

Lecture 8 : Shallow NNs in the feature learning regime

In previous lecture: lazy training: a non-linear model $\Theta \mapsto \phi(\Theta)$ can behave as a linear one when its curvature around initialisation vanishes \rightarrow analysis through TIC / NTK.

\rightarrow For $\phi(\Theta) = \text{NN}(\text{shallow})$ of the form

$$f(\Theta) = \alpha(m) \sum_{j=1}^m \phi(\Theta_j) \quad \Theta_j \sim \mu, \quad \text{Lazy regime}$$

as soon as $m \cdot \alpha(m) \rightarrow \infty$.

\rightarrow Today Active regime $\alpha(m) = 1/m$.

Q1: What is the functional space?

Q2: Learning dynamics?

Shallow NN and Variation-Norm spaces

\rightarrow with choice $\alpha(m) = 1/m$, the model becomes

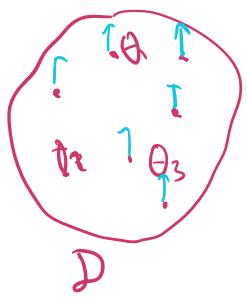
$$f(x; \Theta_1 \dots \Theta_m) = \frac{1}{m} \sum_{j=1}^m \phi(x; \Theta_j)$$

$x \in \mathbb{R}^d$

$$\underline{\phi(x; \theta)} = c \cdot \sigma(\langle x, a \rangle + b) \quad \theta = (a; b; c) \in \mathbb{R}^{d+1+1} \subset \mathbb{D}$$

$\begin{matrix} \downarrow \text{output weight} & \downarrow \text{input weight} & \downarrow \text{bias} \\ c \cdot \sigma(\langle x; 1 \rangle; \bar{\theta}) \end{matrix}$

$\rightarrow f$ is totally characterised by the "point cloud"



$$\theta_1 \dots \theta_m \rightarrow \theta_1, \dots, \theta_m \quad \alpha : \{1, m\} \rightarrow \{1, m\}$$

$$\vdots \quad \vdots \quad \vdots$$

$$\theta_m \quad \theta_{m+1}$$

$\rightarrow \mathcal{I}$ can completely describe f via the empirical measure

$$\mu_m = \frac{1}{m} \sum_{j=1}^m \delta_{\theta_j} \in \mathcal{P}(D) \quad \text{Def: Dirac Mass.}$$

$$\rightarrow f(x; \mu) = \int_D \phi(x; \theta) d\mu_m(\theta)$$

Integral representation of f .

$\rightarrow Q:$ Associated function class? Consider a non-negative regularisation term $V: D \rightarrow \mathbb{R}$ $V(\theta) = \|\theta\|^2$ [Weight decay]

$$V(\theta) = |C| \cdot \|\bar{\theta}\|$$

\rightarrow Define

$$F_{(V)} = \left\{ f(x) = \int_D \phi(x; \theta) d\mu(\theta) \text{ for } \mu \in \mathcal{P}(D) \text{ such that} \right. \\ \left. \int_D V(\theta) d\mu(\theta) < +\infty \right\}$$

\rightarrow Fact: F is a normed (=Banach) space, with

norm

$$\|f\| = \inf \left\{ \|d\mu : f(x) = \int \phi(x; \theta) d\mu\| \right\}$$

\hookrightarrow Ex: For $V(c; \bar{\theta}) = |c| \| \bar{\theta} \|$ this is called
the variation-norm space [Bach' 17]
a.k.a. "Baron" space.

Q: How does this space compare with the RKHS?

A Write $D = \mathbb{R} \times \bar{D}$ and assume $\bar{D} = \mathbb{S}^d$, unit $(d+1)$ -dim. sphere.

\rightarrow Fix an arbitrary measure π on \bar{D} , and consider

$$\mathcal{F}_\pi = \left\{ f(x) = \int_{\bar{D}} g(\bar{\theta}) \Gamma(\langle \tilde{x}, \bar{\theta} \rangle) d\pi(\bar{\theta}) ; \right. \\ \left. g \in L_2(\bar{D}; \pi) \right\} \quad \boxed{\left(\int |g(\bar{\theta})|^2 d\pi(\bar{\theta}) \right)}$$

\hookrightarrow We saw that \mathcal{F}_π is an RKHS, with kernel given by

$$K(x, x') = \left[E_{\bar{\theta} \sim \pi} \left\{ \Gamma(\langle \tilde{x}, \bar{\theta} \rangle) \Gamma(\langle \tilde{x}', \bar{\theta} \rangle) \right\} \right]$$

