

Lecture 7: Optimization in Shallow NNs I:

the "lazy" regime

- Previous weeks: → Discussing approximation, generalization & optimisation aspects separately.
- Always tradeoffs between these sources of error.
- Today: Effect of model parametrization on these tradeoffs.

→ Let $\phi(\theta)$ be a differentiable map $\Theta \xrightarrow{\phi} F$
 θ_0 : initial point

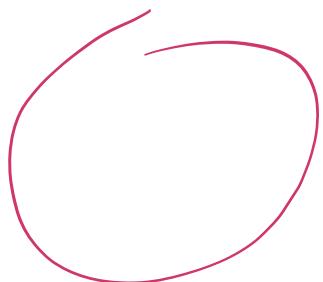
\mathbb{R}^d ↓ parameter space. ↓ function space.

Ex: for a NN: $\phi(\theta) : X \rightarrow \mathbb{R}$

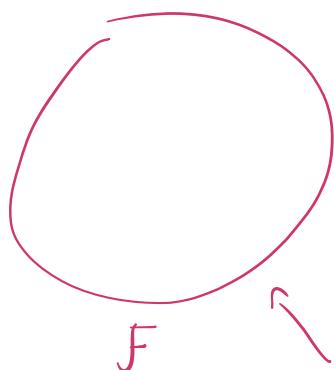
$$x \mapsto \phi(\theta; x) = W_K \sigma(W_{K-1} \sigma \dots W_1 x).$$

→ $R : F \rightarrow \mathbb{R}$ risk functional $R(f) = \mathbb{E}_y [l(f(x), y)]$

→ l is convex w.r.t first argument. ($= \|f(x) - f^*\|_2^2$)
↳ R is convex w.r.t f b for (least-squares).



Θ



F

↳ If ϕ is linear, then the ERM is also convex.

↳ In that case, first-order methods (GD) can find global min efficiently (dim-free iteration complexity).

↳ RICKS example falls in this category

$$\phi(\theta) = \sum_{i=1}^n \alpha_i K(x_i, \cdot) \quad (\text{linear w.r.t parameters}).$$

→ Q: What happens when $\phi(\theta)$ is non-linear?

→ Let us consider gradient-descent with fixed step size γ

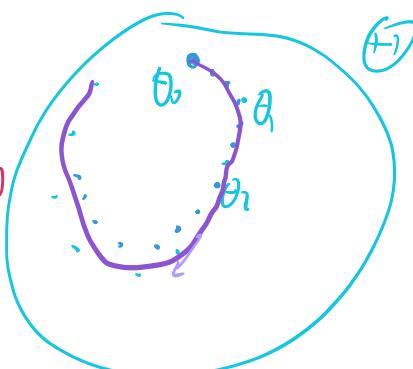
$$\theta_{k+1} = \theta_k - \gamma \nabla R(\phi(\theta_k))$$

$$R(f) \leftarrow \left(\mathbb{E}_{x \sim v} [\ell(f(x), y)] \right)$$

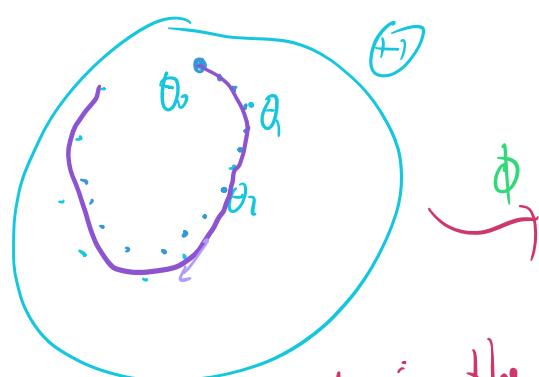
b) empirical over data!

→ As $\gamma \rightarrow 0$, gradient descent is approximating an underlying ODE, the gradient-flow:

$$\dot{\theta}(t) = - \nabla_\theta R(\phi(\theta))$$



→



What is the associated dynamics in function space?



