

Lecture 10: Depth Separation

Recap from last week:

- Optimizing shallow nets in feature learning regime: positive guarantee, BUT in quantitative
- Learning lower Bounds under statistical / "gradient" queries: exponential number of orthogonal possible target functions.

Today: Another source of learning hardness: poor approximation using shallow NN

↳ Better approximation using deeper networks?

Depth Separation: Univariate Case

(Motivating question): Let $f: \mathbb{R} \rightarrow \mathbb{R}$. All approximation schemes we have seen so far (Fourier, Kernels, Shallow NN) are additive:

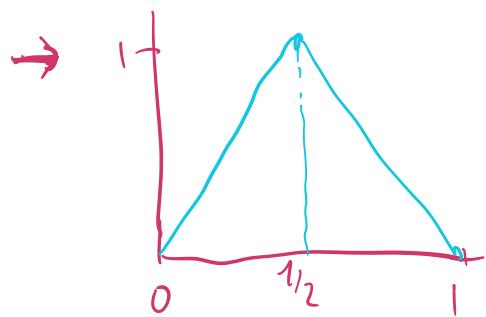
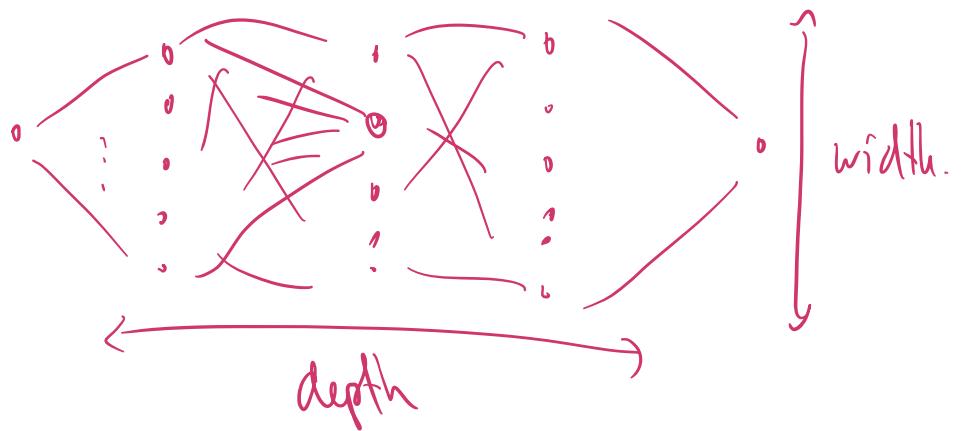
$$f \approx \sum_i \phi_i(x)$$

↳ This is in contrast to Approximation by composition

$$f \approx \phi_L \circ \dots \circ \phi_2 \circ \phi_1 \quad \phi_i: \mathbb{R} \rightarrow \mathbb{R}$$

Q: Which functions can be easily approximated by composition, yet poorly approximated by addition?

In the context of NNs, we will quantify these options in terms of # of neurons.



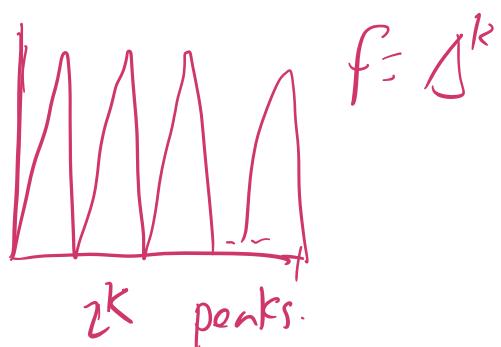
$$\Delta : [0,1] \rightarrow [0,1]$$

$$x \mapsto \begin{cases} 2x & \text{if } x \in [0, \frac{1}{2}] \\ 2-2x & \text{if } x \in [\frac{1}{2}, 1] \end{cases}$$

$$\Delta^2 = \Delta \circ \Delta$$



$$\Delta^k = \underbrace{\Delta \circ \Delta \circ \dots \circ \Delta}_{k \text{ times}}$$



Theorem [Telgarski '15] : Fix $L \geq 2$, and let

$$f = \Delta^M \quad \text{with} \quad M = L^2 + 2. \quad \text{Then}$$

- (i) f can be expressed as a ReLU net with $3M$ neurons and $2M$ layers. (positive approximation)

(ii) f cannot be approximated using shallow nets: any ReLU network g with $\leq 2^L$ neurons and $\leq L$ layers incurs in constant L_1 error: $\|f - g\|_1 \geq \frac{1}{32}$.
 (negative approx).

