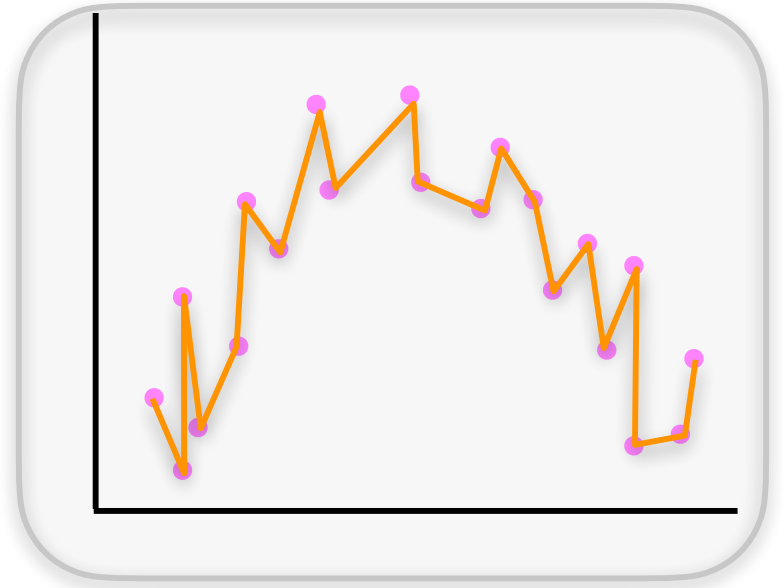
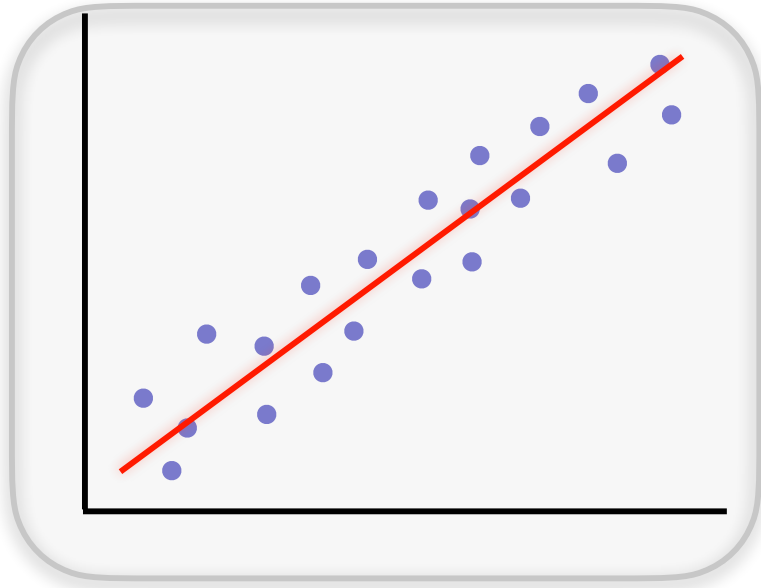


