

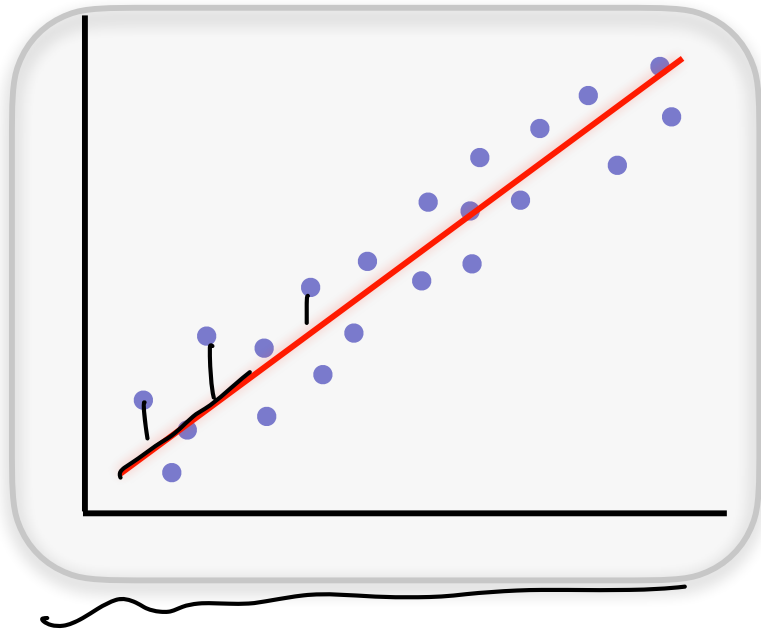
# Approximation Theory

---

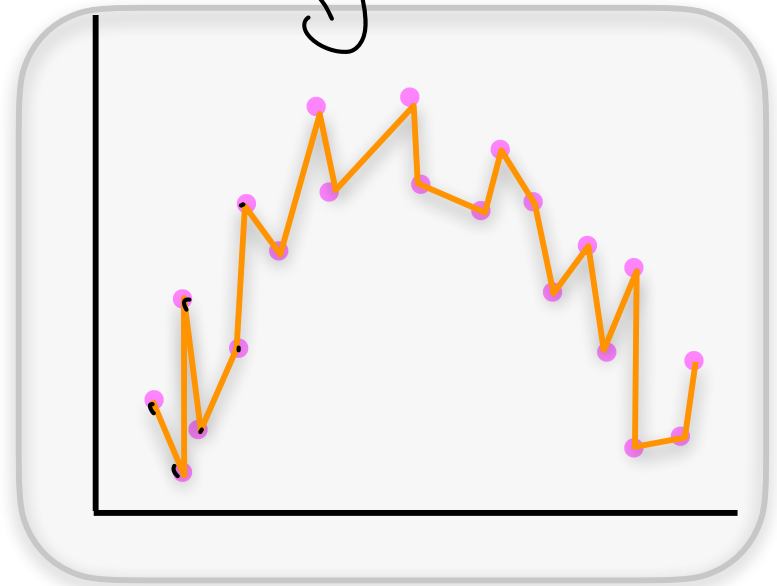


# Expressivity / Representation Power

linear

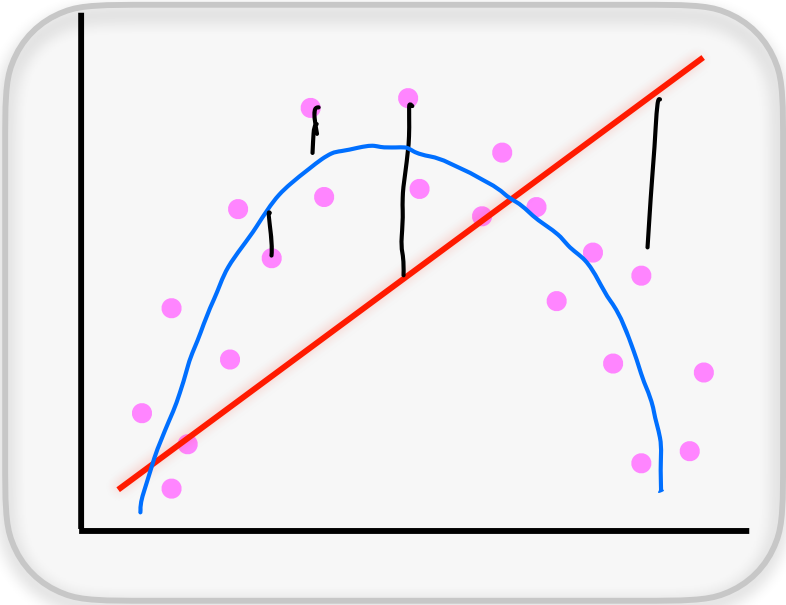