Proof:



(i) Track the number of affine regions, $N_A(f)$, for each ReLU network.

Lemma: Let $f: \mathbb{R} \rightarrow \mathbb{R}$ be a ReLU net with L layers and widths (m_1, \dots, m_L) , with $m = \sum_i m_i$. Then

$$N_A(f) \leq \left(\frac{2m}{L}\right)^L$$

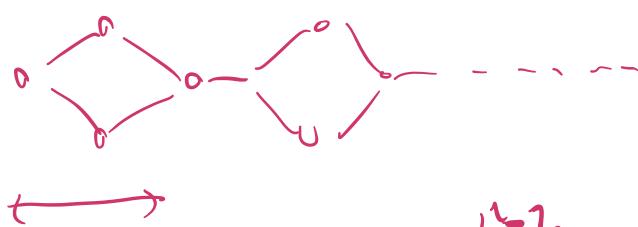
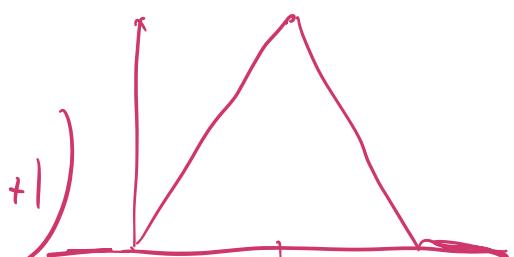
→ Proof of this (Gmn̄a (Talgorski '11)): uses induction over depth.

$$N_A(f + g) \leq N_A(f) + N_A(g)$$

$$N_A(f \circ g) \leq N_A(f) N_A(g) \quad \leftarrow \text{This is tight for } f=g=\text{id.}$$

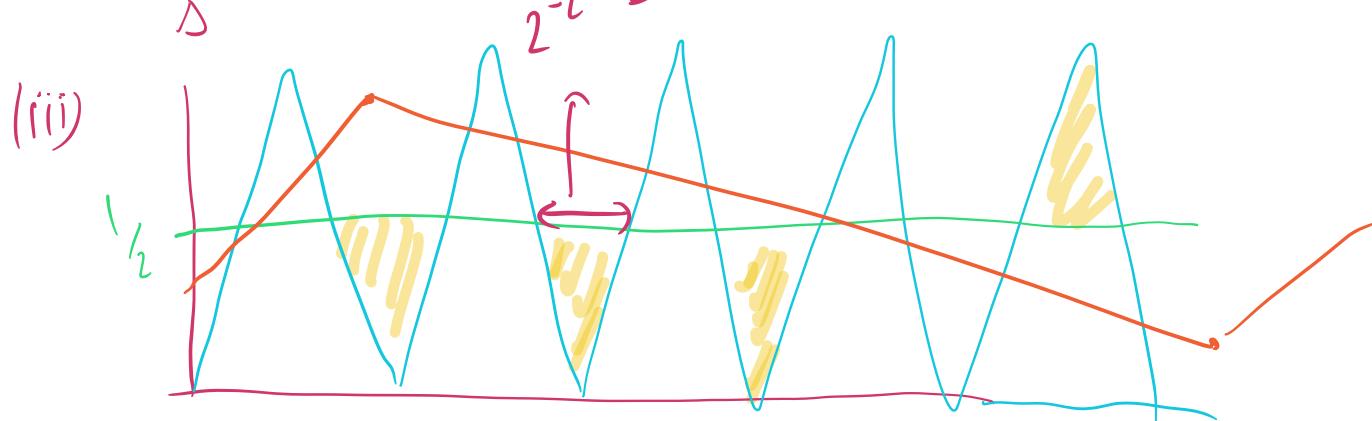
(ii) Observe that $\Delta(x)$ can be implemented with a simple ReLU Net: $\Delta(t) = \max(0, t)$

$$\Delta(x) = \nabla(-2\nabla(x - \frac{1}{2}) - 2\nabla(-x + \frac{1}{2}) + 1)$$



$$\Delta^M$$





↳ we have 2^{L^2+1} copies of $\Delta \rightsquigarrow 2^{L^2+2}-1$ half-triangles. Each halftriangle has area $\frac{1}{4}2^{-L-2} = 2^{-L-4}$

