

Lecture 2: basic norm-based apx.

* Apx goal; norm proxy; apx plan.

* univariate apx: L_n , 1 layer, threshold

* \mathbb{R}^d apx: L_1 , 2 layers, ReLU

* "universal" apx: \mathbb{R}^d , L_n , 1 layer, "any" activation

$$R(f) = \mathbb{E} \ell(f(x)|Y); \quad \hat{R}(f) = \frac{1}{n} \sum_i \ell(f(x_i)|y_i)$$

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

$\bar{f} \in \mathcal{F}$ best in class

$$R(\bar{f}) \approx \inf_{g \in \mathcal{F}} R(g)$$

$$R(\hat{f}) \approx R(\bar{f})$$

(opt + gen)

ex. \approx

$$\inf_{g \in \mathcal{F}} R(g)$$

Small

To argue $\inf_{g \in \mathcal{F}} R(g)$ is small,

lazy/standard/classical approach:

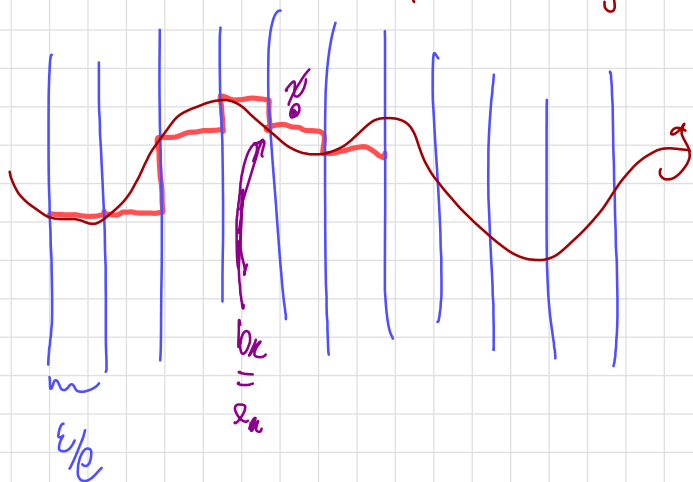
take a b-zy class, e.g., \mathcal{C} all cont. functions.

$\forall h \in \mathcal{C}, \exists g \in \mathcal{F}$ s.t. $\|g - h\|$ small
(for some norm $\|\cdot\|$).

Proposition. Given ϵ -Lipschitz $g: \mathbb{R} \rightarrow \mathbb{R}$, $\epsilon > 0$
 $\exists f(x) = \sum_{j=1}^m a_j \mathbb{1}[x \geq b_j]$, $m = \lceil \frac{e}{\epsilon} \rceil$

s.t. $\forall x \in [0, 1], |f(x) - g(x)| \leq \epsilon$.

Proof.



$$b_0 = 0, \quad a_0 := g(0); \quad b_j := \frac{j\epsilon}{e}, \quad a_j := g(b_j) - g(b_{j-1})$$

Let $x \in (0, 1)$ be given, pick largest $x_k := b_k \leq x$

$$\begin{aligned} |g(x) - f(x)| &\leq |g(x) - g(x_k)| + |g(x_k) - f(x_k)| + |f(x_k) - f(x)| \\ &\leq e \left(\frac{\epsilon}{e} \right) + |g(x_k)| + (g(x_k) - g(x_{k-1})) + g(x_{k-1}) - \dots \\ &= \epsilon. \end{aligned}$$

Remark.

* This is a non-adaptive; e.g., pays for flat regions.

It seems DNNs are good because they effectively model natural phenomena.

* Polynomials pay for flat regions

* Better version in hwk1; but over \mathbb{R}^d open.

