

# Mathematics of Deep Learning

## Course Logistics

Lectures Th 4:55 → 6:55 (5 min break)

Recitation Mondays @ 3pm - 4pm

Carles Domingo (PhD student).

Office Hours Fridays at 9am → 10:30ish.

Evaluation : Final Project 90%  
Participation 10% (OH, Campus Wire, Class).

- Final Project: In-depth survey on a selected topic.  
NOT a Research Project. By groups of 2, 4
- List of topics are in the website (tentative)
- Abstract : due mid-semester

Bibliography: Two major references :  
    { · ML from first principles  
        by Francis Bach  
    { · Lecture Notes by  
        Matus Telgarski

## Main Objective of this course

- When and why can neural nets learn in high-dimension.

- Focus on the simplest learning setup: supervised learning.
  - Focus on the simplest class of NN: shallow nets.
- 50%
- Study role of architecture (depth)
  - Study generative modeling.
  - Discuss old & new results.

## Lecture 1 : The curse of dimensionality ; NNs and approximation.

### Basic Supervised Learning Setup

→ Given data  $\{(x_i, y_i)\}_{i=1 \dots n}$  with  $x_i \in \mathcal{X}$  (high-dim)  
 $y_i \in \mathcal{Y}$  (label target space)

estimate a mapping  $f: \mathcal{X} \rightarrow \mathcal{Y}$  that

generalises on future data "Glorified" interpolation.

→ IID assumption: data is drawn iid from an underlying distribution

✓ on  $\mathcal{X} \times \mathcal{Y}$ .  $(x_i, y_i) \stackrel{\text{iid}}{\sim} \mathcal{D}$

→ Pointwise loss  $l: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$

Given  $f: \mathcal{X} \rightarrow \mathcal{Y}$ , this defines

$$\rightarrow R(f) := \mathbb{E}_y [l(f(x), y)] \quad \text{"population risk".}$$

$$\rightarrow \hat{R}(f) := \frac{1}{n} \sum_{i=1}^n l(f(x_i), y_i) \quad \text{"empirical risk".}$$

$$\mathbb{E}\hat{R} = R \quad (\hat{R} \text{ is an unbiased estimator of } R).$$

$\rightarrow$  Hypothesis Space  $\mathcal{F} = \{ f: \mathcal{X} \rightarrow \mathcal{Y} \}$ . Assume

$\mathcal{F}$  is a normed space (Banach): for each  $f \in \mathcal{F}$

we can assign a "complexity" measure  $(\gamma(f))$   
 $\|f\|_{\mathcal{F}}$

$$\mathcal{F}_{\delta} = \{ f \in \mathcal{F}; \|f\|_{\mathcal{F}} \leq \delta \}.$$



Examples of complexity:

$\rightarrow$  (Euclidean) norm of the weights in a NN

$\rightarrow$  Number of parameters

$\rightarrow$  Number of gradient-descent iterations for models in  $\mathcal{F}$  reached by gradient-descent.

$\rightarrow$  Empirical Risk Minimisation (ERM)

We search for small empirical risk with small complexity

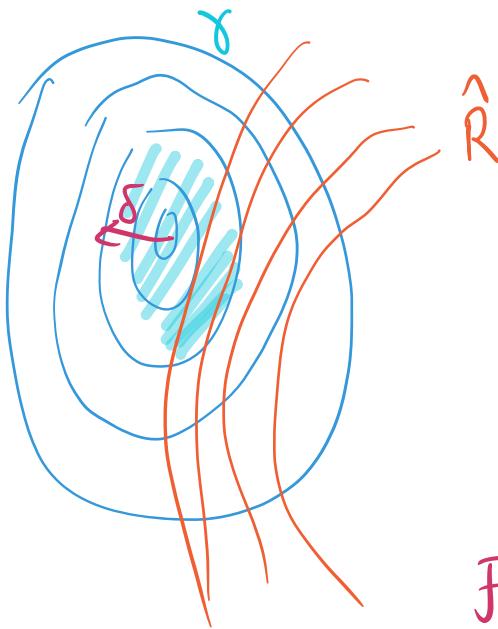
$$(C) \min_{\gamma(f) \leq \delta} \hat{R}(f)$$

$$(P) \min_{f \in \mathcal{F}} \hat{R}(f) + \lambda \gamma(f)$$

$$(I) \min \gamma(f)$$

$$\text{fst. } f(x_i) = y_i \Rightarrow \hat{R}(f) = 0$$

"interpolant" or "memorization" regime.



