

Approximation Theory

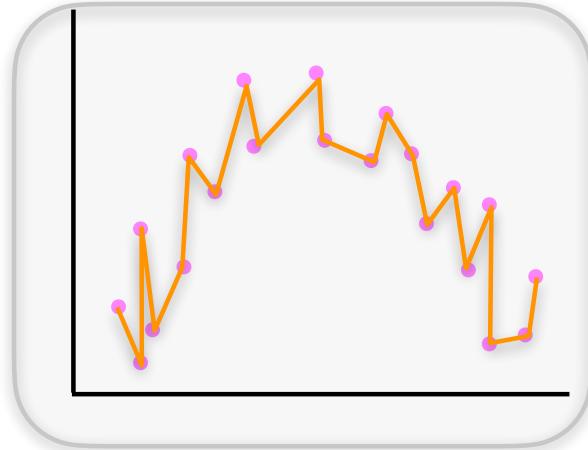
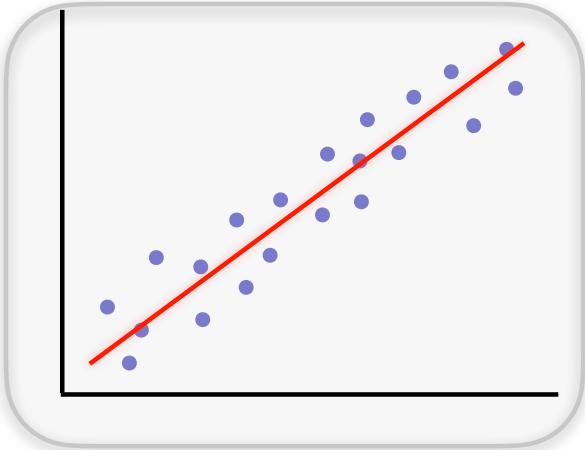


UNIVERSITY *of* WASHINGTON

W

Expressivity / Representation Power

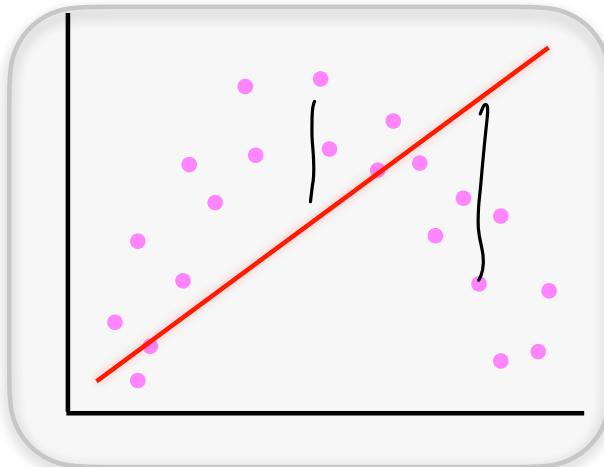
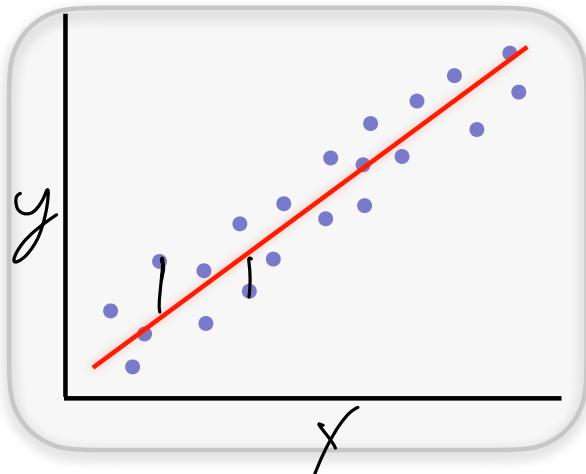
$$f \in \mathcal{F}$$



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

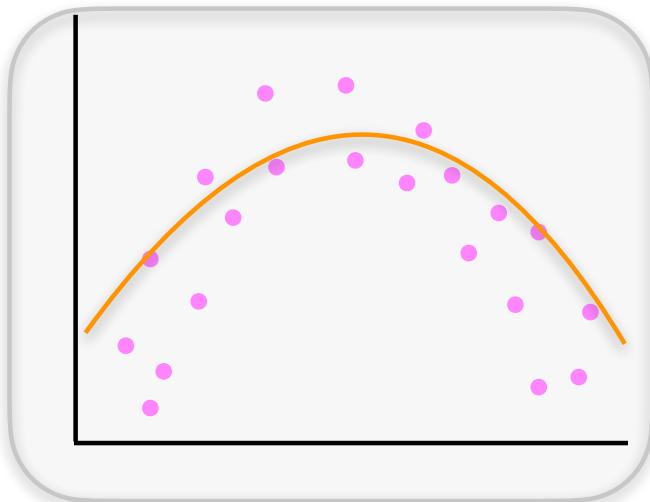
Linear Function

large gaps



best linear fit

Review: generalized linear regression



Transformed data:

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

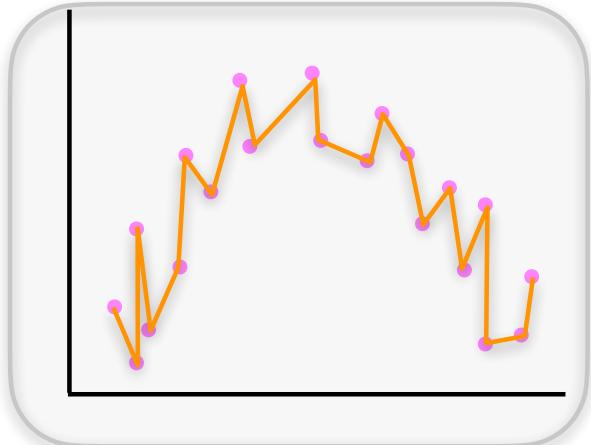
Transformations:

$$h_1(x) = 1$$
$$h_2(x) = x$$
$$h_3(x) = x^2$$

Hypothesis: linear in h

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

Review: Polynomial Regression



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

$$f(x) = \langle w, h(x) \rangle$$
$$w \in \mathbb{R}^n$$

Lagrange's Interpolation Theorem
 $\{(x_i, y_i)\}_{i=1}^n, \exists$ polynomial f
of degree $(n-1)$, $\forall i, y_i = f(x_i)$

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:

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

Example

① $\ell(f(x), y) = \ell(yf(x))$, ρ (lip shift)

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

e.g., hinge loss

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

$$L(f) = \int \ell(yf(x)) dM(x, y)$$

$M(x, y)$: distribution over (x, y)

Decomposition

$$\begin{aligned} & f, g \\ & \int (l(yf(x)) - l(yg(x))) d\mu(x, y) \\ & \leq \int |l(yf(x)) - l(yg(x))| d\mu(x, y) \\ & \leq \int \rho |yf(x) - yg(x)| d\mu(x, y) \\ & \leq \rho \int |f(x) - g(x)| d\mu(x, y) \\ & \text{(assume } |g| \leq 1) \end{aligned}$$

Specific Setups

- “Average” approximation: given a distribution μ

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

- “Everywhere” approximation

$$\underbrace{\|f - g\|_{\infty}}_{x \in S} = \sup |f(x) - g(x)| \geq \|f - g\|_{\mu}$$

$$\begin{aligned}\|f - g\|_{\mu} &= \int_X |f(x) - g(x)| d\mu(x) \\ &\leq \int_X \sup_{\mathcal{X}} |f(x) - g(x)| d\mu(x) = \underbrace{\|f - g\|_{\infty} \cdot \int_X d\mu(x)}_{= \|f - g\|_{\infty}}\end{aligned}$$

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!**

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

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

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

$$w = (f(x_0), f'(x_0), \dots)$$

Kernel Method

$$\begin{cases} x \mapsto \phi(x) \quad , \quad f(x) = \langle w, \phi(x) \rangle \\ \Rightarrow k(x, x') = \langle \phi(x), \phi(x') \rangle \end{cases}$$

Polynomial kernel

d -dimensional x

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

Gaussian Kernel

$$x, x' \in \mathcal{L}$$

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

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

\Rightarrow Kernel (as if representation power)

1D Approximation

Theorem: Let $g : [0,1] \rightarrow R$, 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

$$\text{Pf: } m \triangleq \left\lceil \frac{\rho}{\varepsilon} \right\rceil$$

$$x_i \triangleq \frac{(i-1)\varepsilon}{\rho}$$

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

$$a_0 = g(x_0) \quad a_i = g(x_i) - g(x_{i-1})$$

— if $x < x_1$, $f(x) = a_0 = g(x_0)$

— $x_1 \leq x < x_2$ $f(x) = g(x_1)$

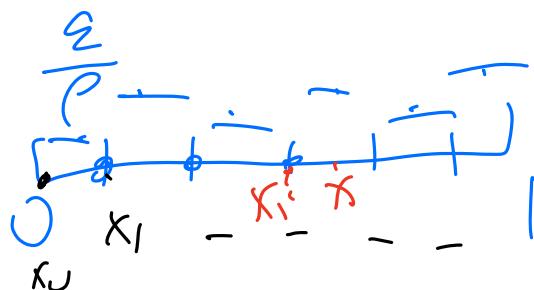
$$|g(x) - f(x)| = |g(x) - f(x_i)|, \quad x_i \leq x, \text{ closest}$$