\rightarrow If $f \in \mathcal{F}_\pi$, let's define $\mu \in P(D)$

$$\mu(c, \bar{\theta}) = \delta(c - g(\bar{\theta})) \cdot \pi(\bar{\theta})$$

$$(i) \quad f(x) = \int_D \phi(x, \theta) \cdot d\mu(\theta) \quad V(c, \bar{\theta}) = |c| \\ D = \mathbb{R} \times \bar{D}$$

$$\text{(ii)} \quad \|f\|_F \leq \int V(c, \bar{\theta}) d\mu = \int_{\bar{D}} |g(\bar{\theta})| d\pi(\bar{\theta}) = \\ = \left| E_{\bar{\theta} \sim \pi} [g(\bar{\theta})] \right|$$

→ B) Jensen's ineq. $|EX|^2 \leq |Ex^2|$

→ This means that $f \in F$

$$\|f\|_F \leq \|f\|_{F_\pi} + \pi$$

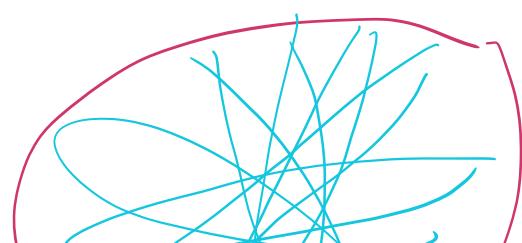
→ Therefore, $F = \bigcup_{\pi \in P(\bar{D})} F_\pi$

→ Ex $f(x) = \sigma(\langle \bar{x}, \bar{\theta}^* \rangle)$ ↗ Realised as $\mu =$
 ↗ $\|\bar{f}\|_F = 1$ $d\theta^*$
 ↘ Not in the RKHS
 F_π for $\pi \ll \text{unif.}_m$

→ Remark: Instance of a "separation" between two different hypothesis spaces

→ Remark: Same thing holds for $V(\theta) = \|\theta\|^2$ under homogeneous activations ($\sigma(tx) = t\sigma(x)$)

Q: Is it still possible to learn efficiently in this bigger



- space?
- ↳ In terms of generalisation, all the Rademacher-based analysis still work (by replacing $\|f\|_{F_\pi}$ by $\underline{\|f\|_F}$).
 - ↳ What about optim?

Shallow NNs (in active regime) as particle interaction systems

→ Focus on least-squares regression:

$$\min_{\theta_1 \dots \theta_m} L(\theta_1 \dots \theta_m) = \| f_\theta - f^* \|_V^2 + \lambda \left(\frac{1}{m} \sum_{j=1}^m V(\theta_j) \right)$$

→ By expanding the square, and using the linear $f_\theta = \frac{1}{m} \sum_{j=1}^m \phi(\theta_j)$, $\langle f, g \rangle_V := \frac{1}{n} \sum_{i=1}^n f(x_i)g(x_i)$ we can write $L \approx$

$$L(\theta_1 \dots \theta_m) = \cancel{\chi} - \frac{1}{m} \sum_{j=1}^m F(\theta_j) + \frac{1}{m^2} \sum_{i,j} K(\theta_i, \theta_j), \text{ with}$$

$$F(\theta) := \langle \phi(\theta), f^* \rangle_V - \frac{1}{2} V(\theta)$$

$$K(\theta, \theta') = \langle \phi(\theta), \phi(\theta') \rangle_V$$

↳ Interpret L as the Hamilton/Energy of a system of m particles?

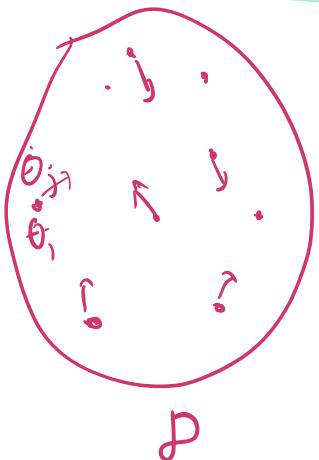
system of m particles, $(\theta_1 \dots \theta_m)$, in \mathcal{D} .

$F \longleftrightarrow$ external field.

$K \longleftrightarrow$ interaction kernel.

→ Gradient flow with respect to $\theta_1 \dots \theta_m$ gives us

$$\dot{\theta}_j = -\frac{m}{2} \nabla_{\theta_j} L(\theta_1 \dots \theta_m) = +\nabla F(\theta_j) - \frac{1}{m} \sum_{j=1}^m \nabla K(\theta_j, \theta_j)$$



↳ This is Lagrangian perspective. (tracking each particle).

↳ This system of m d-dim ODEs is very complex. L is non-convex wrt $\theta_1 \dots \theta_m$.

