

## Lecture 2 : Approximation Main Ingredients

Goals for today: → Notion linear vs non-linear approximation  
→ Get familiar with Fourier represent.  
(→ from Fourier to shallow NNs)

- Signals that we measure as inputs to our ML algorithms (images, video, speech, ...) they are inherently defined over continuous domains
  - Similar functions that we aim to learn with  $f: \mathcal{X} \rightarrow \mathcal{Y}$   $t \in \mathbb{R}$
- $f: \mathcal{X} \rightarrow \mathcal{Y}$   
space of images      label space.
- 
- in both cases  $\mathcal{X}(t), f$  belongs to infinite-dimensional spaces.

Q: How to represent / manipulate signals in infinite-dim spaces?

→ This is what approximation theory tries to answer.

→ General set-up:  $\mathcal{F}$ : normed vector space of signals.

$\mathcal{D} = \{ \phi_\lambda \}_{\lambda \in \Lambda}, \phi_\lambda \in \mathcal{F}$   $\lambda$ : indexing a family/dictionary of functions in  $\mathcal{F}$ .

e.g.:  $\Lambda = \mathbb{N}, \mathbb{Z}, \mathbb{R}$ .

Given  $f \in \mathcal{F}$ , we want to express it as a linear combination of dictionary elements in  $\mathcal{D}$ :

$$\begin{cases} f = \sum_{\lambda \in \Lambda} \alpha_\lambda \cdot \phi_\lambda & (\Lambda \text{ is countable}) \\ f = \int_{\Lambda} g(\lambda) \cdot \phi_\lambda d\lambda & (\text{if } \Lambda \text{ is } \mathbb{R}, \mathbb{R}^d) \end{cases}$$

The information on  $f$  is in the "representation"  $\alpha_\lambda$   
 $g(\lambda)$

Goal: For a target family of functions  $\mathcal{F}^* \subseteq \mathcal{F}$ , obtain efficient approximation for any  $f \in \mathcal{F}^*$  using the dictionary  $\mathcal{D}$ .

→ By "efficient", we mean identifying a finite set  $\mathcal{S}_N$  of  $N$  elements in  $\mathcal{D}$  such that

$$U_N = \{ \tilde{f} \in \mathcal{F}; \tilde{f} = \sum_{\lambda \in \mathcal{S}_N} \alpha_\lambda \phi_\lambda \} \subseteq \mathcal{F}$$

is "close" to  $\mathcal{F}^*$ :

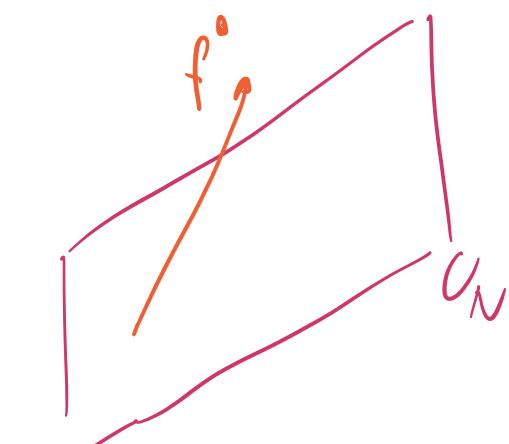
$\forall f \in F^*$ ,  $\inf_{g \in U_N} \|f - g\|_F$  small  $\rightarrow f^*$



→ Observe that we are selecting  $U_N$  that needs to work for any  $f^* \in F^*$ . Once we have chosen  $U_N$ , finding the best approximation of  $f$  in  $U_N$ , amounts to solving a least-squares problem.

$$\min_{g \in U_N} \|f - g\| =$$

$$= \min_{\alpha \in \mathbb{R}^N} \|f - \sum_{k \in U_N} \alpha_k \phi_m\|$$



Example: Suppose  $F = L^2([0, 1])$  squared-integrable functions-

$$\langle f, g \rangle_F = \int_0^1 f(t) g(t) dt.$$

Hilbert space.

$\mathcal{D} = \{\phi_m\}_{m \in \mathbb{N}}$  orthonormal basis of  $F$ .

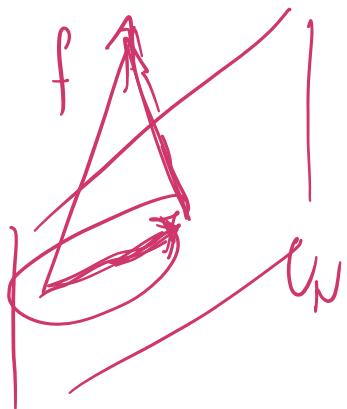
$\{\phi_m\}_{m \in \mathbb{N}_N}$  orthonormal basis of  $\text{span}\{\phi_m\}_{m \in \mathbb{N}_N} = U_N$ .

