# Part 3. Neural Network Approximation

(proof-trivial)

- (Approximation) (Representation)
- Generalization
- Optimization

Goal: $f: \mathbb{R}^d \to \mathbb{R} \quad \to \quad$ neural network, $g: \mathbb{R}^d \to \mathbb{R}$,

Population risk $\quad \int \ell(f(x), y)\, dP(x,y) \quad \longleftrightarrow \quad \int \ell(g(x), y)\, dP(x,y)$

$<$ upper bound: $\ell(\cdot, y)$ 1-Lipschitz: $\int \left( \ell(g(x), y) - \ell(f(x), y) \right) dP(x,y) \leq \int |g(x) - f(x)| \, dP(x,y)$

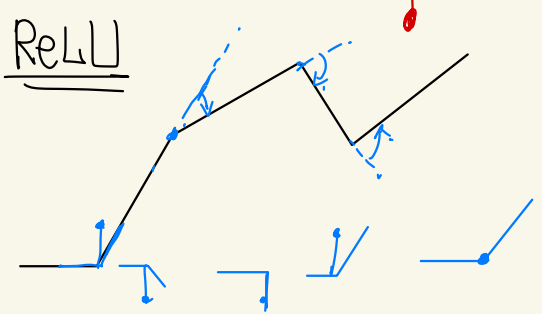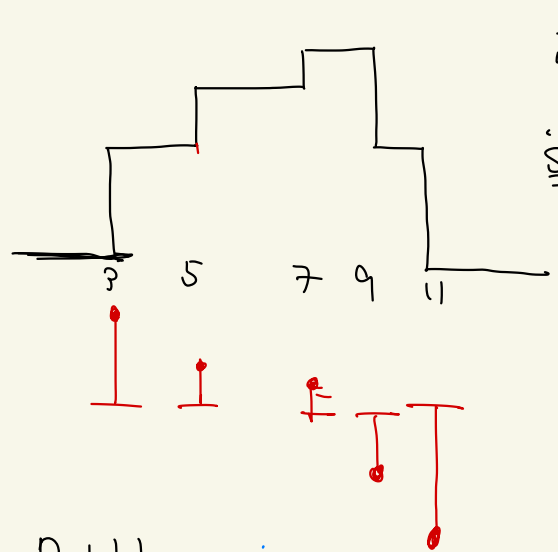lower bound;  w.c. $|g-f|$ large $L_1(P)$, $L_1(\text{Uniform})$

universal, uniform  $L_\infty(P)$ approximation

Deep Network: $\quad x \mapsto A_L \sigma_{L-1}\left( \cdots \sigma_1(A_1 x + b_1) \cdots \right) + b_L$,
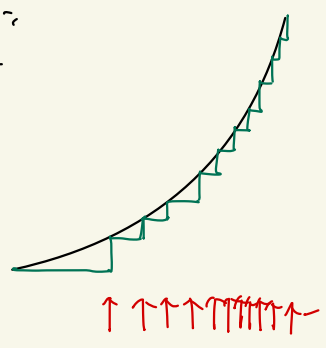
non-linearity / activation: ReLU $\quad z \mapsto \max\{0, z\}$

Univariate

Step Activation

$x \mapsto 2 \cdot \mathbb{1}\{x-3 \geq 0\} + \mathbb{1}\{x-5 \geq 0\} + \mathbb{1}\{x-7 \geq 0\} \cdots$,
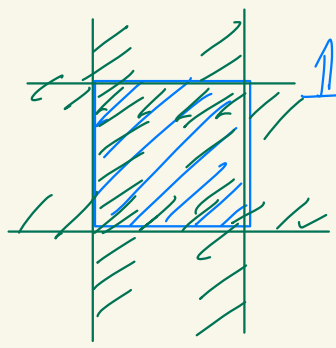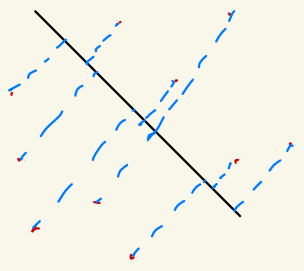
$\xrightarrow{\text{`smoothe'}}$

Lipschitz $L: \dfrac{1}{\varepsilon}$

$f(x) = f(0) + \int_0^x f'(b)\, db$

$= f(0) + \int_0^\infty \mathbb{1}\{x-b \geq 0\} f'(b)\, db$

$\to$ Infinite width network

avg $L / \varepsilon^2$.