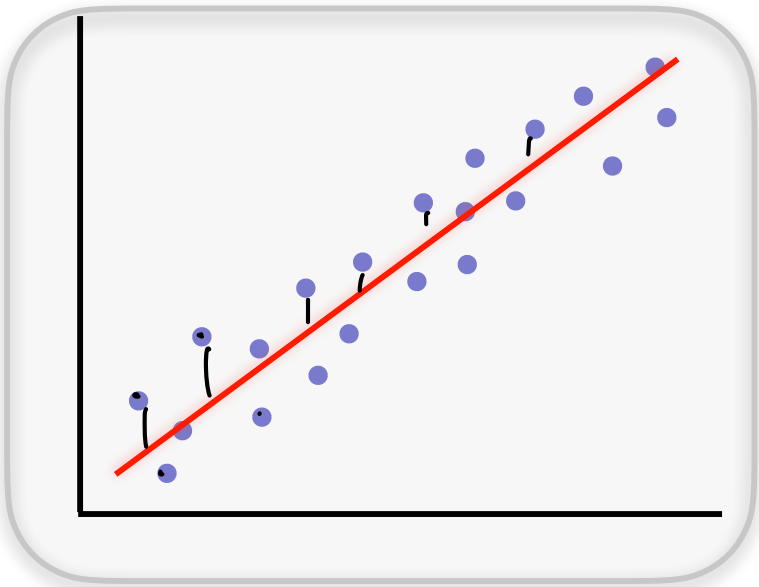


strong



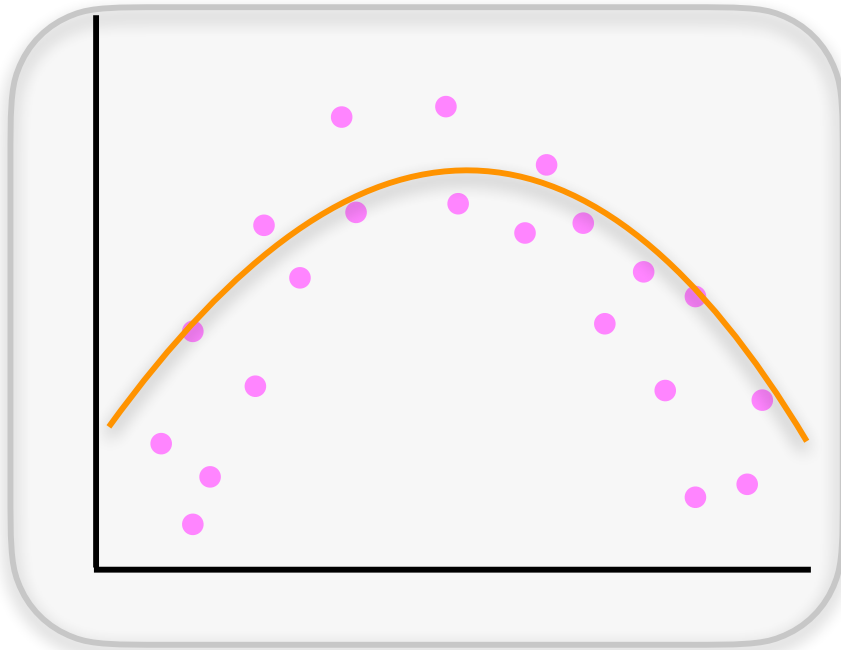
Expressive: Functions in class can represent  
“complicated” functions.

