

## Lecture 3: Shallow Neural Networks and Approximation

- Last week we talked about two basic notions of approximation:

① Linear Approximation: Identifying a single finite-dim subspace (or a nested collection) that approximates uniformly well a whole class of inputs.

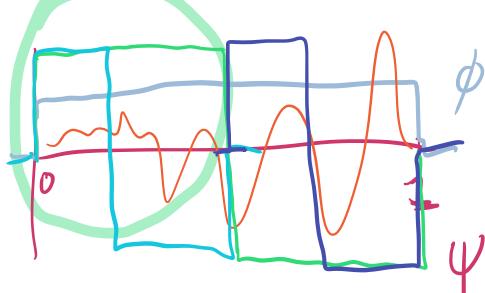
② Non-linear Approx: We allow the finite-dim model to adapt to the input.

Example: Fourier representation. (linear).

$$f(t) = \sum_{k=0}^{\infty} a_k e^{i 2\pi k \cdot t}$$

$$\underbrace{f_N(t)}_{=} = \sum_{k \leq N} a_k e^{i 2\pi k \cdot t}.$$

- Wavelet representations. ; Haar wavelets (1920s)



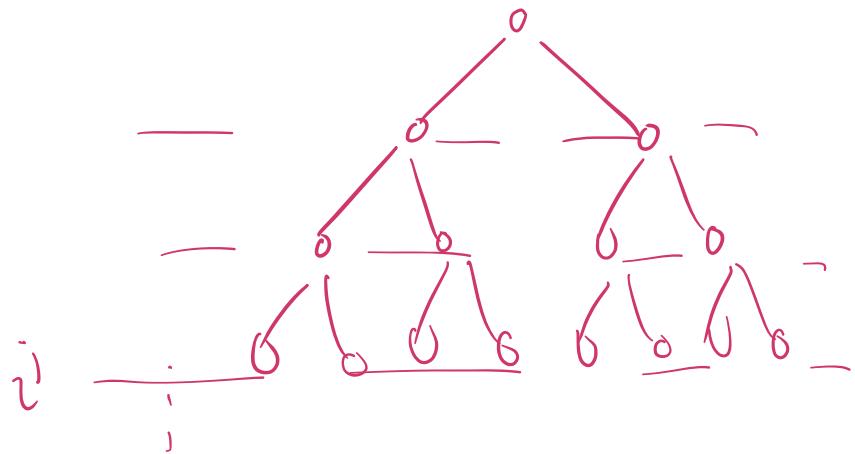
$$\int \phi(t) \psi(t) dt = 0$$

$$\|\phi\| = \|\psi\| = 1$$

$$D = \{ \phi, \psi, \psi \circ \psi \}$$

$$\sim \{1, 1_0, \dots, 1_{j,2}, \dots, 1_{j,k}, \dots, 1_{j+1,1}, \dots\}$$

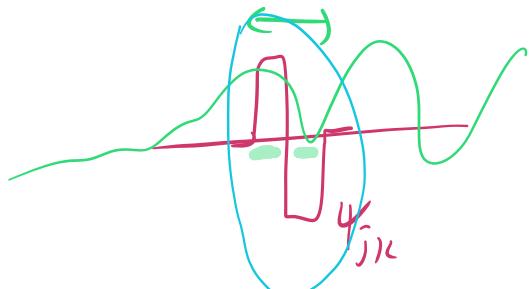
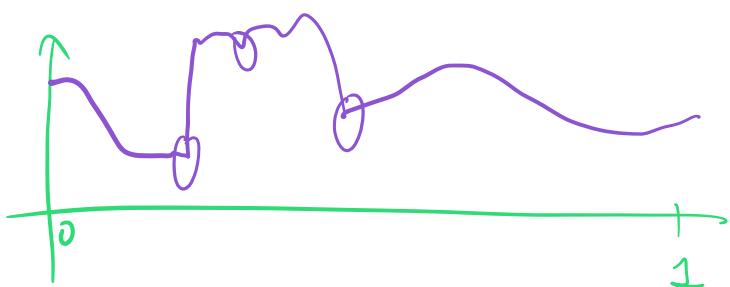
$$\Psi_{j,k}(t) = 2^{-j/2} \Psi_0(2^{-j}t - 2^{-j}k). \quad k=0 \dots 2^{-j}-1$$



Fact:  $\mathcal{D}$  is an orthonormal basis of  $L^2([0,1])$ .