The best approximation

$\inf_{g \in U_N} \|f - g\|_F^2$  can be computed explicitly:

$$\text{if } f = \sum_{m=0}^{\infty} \langle f, \phi_m \rangle \phi_m$$

$$f_N := P_{U_N} f = \sum_{m \in \mathbb{N}_N} \langle f, \phi_m \rangle \phi_m$$



$$f - f_N = \sum_{m \notin \mathbb{N}_N} \langle f, \phi_m \rangle \phi_m, \quad \text{with apprx error}$$

$$\epsilon_\ell(N, f) = \|f - f_N\|^2 = \boxed{\sum_{m \notin \mathbb{N}_N} |\langle f, \phi_m \rangle|^2}$$

- Since  $\|f\|^2 = \int |f(t)|^2 dt = \boxed{\sum_{m=0}^{\infty} |\langle f, \phi_m \rangle|^2 < +\infty}$

We have  $\lim_{N \rightarrow \infty} \epsilon_\ell(N; f) = 0$  if  $f$ .

- Suppose w.l.o.g  $\mathbb{N}_N = \{1, \dots, N\}$ .

Decay of  $\epsilon_{\ell}(N, f)$  as  $N$  increases

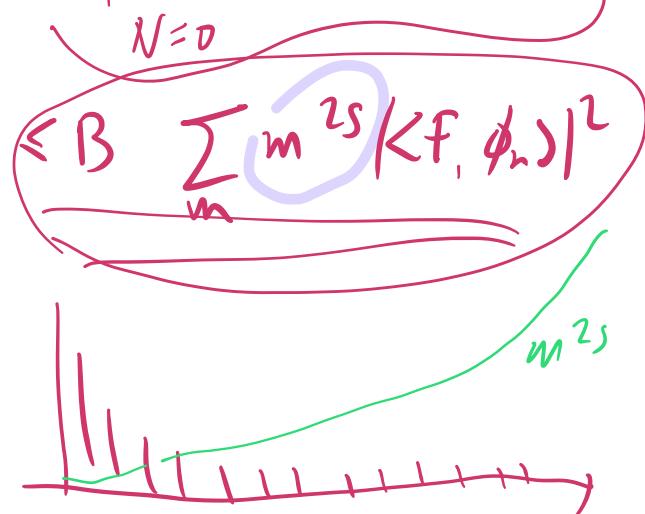
↓  
Decay of  $|\langle f, \phi_m \rangle|$  as  $m$  increases.

Fact: For any  $s > 1/2$ , there exists  $A, B > 0$   
such that if

$$\sum_{m=0}^{\infty} |m|^{2s} |\langle f, \phi_m \rangle|^2 < +\infty$$

then

$$A \sum_m m^{2s} |\langle f, \phi_m \rangle|^2 \leq \sum_{N=0}^{\infty} N^{2s-1} \epsilon_{\ell}(N, f) \leq$$

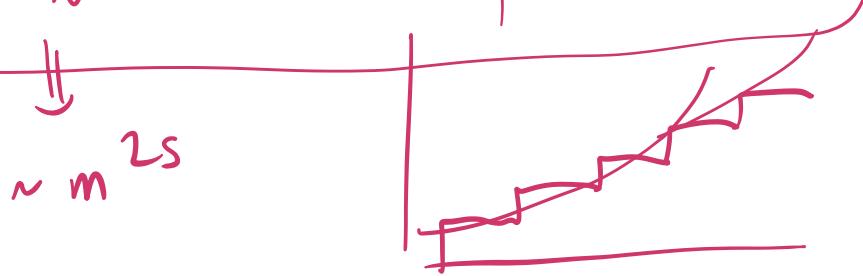


Thus  $\epsilon_{\ell}(N, f) = O(N^{-2s})$ .

Proof: Apply definition  $\epsilon_{\ell}(N, f) = \sum_{m=N}^{\infty} |\langle f, \phi_m \rangle|^2$

$$\sum_{N=0}^{\infty} N^{2s-1} \sum_{m=N}^{\infty} \beta_m^2 = \sum_{m=0}^{\infty} \beta_m^2 \sum_{N=0}^{m-N} N^{2s-1}$$

$$\int_0^m x^{2s-1} dx \leq \sum_{N=0}^m N^{2s-1} \leq \int_0^{m+1} x^{2s-1} dx$$



let's verify  $\epsilon_e(N, f) = o(N^{-2s})$ .

Observe that  $\epsilon_e(N, f) \leq \epsilon_e(m, f)$   $\forall m \leq N$

$$\epsilon_e(N, f) \sum_{m=N/2}^N m^{2s-1} \stackrel{\text{blue hand icon}}{\leq} \sum_{m=N/2}^N \epsilon_e(m, f) m^{2s-1}$$

Since  $\sum_{m=1}^{\infty} \epsilon_e(m, f) m^{2s-1} < +\infty$ , it follows

that  $\lim_{N \rightarrow \infty} \sum_{m=N/2}^N \epsilon_e(m, f) m^{2s-1} = 0$ .