# Linear Function



best linear fit

# Review: generalized linear regression



Transformed data:  $h_0(x) = 1$

$$h_1(x) = x$$

$$h_2(x) = x^2$$

$$h_3(x) = x^3$$

$$\underline{h(x)} = \begin{bmatrix} h_1(x) \\ h_2(x) \\ \vdots \\ h_p(x) \end{bmatrix}$$

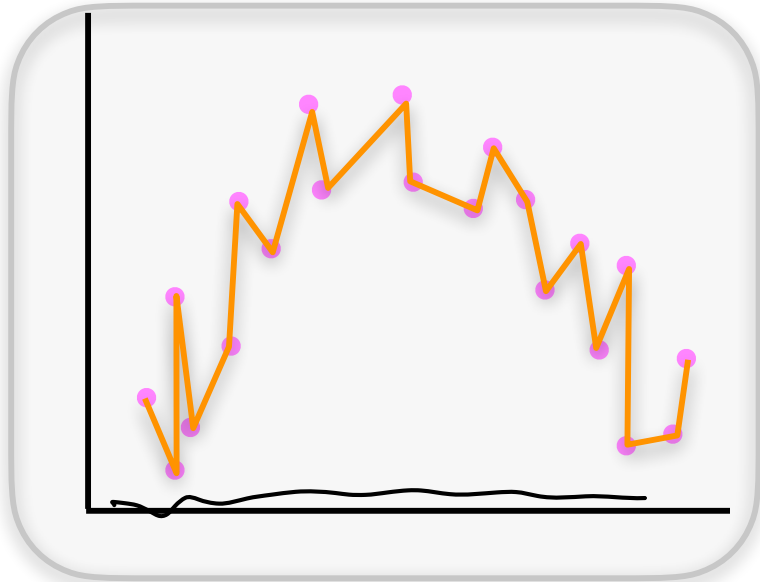
$\vdots$

Hypothesis: linear in  $h$

$$y_i \approx \underbrace{h(x_i)^T}_{} w$$

map powerful  $x \cdot w$

# Review: Polynomial Regression



$$h(x) = \begin{pmatrix} 1 \\ x \\ \vdots \\ x^p \end{pmatrix}$$

$$f(x) = \langle w, h(x) \rangle$$

$w \in \mathbb{R}^{p+1}$

Lagrange's Interpolation Theorem

Given a data set  $\{(x_i, y_i)\}_{i=1}^n$   
 $\exists$  polynomial  $p$  of degree  $n-1$   
s.t.  $y_i = \underline{p}(x_i)$

(condition: no  $(x_i, y_i), (x_i', y_i')$  s.t.  
 $x_i = x_i'$   
 $y_i \neq y_i'$ )

# Approximation Theory Setup

No training

- Goal: to show there exists a neural network that has small error on training / test set.

Sanity

- Set up a natural baseline:

$$\inf_{f \in \mathcal{F}} L(f) \text{ v.s. } \inf_{g \in \text{continuous functions}} L(g)$$

$\mathcal{F} \subseteq \{\text{continuous functions}\}$

Diagram: A wavy line under  $f \in \mathcal{F}$  has an arrow pointing down to "NN". A wavy line under  $g \in \text{continuous functions}$  has an arrow pointing down to the set definition.