$$F_\delta = \{f \in F; \gamma(f) \leq \delta\}$$

### Basic Decomposition of Risk

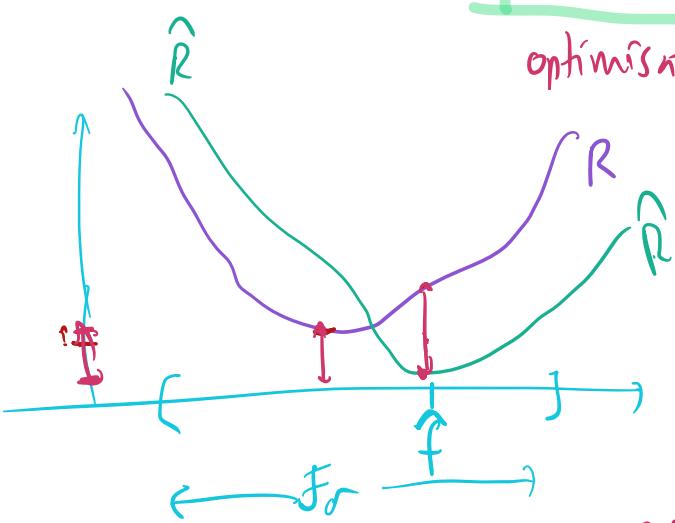
- let  $\hat{f} \in F_\delta$  produced by an arbitrary algorithm.

$$R(\hat{f}) - \inf_{f \in F} R(f) = R(\hat{f}) - \inf_{f \in F_\delta} R(f)$$

$$+ \inf_{f \in F_\delta} R(f) - \inf_{f \in F} R(f)$$

$$= [R(\hat{f}) - \hat{R}(\hat{f})] + [\hat{R}(\hat{f}) - \inf_{f \in F_\delta} \hat{R}(f)] + \inf_{f \in F_\delta} \hat{R}(f) - \inf_{f \in F} R(f) + \text{approx. Err.}$$

optimisation error.



$$\sup_{f \in F_\delta} |R(f) - \hat{R}(f)|$$

$$\left| \inf_{f \in F_\delta} R(f) - \inf_{f \in F_\delta} \hat{R}(f) \right|$$

$$\left| \inf_{f \in F_\delta} R(f) - \inf_{f \in F_\delta} \hat{R}(f) \right| \leq \sup_{f \in F_\delta} |R(f) - \hat{R}(f)|$$

$$f^* = \operatorname{argmin}_f R$$

Suppose  
 $\inf_{\hat{\mathcal{F}}^*} \hat{R} < \inf_{\mathcal{F}^*} R$

$\hat{f}$

$$\left| \inf_{\hat{\mathcal{F}}^*} \hat{R} - \inf_{\mathcal{F}^*} R \right| \leq \left| \hat{R}(f^*) - R(f^*) \right|$$

$$\leq \sup_{f \in \mathcal{F}^*} |R(f) - \hat{R}(f)|$$

$$\leq \underbrace{2 \sup_{f \in \mathcal{F}^*} |R(f) - \hat{R}(f)|}_{\text{(uniform)} \text{ statistical error}} + \varepsilon_{\text{opt}} + \varepsilon_{\text{appx.}}$$

$$\text{"Generalisation" error} \leq \underbrace{\text{stat}(\delta)}_{\text{error}} + \underbrace{\frac{\text{opt}(\delta)}{\text{error}}}_{\text{error}} + \underbrace{\text{appx error}(\delta)}_{\text{error}}$$

$$\begin{bmatrix} \text{As } \delta \rightarrow 0, & \downarrow & ? & \uparrow \\ \text{As } \delta \rightarrow n & \uparrow & ? & \downarrow \end{bmatrix}$$

$$\rightarrow \underbrace{\text{Statistical Error.}}_{\text{Random quantity.}} \quad \sup_{f \in \mathcal{F}^*} |R(f) - \hat{R}(f)|$$

It depends on two basic quantities:

Complexity:  $\delta$  and  $n$ : number of data points.

Fix  $f$ .

$$\hat{R}(f) - R(f) =$$

$$= \frac{1}{n} \sum_{i=1}^n \underbrace{[l(f(x_i), y_i) - E_l(f(x), y)]}_{\text{Random variable.}}$$

$$\frac{1}{n} \sum_{i=1}^n (z_i - \mathbb{E}(z_i))^2 \quad z_1, \dots, z_n \text{ are iid.}$$

$$z \quad \sigma^2 = \mathbb{E}((z - \mathbb{E}z)^2) \quad \mathbb{E}z = m$$

$$\text{Var} \left[ \frac{1}{n} \sum_{i=1}^n (z_i - m) \right] = \frac{\sigma^2}{n}$$

$$\Rightarrow |R(f) - \hat{R}(f)| \sim \frac{\text{std}(\ell(f(x), y))}{\sqrt{n}} \quad \text{Monte-Carlo Estimate.}$$

non-asymptotic; we can quantify with tail probability Bounds.

