

Approximation by ridge functions with weights in a specified set

David Stewart & Palle Jorgensen

University of Iowa
Mathematics

September 11, 2018



Extreme Learning Machines

Extreme Learning Machines are neural networks with one hidden layer where the training is only carried out on the weights in the *output*. The weights before the activation/sigmoidal functions are generated randomly and are left unchanged.

Approximation by ridge functions

We consider approximating functions $f(\mathbf{x})$ for $\mathbf{x} \in [-1, +1]^m = J^m$ by functions in

$$V_{\mathcal{W}} := \text{span} \left\{ \mathbf{x} \mapsto \varphi(\mathbf{w}^T \mathbf{x}) \mid \varphi \in C(\mathbb{R}), \mathbf{w} \in \mathcal{W} \right\}$$

for \mathcal{W} a specific subset in \mathbb{R}^m .

Questions:

- ▶ For what sets \mathcal{W} is $V_{\mathcal{W}}$ dense in $C(J^m)$?
- ▶ If $V_{\mathcal{W}}$ is not dense in $C(J^m)$, how well can we approximate (nice) functions by functions in $V_{\mathcal{W}}$?

General framework

We work in a Banach space X ; the approximating functions form a subspace $V \subset X$.

If $\overline{V} = X$ then every object in X can be approximated (arbitrarily well) by elements of V .

But if $\overline{V} \neq X$, then for every $\epsilon > 0$ there are functions $0 \neq f \in X$ where $\inf_{g \in V} \|f - g\|_X \geq (1 - \epsilon) \|f\|_X$. So

$$\sup_{f: \|f\|_X=1} \inf_{g \in V} \|f - g\|_X \quad \text{is either 0 or 1.}$$

Choose a Banach subspace Z compactly embedded in X and we look determine

$$m(V; Z, X) = \sup_{f: \|f\|_Z=1} \inf_{g \in V} \|f - g\|_X.$$

In our case we use:

- ▶ $X = C(J^m)$
- ▶ $V = V_{\mathcal{W}}$
- ▶ $Z = \text{Lip}(J^m)$, the space of Lipschitz functions with semi-norm

$$\|f\|_Z = \sup_{\mathbf{x}, \mathbf{y} \in J^m} \frac{|f(\mathbf{x}) - f(\mathbf{y})|}{\|\mathbf{x} - \mathbf{y}\|_2}$$

(We can quotient out constant functions since $V_{\mathcal{W}}$ always contains these.)

We say $f \in X$ is *unapproximable* by V if

$$\|f\|_X \leq \|f - g\|_X \quad \text{for all } g \in V.$$

For any $f \in X$ if $h \in V$ is the closest point in V to f then $f - h$ is unapproximable by V .

The existence of a closest point is assured if X is a reflexive Banach space, but generally false otherwise.

Separating Hyperplane Theorem

Theorem: If $C \subset X$ is closed and convex and $y \notin C$, then there is a $\mu \in X'$ and $b \in \mathbb{R}$ where

$$\begin{aligned}\langle y, \mu \rangle + b &> 0 \\ \langle z, \mu \rangle + b &\leq 0 \quad \text{for all } z \in C\end{aligned}$$

Specifically, if C is a closed subspace of X , then ν satisfies

$$\begin{aligned}\langle y, \mu \rangle &> 0 \\ \langle z, \mu \rangle &= 0 \quad \text{for all } z \in C.\end{aligned}$$

Cybenko's universal approximation result

George Cybenko's paper from 1989 shows that if (for example) $\sigma(u) = \tanh(u)$ then

$$\overline{\text{span} \{ \mathbf{x} \mapsto \sigma(\mathbf{w}^T \mathbf{x} + b) \mid \mathbf{w} \in \mathbb{R}^m, b \in \mathbb{R} \}} = C(J^m).$$

The proof uses the Separating Hyperplane Theorem

Note: $C(J^m)' = \mathcal{M}(J^m)$, the space of signed Borel measures on J^m with bounded variation and $\langle g, \mu \rangle = \int g(\mathbf{x}) d\mu(\mathbf{x})$.