## Example

$$\rho: \mathbb{R}_+ \rightarrow \mathbb{R}_+$$

(i) loss  $\ell(f(x), y) = \rho(\underline{y \cdot f(x)})$ ,  $\rho$ -Lipschitz

$$|\ell(z) - \ell(z')| \leq \rho |z - z'|, z, z' \in \mathbb{R}$$

e.g. hinge loss

$$\ell(y f(x)) = \max\{0, 1 - y f(x)\}$$

1-Lipschitz

$$\mathcal{L}(f) = \int \ell(y f(x)) d\mu(x, y)$$

$\mu(x, y)$  distribution over  $(x, y)$

# Decomposition

$f \in \mathcal{F}$   
 $g \in \{\text{continuous functions}\}$

$$\begin{aligned} & L(f) - L(g) \\ &= \int (l(yf(x)) - l(yg(x))) d\mu(x, y) \\ &\leq \int \underbrace{|l(yf(x)) - l(yg(x))|}_{\text{red underline}} d\mu(x, y) \\ &\leq \int \rho |yf(x) - yg(x)| d\mu(x, y) \\ &\leq \int \rho \cdot |y| |f(x) - g(x)| d\mu(x, y) \quad \text{Assume} \\ &\leq \rho \int |f(x) - g(x)| d\mu(x) \quad |y| \leq 1 \end{aligned}$$

red underline



# Specific Setups

$f \in \mathcal{F}$   
 $g \in$  continuous functions

- “Average” approximation: given a distribution  $\mu$

$$\|f - g\|_{\mu} = \int_x |f(x) - g(x)| d\mu(x)$$

- “Everywhere” approximation

$$\|f - g\|_{\infty} = \sup_x |f(x) - g(x)| \geq \|f - g\|_{\mu}$$

$$\begin{aligned} \|f - g\|_{\mu} &= \int_x |f(x) - g(x)| d\mu(x) \\ &\leq \int_x \sup_x |f(x) - g(x)| d\mu(x) \end{aligned}$$

$$= \|f - g\|_{\infty} \int_x d\mu(x) = \|f - g\|_{\infty}$$

# Polynomial Approximation

**Theorem (Stone-Weierstrass):** for any <sup>continuous</sup> function  $f$ , we can **approximate it** on any compact set  $\Omega$  by a sufficiently high degree polynomial: for any  $\epsilon > 0$ , there exists a polynomial  $p$  of sufficient high degree, s.t.,

$$\max_{x \in \Omega} |f(x) - p(x)| \leq \epsilon.$$

Intuition: **Taylor expansion!**

$$f(x) = f(x_0) + f'(x_0)(x-x_0) + \frac{f''(x_0)}{2}(x-x_0)^2 + \dots$$

$$f(x) \approx \langle w, \phi(x) \rangle$$

$$\phi(x) = (1, x-x_0, (x-x_0)^2, \dots)$$

$$w = (f(x_0), f'(x_0), \frac{f''(x_0)}{2}, \dots)$$

# Kernel Method

fixed (can be infinite dimensions)

$$x \mapsto \phi(x), \quad f(x) = \langle w, \phi(x) \rangle \quad \text{for some } w$$

one only needs to evaluate

$$K(x, x') = \langle \phi(x), \phi(x') \rangle$$

## Polynomial kernel

$d$ -dim input

$$\phi(x) = (1, x_1, x_2, \dots, x_d, x_1 x_2, \dots, x_d^2, \dots, x_d^p)$$

## Gaussian Kernel

$1$ -dim

$$K(x, x') = \exp\left(-\frac{|x-x'|^2}{2\sigma^2}\right)$$

$$\phi(x) = e^{-\frac{x^2}{2\sigma^2}} \left(1, \sqrt{\frac{1}{1!}} \frac{x}{\sigma}, \sqrt{\frac{1}{2!}} \left(\frac{x}{\sigma}\right)^2, \dots\right)$$