Approximation Theory



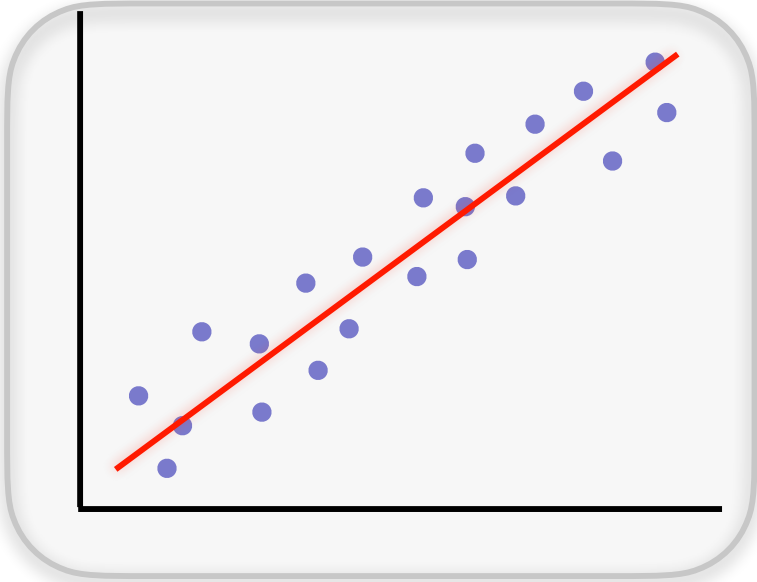
Expressivity / Representation Power

\mathcal{F} : linear, NN, ...
Ziel

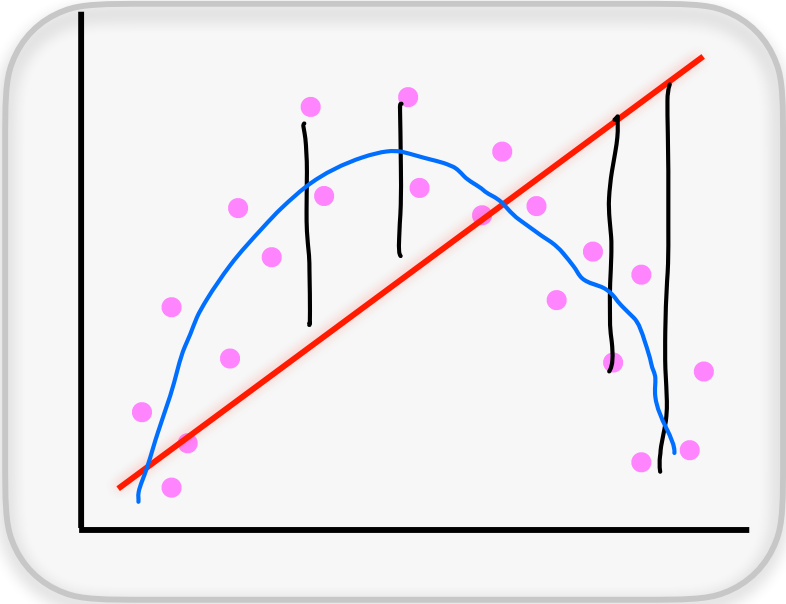


Expressive: Functions in class can represent “complicated” functions.

Linear Function



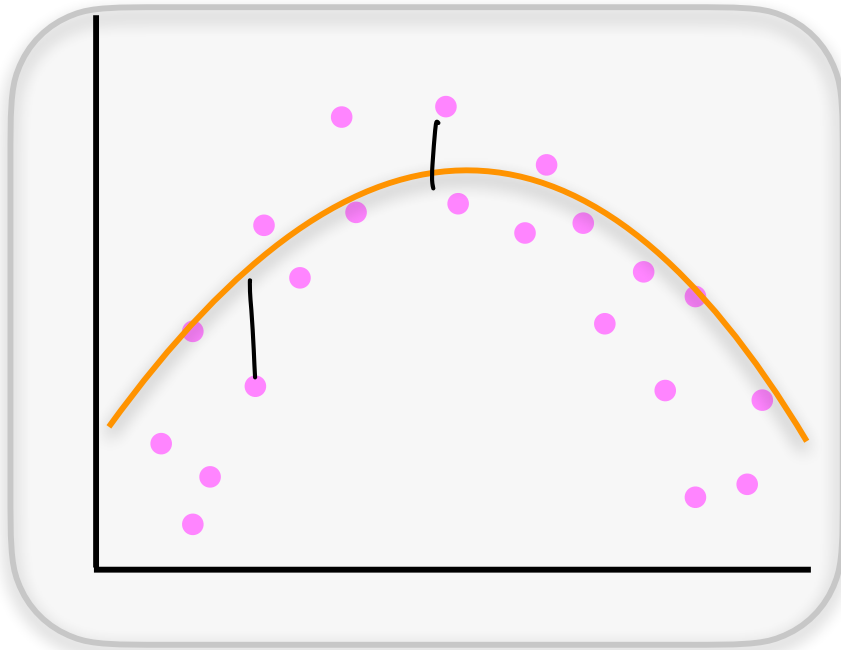
large bias



best linear fit

Review: generalized linear regression

smaller bias



Transformed data:

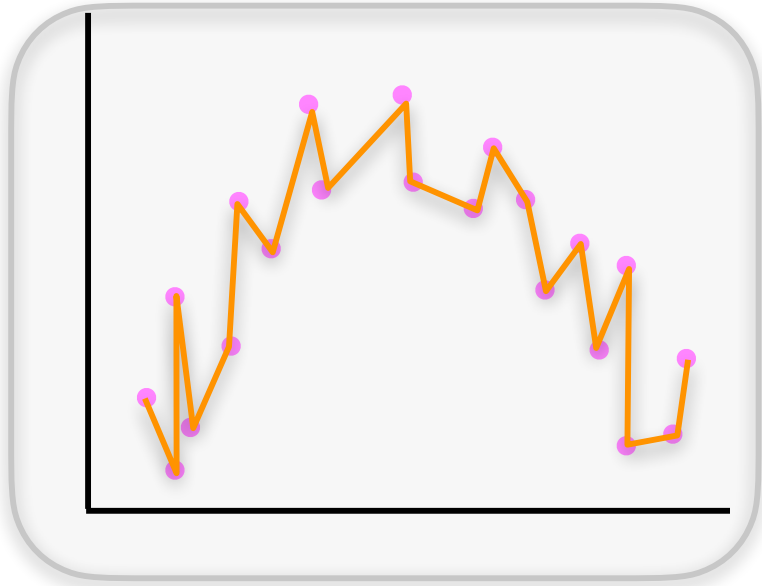
$$\begin{aligned} h_1(x) &= 1 \\ h_2(x) &= x \\ h_3(x) &= x^2 \end{aligned} \quad h(x) = \begin{bmatrix} h_1(x) \\ h_2(x) \\ \vdots \\ h_p(x) \end{bmatrix}$$

Hypothesis: linear in h

$$y_i \approx h(x_i)^T w$$

Review: Polynomial Regression

o bita)



$$h(x) = \begin{pmatrix} 1 \\ x \\ x^2 \\ \vdots \\ x^p \end{pmatrix}$$

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

n data points

$$p \geq n-1$$

Lagrange's Interpolation Theorem

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

Approximation Theory Setup

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

- Set up a natural baseline:

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

L : loss

\mathcal{F} : NN function class

Example

(1) loss $l(f(x), y) = l(y \cdot f(x))$, ρ -Lipshitz

$$|l(z) - l(z')| \leq \rho \cdot |z - z'|, z, z' \in \mathcal{R}$$

e.g. hinge loss

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

1-Lip

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

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

Decomposition

$$\begin{aligned} & L(f) - L(g) \\ &= \int (\ell(yf(x)) - \ell(yg(x))) d\mu(x, y) \\ &\leq \int |\ell(yf(x)) - \ell(yg(x))| d\mu(x, y) \\ &\leq \int \rho \cdot |yf(x) - yg(x)| d\mu(x, y) \end{aligned}$$

(assume $|y| \leq 1$)

$$\leq \rho \cdot \int |f(x) - g(x)| d\mu(x)$$

Specific Setups

- “Average” approximation: given a distribution μ

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

↑ (sample)

- “Everywhere” approximation

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

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

Polynomial Approximation

Theorem (Stone-Weierstrass): for any function f , we can **approximate it** on any compact set Ω 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!**

$$f(x) = f(x_0) + f'(x_0)(x-x_0) + \frac{f''(x_0)}{2}(x-x_0)^2 + \dots$$
$$f(x) = \langle \omega, \phi(x) \rangle$$
$$\phi(x) = (1, x-x_0, (x-x_0)^2, \dots)$$
$$\omega = \left(f(x_0), f'(x_0), \frac{f''(x_0)}{2}, \dots \right)$$