$$(\text{triangle inequality}) \leq |g(x) - g(x_i)| + |g(x_i) - f(x_i)|$$

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

$$\leq \frac{\rho}{\rho} \cdot \frac{\varepsilon}{\rho} = \varepsilon$$

g : target function



□

Multivariate Approximation

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

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

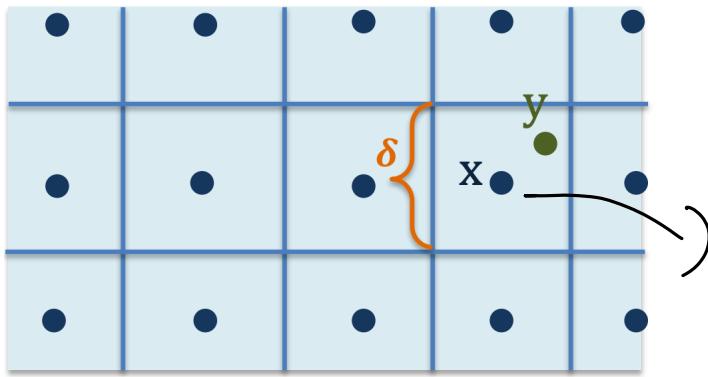
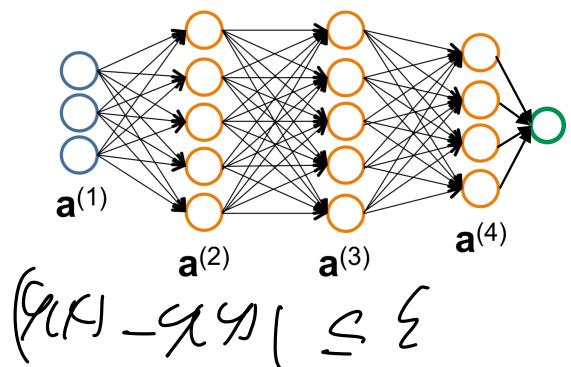


Figure credit to Andrej Risteski



Partition Lemma

generalization of 1d approx

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

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

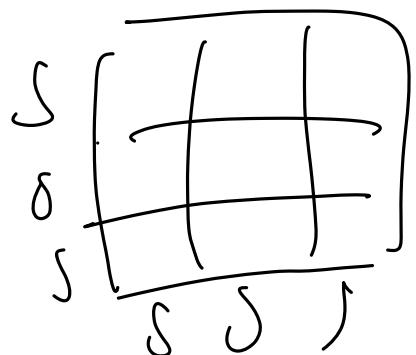
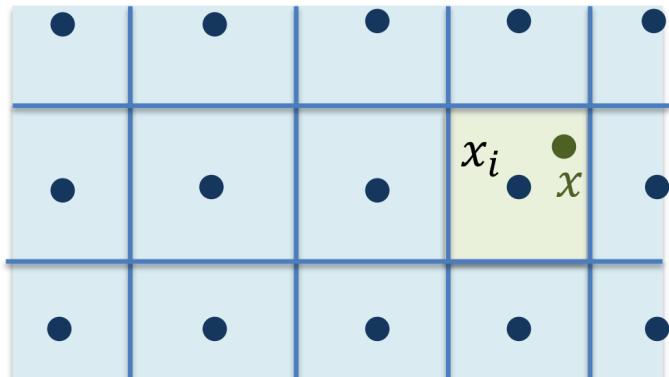


Figure credit to Andrej Risteski

Proof of Partition Lemma

(Pf): For each \mathcal{Q}_i , pick $x_i \in \mathcal{Q}_i$, set
 $\alpha_i \subseteq g(x_i)$ $h(x) = \sum \alpha_i \mathbb{1}_{\{x \in \mathcal{Q}_i\}}$

$$\sup_{x \in [\bar{0}, 1]^d} |g(x) - h(x)| = \sup_{r \in \{1, \dots, N\}} \sup_{x \in \mathcal{Q}_i} (g(x) - h(x))$$

$$\leq \sup_{r \in \{1, \dots, N\}} \sup_{x \in \mathcal{Q}_i} (|g(x) - g(x_i)| + \underline{\frac{|g(x_i) - h(x_i)|}{\partial}})$$

$$|(x - x_i)|/\omega \leq \delta$$

$$\rightarrow |g(x) - g(x_i)| \leq \varepsilon \leq \varepsilon$$

(D)