→ Let's write down the dynamics in F :

$$\phi_t = \phi(\theta_t)$$

$$\frac{d}{dt} \theta_t = -\nabla_{\theta} R(\phi(t))$$

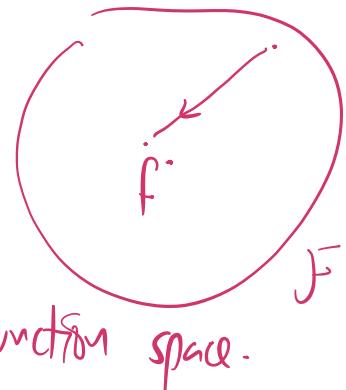
$$\frac{d}{dt} \phi_t = D\phi(\theta_t)^T \cdot \frac{d}{dt} \theta_t$$

$$= -\underbrace{D\phi(\theta_t)^T \cdot D\phi(\theta_t)}_{K_t} \cdot R'(\phi(\theta_t)).$$

R' : functional derivative / first-variation of R in function space. eg $R(f) = \frac{1}{2} \|f - f^*\|_2^2$ $R'(f) = f - f^*$

$K_t = D\phi(\theta(t))^T \cdot D\phi(\theta(t))$, is the tangent kernel.

$\langle K_t f \rangle: \mathcal{F} \rightarrow \mathcal{F}$ maps the "ambient" derivative $R'(\phi_t)$ to a "parameter curve" update in function space.



- In coordinates,

$$(\langle K_t f \rangle)(x) = \int K_t(x, x') f(x') d\nu(x')$$

$$K_t(x, x') = \left\langle \nabla_{\theta} \phi(\theta_t, x), \nabla_{\theta} \phi(\theta_t, x') \right\rangle_{L^2(\theta)}$$

- Let's compute this kernel when $\phi(\theta)$ is a shallow NN of the form $\phi(\theta; x) = \frac{1}{\sqrt{m}} \sum_{i=1}^m g(\theta_i, x)$; $\theta_i \sim \mu_0$ initial distnb.

$$K(x, x') = \frac{1}{m} \sum_{i=1}^m \left\langle \nabla_{\theta} g(\theta, x) \Big|_{\theta=\theta_i}, \nabla_{\theta} g(\theta, x') \Big|_{\theta=\theta_i} \right\rangle$$

$| g(\theta, x) = a \Gamma(x, \theta)$

at fixed x, x' this is an average of
m iid random variables.

$m \rightarrow \infty$

$$\downarrow \mathbb{E}_{\theta \sim p_{\text{true}}} [\langle \nabla_{\theta} g(\theta, x), \nabla_{\theta} g(\theta, x') \rangle] := \bar{K}(x, x').$$

↳ By choosing the scaling in $1/m$, the tangent kernel
initial

has a well-defined limit as $m \rightarrow \infty$. at fine $t=0$.

↳ Q: What happens as we start learning using
gradient dynamics in this regime?

Theorem: [Jacot et al, 18, NTK] For all $T > 0$,
as $m \rightarrow \infty$ uniformly in $[0, T]$ and on compacts in
 \mathbb{R}^N , $K_t(x, x') \rightarrow \bar{K}(x, x')$. $\forall x, x' \in \mathcal{X}$
 $\forall t \in [0, T]$.

$g(\theta, x)$ | Consequence: Learning happens in a linear
space \rightarrow The RKHS whose kernel is \bar{K} .

$$\frac{d}{dt} \phi_t = - K_t \cdot R'(\phi_t) \quad \text{becomes}$$

when $\bar{K} \succ 0$

$$\frac{d}{dt} \phi_t = - \bar{K} \cdot (\phi_t - f^*) \quad \Rightarrow \text{global convergence at}$$

exponential rate!