→ Gaussian Kernel has <sup>inf</sup> strong power representation

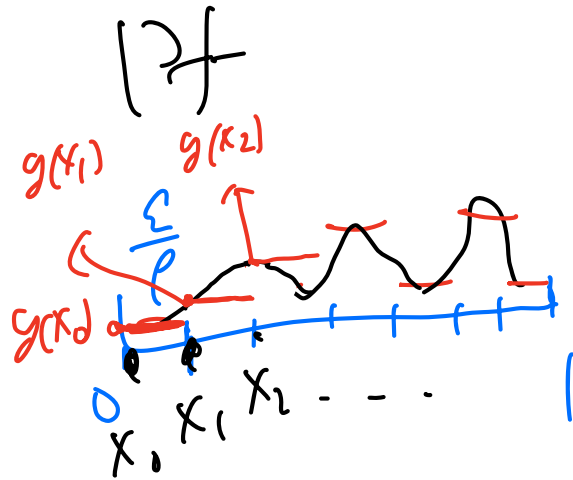
# 1D Approximation



**Theorem:** Let  $g : [0,1] \rightarrow \mathbb{R}$ , and  $\rho$ -Lipschitz. For any  $\epsilon > 0$ ,  $\exists$  2-layer neural network  $f$  with  $\lceil \frac{\rho}{\epsilon} \rceil$  nodes, threshold activation:  $\sigma(z) : z \mapsto \mathbf{1}\{z \geq 0\}$  such that

$$\sup_{x \in [0,1]} |f(x) - g(x)| \leq \epsilon.$$

# Proof of 1D Approximation



$$\text{Let } m \triangleq \left\lceil \frac{\rho}{\epsilon} \right\rceil, \quad x_i \triangleq \frac{(i-1)\epsilon}{\rho}$$

$$f(x) = \sum_{i=1}^m a_i \cdot \mathbb{1}_{\{x - x_i \geq 0\}}$$

$$a_1 = g(x_0), \quad a_i = g(x_i) - g(x_{i-1}), \quad i=1, \dots, m$$

• if  $x < x_1$ ,  $\mathbb{1}_{\{x - x_i \geq 0\}} = 0, i=1, \dots, m$

$$\underline{f(x) = g(x_0)}$$

• if  $x_1 \leq x < x_2$ ,  $\mathbb{1}_{\{x - x_i \geq 0\}} = 0, i=2, \dots, m$

$$\underline{f(x) = g(x_0) + g(x_1) - g(x_0) = g(x_1)}$$

$$|g(x) - f(x)| = |g(x) - f(x_i)|$$

$$\leq |g(x) - g(x_i)| + \underbrace{(g(x_i) - f(x_i))}_0$$

$$\leq \rho \cdot |x - x_i|$$

$$\leq \rho \cdot \frac{\epsilon}{\rho} = \epsilon \quad \square$$

$$x_i \leq x \leq x_{i+1}$$

$$\Rightarrow f(x) = g(x_{i+1})$$

$\forall x,$

say  $x_i \leq x$ , closest  
 $i=0, \dots, m$

# Multivariate Approximation

**Theorem:** Let  $g$  be a continuous function that satisfies  $\left(\frac{\epsilon}{\rho}\right)$ -Lip  
 $\|x - x'\|_\infty \leq \delta \Rightarrow |g(x) - g(x')| \leq \epsilon$  (Lipschitzness).

Then there exists a **3-layer ReLU neural network** with  $O\left(\frac{1}{\delta^d}\right)$  nodes that satisfy