→ From point-wise control (fixing  $f$  before looking at the data) to uniform control.

Requires concentration tools from empirical process (e.g. Rademacher complexity); but final bound

$$\epsilon_{\text{stat}} \sim \frac{g(f, \delta)}{\sqrt{n}}$$

## The Curse of Dimensionality

Q: How do approximation and statistical errors behave as the input dimension increases?

Statistical curse

We observe  $\{(x_i, f^*(x_i))\}_{i=1 \dots n}$   $x_i \sim N(0, I_d)$   
 $f^*$  is the target function.

Q: How many samples are needed to estimate  $f^*$  up to error  $\epsilon$ , ie  $|E_x |f(x) - f^*(x)|^2 \leq \epsilon$ ?

↳ suppose first that  $f^*$  is linear;  $f^*(x) = \langle x, \theta^* \rangle$

$$f = \left\{ f: \mathbb{R}^d \rightarrow \mathbb{R} ; f(x) = \langle x; \theta \rangle \right\} \cong \mathbb{R}^d$$

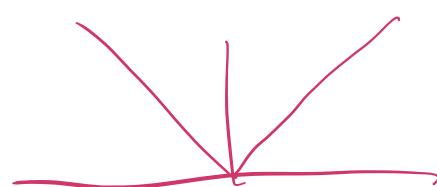
$\theta \in \mathbb{R}^d$   
 $x \in \mathbb{R}^d$

A:  $n = d$  (we just need to solve a linear system of equations).

Remark:  $f^* = \hat{\varphi}(\langle x, \theta^* \rangle)$  ( $\varphi$  even and smooth).

A:  $n = \underbrace{d+1}_{\text{is sufficient for exact recovery}}$

$$\varphi(t) = |t|$$



↳ Suppose now that  $f^*$  is only locally linear, ie

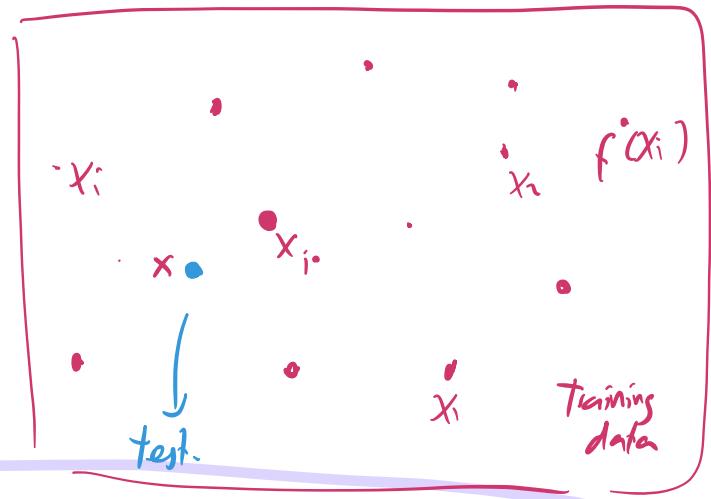
$f^*$  is  $\beta$ -Lipschitz:  $|f(x) - f(x')| \leq \beta \cdot \|x - x'\|$

$$f = \left\{ f: \mathbb{R}^d \rightarrow \mathbb{R} ; f \text{ is Lipschitz + Bounded} \right\}$$

(now  $f$  is a Banach space;  $\|f\| = \text{Lip}(f) + \|g\|_\infty$ .

$$\text{Lip}(f) = \inf \left\{ \beta ; |f(x) - f(x')| \leq \beta \|x - x'\| \right\}$$

$\rightarrow$



$$\hat{f} = \underset{f \in F}{\operatorname{argmin}} \left\{ \text{Lip}(f) ; f(x_i) = \hat{f}(x_i) \forall i \right\}$$