Kernel Method

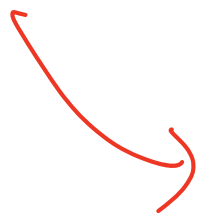
(can be inf-dim, fixed)
 $x \mapsto \phi(x), f(x) = \langle w, \phi(x) \rangle$

only need evaluate

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

Polynomial kernel d -dim

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



*if p is large
→ strong approx power*

Gaussian Kernel

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

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



kernels have strong approximation power

1D Approximation

$g \in C_1 : \rho$ -Lipschitz

Theorem: Let $g : [0,1] \rightarrow \mathbb{R}$, and ρ -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

Pf :

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

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

$$a_1 = g(0), \quad a_i = g(x_i) - g(x_{i-1}), \quad i=2, \dots, m$$

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

$$f(x) = g(0)$$

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

$$f(x) = g(x_0) + g(x_1) - g(x_0) = g(x_1)$$

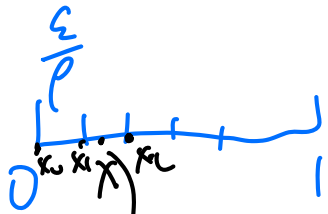
$$|g(x) - f(x)| = |g(x) - f(x_i)|, \quad x_i \in x \text{ does not}$$

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

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

$$= \rho \cdot \frac{\epsilon}{\rho} = \epsilon$$

□



$g(x)$

$$f(x) = g(x_1)$$

Multivariate Approximation

$$x \in \mathcal{D}^d, \quad g(x) \in \mathcal{R}$$

Theorem: Let g be a continuous function that satisfies $\|x - x'\|_\infty \leq \delta \Rightarrow |g(x) - g(x')| \leq \epsilon$ (Lipschitzness).

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

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

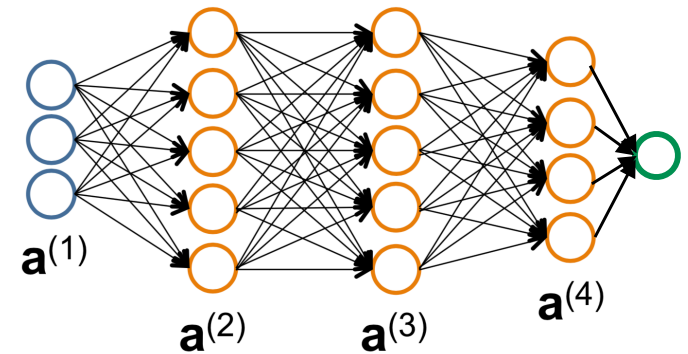
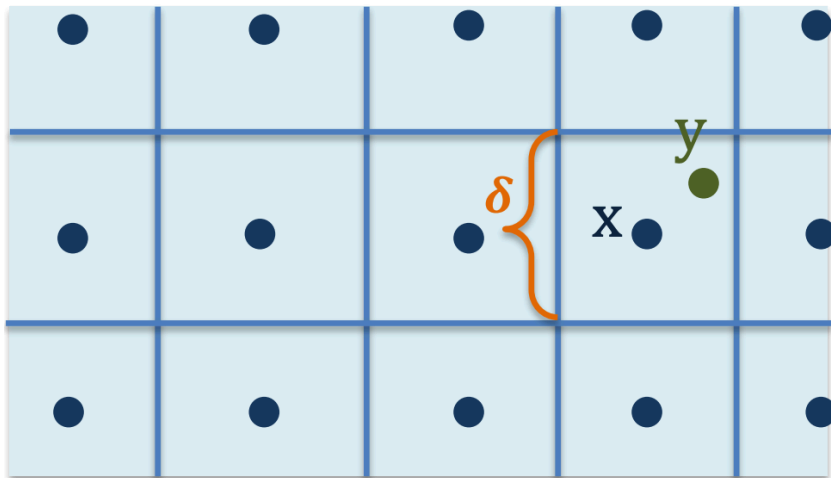