*uniform average approximation*

*(curse of dimensionality)*  
*right  $\Omega\left(\frac{1}{\delta^d}\right)$*

$$\int_{[0,1]^d} |f(x) - g(x)| dx = \|f - g\|_1 \leq \epsilon$$

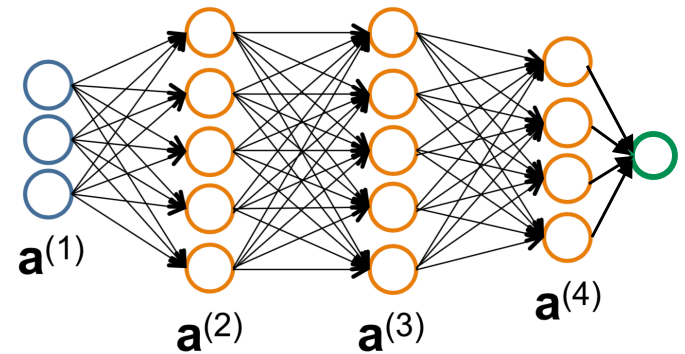
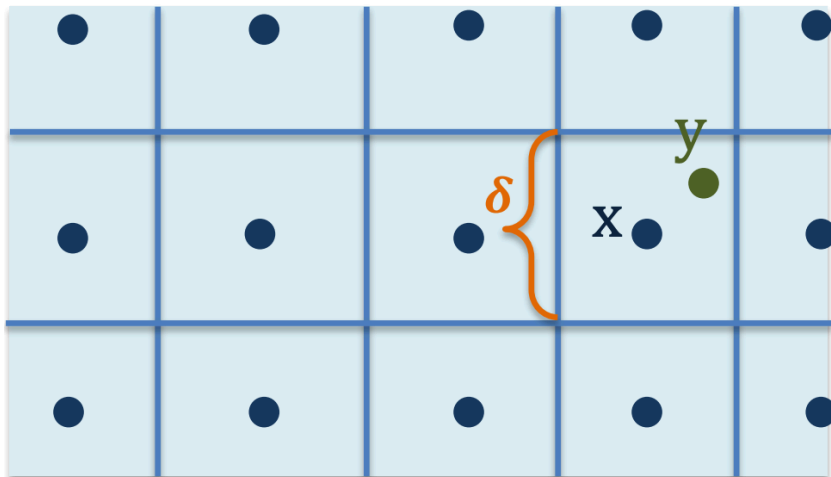


Figure credit to Andrej Risteski

# Partition Lemma

**Lemma:** let  $g, \delta, \epsilon$  be given. For any partition  $P$  of  $[0,1]^d$ ,  $P = (R_1, \dots, R_N)$  with all side length smaller than  $\delta$ , there exists  $(\alpha_1, \dots, \alpha_N) \in \mathbb{R}^N$  such that

*everywhere approximately*

$$\sup_{x \in [0,1]^d}$$

$$|g(x) - h(x)| \leq \epsilon \text{ with } h(x) :=$$

$$\sum_{i=1}^N \alpha_i \mathbf{1}_{R_i}(x).$$

$\mathbf{1}_{R_i}(x) = \mathbb{1}\{x \in R_i\}$

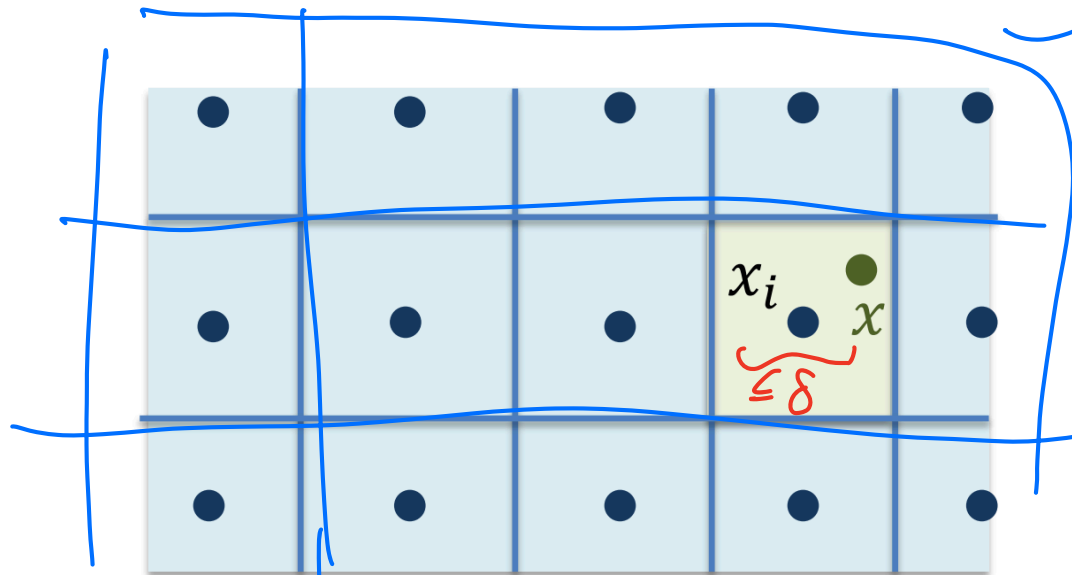
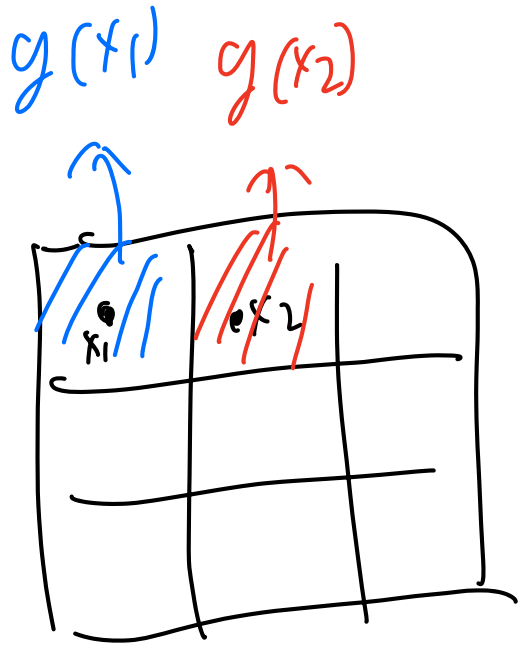