ReLU

change of slope.

$f(x) = f(0) + r(x) \cdot f'(0) + \int_0^\infty r(x-b) f''(b) \cdot db$, $\quad \dfrac{\text{avg } L}{\varepsilon^2}$.

**Multivariate**    **box**    $\mathbb{1}\{$ inside of the box $\}$

$$\begin{array}{|c|c|c|} \hline 2 & 3 & 2 \\ \hline 3 & 4 & 3 \\ \hline 2 & 3 & 2 \\ \hline \end{array}$$

① product ★★

② threshold at 3.5

$$\begin{array}{|c|c|c|} \hline 0 & 0 & 0 \\ \hline 0 & 1 & 0 \\ \hline 0 & 0 & 0 \\ \hline \end{array} \longleftarrow$$ add a layer

2 Layer

**ball** ("bump")

$\rightarrow$ radial function    $\underline{2^d}$ nodes

"bump"

**RBF:** convolve with $f$

$$\underline{\left| \underline{f(x)} - \int \underline{f(z)} \cdot \underline{p(x-z)} \, dz \right|}$$

$$= \left| \underline{f(x)} - \int f(x-z) \cdot \underline{p(z)} \, dz \right|$$

$$= \left| \int f(x) \, p(z) \, dz - \int f(x-z) \, p(z) \, dz \right| \leq \int |f(x) - f(x-z)| \, p(z) \, dz$$

$$\left( \frac{d \cdot L}{\varepsilon} \right)^{O(d)}$$    (Mhaskar-Michelli '92)

$\begin{cases} \text{box:} & \text{3 layer} \quad \left( \frac{L}{\varepsilon} \right)^{O(d)} \longrightarrow \infty \\ \text{ball} & \text{; 2 layer RBF} \quad \left( \frac{d \cdot L}{\varepsilon} \right)^{O(d)} \end{cases}$

---

**Univariate bump:** $\underline{\cos x^P}$    $\underline{\mathbb{1}\{\|x\|_\infty \leq 1\} = \prod_{i=1}^{d} \mathbb{1}\{|x_i| \leq 1\}}$

$$\cos x \cdot \cos x = \frac{1}{2} \cos 2x + 1$$

$$2 \cos x_1 \cos x_2 = \cos(x_1 + x_2) + \cos(x_1 - x_2)$$

↓ Polynomial

1885

**Weierstrass Approximation Thm** : Polynomials can uniformly approximate continuous functions over compact sets.
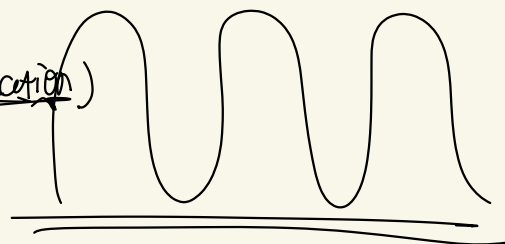
# Bernstein's proof

$\Downarrow$

$$\sum \binom{n}{x} p^x (1-p)^{n-x}$$

## Stone-Weierstrass thm : polynomial-like functions
approximate ctx functions ( closed under multiplication )

(Thm) ( Hornik-Stinchcombe-White '89 )

$\sigma: \mathbb{R} \to \mathbb{R}$,    $\lim\limits_{z \to \infty} \sigma(z) = 0$,  $\lim\limits_{z \to \infty} \sigma(z) = 1$

$$H_\sigma = \{ x \mapsto \sigma(a^T x - b) : (a,b) \in \mathbb{R}^{d+1} \}$$

$\text{span}\{H_\sigma\}$ uniformly approximates ctx functions on compact sets

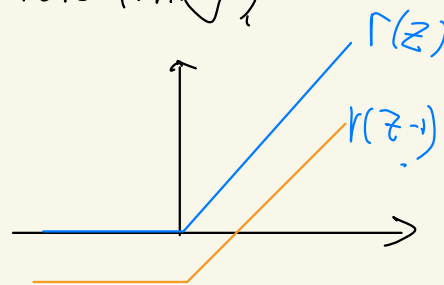proof sketch :    $H_{\cos}$ is closed.    $2\cos a \cos b = \cos(a+b) + \cos(a-b)$

$H_{\cos}$  with  $\text{span}\{H_\sigma\}$
     uniformly approximate  ( univariate fitting )

\* $H_{\exp}$ :  $e^a \cdot e^b = e^{a+b}$     $\square$

\* Leshno-Lin-Pinkus-Schocken '93.