For the μ in the Separating Hyperplane Theorem

$$\int \sigma(a\mathbf{w}^T \mathbf{x} + b) d\mu(\mathbf{x}) = 0 \quad \text{for all } a, b \in \mathbb{R}$$

so we can show that

$$0 = \mu^{\mathbf{w}}(F) := \mu \left(\left\{ \mathbf{x} \mid \mathbf{w}^T \mathbf{x} \in F \right\} \right) \quad \text{for all Borel } F \subset \mathbb{R}.$$

Fourier Transforms

A Borel measure μ with bounded variation has a Fourier Transform

$$\widehat{\mu}(\boldsymbol{\xi}) = \int_{\mathbb{R}^m} e^{-i\boldsymbol{\xi}^T \mathbf{x}} d\mu(\mathbf{x}).$$

Note that for the μ from the SHT

$$\begin{aligned}\widehat{\mu}(\mathbf{s}\mathbf{w}) &= \int_{\mathbb{R}^m} e^{-i\mathbf{s}\mathbf{w}^T \mathbf{x}} d\mu(\mathbf{x}) \\ &= \int_{\mathbb{R}} e^{-ist} d\mu^{\mathbf{w}}(t) = \widehat{\mu^{\mathbf{w}}}(\mathbf{s}) = 0.\end{aligned}$$

For Cybenko's result, this is true for all $\mathbf{w} \in \mathbb{R}^m$ so $\widehat{\mu}(\boldsymbol{\xi}) = 0$ for all $\boldsymbol{\xi}$, and so $\mu = 0$ contradicting the SHT.

Thus there is no f in $C(J^m)$ that is not in

$$\text{span} \{ \mathbf{x} \mapsto \sigma(\mathbf{w}^T \mathbf{x} + b) \mid \mathbf{w} \in \mathbb{R}^m, b \in \mathbb{R} \}$$

What about specific (finite) \mathcal{W} ?

What about spans of ridge functions

$$V_{\mathcal{W}} := \text{span} \left\{ \mathbf{x} \mapsto \varphi(\mathbf{w}^T \mathbf{x}) \mid \varphi \in C(\mathbb{R}), \mathbf{w} \in \mathcal{W} \right\} ?$$

We can get *lower bounds* on how badly a Lipschitz function f can be approximated by $V_{\mathcal{W}}$ as follows: Pick a measure μ with support in J^m where $\mu^{\mathbf{w}} = 0$ for every $\mathbf{w} \in \mathcal{W}$. Then look for a function f where $\langle f, \mu \rangle = \|f\|_{\infty} \|\mu\|_{\mathcal{M}} \neq 0$.

The measure μ has to satisfy $\hat{\mu}(t\mathbf{w}) = 0$ for all $t \in \mathbb{R}$ and $\mathbf{w} \in \mathcal{W}$.

Example: $\mathcal{W} = \{\mathbf{e}_1, \mathbf{e}_2\}$.

The Fourier transform $\widehat{\mu}(t\mathbf{e}_1) = \widehat{\mu}(t\mathbf{e}_2) = 0$ for all $t \in \mathbb{R}$.

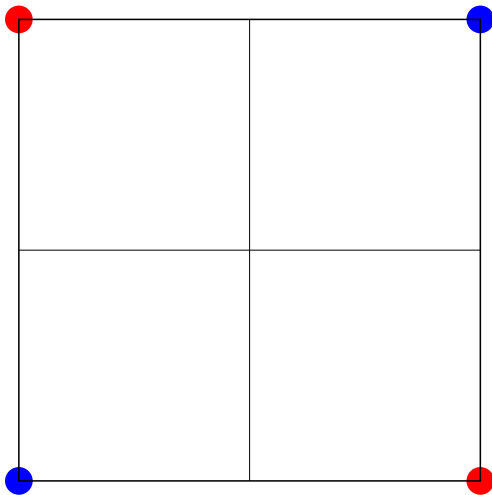
Since the Fourier transform of $\delta_{\mathbf{v}} = \delta(\cdot - \mathbf{v})$ is $\widehat{\delta}_{\mathbf{v}}(\boldsymbol{\xi}) = \exp(-i\boldsymbol{\xi}^T \mathbf{v})$, so we look for $\widehat{\mu}(\boldsymbol{\xi})$ that involves complex exponentials $\exp(-i\boldsymbol{\xi}^T \mathbf{v})$ for $\mathbf{v} \in \mathcal{J}^m$.