↳ if $\mu_A(g) \leq \left(2 \cdot \frac{2^L}{2}\right)^L \leq 2^{L^2}$ linear regions.

$$\begin{aligned} \int_0^1 |f \cdot g| dx &\geq (\# \text{ untouched triangles}) \cdot (\text{triangle area}) \\ &\geq \frac{1}{2} (\text{Total triangles} - 2 \cdot N_A(g)) \cdot 2^{-L-4} \\ &\geq \frac{1}{2} (2^{L^2+2}-1 - 2 \cdot 2^{L^2}) \cdot 2^{-L-4} \geq 2^{-5} \quad \square \end{aligned}$$

- Remarks:
- Separates $\sim L^2$ depth from $\sim L$ depth.
 - Very tailored to Reh/piecewise linear structures \rightsquigarrow NOT too "intrinsic".

Q: $\xrightarrow{\hspace{1cm}}$ High-dimensional perspective?

Depth Separation in the High-Dimensional Regime

→ From UAT, shallow is sufficient in a qualitative sense.
Now we will be interested in finite guarantees.

→ Focus on simplest instance of Depth separation: Shallow vs
2 hidden layer. (1 hidden
layer)

→ Problem setup: Let $\Omega \subseteq \mathbb{R}^d$ input domain; let
 ν be a probability distribution (data) on Ω .

Let $\mathcal{F}_2(m) := \left\{ f: \Omega \rightarrow \mathbb{R}; f(x) = \sum_{j=1}^m \phi_j(\langle x, w_j \rangle) \right. \begin{array}{l} \phi_j \text{ is Lipschitz} \\ : \mathbb{R} \rightarrow \mathbb{R} \end{array} \right\}$

(ex $\phi_j(t) = a_j \cdot \sigma(t - b_j)$ recovers $2MN$ "classes".

$\mathcal{F}_3(m) = \left\{ f: 3 \text{ layer NN of widths } m_1, m_2 \right. \begin{array}{l} m = m_1 + m_2 \end{array} \right\}$

Fix $\varepsilon > 0$, and consider a target function $f^*: \Omega \rightarrow \mathbb{R}$

The rate of approximation of f^* on each $\mathcal{F}_2(m)$

$$M_2(\varepsilon, d) = \inf \left\{ m; \inf_{f \in \mathcal{F}_2(m)} \|f^* - f\|_\nu \leq \varepsilon \right\}$$

In words, width required to approximate f^* to error ε
in dimension d .

$$\left\| \mathbb{E}_{x \sim \nu} (f^*(x) - f(x))^2 \right\|$$

Punchline: We can find f^* such that

$$m_3(\varepsilon, d) = \text{poly}(d, 1/\varepsilon), \quad \underline{\text{but}}$$

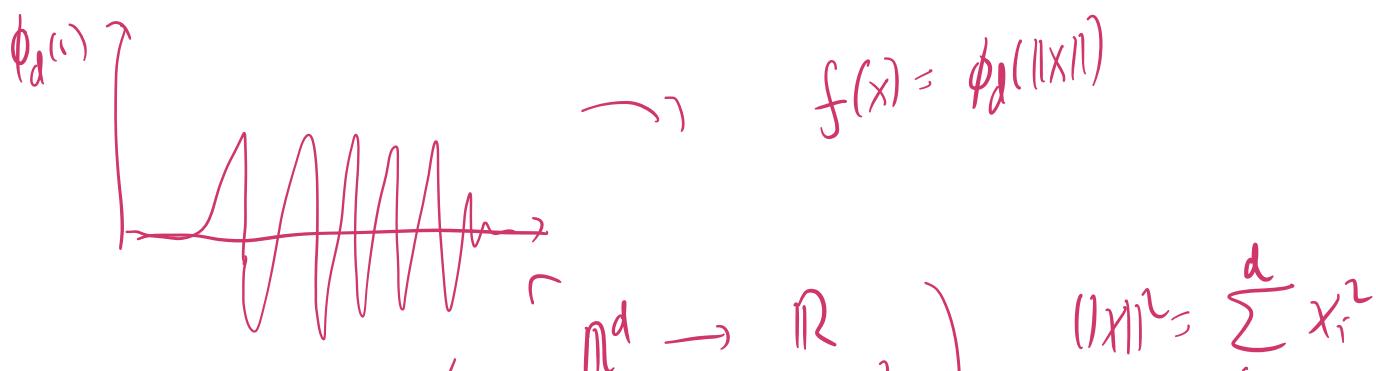
$$m_2(\varepsilon, d) = \exp(d) \quad (\text{for any fixed } \varepsilon)$$