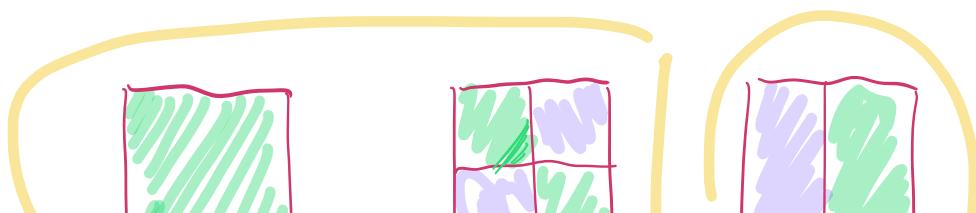
Q: What are the functions that can be well approximated by few (Haar) wavelet coefficients?

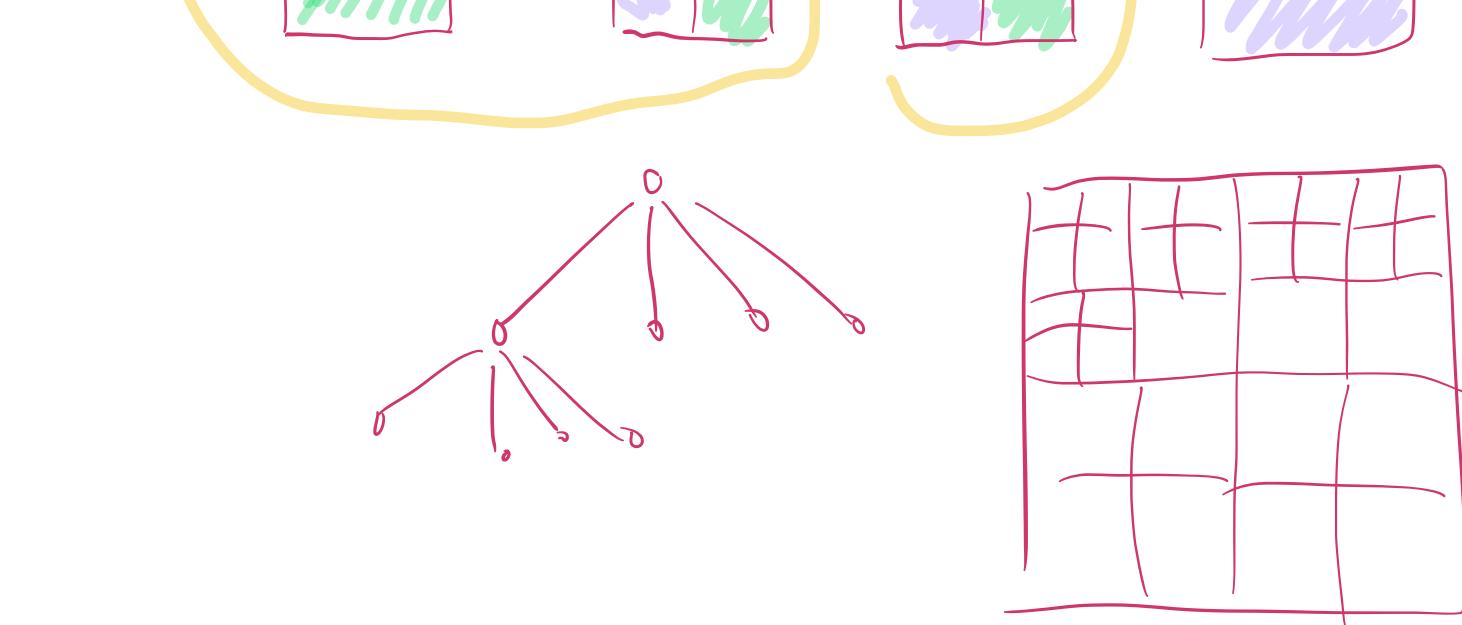
$$\boxed{(f, \Psi_{j,k})} \approx 0$$



$$\begin{aligned} & \int f(t) \Psi_{j,k}(t) dt \\ &= 2^{-jk} \left[ \int_{2^{-j}k}^{2^{-j}(k+1)} f(t) dt - \int_{2^{-j}(k+1)}^{2^{-j}(k+2)} f(t) dt \right] \end{aligned}$$

$$\rightarrow f(u,v) \quad (u,v) \in [0,1]^2$$





→ In general, in high-dimensions, this wavelet construction increases exponentially in the input dimension!!