Figure credit to Andrej Risteski

Partition Lemma

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

$$\sup_{x \in [0,1]^d} |g(x) - h(x)| \leq \epsilon \text{ with } h(x) := \sum_{i=1}^N \alpha_i \mathbf{1}_{R_i}(x).$$

Generalization of

1d theorem



non-parametric regression

Hölder space

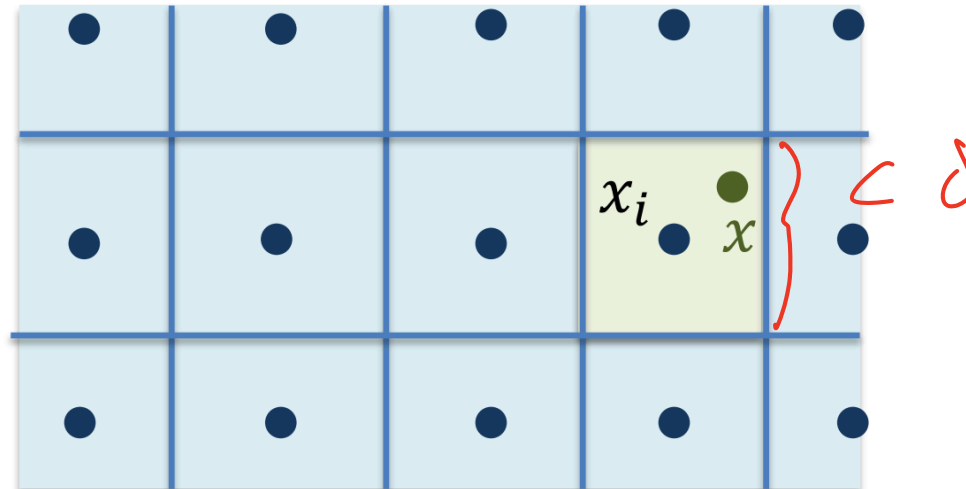
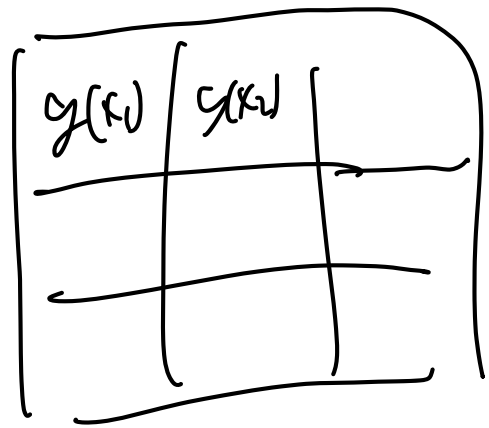


Figure credit to Andrej Risteski

Proof of Partition Lemma

Pf: For each \mathcal{P}_i , pick $x_i \in \mathcal{P}_i$, set $\alpha_i \triangleq g(x_i)$



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

$$|g(x) - h(x)| = \sup_{i \in \{1, \dots, n\}} \sup_{x \in \mathcal{P}_i} |g(x) - h(x)|$$

$$\leq \sup_{i \in \{1, \dots, n\}} \sup_{x \in \mathcal{P}_i} \left(\underbrace{|g(x) - g(x_i)|}_{\leq \epsilon} + \underbrace{|g(x_i) - h(x_i)|}_{\leq \epsilon} \right)$$

$$\leq \epsilon + \epsilon$$

$$= 2\epsilon$$

□

Proof of Multivariate Approximation Theorem

Idea: $h(x) = \sum_i \alpha_i \mathbb{1}_{\mathcal{R}_i}(x)$ in the learning

- 1) use 2-layer NN to approximate $x \mapsto \mathbb{1}_{\mathcal{R}_i}(x)$
- 2) find a linear combination to represent h

$\Rightarrow \|f - g\|_1 \leq \|f - h\|_1 + \underbrace{\|h - g\|_1}_{\leq \epsilon}$

