all } i$

Idea of the proof: Exploit the exponential tails of the logistic loss (+the Laplace method)

(-) Implicit bias of direction  $\frac{\theta^{(k)}}{\|\theta^{(k)}\|}$  rather than vector itself is inherent to logistic losses.

Q: Extension to non-linear models?

A: Yes! In [Chizat, Bach '20], extension to shallow NN.

Margin of a predictor  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ :

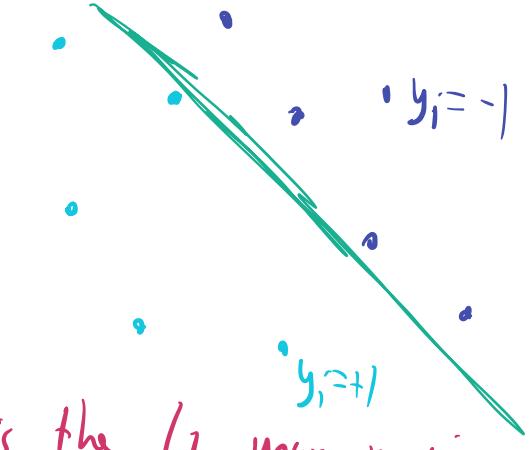
$$\min_{i \in \{1, \dots, n\}} y_i f(x_i) \quad \left( = \min_i \langle \theta, x_i y_i \rangle \text{ for linear } f(x) = \langle x, \theta \rangle \right)$$

Max-Margin predictor in a function class  $\mathcal{F}$ :

$$\max_{\|f\|_{\mathcal{F}} \leq 1} \min_{i \in \{1, \dots, n\}} y_i f(x_i)$$

Theorem [Chizat, Bach]: Consider a two-layer ReLU network in

the mean-field scaling regime:  $f(x; \theta_1, \dots, \theta_m) = \frac{1}{m} \sum_i \phi(\theta_i; x)$ .



Then, under mild assumptions (non-quantitative on  $m!$ ), gradient flow over logistic loss converges (in direction) to the  $\mathcal{F}_1$ -max-margin solution, when  $\mathcal{F}_1$  is the variation-norm (a.k.a. Barron) space.

### Other examples:

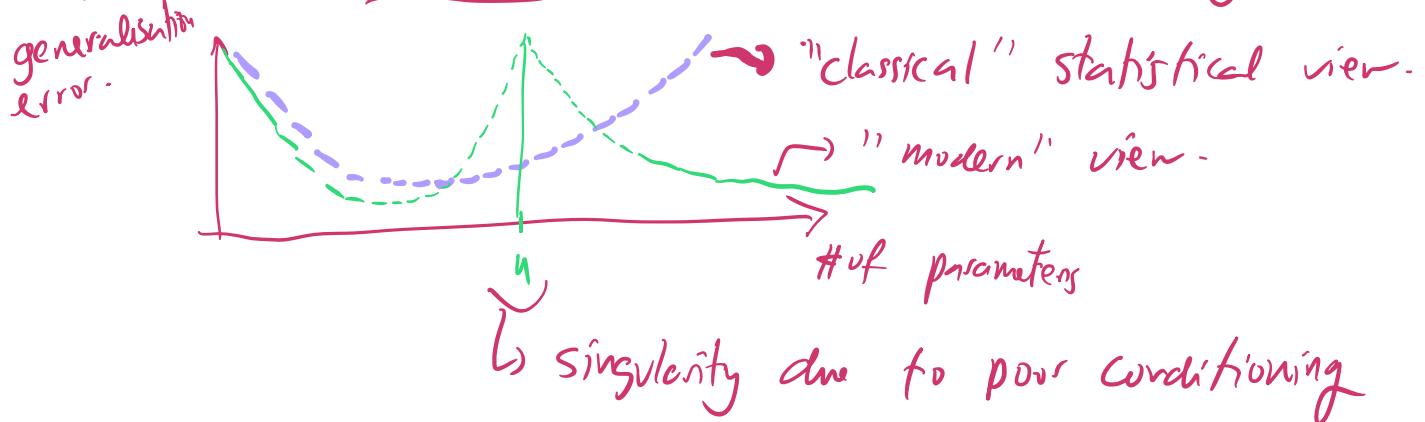
[Ji & Telgarsky '19] : Extend to more general losses, with rates.

[Gunasekar et al.]      { Matrix / Tensor factorization  
 [Li, Ma et al.]            ↳ GD finds low-rank solution.  
 [Arora et al.]

### Takeaways so far:

- ↳ GD has some "built-in" capacity-control, sometimes encoded in its initialization + underlying choice of metric.
- ↳ While in some cases this implicit bias leads to strong statistical guarantees, it is not always clear what are the advantages over explicit regularisation.

↳ Example: Double descent [Belkin, Hsu, et al. + many others..]



