## Universal approximation of one-hidden layer feedforward NN
## A simple proof.

Sébastien GERCHINOVITZ

The results and proofs below are combinations of results and proofs from Cybenko (1989) and Barron (1993).
(See also Jones 1992 and Leshno et al. 1993.)

**Definition:** $\sigma : \mathbb{R} \to \mathbb{R}$ is "universal" iff, for all $a < b$, any $\mathcal{C}^\infty$ function $f : [a,b] \to \mathbb{R}$ can be arbitrarily well approximated by a 1-hidden layer NN with activation function $\sigma$, i.e.: $\forall \varepsilon > 0, \exists N \geq 1, \exists v_i, w_i, b_i \in \mathbb{R}$ s.t.

$$\forall x \in [a,b], \quad \left| f(x) - \sum_{i=1}^{N} v_i \, \sigma(w_i x + b_i) \right| \leq \varepsilon.$$

See Eldan and Shamir (2016) for a related definition.

**Ex:** $\sigma = \text{ReLU}$ or Heaviside (see later for more examples)

**Theorem:** Let $K \subset \mathbb{R}^d$ compact and $\sigma : \mathbb{R} \to \mathbb{R}$ universal. Then, any continuous function $f : K \to \mathbb{R}$ can be arbitrarily well approximated by a 1-hidden layer NN with activation function $\sigma$, i.e.: $\forall \varepsilon > 0, \exists N \geq 1, \exists v_i \in \mathbb{R}, w_i \in \mathbb{R}^d, b_i \in \mathbb{R}$ s.t.

$$\forall x \in K, \quad \left| f(x) - \sum_{i=1}^{N} v_i \, \sigma(\langle w_i, x \rangle + b_i) \right| \leq \varepsilon.$$

# Proof scheme

① Reduce the problem to $f = g_{|K}$, with $g \in \mathcal{C}_c^\infty(\mathbb{R}^d, \mathbb{R})$.

② Use the Fourier decomposition

$$g(x) = \int_{\mathbb{R}^d} e^{i \langle \omega, x \rangle} \hat{g}(\omega) \, d\omega \qquad (1)$$

$$= \int_{\mathbb{R}^d} e^{i \left( \langle \omega, x \rangle + \theta(\omega) \right)} |\hat{g}(\omega)| \, d\omega \qquad (2)$$

and approximate the first integral (uniformly in $x \in K$) by

$$\tilde{g}(x) = \sum_{k=1}^{K} e^{i \langle \omega_k, x \rangle} \hat{g}(\omega_k) \lambda(A_k)$$

where the $A_k$ are "small", pairwise disjoint, and s.t.

$$\int_{\mathbb{R}^d \setminus \bigcup_{k=1}^{K} A_k} |\hat{g}(\omega)| \, d\omega \le \frac{\varepsilon}{2} . \qquad \text{with } \operatorname{diam}\left( \bigcup_{k=1}^{K} A_k \right) < +\infty .$$

This is possible since $g$ and thus $\hat{g}$ belong to the Schwartz space.

Take the real part to see that we can approximate $g$ with

$$\tilde{g} \in \operatorname{span}\left\{ \cos(\langle \omega, \cdot \rangle + \varphi) : \omega \in \mathbb{R}^d, \varphi \in [0, 2\pi] \right\}$$

③ Approximate $t \mapsto \cos(t)$ by $t \mapsto \sum_{i=1}^{N} \sigma_i \, \tau(\nu_i t + \beta_i)$
uniformly on $\left[ -R_\omega R_x, R_\omega R_x + 2\pi \right]$, where $R_\omega := \sup\left\{ \|\omega\| : \omega \in \bigcup_{k=1}^{K} A_k \right\}$
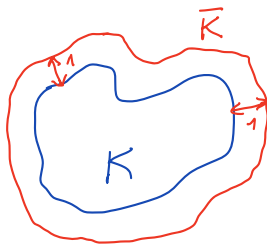
$$R_x := \sup\left\{ \|x\| : x \in K \right\}$$

Conclude.

**Details:** Let $\varepsilon > 0$.

(1) We show that there exists $g \in \mathcal{C}_c^\infty(\mathbb{R}^d, \mathbb{R})$ such that
$$\forall x \in K, \quad |f(x) - g(x)| \le \varepsilon$$

First, by the Tietze extension theorem, there exists $\bar{f} : \bar{K} \to \mathbb{R}$ continuous on $\bar{K} = \bigcup_{x \in K} \bar{B}(x, 1)$ and such that $\bar{f}(x) = f(x)$ for all $x \in K$.



Since $\bar{f}$ is uniformly continuous on the compact $\bar{K}$, let $\delta \in (0, 1)$ be s.t.
$$\forall x, x' \in \bar{K}, \quad \|x - x'\| \le \delta \implies |\bar{f}(x) - \bar{f}(x')| \le \varepsilon$$

Define $\varphi_\sigma(h) = \dfrac{1}{C_\sigma} \exp\left(-\dfrac{1}{1 - \left(\frac{\|h\|}{\sigma}\right)^2}\right)$, density function over $\mathbb{R}^d$

Note that $\varphi_\sigma \in \mathcal{C}^\infty(\mathbb{R}^d, \mathbb{R})$ and $\operatorname{supp}(\varphi_\sigma) = B(0, \sigma)$