Moreover, since  $\exists C > 0$   $\sum_{m=N/2}^{N-1} m^{2s-1} \geq C \cdot N^{2s}$

$$\lim_{N \rightarrow \infty} \epsilon_e(N, f) \cdot N^{2s} = 0. \quad \square$$

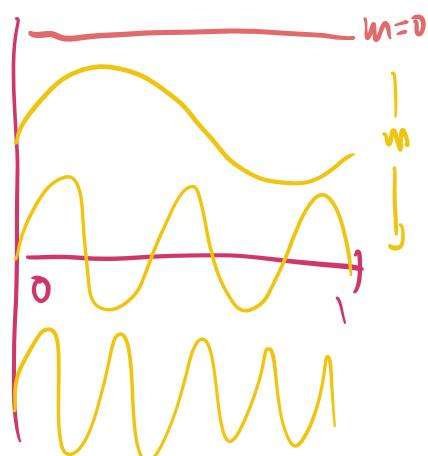
In summary, the linear approximation error of  $f$  in  $D$  decays at least as  $N^{-2s}$  if  $f$  belongs to the set

$$W_{D,s} := \left\{ f \in \mathcal{F}; \quad \sum m^{2s} |(f, g_m)|^2 < \infty \right\}$$

Q: What is a "canonical" example of orthonormal basis of  $L^2([0,1])$  where this linear approx. scheme is appropriate?

The Fourier Basis:  $\{e^{i2\pi mt}\}_{m \in \mathbb{Z}}$  is an orthonormal basis of  $L^2([0,1])$ .

$$f(t) = \sum_{m=-\infty}^{\infty} \underbrace{\langle f, \phi_m \rangle}_{\text{Fourier coefficients of } f} \phi_m \leftarrow (*)$$



Next week: More flavors of Fourier transform.

Recall:

Decay of Fourier coefficients  $(f_m)_m$

↑  
smoothness of  $f$

- Linear approximation of  $f$  in Fourier: keep the  $N$  lowest frequencies and discard the rest.

$$f_N(t) = \sum_{|m| \leq N/2} \hat{f}_m \phi_m \quad \begin{array}{l} \text{low-pass filter.} \\ \text{N: Sampling rate.} \end{array}$$

- Recall: The Fourier transform of  $f'$  (derivative of  $f$ )

is  $\int |f(w)|^2 \hat{f}(w) dw$  (where  $\hat{f}(w)$  is the F.T. of  $f$ ).

so  $\sum_m |m|^{2s} |\hat{f}_m|^2 < +\infty$  roughly measures

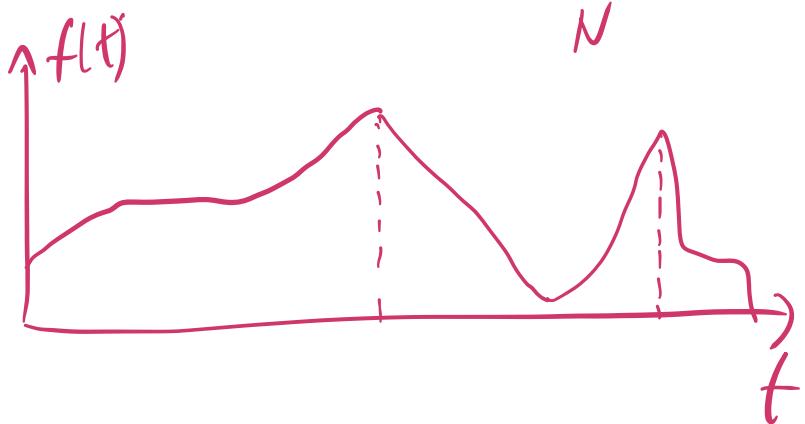
if  $f$  admits  $s$  derivatives (Sobolev Space  $W^s([0,1])$ )

↳  $\|f\|_2^2 = \sum_m |\hat{f}_m|^2$  (Parseval Identity).

( $s=1$ )  $\sim \|f'\|_2^2 = \sum_m |m^2 |\hat{f}'(m)|^2$

In conclusion, we can characterize  $f \in W^s([0,1])$

if and only if  $\sum_N N^{2s-1} \epsilon_\ell(N, f) < +\infty$



↳ When  $f$  is discontinuous, (or  $f'$  discontinuous)  
rates of approximation cannot leverage local  
regularity of  $f$ .

→ Linear approximation exhibits no adaptivity.

The subspace  $U_1 \subset L^2([0,1])$  is chosen one and

for all target function in  $\mathcal{F}$ .

Q: Can we do better?