→ Alternative?

→ Go from "grid"-like dictionaries to  
"off-the-grid" representations!

| Shallow Neural Networks as off-the-grid  
representations |

Consider  $\Gamma: \mathbb{R} \rightarrow \mathbb{R}$  Lipschitz activation function.

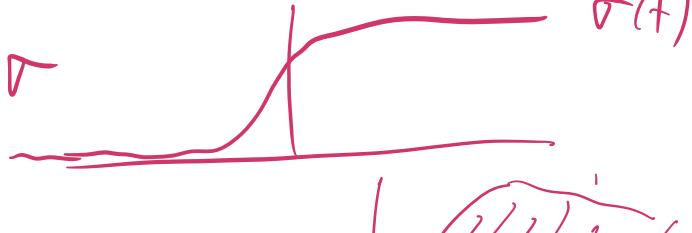
$$(\phi(x; \theta)) = \underbrace{\Gamma((x, w) + b)}$$

activation-

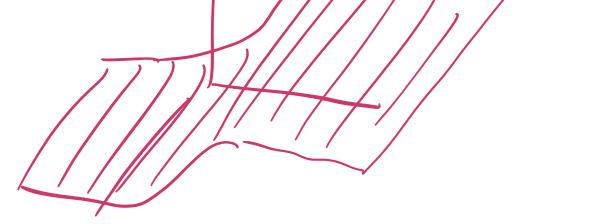
γ - Ridge function

$$\theta = (w, b) \in \mathbb{R}^{d+1}$$

affine function from  $\mathbb{R}^d \rightarrow \mathbb{R}$



$$\phi: \mathbb{R}^d \rightarrow \mathbb{R}$$

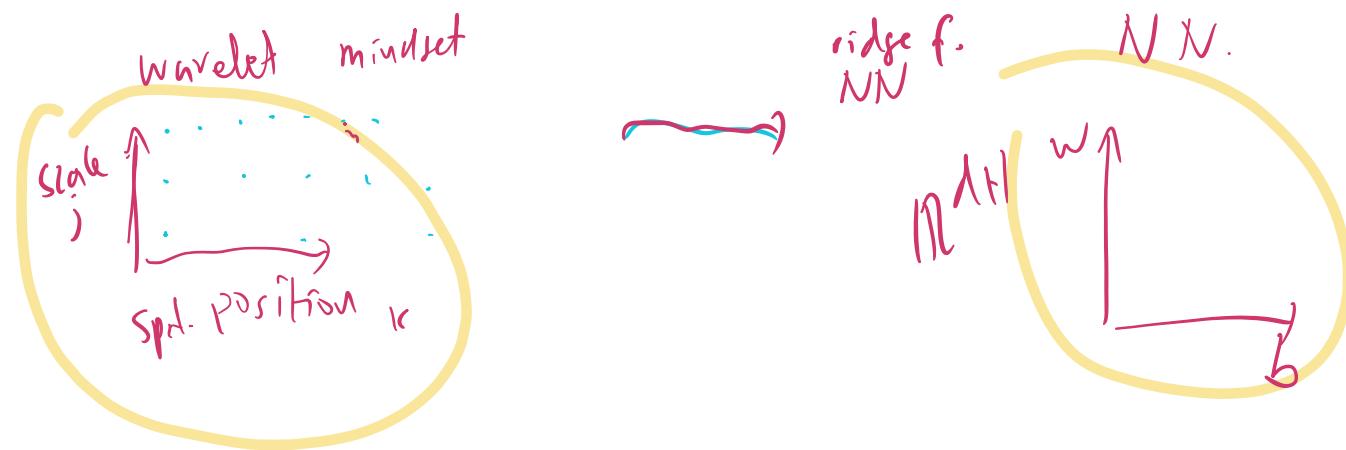


$$\rightarrow f_K(x; \theta) = \sum_{k=1}^K \phi(x; \theta_k)$$

linear combination  
of ridge functions.