   HSW holds __iff__ $\sigma$ is not a polynomial

(Barron '93)    $x \mapsto \exp(i a^T x)$ $\rightrightarrows$ $\int \exp(i a^T x) \tilde{f}(a) \, da$

$\Gamma(z)$
$r(z-1)$
$\sigma(z) \triangleq r(z) - r(z-1)$

---

[Depths]  → (Radial functions)  2 layer ReLU

$f(\|x\|^2)$ with Lipschitz constant $L$.

• $h(x) \underset{\varepsilon}{\approx} \|x\|_2^2 = \sum x_i^2$   → Layer 1 →  $d \cdot \frac{1}{\varepsilon}$. ReLU

• $g \underset{\varepsilon}{\approx} f$.   → Layer 2 →  $\frac{L}{\varepsilon}$ ReLU

$$\left| f(\|x\|^2) - g(h(x)) \right| \leq \left| f(\|x\|^2) - f(h(x)) \right| + \left| f(h(x)) - g(h(x)) \right|$$

$$\leq L \cdot \left| \|x\|^2 - h(x) \right| + \varepsilon \lesssim O(\varepsilon)$$

(g ∘ h) → size: poly$(L, d \cdot \frac{1}{\varepsilon})$

However, (Thm) $\exists$ radial function $f$ —

(Eldan-Shamir '15)     expressible with two layer ReLU of width poly($d$)

s.t. every $g$ with a single ReLU layer of width $2^{O(d)}$

satisfies   $\int (f(x) - g(x))^2 \, dP(x) \geq \Omega(1)$

$\uparrow$ A probability measure $P$

(Thm) (Daniely, 17')   $(x, x') \sim P = \text{Uniform}(S^2)$
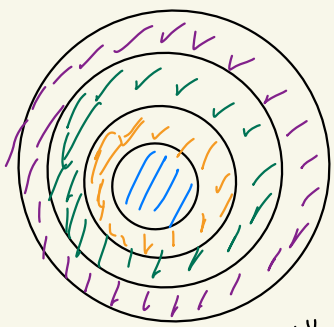
$$h(x, x') = \sin(\pi d^3 x^T x')$$

$\forall g$ with a single ReLU layer of width

$d^{O(d)}$ and weight magnitude $O(2^d)$

$$\int (h(x,x') - g(x,x'))^2 \, dP(x,x') \geq \Omega(1)$$

$h$ can be approximated to accuracy $\varepsilon$ by $f$ with 2 layer ReLU of size poly($d, \frac{1}{\varepsilon}$)
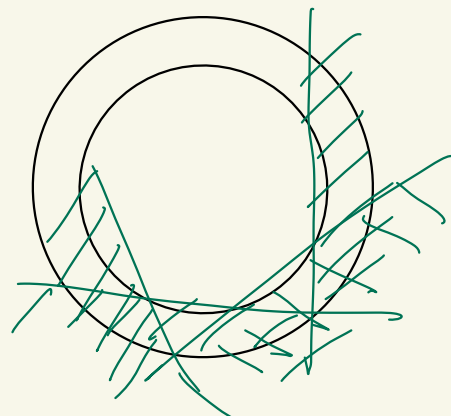


each shell approximation
$\downarrow$
overall function

"shell"

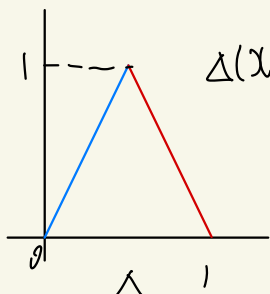(Depth) $\leftarrow$ (Benefits of depth.)