Theorem [Eldan & Shamir '16] Suppose Γ activation function satisfies mild growth conditions and \neq polynomial. Then for each $d \in \mathbb{N}$, there exists probability measure ν in \mathbb{R}^d and $f^\circ: \mathbb{R}^d \rightarrow \mathbb{R}$ such that $m_3(\varepsilon, d) = \text{poly}(d, \varepsilon^{-1})$, but f° cannot be ε -approximated for $\varepsilon = \varepsilon(d)$ by any shallow NN of width $\exp(o(d))$.

Q: What is the nature of f° ?

f° is chosen as a radial function:

$f^\circ(x) = \phi_d(\|x\|)$, for $\phi_d: \mathbb{R} \rightarrow \mathbb{R}$, and when ϕ_d will oscillate, $\text{Lip}(\phi_d) = \text{poly}(d)$



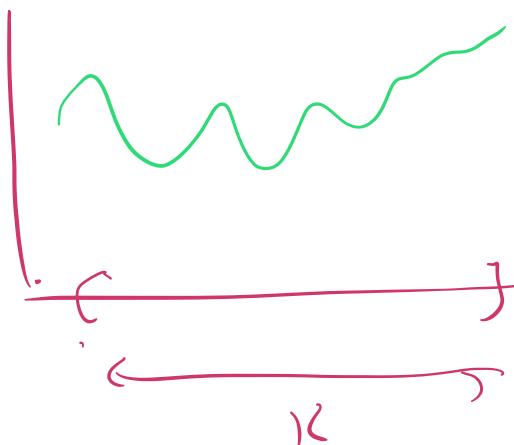
(+) positive result: $(F_1: \mathbb{R}^d \rightarrow \mathbb{R}, x \mapsto \|x\|^2)$ $\|x\|^2 = \sum_{i=1}^d x_i^2$

$(F_2: \mathbb{R} \rightarrow \mathbb{R}, t \mapsto \phi_d(\sqrt{t})) \leftarrow f^\circ = F_2 \circ F_1$

Idea:

"spend" 1st hidden layer to approximate F_1 .
 "spend" 2nd hidden layer to approximate F_2 .

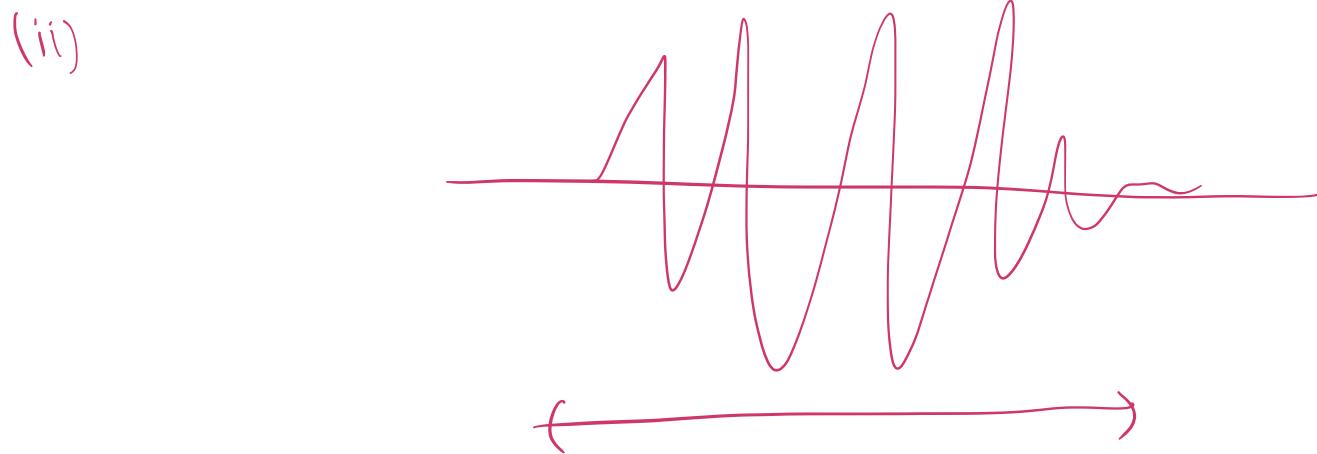
Given a compact domain $K \subset \mathbb{R}$ and $F: K \rightarrow \mathbb{R}$ with $\text{Lip}(F) \leq \beta$, we