$\theta = \{\theta_k\}_{k=1}^K$

$H_\Gamma = \{f_K ; K \in \mathbb{N}\}$  : space of functions represented as shallow



Q: How expressive is the set  $H_\Gamma$ ?

### Universal Approximation Theorems

$\rightarrow$  We want to understand whether  $H_\Gamma$  can approximate arbitrarily well continuous functions  $f: \mathbb{R}^d \rightarrow \mathbb{R}$ .

For a given error metric  $d$ ,  $\forall f \in C(\mathbb{R}^d)$ ,  $\exists \varepsilon_0$   
 $\exists \tilde{f} \in H_\Gamma$  with  $d(f, \tilde{f}) \leq \varepsilon$ .

Let's first establish UAT using 3-layer NN.

Theorem: Let  $x: \mathbb{R}^d \rightarrow \mathbb{R}$  be continuous  $\varepsilon > 0$ .

Theorem: Let  $g: \mathbb{R}^d \rightarrow \mathbb{R}$  continuous,

Let  $\delta > 0$  so that  $\forall x, x' \in [0,1]^d$ ,  $\|x - x'\|_\infty \leq \delta$

then  $|g(x) - g(x')| \leq \varepsilon$ . Then there exists a 3-layer Neural Network  $f$  with ReLU activation and  $\Omega(\delta^{-d})$  such that  $\|f - g\|_1 \leq 2\varepsilon$ .

Proof: [ Telgarski '20 ]

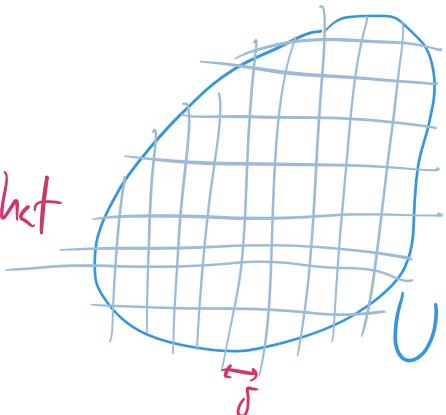
Lemma: Let any  $U \subset \mathbb{R}^d$ , along with a partition  $P$  of  $U$  into rectangles of diameter  $\delta$

$$P = \{R_1, \dots, R_n\}. \text{ Then there}$$

exists scalars  $\alpha_1, \dots, \alpha_N$  such that

$$\sup_{x \in U} |g(x) - h(x)| \leq \varepsilon,$$

$$h(x) = \sum_i \alpha_i \mathbf{1}_{R_i}(x).$$

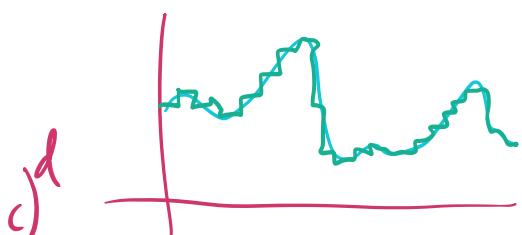


in other words, we are building a piece-wise constant approximation of a continuous function.

Proof of lemma: exercise.

Let  $P$  be a partition of  $[0, 1]^d$

$$\text{into rectangles } \prod_j [a_j, b_j]$$



with diameter  $\leq \delta$ .  
 $b_j - a_j \leq \delta$ .

→ Let  $h = \sum_i \alpha_i \mathbb{1}_{R_i}$  from Lemma, so we know  
 $\|g - h\|_1 \leq \varepsilon$ .

→ We will build  $f(x) = \sum_i \alpha_i (g_i(x))$   
 where each  $g_i$  is a ReLU Net that approximates  
 $\mathbb{1}_{R_i}$  with  $O(d)$  nodes.