For $\sigma \in (0, \delta)$ and $x \in \mathbb{R}^d$, we have
$$(\bar{f} * \varphi_\sigma)(x) = \int_{\mathbb{R}^d} \bar{f}(x-h)\, \varphi_\sigma(h)\, dh = \mathbb{E}_{Z \sim \varphi_\sigma}\left[\bar{f}(x - Z)\right]$$
with $\bar{f}(x') \overset{conv}{=} 0$ if $x' \notin \bar{K}$
(never happens if $x \in K$, because $\operatorname{supp}(\varphi_\sigma) = B(0,\sigma)$ and $K + B(0,\sigma) \subset \bar{K}$)

so that $\left|(\bar{f} * \varphi_\sigma)(x) - \bar{f}(x)\right| \le \varepsilon$ for all $x \in K$ $\Big\}$ This proves the result with $g = \bar{f} * \varphi_\sigma$

and $\bar{f} * \varphi_\sigma \in \mathcal{C}_c^\infty(\mathbb{R}^d; \mathbb{R})$

(2) We now approximate $g$ with a 1-hidden layer NN.

Since $g \in \mathcal{S}_c^\infty(\mathbb{R}^d, \mathbb{R})$ is "rapidly decreasing" ($g$ lies in the Schwartz space), the inversion formula for the Fourier transform holds: $\forall x \in K$,

$$g(x) = \int_{\mathbb{R}^d} e^{i\langle w, x\rangle} \, \hat{g}(w) \, dw$$

$$\overset{\pm \varepsilon}{\approx} \int_C e^{i\langle w, x\rangle} \hat{g}(w) \, dw \quad \text{for some hypercube } C \subset \mathbb{R}^d \text{ of sidelength } \rho.$$

$$\left( \text{we write } a \overset{\pm \varepsilon}{\approx} b \text{ to mean that } |a-b| \leq \varepsilon \right)$$

This is because $\int_{\mathbb{R}^d} |\hat{g}(w)| \, dw < +\infty$ ($\hat{g}$ is also rapidly decreasing)

$$\overset{\pm \varepsilon}{\approx} \sum_{k=1}^{m^d} e^{i\langle w_k, x\rangle} \hat{g}(w_k) \, \text{vol}(A_k) \quad \text{for } A_k : \text{subcubes of } C \text{ with sidelength } \frac{\rho}{m}.$$

We can choose $m \geq 1$ large enough and **independently of $x \in K$** such that the integral-sum $\varepsilon$-approximation above holds, because:

$$\exists \rho' > 0 : \| \, - w'\|_\infty \leq \rho' \Rightarrow \forall x \in K, \, \left| e^{i\langle w, x\rangle} \hat{g}(w) - e^{i\langle w', x\rangle} \hat{g}(w') \right| \leq \frac{\varepsilon}{\text{vol}(C)}$$

(indeed, $K$ is bounded and both $t \mapsto e^{it}$ and $\hat{g}$ are uniformly continuous)

Therefore, for all $x \in K$,

$$\left| g(x) - \underbrace{\text{Re} \sum_{k=1}^{m^d} e^{i\langle w_k, x\rangle} \hat{g}(w_k) \, \text{vol}(A_k)}_{\text{of the form } \sum_{k=1}^{m^d} \alpha_k \cos\left(\langle w_k, x\rangle + \theta_k\right) \text{ with } \alpha_k \in \mathbb{R} \text{ and } \theta_k \in [0, 2\pi]} \right| \leq \left| g(x) - \sum_{k=1}^{m^d} e^{i\langle w_k, x\rangle} \hat{g}(w_k) \, \text{vol}(A_k) \right| \leq 2\varepsilon$$

$$\overset{\pm \varepsilon}{\approx} \sum_{k=1}^{m^d} \alpha_k \left( \sum_{i=1}^N a_i \sigma \left( b_i [\langle w_k, x\rangle + \theta_k] + c_i \right) \right)$$

$$= \sum_{k=1}^{m^d} \sum_{i=1}^N \alpha_k a_i \sigma \left( \langle b_i w_k, x\rangle + b_i \theta_k + c_i \right). \qquad \blacksquare$$

# Examples of universal activation functions $\sigma$

(a) $\sigma(x) = \max\{x, 0\}$

→ approximate indicator function ⎍

+ convex combination of two such bumps at the boundary between two cells in a partition of $[a, b]$.

(b) $\sigma \in \mathscr{C}^\infty(\mathbb{R}, \mathbb{R})$ s.t. $\sigma$ is not a polynomial.

→ Leshno et al. (1993), Step 3

→ further extensions in the same paper.

# Dropping the assumption that $K$ is compact?

Under mild assumptions on $f: \mathbb{R}^d \to \mathbb{R}$, the approximation result holds in $L^\infty_{loc}$-norm, with the Lebesgue measure. loc ↪ local K

→ see Leshno et al. (1993), Theorem 1

→ csq, same paper, Proposition 1: $L^p(\mu)$, $1 \leq p < +\infty$, $\mu$ proba with compact support, $\mu \ll \lambda$.

In $L^2(\mu)$-norm, there is a quantitative approximation result via an argument due to Maurey ($\overline{cv(G)}$, $G \subset$ Hilbert space)

→ see Barron (1993), Lemma 1 in particular

→ the integral in (2) above is reinterpreted as a convex combination

NB: Barron assumes $\int \|w\| |\hat{f}(w)| \, dw < +\infty$, but $\int |\hat{f}(w)| \, dw < +\infty$

seems to work too (weaker assumption provided $f \in L^1(\mathbb{R}^d)$).