Figure credit to Andrej Risteski

# Proof of Partition Lemma

Pf: For each  $R_i$ , pick  $\forall x_i \in R_i$ , set  $\alpha_i \stackrel{\circ}{=} g(x_i)$



$$\sup_{x \in [0,1]^d} |g(x) - h(x)|$$

$$= \sup_{i \in \{1, \dots, N\}} \sup_{x \in R_i} |g(x) - h(x)|$$

$$\leq \sup_{i \in \{1, \dots, N\}} \sup_{x \in R_i} \left( \underbrace{|g(x) - g(x_i)|}_{\circ} + \underbrace{|g(x_i) - h(x_i)|} \right)$$

$$\leq \sup_{i \in \{1, \dots, N\}} \sup_{x \in R_i} (\epsilon + 0)$$

$$= \epsilon$$

□



# Proof of Multivariate Approximation Theorem

Idea:  $h(x) = \sum_{i=1}^N \alpha_i \mathbb{1}_{D_i}(x)$  in the lemma

1) use a 2-layer NN to approximate  
 $x \mapsto \mathbb{1}_{D_i}(x)$

2) find a linear combination to represent  $h$

$$\Rightarrow \|f - g\|_1 \leq \|f - h\|_1 + \|h - g\|_1 \leq \varepsilon$$

Let  $f(x) = \sum_{i=1}^N \alpha_i f_i(x)$      $\alpha_i = y(x_i)$ ,  $f_i$  to approximate  $\mathbb{1}_{D_i}(x)$

$$\|f - h\|_1 = \left\| \sum_{i=1}^N \alpha_i (\mathbb{1}_{D_i} - f_i) \right\|_1$$

$$\leq \sum_{i=1}^N |\alpha_i| \|\mathbb{1}_{D_i} - f_i\|_1 \leq \varepsilon$$

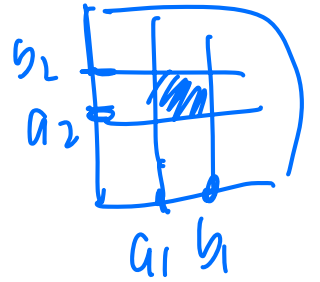
want to show  $\|\mathbb{1}_{D_i} - f_i\|_1 \leq \frac{\varepsilon}{\sum_{i=1}^N |\alpha_i|}$

what if  $\sum_{i=1}^N |\alpha_i| = 0 \rightarrow y(x_i) = 0, |y(x)| \leq \varepsilon$  use 0-function