Recall MSE:  $E_{x \sim \nu} |f(x) - \hat{f}(x)|^2$

$$x \sim \nu \quad |f(x) - \hat{f}(x)| \leq |\hat{f}(x) - \hat{f}(x_{i^*})| + |\hat{f}(x_{i^*}) - \hat{f}(x_i)| + |\hat{f}(x_i) - f(x_i)|$$

$$|\hat{f}(x_{i^*}) - f(x_i)| \leq \beta \|x_{i^*} - x_i\|$$

$$|\hat{f}(x_{i^*}) - \hat{f}(x)| \leq \text{Lip}(\hat{f}) \|x_{i^*} - x\|$$

$$|\hat{f}(x) - f(x)| \leq 2\beta \|x_{i^*} - x\|$$

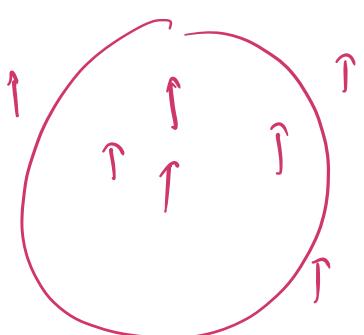
$$\Rightarrow E |\hat{f}(x) - f(x)|^2 \leq 4\beta^2 E_x \|x - x_{i^*}\|^2$$

$$x \sim N(0, I)$$

$$W_2^2(\nu, \hat{\nu}_n) \simeq \boxed{n^{-1/d}}$$

$$\nu = N(0, I_d)$$

$$\hat{\nu}_n = \frac{1}{n} \sum_{i=1}^n \delta_{x_i}$$

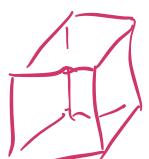


If we want  $\|f - f^*\|^2 \leq \epsilon$

$$\epsilon \simeq n^{-1/d} \rightarrow n \sim \epsilon^{-d}$$

This sample complexity is cursed by dimension.

→ Is this necessary? Can we learn without exponential dependency on dimension?



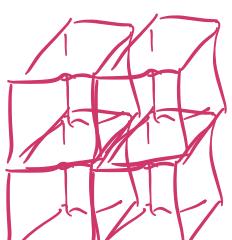
$$B = [-\frac{1}{2}, \frac{1}{2}]^d$$



$$\psi : B \rightarrow \mathbb{R}$$

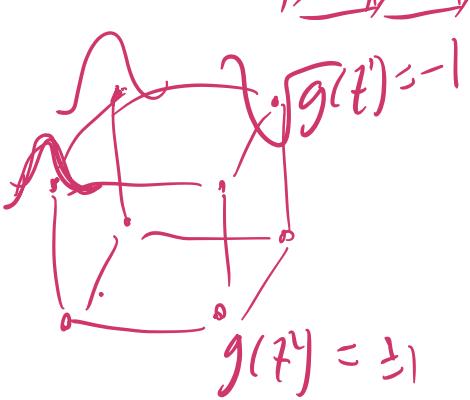
$$\underline{\psi(x) = \text{dist}(x; \partial B)}$$

Fact:  $\psi$  is 1-Lipschitz and supported in  $B$ .



$2^d$  copies of  $B$ .

$$z = z_1 \dots z_d \quad z_i = \pm 1$$

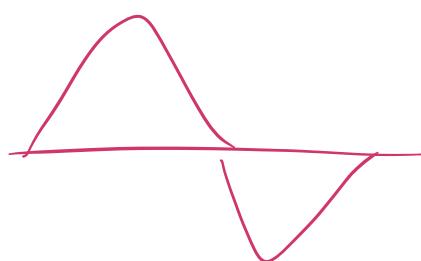
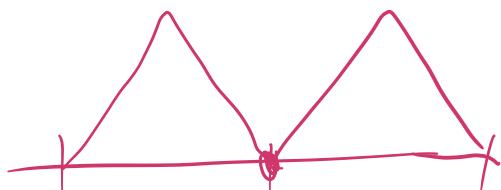


$$z \in \{-1, +1\}^d$$

$$f^*(x) = \sum_{z \in \{\pm 1\}^d} g(z) \psi(x-z)$$

Claim:  $f^*$  is also 1-Lipschitz

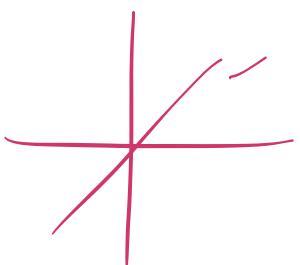
$$f_1, f_2$$



→ If  $n < 2^d$  then any estimator will incur in a relative error

$$\mathbb{E}_x |\hat{f}(x) - f^*(x)|^2 = \Theta(1) \cdot \mathbb{E} |f^*(x)|^2$$

→



most quadrants will be empty of data.

→ To summarize: linear functions (global reg)  $\rightarrow n=d$  O(n)

Lipschitz functions (local reg)  $\rightarrow n \approx \epsilon^{-d}$

↳ Too hard to learn

Next lectures → Can <sup>Shallow</sup> NNs define "interesting" hypothesis spaces somewhere in between these two extremes?