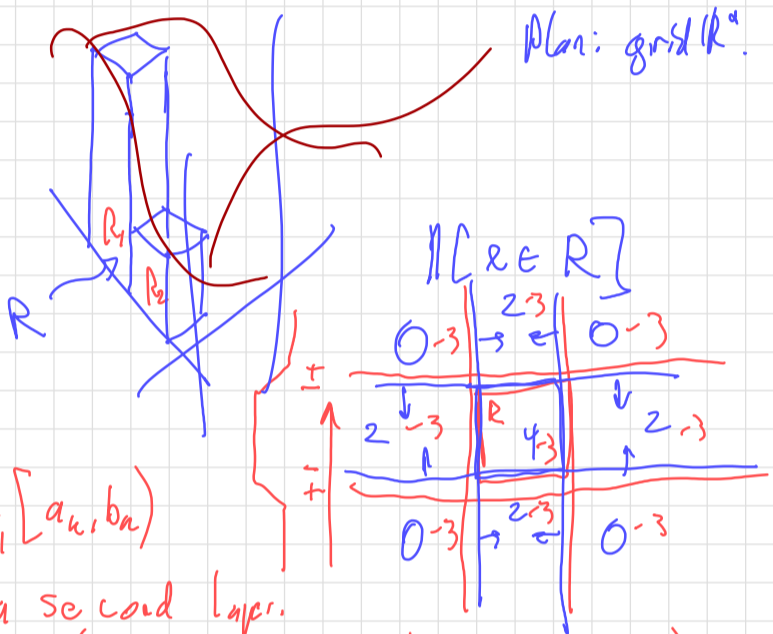
Theorem. \forall cont $g: \mathbb{R}^d \rightarrow \mathbb{R}$ $\varepsilon > 0 \in [0, 1]^d$
 given $\delta > 0$ s.t. $\|x-x'\|_\infty \geq \delta \Rightarrow |g(x)-g(x')| \leq \varepsilon$
 $\exists f: \mathbb{R}^d \rightarrow \mathbb{R}$ 2-hidden layer ReLU
 network with $\geq (\frac{1}{\delta})^d$ nodes
 s.t. $\int_{[0,1]^d} |f(x) - g(x)| dx \leq \varepsilon$.

Remark.

* L_1 is bad; next theorem uses L_∞
 * $\frac{1}{\delta}^d$ is a disaster; e.g., CIFAR dim = 3072.

Also matching lower bound.

Proof.



$$R = \prod_{k=1}^d [a_k, b_k]$$

Use a second layer.

$$g_\gamma(x) = \sigma \left(\sum_{k=1}^d \left[\sigma \left(\frac{a_k - (x_k - \gamma)}{\gamma} \right) - \sigma \left(\frac{a_k - x_k}{\gamma} \right) - \sigma \left(\frac{b_k - x_k}{\gamma} \right) + \sigma \left(\frac{b_k - (x_k + \gamma)}{\gamma} \right) \right] \right) - (d-1)$$

Note $g_\gamma \rightarrow \mathbb{1}[\cdot \in R]$ pointwise as $\gamma \rightarrow 0$

Analysis part: $\exists h(x) := \sum_{j=1}^m a_j \mathbb{1}[x \in R_j]$

s.t. $\sup_{x \in [0,1]^d} |h(x) - g(x)| \leq \varepsilon/2$.

Define $f(x) := \sum_{j=1}^m a_j g_{\gamma_j}(x)$.

Then $\lim_{\gamma \rightarrow 0} \int |f - g| \leq \frac{\varepsilon}{2} + \lim_{\gamma \rightarrow 0} \int |f - h|$
 $\leq \frac{\varepsilon}{2} + \sum_{j=1}^m |a_j| \lim_{\gamma \rightarrow 0} \int_{[0,1]^d} |\mathbb{1}[x \in R_j] - g_{\gamma_j}|$
 $\rightarrow 0$.

- * "universal approximation"
- * Barron class / Fourier.
- * NTK / "near initialization" ↙ related
- * "Benefits of depth" (non-algorithms)

Definition. \mathcal{F} is a "universal approximator" if \forall cont $g: [0,1]^d \rightarrow \mathbb{R}$, $\forall \epsilon > 0$
 $\exists f \in \mathcal{F}$, $\sup_{x \in [0,1]^d} |f(x) - g(x)| \leq \epsilon$.

Theorem. Let "sigmoidal" σ be given:
 σ continuous, $\lim_{z \rightarrow -\infty} \sigma(z) = 0$, $\lim_{z \rightarrow \infty} \sigma(z) = 1$

Let \mathcal{F} denote 1-hidden layer networks of unbounded width $(\sum_{j=1}^m a_j \sigma(w_j^T x + b_j))$ $\forall m \geq 0$
 $\forall a_j, w_j, b_j$.
 Then \mathcal{F} is a universal approximator.

Remark:

- * Found concurrently in ~1988
 Cybenko, Hornik - Stinchcombe - White, Funkeleski,
 Barron. (Proofs differ)
- * Skill you!

Proof idea: (Hornik - Stinchcombe - White)
 $\mathbb{1} \left[x \in \bigcap_{k=1}^d [a_k, b_k] \right] = \prod_{k=1}^d \mathbb{1} \left[x_k \in [a_k, b_k] \right]$

Weierstrass: polynomials of unbounded degree are universal approximator

Stone-Weierstrass: polynomial-like \mathcal{F} are uni. appx

Why are shallow networks \rightarrow

Note $e^{w^T x + b} e^{u^T x + c} = e^{(w+u)^T x + (b+c)}$

Expl.) activations are closed under products;
 Arbitrary activations okap.