→ Analysis is tedious; it needs to be done in a case-by-case

basis.; currently we don't have a good description for implicit bias on complex NN architectures.

→ Negative results: for certain non-linear models, implicit bias of GD probably cannot be characterized by any norm.

[Vardi, Shamir] [Razin & Cohen]

→ We know that vanilla GD is sometimes brittle (trapped in bad local minima).

↳ Q: other forms of algorithmic regularization?

## Part II: Noisy Gradient Dynamics

→ Vanilla GD step  $\theta^{(k+1)} = \theta^{(k)} - \eta \nabla L(\theta^{(k)})$

↳ Each update requires all training samples!  $\frac{1}{n} \sum_{i=1}^n \nabla_{\theta} l(\phi(\theta^{(k)})(x_i), y_i)$

→ Since  $\nabla L(\theta^{(k)}) = \hat{\mathbb{E}}_j [\nabla l(\phi(\theta^{(k)})(x_j), y_j)]$ , a natural (stochastic) approximation is to do a MC estimate:

$$\hat{\nabla} L(\theta^{(k)}) = \nabla L(\phi(\theta^{(k)})(\tilde{x}), \tilde{y}) \text{ with } (\tilde{x}, \tilde{y}) \sim \nu$$

↳ Stochastic Approximation / Optimization [Robbins & Munro, 51]:

$$\theta^{(k+1)} = \theta^{(k)} - \eta_k \cdot \underbrace{g(\theta^{(k)}, z_k)}_{\text{stochastic gradient, } \mathbb{E} g(\theta^{(k)}, z) = \nabla L}$$

with

Folklore: Under mild assumptions on the variance

$\mathbb{E} \|g(\theta^{(k)})\|^2 \leq \eta_k^{-1} \|f'\|^2$  and for appropriate learning rate

If  $\|\theta^{(k+1)} - \theta^*\|$  and  $\|\nabla L\|$  are small enough, then SGD "recovers" same convergence guarantees as GD (global convergence for convex losses, local convergence for non-convex losses).

Q: What about implicit bias, in overparametrised settings?  
Is it the same as GD?

↳ This SGD is viewed as adding noise to the gradients.  
Next week: we are going to study precisely the convergence of a (simplified) version of this algorithm.

TLDL: (next week): solutions of SGD "concentrate" around local minimisers of the loss.

Q: Other forms of noisy gradient descent?

### Label Noise as Implicit Regularization

Idea: Rather than adding noise into the gradients, we instead add noise to the targets:

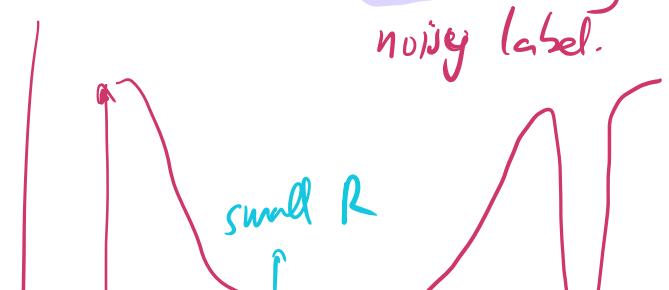
→ For  $k=0, 1, \dots$  do:

- Sample a data-point  $x^{(k)}$  uniformly from  $\{(x_i, y_i)\}_{i=1..n}$
- sample  $\epsilon \sim \{-\sigma, \sigma\}$ .
- $\theta^{(k+1)} = \theta^{(k)} - \eta \nabla_{\theta} \left[ (f(\theta^{(k)})(x^{(k)}) - y^{(k)} - \epsilon)^2 \right]$

noisy label.

→ Consider the term

$$R(\theta) = \frac{1}{2n} \text{tr} \log \left( I - \frac{\eta}{2} \nabla^2 L(\theta) \right)$$



$$(\text{notice for } \eta \rightarrow 0, R(\theta) \approx \frac{1}{4} \text{tr } \nabla^2 L(\theta))$$

measure of  
"total" curvature of the loss.

"noise robust"  
landscape "perspective!"  
large  $R$

Theorem [Damian, Lee & Ma, follow up from Blanc et al]. Under smoothness and Polyak-Lojasiewicz sharpness assumptions, then GD with label noise finds an approximate stationary point of  $\tilde{L}(\theta) = L(\theta) + \lambda \cdot R(\theta)$  ( $\lambda$  is explicit)

In other words, label noise favors "flat" solutions over non-flat ones.

Summarise:

- Algorithmic regularization: powerful method to
  - (+) obtain capacity-control and achieve statistical learning guarantees.
  - (-) sometimes it is the "wrong" way to achieve a desirable effect (least-squares), whereas sometimes leads to efficient algorithms (logistic regression in F, label noise).