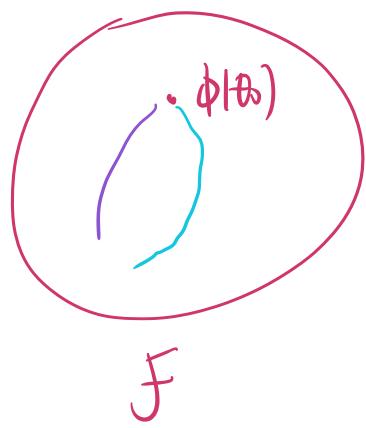
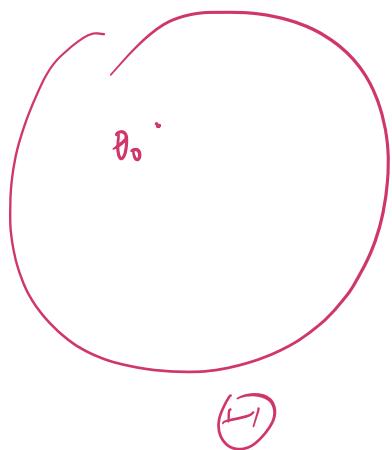
→ How come a non-linear model (shallow NN) behaves like a linear one?

Lazy Dynamics [Chizat, Oyallon, Bach'19]

→ Let $\phi(\theta)$ be a differentiable map $\Theta \in \mathbb{D} \xrightarrow{\phi} \mathcal{F}$
 θ_0 : initial point.

→ Consider the linearized tangent model at θ_0 :

$$\bar{\phi}(\theta) = \phi(\theta_0) + \underline{D\phi(\theta_0)} \cdot (\theta - \theta_0)$$



Two dynamics in function space:

(i) $\dot{\theta}_t = -\nabla R(\phi(\theta_t)) ; \quad \boxed{\phi_t = \phi(\theta_t)} \quad \theta_0$

(ii) $\dot{\bar{\theta}}_t = -\nabla R(\bar{\phi}(\bar{\theta}_t)) ; \quad \boxed{\bar{\phi}_t = \bar{\phi}(\bar{\theta}_t)} \quad \bar{\theta}_0 = \theta_0$

Question: When are the learning dynamics under these two models similar?

→ Denote $F(\theta) = R(\phi(\theta))$ and assume θ_0 s.t $F(\theta_0) > 0$.

→ A gradient step is $\theta_1 = \theta_0 - \gamma \nabla F(\theta_0)$

→ Relative change in objective function: $F(\theta) \approx F(\theta_0) + \langle \nabla F(\theta_0), \theta_1 - \theta_0 \rangle$

$$\Delta(F) = \frac{F(\theta_0) - F(\theta_1)}{F(\theta_0)} \approx \gamma \frac{\|\nabla F(\theta_0)\|^2}{F(\theta_0)}$$

→ Relative change in tangent space: $\Delta(D\phi) = \frac{\|D\phi(\theta_1) - D\phi(\theta_0)\|}{\|D\phi(\theta_0)\|}$

C-S.

$$\leq \gamma \frac{\|D^2\phi(\theta_0)\| \cdot \|\nabla F(\theta_0)\|}{\|D\phi(\theta_0)\|}$$

→ Lazy Regime: $\boxed{\Delta(D\phi) \ll \Delta(F)}$

tangent space moves much slower than the loss itself.

→ If $R(f) = \|f - f^*\|^2$, $F(\theta) = \|\phi(\theta) - f^*\|^2$

$$\begin{aligned} \|\nabla F(\theta_0)\| &= \|D\phi(\theta_0)^\top (\phi(\theta_0) - f^*)\| \\ &\leq \|D\phi(\theta_0)\| \cdot \|\phi(\theta_0) - f^*\| \end{aligned}$$

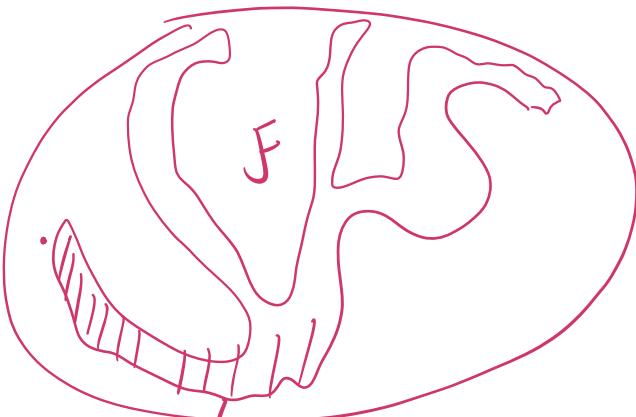
→ If we define the relative scale of ϕ at point θ_0 as

$$K_\phi(\theta) := \|\phi(\theta) - f^*\| \cdot \frac{\|D^2\phi(\theta)\|}{\|D\phi(\theta)\|^2}, \text{ the lazy regime}$$

is equivalent to

$$K_\phi(\theta) \ll 1$$

→ Why the name "lazy"? parameters almost do not move, because loss already converges at linear rate.