$\searrow$ What do shallow representations
      do exceptionally badly?

$$x \mapsto \mathbb{1}\left\{ \|x\| \in [1 - \tfrac{1}{d}, 1] \right\}^2$$

fraction
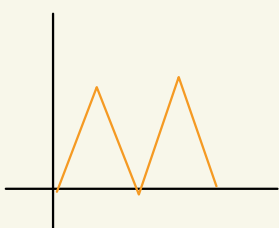


$\Delta(x) = r(2x) - r(4x - 2)$

$= \begin{cases} 2x, & x \in [0, \tfrac{1}{2}] \\ 2(1-x), & x \in [\tfrac{1}{2}, 1] \end{cases}$
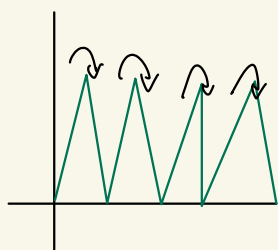
Composition:

$f(\Delta(x)) = \begin{cases} x \in [0, \tfrac{1}{2}] \Rightarrow f(2x) = f \text{ squeezed into } [0, \tfrac{1}{2}] \\ x \in [\tfrac{1}{2}, 1] \Rightarrow f(2(1-x)) = f \text{ reversed, squeeze} \end{cases}$

$\underline{\Delta^2 = \Delta \circ \Delta}$

$\underline{\Delta^K}$ : $O(K)$ layer & nodes, $\Theta(2^K)$ bumps (oscillation)

(Thm) (Telgarsky '15)   K: # of layer,

∃ ReLU network $f: [0,1] \to [0,1]$ with { 4 distinct parameters

$3K^2+9$ nodes

$2K^2+9$ layers

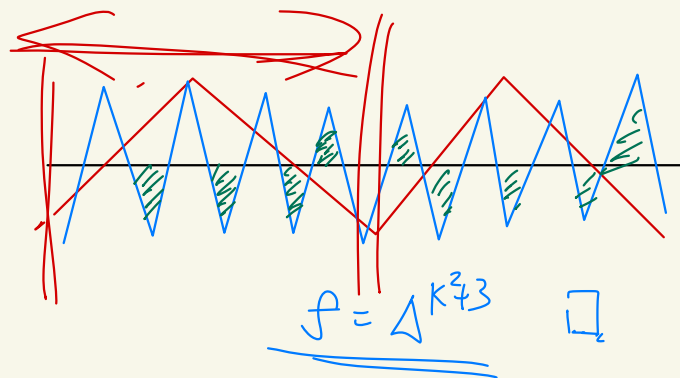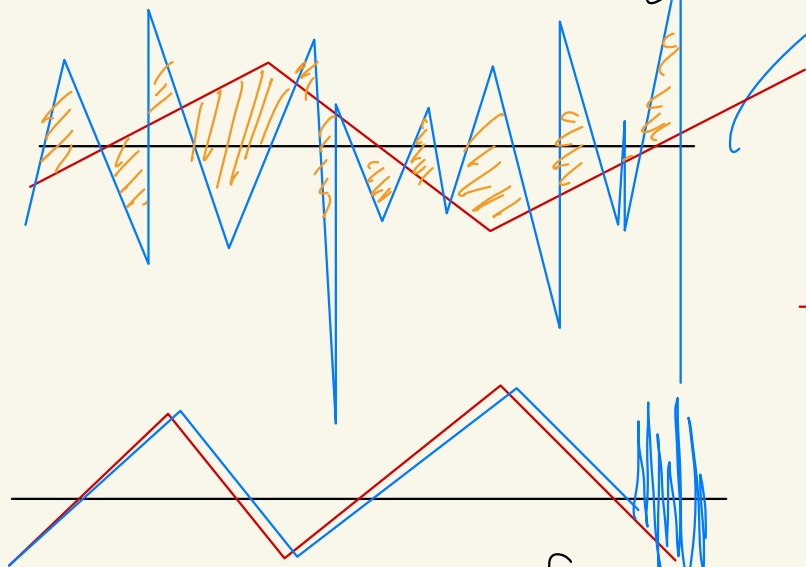s.t. ∀ ReLU network $g: \mathbb{R}^a \to \mathbb{R}$

with $\leq K$ layer. $\leq 2^K$ nodes

$\Rightarrow \int_{[0,1]} |f(x) - g(x)| \, dx \geq \boxed{\frac{1}{32}}$

**Proof:**

1. $g$ with few oscillations   ✗

2. $f = \triangle^{K^2+3}$ : regular oscillatory f.

3. width $m$, depth $L$ } $< o(m^L)$



$f = \triangle^{K^2+3}$   □

Depth { Radial function.

$\triangle^{K^2+3}$ function

$O(K^2)$  $O(1)$

if $O(K)$   $O(2^K)$

depth   width

$\int_{[0,1]} |f-g|$  ✗

$L_\infty$  ✓

$L_1$  ✓

$\triangle^K : 2^K$-Lipschitz.

non-realistic

Sobolev Spaces