# Proof of Multivariate Approximation Theorem

(2) Construct  $f_i$

Recall:  $D_i \subseteq [a_1, b_1] \times [a_2, b_2] \cdots [a_d, b_d]$



★ bump function

Given  $\gamma > 0$ , define  $\{b: D \in U\}$

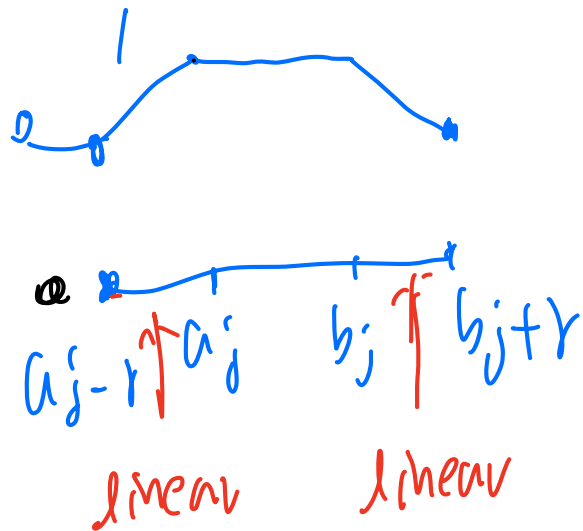
$$g_{r,i,j}(z) = b\left(\frac{z - (a_j - r)}{\gamma}\right) - b\left(\frac{z - a_j}{\gamma}\right)$$

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

$$z \in [a_j, b_j], g_{r,i,j}(z) = 1$$

$$z \in [a_j - r, b_j + r], g_{r,i,j}(z) = 0$$

$$\star \Leftrightarrow 0, g_{r,i,j} \rightarrow \perp [a_j, b_j]$$

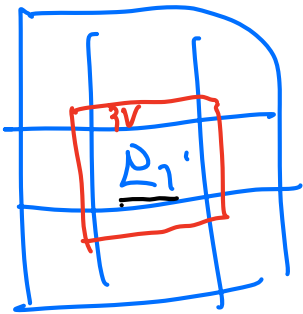


# Proof of Multivariate Approximation Theorem

Define  $g_r(x) = 6 \left( \frac{d}{\sum_{j=1}^d g_{r,j}(x^j)} - (d-1) \right)$

$$x = \begin{pmatrix} x^1 \\ x^2 \\ \vdots \\ x^d \end{pmatrix}$$

$$g_r(x) = \begin{cases} 1 & \text{if } x \in Q_i \\ 0 & \text{if } x \notin [a_1-r, b_1+r] \times \dots \times [a_d-r, b_d+r] \end{cases}$$



Since  $r \rightarrow 0$ ,  $g_{r,j} \rightarrow \mathbb{1}_{[a_j, b_j]}$

$$\rightarrow g_r \rightarrow \mathbb{1}_{Q_i}$$

$\exists r$  with  $\|g_r - \mathbb{1}_{Q_i}\|_1 \leq \frac{\epsilon}{\sum_{i=1}^d \alpha_i}$

Let  $f_i = g_r$

$$f = \sum_{i=1}^d \alpha_i f_i$$

□

# Universal Approximation

---

**Definition:** A class of functions  $\mathcal{F}$  is **universal approximator** over a compact set  $S$  (e.g.,  $[0,1]^d$ ), if for every continuous function  $g$  and a target accuracy  $\epsilon > 0$ , there exists  $f \in \mathcal{F}$  such that

$$\sup_{x \in S} |f(x) - g(x)| \leq \epsilon$$

# Stone-Weierstrass Theorem

---

**Theorem:** If  $\mathcal{F}$  satisfies

1. Each  $f \in \mathcal{F}$  is continuous.
2.  $\forall x, \exists f \in \mathcal{F}, f(x) \neq 0$
3.  $\forall x \neq x', \exists f \in \mathcal{F}, f(x) \neq f(x')$
4.  $\mathcal{F}$  is closed under multiplication and vector space operations,

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

$$\forall g : S \rightarrow R, \epsilon > 0, \exists f \in \mathcal{F}, \|f - g\|_{\infty} \leq \epsilon.$$

# Example: cos activation

---

# Example: cos activation

---

# Other Examples

---

**Exponential activation**

**ReLU activation**





# Recent Advances in Representation Power

---

- Depth separation
- Analyses of different architectures
  - Graph neural network
  - Attention-based neural network
- Finite data approximation
- ...