Proof of Multivariate Approximation Theorem

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

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

2) then find a linear configuration to represent

$$\Rightarrow \|f - g\|_1 \leq \underbrace{\|f - h\|_1}_{\text{want}} + \underbrace{\|h - g\|_1}_{\leq \varepsilon} \leq \varepsilon$$

$$f(x) = \sum_{i=1}^N \alpha_i f_i(x)$$

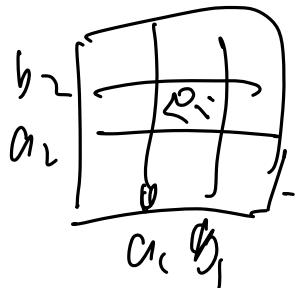
$$\begin{aligned} \|f - h\|_1 &= \left\| \sum_{i=1}^N \alpha_i (\mathbb{1}_{R_i} - f_i) \right\|_1 \\ &\leq \sum_{i=1}^N |\alpha_i| \cdot \|\mathbb{1}_{R_i} - f_i\|_1 \end{aligned}$$

$$\text{Want } \|\mathbb{1}_{R_i} - f_i\|_1 \leq \frac{\varepsilon}{\sum_{i=1}^N |\alpha_i|}$$

Proof of Multivariate Approximation Theorem

② Construct f_j

$$\Omega_j \stackrel{\text{def}}{=} [\bar{a}_1, \bar{b}_1] \times [\bar{a}_2, \bar{b}_2] \dots \times [\bar{a}_d, \bar{b}_d]$$

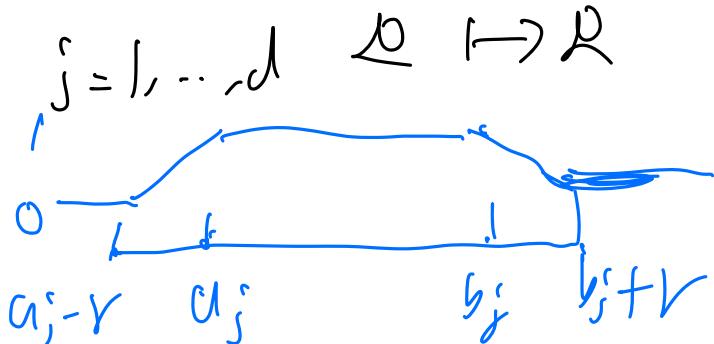


☆ bump function

Given $r > 0$,

Define $g_{r,j}(z) =$

$$6 \left(\frac{z - (a_j - r)}{r} \right) - 6 \left(\frac{z - a_j}{r} \right) - 6 \left(\frac{z - b_j}{r} \right) + 6 \left(\frac{z - (b_j + r)}{r} \right)$$



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

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

$\rightarrow 0$, $g_{r,j}(z) = 2[a_j, b_j]$

Proof of Multivariate Approximation Theorem

Define $g_r(x) = \sum_{j=1}^d g_{r,j}(x_j) - (d-1)$

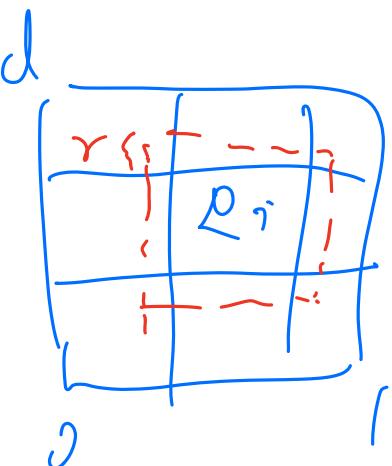
$$x = \begin{pmatrix} x_1 \\ \vdots \\ x_d \end{pmatrix}$$

$$g_r(x) = \begin{cases} 1 & \text{if } x \in R; \\ 0 & \text{if } x \notin [a_i-r, b_i+r] \text{ for all } i \\ [0, 1] & \text{o.w.} \end{cases}$$

Since $r \rightarrow 0$, $g_r \rightarrow 1_{\mathbb{R}^d}$ ϵ
 $\exists r \text{ s.t. } \| (g_r - 1_{\mathbb{R}^d}) \|_1 \leq \frac{\epsilon}{2d}$

Define $f_i = g_r$

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



□

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

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

$$h(x) = \sum \alpha_i \mathbb{1}_{\mathcal{D}_i}$$

Stone-Weierstrass Theorem

Theorem: If \mathcal{F} satisfies

1. Each $f \in \mathcal{F}$ is continuous.
2. $\forall x, \exists f \in \mathcal{F}, f(x) \neq 0$ $g(x) \neq g(x')$
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, $f_1, f_2 \in \mathcal{F}, f_1 \cdot f_2 \in \mathcal{F}$

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

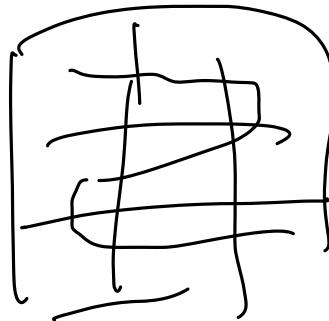
Exponential activation

ReLU activation

Curse of Dimensionality

$(\frac{1}{\delta^d})$ modes

- Unavoidable in the worse case



- Barron's theory

Equiv^l

Recent Advances in Representation Power

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