Note: $\widehat{\mu}(\boldsymbol{\xi})$ is complex analytic everywhere (entire) so we can look for Taylor series. So...

$$\widehat{\mu}(\boldsymbol{\xi}) = c \xi_1 \xi_2 + \dots$$

We can put $\widehat{\mu}(\boldsymbol{\xi}) = (e^{-i\xi_1} - e^{+i\xi_1}) (e^{-i\xi_2} - e^{+i\xi_2})$, ($c = (-2i)^2$).

Here μ is a sum of δ -functions at (x_1, x_2) with each $x_i = \pm 1$, and the weights at each of these points is $x_1 x_2$.



● = $\text{supp } \mu_-$, ● = $\text{supp } \mu_+$

Now we need to find a function f with, say, $\|f\|_\infty = 1$ and $\langle f, \mu \rangle = \|f\|_\infty \|\mu\|_{\mathcal{M}}$.

Since μ is a sum of (scaled) δ -functions, we can choose $f(\mathbf{v}) = \pm 1$ at each of these points, choosing the sign of $f(\mathbf{v})$ to match the sign of the scaling of the associated δ -function.

We can put $f(x_1, x_2) = x_1 x_2$

Note that this f is Lipschitz with Lipschitz constant $\sqrt{2}$.

In general: if $A = \text{supp } \mu_+$ and $B = \text{supp } \mu_-$ we have a Lipschitz function f where $f(\mathbf{x}) = +1$ for $\mathbf{x} \in \text{supp } \mu_+$ and $f(\mathbf{x}) = -1$ for $\mathbf{x} \in \text{supp } \mu_-$:

$$f(\mathbf{x}) = \frac{d(\mathbf{x}, B) - d(\mathbf{x}, A)}{d(\mathbf{x}, B) + d(\mathbf{x}, A)}$$

$$\text{Lip } f = \frac{2}{\min_{\mathbf{a} \in A, \mathbf{b} \in B} \|\mathbf{a} - \mathbf{b}\|}$$

Non-trivial lower bounds

What if there are many vectors in \mathcal{W} ? How many do we need to get a reasonable approximation?

Choose

$$\mathbf{z}_{k+1} \perp \{\mathbf{z}_1, \dots, \mathbf{z}_k, \mathbf{w}_{s(k)}, \dots, \mathbf{w}_{s(k+1)-1}\}$$

where $s(k+1) = s(k) + m - 1 - k$ for $k = 1, 2, \dots, m-2$. Put

$$\hat{\mu}(\boldsymbol{\xi}) = c \prod_{k=1}^{m-1} \left(\exp(-i\mathbf{z}_k^T \boldsymbol{\xi}) - \exp(+i\mathbf{z}_k^T \boldsymbol{\xi}) \right) \quad \text{so}$$

$$\mu = c \sum_{\mathbf{u} \in \{\pm 1\}^{m-1}} \left(\prod_{k=1}^{m-1} u_k \right) \delta_{\sum_{k=1}^{m-1} u_k \mathbf{z}_k}$$

Choose $\|\mathbf{z}_k\|_2 = 1/\sqrt{m-1}$ so that $\|\sum_k u_k \mathbf{z}_k\|_\infty \leq 1$ for all $\mathbf{u} \in \{\pm 1\}^{m-1}$.

$$\text{supp } \mu_+ = \left\{ \sum_k u_k \mathbf{z}_k \mid \mathbf{u} \in \{\pm 1\}^{m-1} \text{ \& \#}\{k \mid u_k > 0\} \text{ is even} \right\}$$

$$\text{supp } \mu_- = \left\{ \sum_k u_k \mathbf{z}_k \mid \mathbf{u} \in \{\pm 1\}^{m-1} \text{ \& \#}\{k \mid u_k > 0\} \text{ is odd} \right\}$$

Thus given \mathcal{W} with $|\mathcal{W}| \leq \frac{1}{2}m(m-1)$ there is a function f of Lipschitz constant $\sqrt{m-1}$ with $\|f\|_\infty = 1$ that is unapproximable by $\text{span} \{ \mathbf{x} \mapsto \varphi(\mathbf{w}^T \mathbf{x}) \mid \varphi \in \mathcal{C}(\mathbb{R}), \mathbf{w} \in \mathcal{W} \}$.