→ Can we aim to obtain a "simple" collective behavior?

→ Let's consider the Eulerian perspective instead.

$$\vec{\theta} = (\theta_1 \dots \theta_m) \in \mathcal{D}^m \iff \mu = \frac{1}{m} \sum_{j=1}^m \delta_{\theta_j} \in \mathcal{P}(\mathcal{D})$$

$$f(x) = \int_D \phi(x; \theta) d\mu$$

→ The Energy/loss now becomes-

$$L(\theta_1 \dots \theta_m) = \cancel{E} - \frac{1}{m} \sum_{j=1}^m F(\theta_j) + \frac{1}{m^2} \sum_{i,j} K(\theta_i, \theta_j), \text{ with}$$

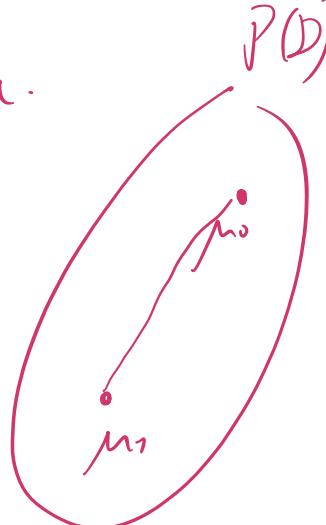
$$\mathcal{L}[\mu] = \underbrace{-2 \int_D F(\theta) d\mu}_{\mu \in \mathcal{P}(D)} + \iint_{D \times D} K(\theta, \theta') d\mu(\theta) d\mu(\theta')$$

$\| \int \phi d\mu - f^* \|_V^2 + \int V(\theta) d\mu$

→ \mathcal{L} is convex with respect to μ .

↳ If $\mu_0, \mu_1 \in \mathcal{P}(D)$ and

$$\mu_t = t\mu_1 + (1-t)\mu_0$$



then $\mathcal{L}[\mu_t] \leq t\mathcal{L}[\mu_1] + (1-t)\mathcal{L}[\mu_0]$

→ Is this convexity "felt" by our dynamics?

Relationship between GD/CF and a choice of metric?

Proximal View point of GD

→ The standard (Euclidean) GD step

$$\theta_{t+1} = \theta_t - \gamma \nabla f(\theta_t)$$

$$= \underset{\theta}{\operatorname{argmin}} \left\{ f(\theta_t) + \langle \nabla f(\theta_t), \theta - \theta_t \rangle + \frac{1}{2\gamma} \|\theta - \theta_t\|^2 \right\}$$

Linear approx of f at θ_t proximity term.

→ Proximity can be measured by different metrics!

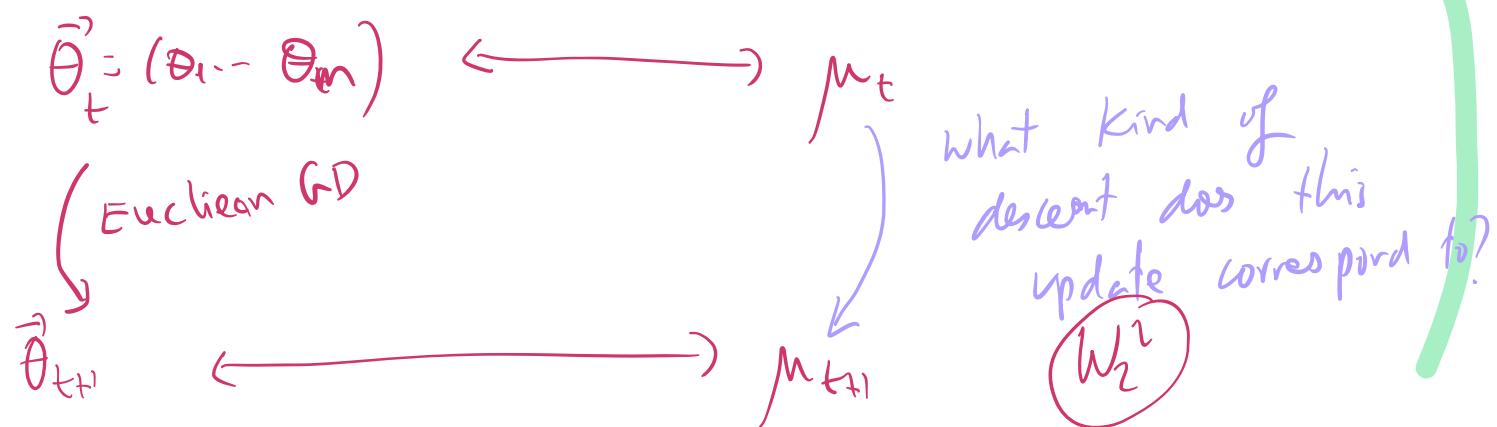
→ The extension of GD to general (non-Euclidean) metrics is given by Mirror Descent [Nemirovski & Yudin 83]

$$\theta_{t+1} = \arg \min_{\theta \in \Theta} \left\{ \langle \nabla f(\theta_t), \theta - \theta_t \rangle + \frac{1}{\eta_t} D(\theta; \theta_t) \right\}$$

$$\approx \arg \min_{\theta \in \Theta} \left\{ f(\theta) + \frac{1}{\eta_t} D(\theta; \theta_t) \right\}$$