Since  $|P| = \Theta(\delta^{-d})$ , then

$$O(d) \cdot \Theta(\delta^{-d}) = \underbrace{\Omega(\delta^{-d})}.$$

$$\begin{aligned} \|g - f\|_1 &\leq \|f - h\|_1 + \|h - g\|_1 \\ &\leq \left\| \sum_i \alpha_i (\mathbb{1}_{R_i} - g_i) \right\|_1 + \varepsilon \\ &\leq \sum_i |\alpha_i| \|\mathbb{1}_{R_i} - g_i\|_1 + \varepsilon \end{aligned}$$

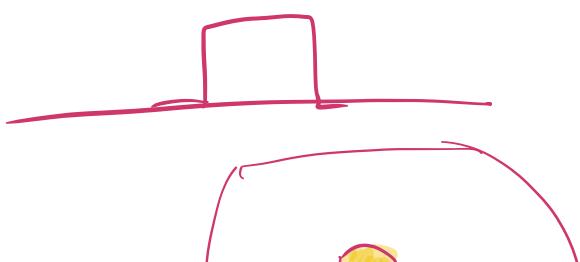
→ If we can build  $g_i$  so that

$$\|\mathbb{1}_{R_i} - g_i\|_1 \leq \frac{\varepsilon}{\|\alpha\|_1}, \text{ then}$$

$$\|f - g\|_1 \leq 2\varepsilon.$$

• Fix  $i$ , and let

$$\mathbb{1}_{R_i}^d = \dots$$



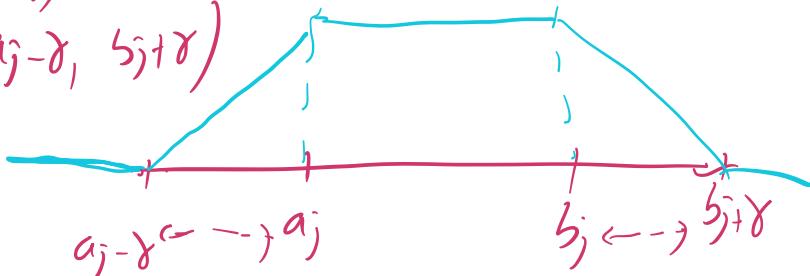
$$K_i = \bigotimes_{j=1}^d [a_{i,j}, b_{i,j}]$$

For  $\gamma > 0$ , for each  $j \in \{1, \dots, d\}$  we define

$$g_{j,\gamma}(z) = \Gamma\left(\frac{z - (a_j - \gamma)}{\gamma}\right) - \Gamma\left(\frac{z - a_j}{\gamma}\right)$$

$$- \Gamma\left(\frac{z - b_j}{\gamma}\right) + \Gamma\left(\frac{z - (b_j + \gamma)}{\gamma}\right)$$

$$g_{j,\gamma}(z) = \begin{cases} 1 & \text{if } z \in (a_j, b_j) \\ 0 & \text{if } z \notin (a_j - \gamma, b_j + \gamma) \\ 0,1 & \text{otherwise.} \end{cases}$$

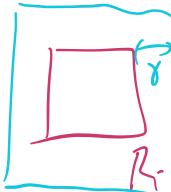
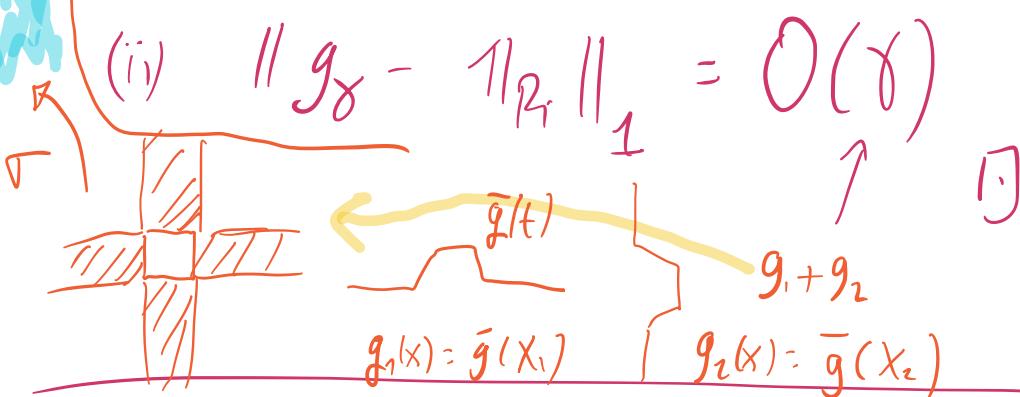


$$g_\gamma(x) = \Gamma\left(\frac{1}{d} \sum_{j=1}^d g_{j,\gamma}(x_j) - (d-1)\right)$$

