

Lecture 2: shallow apx

Ann.

* project draft up

* lecture 2-6 fixed by 11:50pm Monday

Index

* shallow constructive apx

* univariate apx

* multivariate apx

* universal apx

* infinite width

* initialization & overparameterization

* apx - opt - gen - other topics

1 $R =$ abstract measure of future performance

$\mathcal{F} =$ deep network architecture

$\hat{F} =$ output of \mathcal{M}_g

$\bar{F} =$ good choice for R in \mathcal{F}

$R(\hat{F}) \approx R(\bar{F})$ (opt/gen)

now most $R(\hat{F})$ small (apx)

want to formalize "is inf $R(\hat{F})$ small?"

(A) IF: (i) know what future data looks like
(ii) believe in a perf criterion for datum, $\mathcal{L}(f(x), y)$

\Rightarrow measure inf $\int_{\text{test}} \mathcal{L}(f(x), y)$ small.

(B) IF R satisfies "some regularity", we can say

$\inf_{\text{test}} R(\hat{F}) \approx \inf_{\text{test}} R(\bar{F})$
"every function"

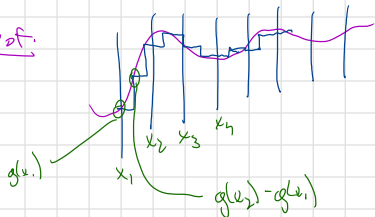
$\Rightarrow \forall g \exists f, \forall x | f(x) - g(x) | \leq \epsilon.$

Proposition. $\forall \epsilon$ -Lip $g: \mathbb{R} \rightarrow \mathbb{R}$, $\epsilon > 0$,

3 f 2-layer network with biases
activation $\sigma(z) = \mathbb{1}[z \geq 0]$ $m = \lceil \frac{e}{\epsilon} \rceil$

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

Proof.



$$b_j := \frac{(j-1)\epsilon}{e}, \quad a_j := g(b_j)$$

$$a_{j+1} := g(b_{j+1}) - g(b_j)$$

$$f(x) = \sum_j a_j \sigma(x - b_j)$$

Let x given, $b_n = x$ be largest

$$|f(x) - g(x)| \leq |f(x) - f(b_n)| + |f(b_n) - g(b_n)| + |g(b_n) - g(x)|$$

$$\leq 0 + 0 + e \left(\frac{e}{e} \right)$$

$$f(b_n) = \sum_{j=1}^n a_j \sigma(b_n - b_j) = \sum_{j=1}^n a_j$$

$$= g(b_1) + \sum_{j=2}^n (g(b_j) - g(b_{j-1})) = g(b_n).$$

Remarks.

* pays for large flat regions,
limitation of the proof.

* polynomials pay for the flat part

Theorem. $\forall \epsilon$ -Lipschitz $g: \mathbb{R}^d \rightarrow \mathbb{R}$ & $\Sigma \times \Omega$
 \exists f 3-layer biased network with
 $\sigma(z) = \text{ReLU}(z)$, $m = O\left(\frac{d}{\epsilon}\right)^d$

$$\int_{\Omega \times \Sigma} |f(x) - g(x)| dx \leq \epsilon.$$

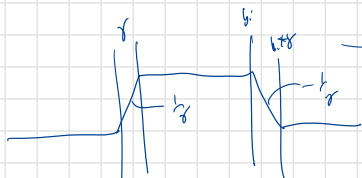
Proof.

Claim: done if p_i , a 2-layer ReLU net, with

$$\int_{\Omega \times \Sigma} |p_i(x) - \mathbb{1}_{x \in S_i}| \leq \frac{\epsilon}{\sum |g_j|}$$

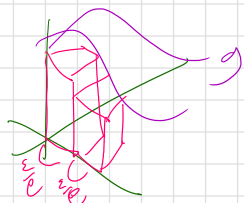
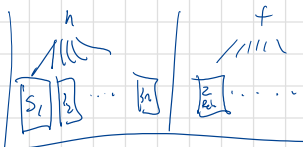
Proof. Using $f(x) = \sum g_j(x) p_j(x)$,

$$\int |h(x) - f(x)| dx \leq \sum_{j=1}^m |g_j(x)| \int |p_j(x) - \mathbb{1}_{x \in S_j}| dx \leq \epsilon.$$



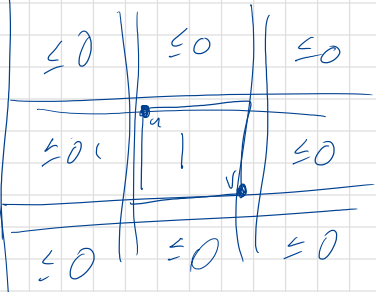
$$g_i(r) = \sigma\left(\frac{r - (a_i - \delta)}{\gamma}\right) - \sigma\left(\frac{r - (a_i)}{\gamma}\right) - \sigma\left(\frac{r - b_i}{\gamma}\right) + \sigma\left(\frac{r - (b_i + \delta)}{\gamma}\right).$$

$$x \mapsto \sigma\left(\sum g_i(x_i)\right) - (d-1)$$



(S_1, \dots, S_m) partition of $[0,1]^d$ into cubes, $S_j := \prod_{i=1}^d [a_{ij}, b_{ij})$

$$h(x) = \sum_{j=1}^m g_j(x) \mathbb{1}_{x \in S_j} = \sum_{j=1}^m g_j(x) \prod_{i=1}^d \mathbb{1}_{x_i \in [a_{ij}, b_{ij})}$$



Standard universal approximation.

Theorem. [Hornik - Stinchcombe - White '89, Leshno '43]

Suppose $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ is \wedge not a polynomial,
continuous, but

$\forall g: \mathbb{R}^d \rightarrow \mathbb{R}$ continuous, $\forall \epsilon$

$\exists f: \mathbb{R}^d \rightarrow \mathbb{R}$ 2-layer biased σ network

$\forall x \in [0, 1]^d, |f(x) - g(x)| \leq \epsilon.$ [supremum norm
uniform norm]

Proof for.

* We'll invoke lemma for polynomial-linear
kernels.

* Consider $\sigma(r) = \exp(r),$

where $\exp(r+s) = \exp(r)\exp(s).$

Remarks.

- * Not unique to neural networks (e.g., SVM + RBF).
- * implicit exponentially large set.
- * proofs.
- * not polynomial is necessary