

Lecture 3 Universal approx continued

<ul style="list-style-type: none"> * hw tomorrow * typed notes... * tablet tweaks * hybrid 	Plan for next lectures: 1-4 shallow constructive appx 5-8 initialization / overparameterization 9-11 deeper topics / RNN transfer / distribution modeling 12- opt 20 - ??
--	--

Last week: folklore constructive appx over $\mathbb{R}^d, \mathbb{R}^d$

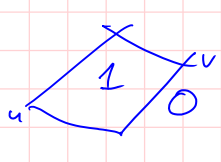
Theorem ("universal approximation", Hornik-Stinchcombe-Ullate '89, Leschao '93).
 Let any $\sigma: \mathbb{R} \rightarrow \mathbb{R}$ continuous & not a polynomial.
 For any cont. $g: \mathbb{R}^d \rightarrow \mathbb{R}$, any $\varepsilon > 0$
 \exists 2-layer σ -network $f: \mathbb{R}^d \rightarrow \mathbb{R}$ s.t. $|f(x) - g(x)| \leq \varepsilon \forall x \in \mathbb{R}^d$.
↳ choice of m

Remarks -

- * C-Lip, 3-layer, avg dist.
- * "celebrated", even though true for all MC models
- * Bahes in exponential dependence on dimension
- * false if m fixed or σ a polynomial

Lemma. Same statement, except $\sigma(r) = \exp(r)$.

Remark on proof. Recall our folklore proof need to appx $\prod_{i=1}^d \mathbb{1}[x_i \geq u_i] \mathbb{1}[x_i < v_i]$
 \Rightarrow products are helpful
 \Rightarrow proof technique is \rightarrow reduction to polynomials.



Proof. Define $\mathcal{F} := \{x \mapsto a^T \exp(Vx) : m \geq 0, a \in \mathbb{R}^m, V \in \mathbb{R}^{m \times d}\}$.

Proof is complete if we can show \mathcal{F} satisfies conditions of the Stone-Weierstrass (aka \mathcal{F} is polynomial-like).

- $\mathcal{F} \ni f$ is continuous \checkmark
- $\forall x \in \mathbb{R}^d, \exists f \in \mathcal{F}, f(x) \neq 0$ (easy: $x \mapsto \exp(0^T x) = 1$)
- $\forall x \neq x', \exists f, f(x) \neq f(x')$ (easy: $z \mapsto \exp(\langle -x', x-x' \rangle) \exp(\langle z, x-x' \rangle)$
 note $f(x) = \exp(0) = 1 \neq \exp(\langle x-x', x-x' \rangle)$ \checkmark)

(4) \mathcal{F} "closed under VS & poly operations": given $a^T \exp(Vx)$ & $u^T \exp(Wx)$, $b \in \mathbb{R}, c \in \mathbb{R}$

$$b a^T \exp(Vx) + c u^T \exp(Wx) = \begin{bmatrix} \mathbb{R}^m \\ \mathbb{R}^n \end{bmatrix} \begin{bmatrix} b a \\ c u \end{bmatrix} \exp \left(\begin{bmatrix} V \\ W \end{bmatrix} x \right);$$

$\mathbb{R}^{m+n}; m \times n$

for products

$$\left[\sum_{i=1}^m a_i \exp(v_i^T x) \right] \left[\sum_{j=1}^n u_j \exp(w_j^T x) \right] = \sum_{ij} a_i u_j \exp(\langle v_i + w_j, x \rangle) \checkmark$$

new outer weights

Remarks of (4) fails for polynomials of fixed degree.

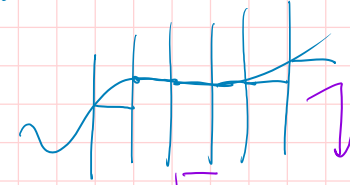
* Going through proof in detail reveals exponential dependence on dim

note to self: adaptive open problem empirical part & math part.

To handle other activations:

- Given $\varepsilon > 0, g$, pick $(a, v) \in \mathbb{R}^m \times \mathbb{R}^{m \times d}$ s.t. $\forall x \in \mathbb{C}, \mathbb{I}^d \quad |a^T \exp(Vx) - g(x)| \leq \varepsilon/2$.
- Write $\exp(r) = \sum_{j=1}^n u_j \sigma(w_j r + b_j)$, define $f(x) = \sum_{i=1}^m a_i \sum_{j=1}^n u_j \sigma(w_j v_i^T x + b_j)$.
univariate approximation, easier, maybe hard problem.

Appx. with infinite width
 * popular in recent years
 * can always convert to finite width (e.g., via sampling)
 * can often write target with equality
 * maybe helps with adaptivity



Recall our univariate approximation proof

Proposition. Let $g: \mathbb{R} \rightarrow \mathbb{R}$ be continuously differentiable $g(0) = 0$.
 Then, $\forall x \in \mathbb{C}, \mathbb{I}^d$

$$g(x) = \int_0^x g'(b) \mathbb{1}[x \geq b] db.$$

"infinite width network"

$$\sum_{j=1}^m a_j \mathbb{1}[x \geq b_j]$$

$m = \lceil \frac{\varepsilon}{\varepsilon} \rceil$

Proof. By FTC $g(x) - g(0) = \int_0^x g'(b) db = \int_0^x g'(b) \mathbb{1}[x \geq b] db$.

Remarks:

- * Sampling. Define $Z := \int_0^1 |g'(b)| db$, and note $\frac{|g'(b)|}{Z}$ is a probability dist over \mathbb{C}, \mathbb{I}^d .
- * sample $b_j \sim (\cdot)$, define $a_j := Z \operatorname{sgn}(g'(b_j))$;

note via proposition that

$$\begin{aligned} \mathbb{E} a_j \mathbb{1}[x \geq b_j] &= \int_0^1 \underbrace{[Z \operatorname{sgn}(g'(b))]}_{a_j} \mathbb{1}[x \geq b] \underbrace{\frac{|g'(b)|}{Z}}_{\text{probability dist}} db \\ &= \int_0^1 (\operatorname{sgn}(g'(b)) \cdot |g'(b)|) \mathbb{1}[x \geq b] \frac{Z}{Z} db \\ &= \int_0^1 g'(b) \mathbb{1}[x \geq b] db = g(x), \end{aligned}$$

so this single node is an unbiased estimator of $f(x)$ $\forall x \in \mathbb{C}, \mathbb{I}^d$

* Sample m such nodes, define $f(x) := \frac{1}{m} \sum_{j=1}^m a_j \mathbb{1}[x \geq b_j]$.

Sampling theorem (in nodes) say $\mathbb{E} (f(x) - g(x))^2 \leq \frac{1}{m} \int \frac{Z^2}{b}$

Define $f(x) := \frac{1}{m} \sum_{j=1}^m a_j \mathbb{1}[x \geq b_j]$, then $\mathbb{E} (f(x) - g(x))^2 \leq \frac{Z^2}{b}$;

width does not pay for large flat regions (somewhat adaptivity).

Next time: we'll write multivariate continuous functions using Fourier functions, # nodes will scale some function of the function ("Barron norm").