*cheat!*

we verify that

$$g_\gamma(x) = \begin{cases} 1 & \text{if } x \in R; \\ 0 & \text{if } x \notin (x, (a_j - \gamma, b_j + \gamma)) \\ 0,1 & \text{otherwise.} \end{cases}$$



Q: Are there more powerful characterizations of universal approximation?

Thm: [Stone-Weierstrass] Let  $\mathcal{I} = [0, 1]^d$ . Let

$\mathcal{F}$  function class with:

- (i) each  $f \in \mathcal{F}$  is continuous.
- (ii)  $\forall x \in \mathcal{I}, \exists f \in \mathcal{F}$  with  $f(x) \neq 0$ .
- (iii)  $\forall x \neq x'$ ,  $\exists f \in \mathcal{F}$  with  $f(x) \neq f(x')$
- (iv)  $\mathcal{F}$  is closed under multiplication and  
a vector space operations (ie  $\mathcal{F}$  is an algebra).

Then  $\mathcal{F}$  is a universal approximator: For any

$g: \mathbb{R}^d \rightarrow \mathbb{R}$  continuous and  $\varepsilon > 0$ ,  $\exists f \in \mathcal{F}$  with

$$\sup_{x \in \mathcal{I}} |f(x) - g(x)| \leq \varepsilon.$$

- $\sigma$  a "sigmoid"  $\lim_{t \rightarrow -\infty} \sigma(t) = 0$ ,  $\lim_{t \rightarrow +\infty} \sigma(t) = 1$ ,  
[Hornik, Stinchcombe, White '89]
  - $\sigma \neq$  polynomial [Leshno' 93]
- Both results use S-A.

Q: Is this VAT "surprising" / different from last week?

### The Fourier Perspective

$\rightarrow$  let  $f \in C(\mathbb{R}^d)$  and consider its restriction

in a compact set  $\Omega$ . ( $\Omega = [0, 1]^d$ ).

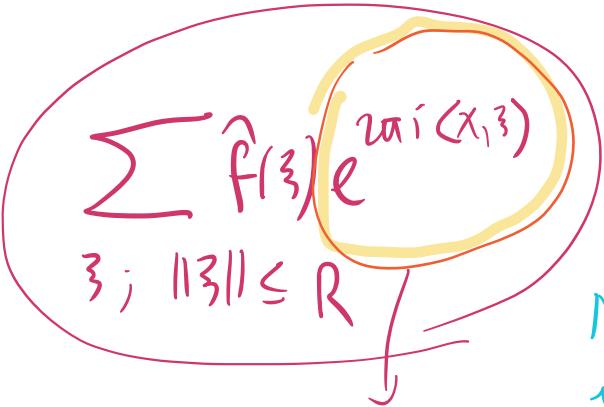
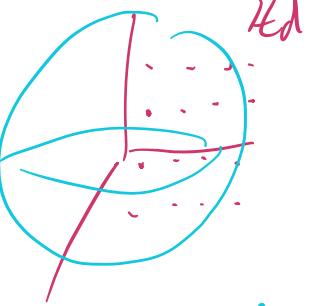
Fourier decomposition:  $\forall \vec{\zeta} \in \mathbb{Z}^d$ , we consider

$$\hat{f}(\vec{\zeta}) = \left\langle f, e^{i2\pi \langle x, \vec{\zeta} \rangle} \right\rangle_{L^2(\Omega)} =$$

$$= \int_{\Omega} f(x) e^{-i2\pi \langle x, \vec{\zeta} \rangle} dx. \boxed{f \approx \sum \alpha_i \frac{1}{R_i} e^{i2\pi \langle x, \vec{\zeta} \rangle} \text{ before.}} \\ \boxed{f \approx \sum \alpha_i \sigma(x, w_i)}$$

Fourier Inversion lemma: If we consider the

partial sums

$$f_M(x) := \sum_{\vec{\zeta}; \|\vec{\zeta}\| \leq R} \hat{f}(\vec{\zeta}) e^{i2\pi \langle x, \vec{\zeta} \rangle}$$



M: number of integer points in  $\{ \|x\| \leq R \}$ .