① Let $f = \sum_{i=1}^N \alpha_i f_i$, $\alpha_i = g(x_i)$

goal: f_i to approximate $\mathbb{1}_{\mathcal{R}_i}(x)$

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

$$\leq \sum_{i=1}^N |\alpha_i| \cdot \|\mathbb{1}_{\mathcal{R}_i} - f_i\|_1 \stackrel{\epsilon}{\leq}$$

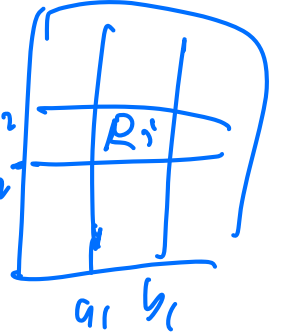
want to show $\|\mathbb{1}_{\mathcal{R}_i} - f_i\|_1 \leq \frac{\epsilon}{\sum_{i=1}^N |\alpha_i|}$

what if $\sum_{i=1}^N |\alpha_i| = 0$? $g(x) \equiv 0$, use 0-function, $|g(x)| \leq \epsilon$

Proof of Multivariate Approximation Theorem

★ bump function goal: construct f_i

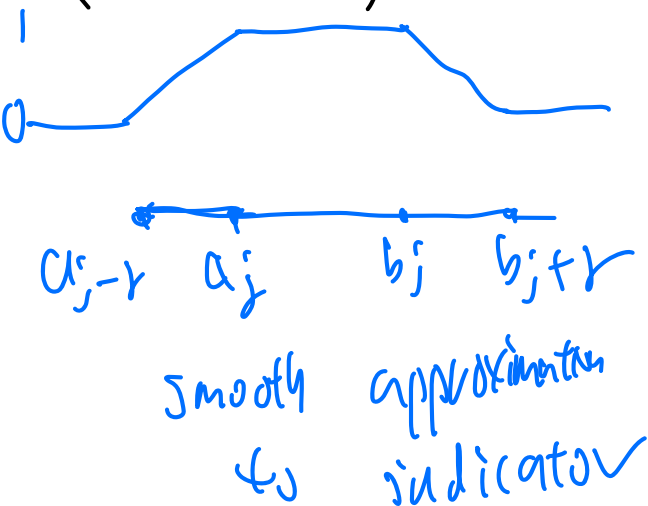
Recall $\mathcal{P}_r = \mathcal{O} \left[\underbrace{[a_1, b_1] \times [a_2, b_2] \times \dots \times [a_d, b_d]}_{\substack{b_2 \\ a_2}} \right]$



Given $r > 0$, define $\phi: \mathbb{R} \rightarrow [0, 1]$

$$g_{r,j}(z) = \phi\left(\frac{z - (a_j - r)}{r}\right) - \phi\left(\frac{z - a_j}{r}\right) - \phi\left(\frac{z - b_j}{r}\right) + \phi\left(\frac{z - (b_j + r)}{r}\right)$$

($j = 1, \dots, d$)



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

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

$$r \rightarrow 0, \quad \underbrace{g_{r,j} \rightarrow \mathbb{1}_{[a_j, b_j]}}_{\text{in } \mathcal{L}}$$

Proof of Multivariate Approximation Theorem

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

$$x = \begin{pmatrix} \vdots \\ x^d \end{pmatrix} \quad g_r(x) = \begin{cases} 1 & \text{if } x \in \mathcal{Q}_i \\ 0 & \text{if } x \notin [\underline{a}_i - r, \overline{a}_i + r] \times [\underline{a}_2 - r, \overline{a}_2 + r] \cdots [\underline{a}_d - r, \overline{a}_d + r] \\ [0, 1] & \text{otherwise} \end{cases}$$

Since $r \rightarrow 0$, $g_r \rightarrow \underline{\mathbb{1}}_{\mathcal{Q}_i}$

$$\exists r \text{ with } \|g_r - \underline{\mathbb{1}}_{\mathcal{Q}_i}\|_1 \leq \frac{\epsilon}{\sum_j |d_j|}$$


Let $f_i = g_r$

$$f \stackrel{\circ}{=} \sum_{i=1}^N \alpha_i f_i$$

□

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
- ...