can approximate F as

$$\tilde{F}(t) = \sum_{i=1}^N \alpha_i \sigma(t - t_i)$$

$|K|, \beta, \varepsilon^{-1}$
 poly(d)



$$\mathbb{E}_{x \sim \nu} |F_1(x) - \tilde{F}_1(x)|^2 \leq \varepsilon$$

(ii) negative result: Suppose data distribution ν admits a density
 ↗ since density
 with respect to Lebesgue: $\nu(dx) = \varphi^2(x) dx$. Then

$$\|f - g\|_{\nu}^2 = \int |f(x) - g(x)|^2 \varphi^2(x) dx = \int (f\varphi - g\varphi)^2 dx = \|f\varphi - g\varphi\|_2^2$$

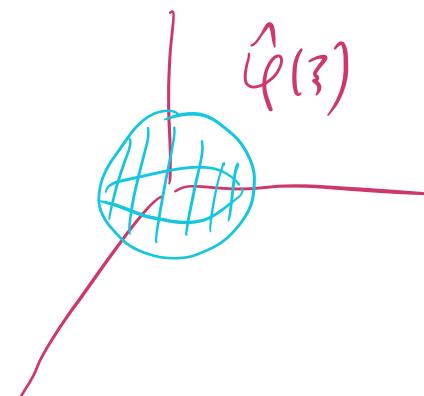
- By Parseval identity, we can express this norm in Fourier:

$$\|f\varphi - g\varphi\|_2^2 = \|\widehat{f\varphi} - \widehat{g\varphi}\|_2^2$$

- Recall $\widehat{f\varphi} = \widehat{f} * \widehat{\varphi}$. We construct data distribution $\nu(x) = \varphi^*(x) dx$ such that $\widehat{\varphi}$ is well localised in freq;
 $\widehat{\varphi}(\xi) = \underbrace{\mathbb{1}}_{B} \{ \|\xi\| \leq 1 \}$ (indicator of unit ball in d dimensions)

↳ we verify that $\varphi^*(x)$ is indeed a density. $\varphi^*(x) \geq 0$

$$\int \varphi^*(x) dx = \int |\widehat{\varphi}|^2(\xi) d\xi = 1.$$



- Now we compute Fourier transform of

$$f(x) = \sum_{j=1}^m \phi_j(\langle x, w_j \rangle).$$

"ridge" function.

↳ Focus on a single term $f_0(x) = \phi(\langle x, w \rangle)$

$$\widehat{f}_0(\xi) = \int f_0(x) e^{-i\langle x, \xi \rangle} dx$$

; it is well-defined in
the sense of tempered distributions, thanks to
polynomial growth assumptions
on ϕ .

$$\widehat{f}_0(\xi) = \left(\phi(\langle x, w \rangle) e^{-i\langle x, \xi \rangle} dx \right)$$

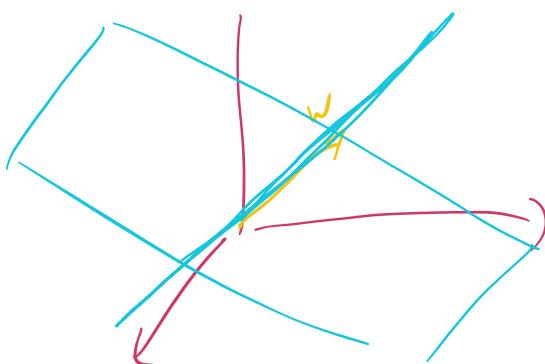
change of
variables
 $w \rightarrow \beta_1$

$$= \int \phi(x_1) e^{-i x_1 \beta_1} e^{-i \langle x_1, \beta_{-1} \rangle} dx_1 \cdot dx_{-1}$$

$$= \left(\int \phi(x_1) e^{-i x_1 \beta_1} dx_1 \right) \cdot \left(\int e^{-i \langle x_1, \beta_{-1} \rangle} dx_{-1} \right)$$