## Non-linear approximation

In the previous linear approx, we can characterise the best approximate  $N$ -dim subspace as

$$\inf_{\substack{U; \dim(U) \leq N}} \sup_{f \in \mathcal{F}^+} \|f - P_{U_N} f\|^2 = \varepsilon_e(N; \mathcal{F}^+) \quad \text{(N-width class)}$$

$$\sup_{f \in \mathcal{F}^+} \inf_{\substack{U; \dim(U) \leq N}} \|f - P_{U_N} f\|^2 = \varepsilon_n(N; \mathcal{F})$$

Non-linear approx

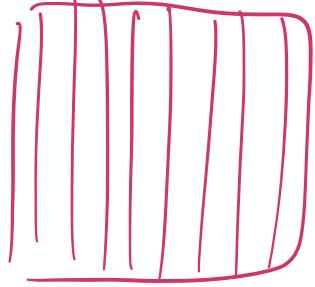
$$\inf_x \sup_y F(x, y) \quad \text{vs} \quad \sup_y \inf_x F(x, y)$$

Fix  $(x, y)$

$$\left[ \inf_{x'} \left[ \inf_{y'} F(x', y') \right] \right] \leq F(x, y) \leq \left[ \sup_{y'} \left[ \sup_{x'} F(x', y') \right] \right]$$

$$f(x, y) \quad g(y) \leq h(x)$$

$$\Rightarrow \sup_y g(y) \leq \inf_x H(x)$$



→ First instance of a fundamental tradeoff that we will encounter in DL / ML.

Better mathematical / Statistical performance

vs

More computationally challenging.

→ How does non-linear approximation work when  $\mathcal{D}$  is still an orthonormal basis?

→  $f \sim \beta_m = \langle f, \phi_m \rangle_m, m \in \mathbb{N}.$

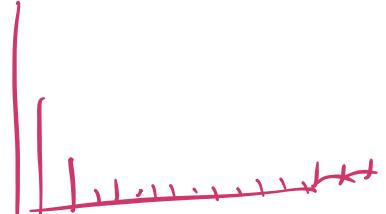
→ let us sort the coefficients  $|\beta_m| = |\langle f, \phi_m \rangle|$

$\beta_{s(m)} = \bar{\beta}_m$  with  $\bar{\beta}_{m+1} > \bar{\beta}_m \neq 0$ .

→ The best  $N$ -term non-linear approximation is

$f_N = \sum_{m=1}^N \beta_{s(m)} \phi_{s(m)},$  with error

$$\varepsilon_n(N, f) = \sum_{m=N+1}^{\infty} |\beta_{s(m)}|^2$$



Similarly as before, the rate of approximation  $\epsilon_n(N, f)$  is driven by the decay of  $\beta_m$ .

Q: When is non-linear approximation provably better than linear approximation?

A: If we work with Fourier representation, gains are not relevant for "natural" classes of signals

Lacunary Series



However, if we go beyond Fourier dictionaries, and into wavelet representations, then we obtain substantial separation between linear and non-linear approximation.

Theorem: If  $f$  has  $K$  discontinuities in  $[0,1]$

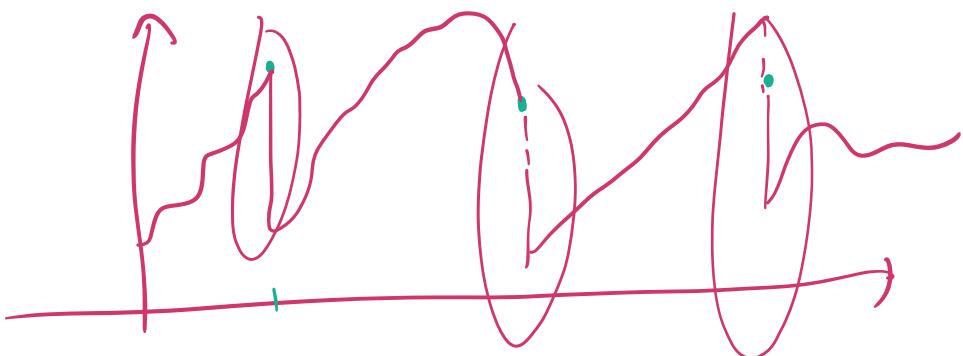
and is uniformly Lipschitz  $\alpha$  between these discontinuities, then  $|f(x) - f(x')| \leq C \cdot |x - x'|^\alpha$

$$\epsilon_Q(N, f) = \Theta\left(K \cdot \|f\|_{C^\alpha}^2 N^{-1}\right) \begin{matrix} \text{vs} \\ \text{non-adapt.} \end{matrix}$$

$$\epsilon_n(N, f) = \Theta\left(\|f\|_n^2 \cdot n^{-2\alpha}\right) \leftarrow$$

(using an appropriate dictionary of wavelets)

adaptive  
to the local  
regularity off



Next week: Take the two concepts from today to shallow NNs.