↳ Level set of training loss.

Theorem [COB, Thm 2.3 Simplified] Assume $\phi, D\phi$ are Lipschitz in a neighborhood of θ_0 . Let $\theta(t), \bar{\theta}(t)$ be the gradient flows associated with $\phi, \bar{\phi}$ respectively. Then for $t \leq C_\phi$, it holds

$$\frac{\|\phi(\theta(t)) - \bar{\phi}(\bar{\theta}(t))\|}{\|\phi(\theta_0) - f^*\|} \lesssim t^2 \underbrace{K_\phi(\theta_0)}_{\text{by relative scaling.}}$$

- Remarks:
- Constant C_ϕ is explicit
 - In the lazy regime, ($K_\phi(\theta_0) \ll 1$), the non-linear trajectory tracks the linearized traj.
 - In the over-parametrised regime, then the result can be extended uniformly in time.

(Theorem 2.4).

Proof: See the paper \rightarrow elementary.

Q: When does the lazy regime happen?

(i) Scaled Models: $\phi_\alpha(\theta) = \alpha \cdot \phi(\theta)$

$$K_{\phi_\alpha}(\theta_0) = \frac{1}{\alpha} \|\alpha \phi(\theta_0) - f^*\| \cdot \frac{\|D^2\phi(\theta_0)\|}{\|D\phi(\theta_0)\|^2}$$

in particular, when $\phi(\theta_0) = 0$, then

$$K_{\phi_\alpha}(\theta_0) = \frac{1}{\alpha} K_\phi(\theta_0).$$

(ii) Homogeneous Models: $\phi(\lambda \theta) = \lambda^r \phi(\theta) \quad \forall \theta \neq 0$

same as before: $K_\phi(\lambda \cdot \theta_0) = \lambda^{-r} K_\phi(\theta_0)$. (for 0 init).

(iii) Single hidden layer NN

$$\phi_m(\theta) = \alpha(m) \sum_{i=1}^m g(\theta_i), \quad \theta_i \sim \mu \text{ iid}$$

with $E_{\mu}[g(\theta)] = 0$, Dg Lipschitz.

$$K_m = E K_{\phi_m}(\theta) ?$$

$$\cdot |E \| \phi_m(\theta) \|^2| = m \cdot \alpha(m)^2 |E \|g(\theta)\|^2|$$

$\hookrightarrow \phi_m$ is a sum of iid terms 0-mean.

$$\cdot \quad D\phi_m(\theta) = \alpha(m) [Dg(\theta_1), \dots, Dg(\theta_m)]$$

$$\frac{D\phi_m(\theta) D\phi_m(\theta)^T}{m \cdot \alpha(m)^2} = \frac{1}{m} \sum_{i=1}^m Dg(\theta_i) Dg(\theta_i)^T$$

$\downarrow m \rightarrow \infty$

$$E[Dg(\theta) Dg(\theta)^T]$$

$$\Rightarrow E[\|D\phi_m(\theta)\|^2] = E[\|D\phi_m(\theta) D\phi_m(\theta)^T\|]$$

$$\sim m \alpha(m)^2 E[Dg(\theta) Dg(\theta)^T]$$

$$= m \alpha(m)^2 E[\|Dg(\theta)\|^2]$$

$$\cdot \|D^2\phi_m(\theta)\| = \sup_{u \in \mathbb{R}^{dxm}} \alpha(m) \sum_{i=1}^m u_i^T D^2g(\theta_i) u_i$$