$R \rightarrow \infty$

$f$

$\sigma(t) = e^{i\pi t}$   $\sigma: \mathbb{R} \rightarrow \mathbb{C}$

$\Gamma(\langle x, \vec{\zeta} \rangle)$

This Fourier Inversion  $\hookrightarrow$  VAT.

Moreover, we have tight control over approximation rates.

Smoothness of  $f \leftrightarrow$  decay of  $\hat{f}(\vec{\zeta})$ .

$f \in H^s$

$$\Rightarrow \boxed{\| f - f_M \| = O(M^{-s/d})}$$

Sobolev class

$\| \cdot \|_s$  Sobolev

$$\| f \|_s = (\int_{\Omega} (1 + |\nabla f|^2)^{s/2})^{1/s}$$

$$f(t) = \cos(\omega t) + i \cdot \sin(\omega t)$$

↳ curse of dimensionality.

→ Q: How to overcome this curse?

A: look for another class of functions with "sparse" Fourier decomposition.

Idea: Suppose that our target function

$$f(x) = \int \tau(\langle x, \theta \rangle) \cdot g(\theta) d\theta$$

with  $\rightarrow \int |g(\theta)| d\theta < +\infty$  ( $\Rightarrow g \in L^1$ )

(!!)  $g(\theta) = \text{sign}(g(\theta)) \cdot \|g\|_1 \cdot \frac{|g(\theta)|}{\|g\|_1} := q(\theta)$

$$q \geq 0 \quad \int q(\theta) d\theta = 1.$$

$$f(x) = \int \tau(\langle x, \theta \rangle) \text{Sign}(g(\theta)) \cdot \|g\|_1 \cdot q(\theta) d\theta.$$

$$= \|g\|_1 \cdot \underbrace{\mathbb{E}_{\theta \sim q} \left[ \tau(\langle x, \theta \rangle) \cdot \text{Sign}(g(\theta)) \right]}_{\phi(x, \theta)}.$$

So we can consider a Monte-Carlo approximation!

$$f_M(x) = \frac{1}{M} \sum_{i=1}^M \phi(x, \theta_i), \quad \theta_i \stackrel{iid}{\sim} q.$$

$$\|f - f_M\|_2^2 = \mathbb{E}_{\theta \sim q} [\|\tau(\langle x, \theta \rangle) - \mathbb{E}_{\theta \sim q}[\tau(\langle x, \theta \rangle)]\|^2]$$

$$\mathbb{E} \|\hat{f} - f_M\| = \frac{1}{M} \text{Var}(\phi(x, \theta)) \underbrace{\left[ n \varphi(x, \theta) - \bar{x} \cdot \nabla \right]}_{\text{Bias}}$$

$$\leq \frac{1}{n} \mathbb{E} \|\phi(x, \theta)\|^2 \leq \left( \frac{1}{M} \sup_{\theta} \mathbb{E}_x \phi^2 \right) \|g\|_1^2$$

There is no curse! (statistical rate of approximation).

Theorem: [Barron '93] Suppose

$$C = \int \|\widehat{\nabla f}(\zeta)\| d\zeta < +\infty$$

with  $f, \widehat{f} \in L_1$ . Then we can

$$\widehat{\nabla f}(\zeta) = \zeta \cdot \widehat{f}'(\zeta)$$

(a weighted  $L_1$  norm in Fourier).

approximate  $f$  to accuracy  $\epsilon$  with a ReLU/sigmoid.

shallow NN with  $\sim \mathcal{O}_{\epsilon^2}$  hidden units.

(no curse in this function class) "Barron Space!"