univariate F. Transform of
 ϕ $\hat{\phi}(\beta_1)$

F. Transform of constant f. $\sim \delta(\beta_{-1})$



The support of \hat{f}_0 is the ray $\text{span}\{w\}$

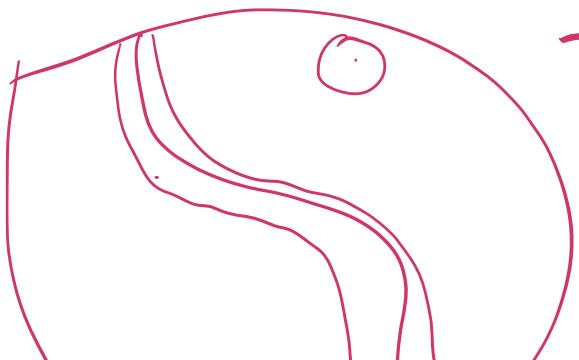
→ The support of the Fourier transform of

$$f(x) = \sum_{j=1}^m \phi_j(\langle x, w_j \rangle)$$

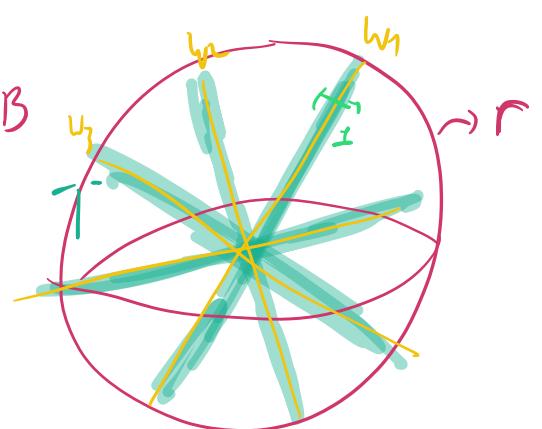
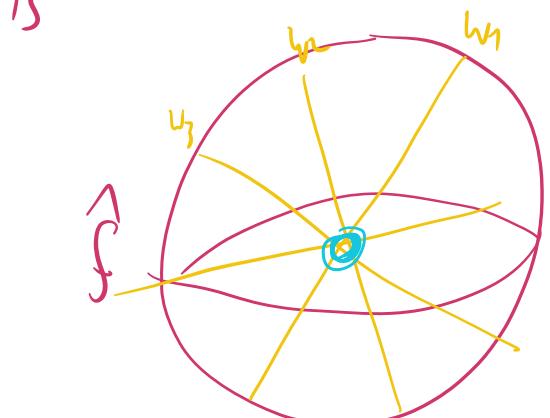
the union of rays $\bigcup_{j=1}^m \text{span}\{w_j\}$

→ The support of $\hat{f} * \hat{\varphi}$?

Recall that $\hat{\varphi} = \mathbb{1}_{\{\beta \in B\}}$



$$T = \bigcup_{j=1}^m \text{span}\{w_j\} + B$$



→ What fraction of the surface of the d-dim sphere of radius r intersects T ?

$$\frac{\text{Vol}_{d-1}(T \cap rS^{d-1})}{\text{Vol}_{d-1}(rS^{d-1})} \lesssim m \cdot e^{-d}$$

when r is large
wrt radius of B .

→ Can we find g (target) such that $\hat{g} * \varphi$ has energy "spread" out in $r S^{d-1}$?

→ A: Radial Function! g such that
 $g(\theta x) = g(x) \quad \forall x \quad \forall \varphi^T \varphi = \text{Id.}$

$$\hat{g}(\varphi z) = \int g(x) e^{-i \langle x, \varphi z \rangle} dx$$

$$= \int g(x) e^{-i \langle \varphi^T x, z \rangle} dx$$

$$= \int g(\varphi x) e^{-i \langle x, z \rangle} dx$$

$$= \hat{g}(z)$$

↳ So we need to construct g such that:

- (i) g is radial $\left(\Rightarrow g\varphi$ radial $\Rightarrow \hat{g} + \hat{\varphi}$ radial $\right)$
 $g(x) = \phi_d(\|x\|)$
- (ii) ϕ_d to oscillate. j energy in the high-frequencies.

Remarks

- The data distribution is heavy-tailed; what happens when we replace it for ν more concentrated; $\nu = \text{Unif}(S^d)$ [Daniely '17] gets also exponential lower bounds (but conditional on weights not too large).
- Extension beyond radial functions [VBJO'22] using similar prof technique.