$$(\|A\| = \sup_{\|x\| \leq 1} x^T A x)$$

$$\|u\| \leq$$

$$\leq \alpha(m) \sup_{\theta_i} \|D^2g(\theta_i)\|$$

$$\leq \alpha(m) \cdot \text{Lip}(Dg)$$

$$\hookrightarrow K_m = E[\|\phi(\theta) - f^*\|] \cdot \frac{\|D^2\phi_m(\theta)\|}{\|D\phi_m(\theta)\|^2}$$

$$\leq (E[\|\phi(\theta)\|] + \|f^*\|) \frac{E[\|D^2\phi_m(\theta)\|]}{E[\|D\phi_m(\theta)\|^2]}$$

$$\frac{(\sqrt{m} \cdot \alpha(m) \cdot A + B)}{m \cdot \alpha(m)^n \cdot D} \leq \frac{\alpha(m) \cdot C}{m \cdot \alpha(m)^n \cdot D}$$

$$\leq \frac{C_1}{\sqrt{m}} + \frac{C_2}{m \cdot \alpha(m)}$$

Conclusion: if $m \cdot \alpha(m) \rightarrow \infty$ as $m \rightarrow \infty$, then
the NN enters the lazy regime!

(in particular, NTIC scaling $\alpha(m) = \frac{1}{\sqrt{m}}$).

Is this good or bad?

- ⊕ It guarantees global convergence; optimisation is not the bottleneck.
- ⊖ We are doing linear learning: No features are learnt; RICS learning \rightarrow approximation becomes the bottleneck!

Q: What happens when $m \cdot \alpha(m) = \Theta(1)$, ie

$$\underline{\alpha(m) = 1/m} ?$$

Shallow NN as particle interaction systems

→ Now our model becomes $\phi(x; \overbrace{\theta_1 \dots \theta_m}^{\theta}) = \frac{1}{m} \sum_{i=1}^m g(x; \theta_i)$

$$g(x; \theta) = c \cdot V(\langle x, \omega \rangle + b), \theta = (a; b; c) \in \mathbb{R}^3$$

→ Integral Representation:

$$\phi(x; \theta) = \int_D g(x; \theta) d\mu_m(\theta)$$

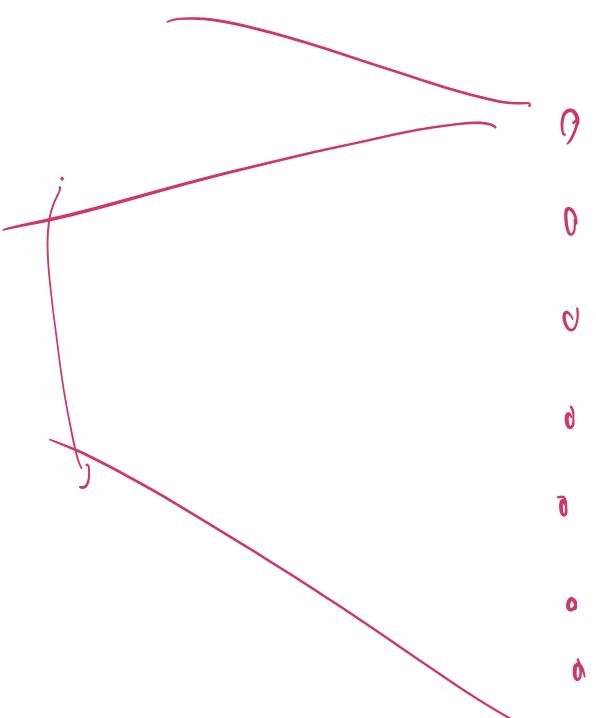
$$\left(\int f(x) d\delta_y(x) = f(y) \right)$$

$$\mu_m = \frac{1}{m} \sum_{i=1}^m \delta_{\theta_i}$$

the empirical
measure associated
with $\theta = (\theta_1, \dots, \theta_m)$

↳ In next lecture, we will study:

- (i) How the associated function is how bigger than the RICS associated with lazy training
- (ii) training dynamics in measure space
- (iii) Natural space where we have representation learning.



0