Bregman
Divergence.

→ Back to our setting in measure space.



→ Recall that $\dot{\theta}_j = \nabla F(\theta_j) - \frac{1}{m} \sum_{j'=1}^m \nabla K(\theta_j, \theta_{j'})$
using $\mu_t = \frac{1}{m} \sum_{j=1}^m \delta_{\theta_j(t)}$, we have

→ $\dot{\theta}_j(t) = \nabla F(\theta_j(t)) - \int_D \nabla K(\theta_j^{(t)}, \theta) d\mu_t(\theta)$
 $:= -\nabla G(\theta_j(t); \mu_t)$, with

$$G(\theta; \mu) = -F(\theta) + \int K(\theta, \theta') d\mu(\theta')$$

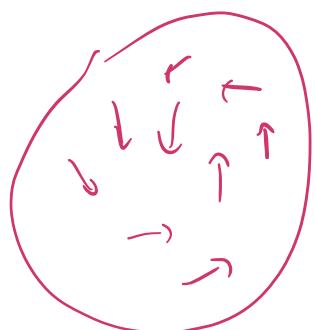
$$\text{Recall } \mathcal{L}[\mu] = -\int_{\mathbb{P}} F(\theta) d\mu + \iint K(\theta, \theta') d\mu d\mu(\theta')$$

$$\frac{\delta \mathcal{L}}{\delta \mu(\theta)} \Rightarrow -2F(\theta) + 2 \int_{\mathbb{P}} K(\theta, \theta') d\mu(\theta')$$

G is the first variation of \mathcal{L} w.r.t. μ

$$G = \frac{\delta \mathcal{L}}{\delta \mu}.$$

→ In physics, G is the instantaneous potential of the system.



→ Consider a (smooth) test function

$$\chi : \mathbb{D} \rightarrow \mathbb{R}$$

$$\begin{aligned} \left(\text{time derivative} \right) \int_{\mathbb{D}} \chi(\theta) d\mu_t(\theta) &= \frac{1}{m} \sum_{j=1}^m \chi(\theta_j(t)) \\ \int_{\mathbb{D}} \chi(\theta) \partial_t d\mu_t(\theta) &= \frac{1}{m} \sum_{j=1}^m \langle \nabla \chi(\theta_j(t)), \dot{\theta}_j(t) \rangle \quad \left. \right) \text{time derivative} \\ &= -\frac{1}{m} \sum_{j=1}^m \langle \nabla \chi(\theta_j(t)), \nabla G(\theta_j(t), \mu_t) \rangle \end{aligned}$$

$$= - \int_D \langle \nabla \chi(\theta), \nabla G(\theta; \mu_t) \rangle d\mu_t(\theta)$$

↳ This is the continuity equation / transport equation /

$$\partial_t \mu_t = \operatorname{div}(\underline{\nabla G(\theta; \mu_t)} \cdot \mu_t) \quad \text{Lagrangian equation.}$$

↳ Time-dependent velocity field.

PDE in the space of measures.

Q: Is this PDE a gradient flow? wrt which metric?

↳ The proximal interpretation of this dynamics is

$$\mu_{t+1} = \underset{\mu \in P(D)}{\operatorname{argmin}} \quad \left\{ L[\mu] + \frac{1}{2\eta} W_2^2(\mu_t, \mu) \right\}$$

$$D \xrightarrow{\text{metric space}} P(D)$$

Remarks: This formulation of SNN dynamics was developed in [Chizat, Bach' 18]

[Mei, Nguyen, Montanari' 18]

[Rotskoff, EPE' 18]

[Sriram, Spiliopoulos' 18]

- Description is exact

• Adding noise to the dynamics ($\text{GF} \rightarrow$ Langevin dynamics)

↓
McKean-Vlasov
equation.

→ Two main questions

- (i) Under what conditions this PDE converges to the global optimum of \mathcal{L} (convergence in t)
- (ii) Effect of over-parametrisation (convergence in m)?

.