

A Organization of the appendices

This paper is a contribution to the mathematical foundations of machine learning, and our results are motivated by expanding the applicability and performance of neural networks. At the same time, we give precise mathematical formulations of our results and proofs. The purposes of these appendices are several:

1. To clarify the mathematical conventions and terminology, thus making the paper more accessible.
2. To provide full proofs of the main results.
3. To develop context around various construction appearing in the main text.
4. To discuss in detail examples, special cases, and generalizations of our results.

We now give a summary of the contents of the appendices.

Appendix B contains proofs the universal approximation results (Theorems 3 and 5) stated in Section 4 of the main text, as well as proofs of additional bounded width results. The proofs use notation given in Appendix B.1, and rely on preliminary topological considerations given in Appendix B.2.

In Appendix C, we give a proof of the model compression result given in Theorem 6, which appears in Section 5. For clarity and background we begin the appendix with a discussion of the version of the QR decomposition relevant for our purposes (Appendix C.1). We also establish elementary properties of radial rescaling activations (Appendix C.2).

The focus of Appendix D is projected gradient descent, elaborating on Section 6. We first prove a result on the interaction of gradient descent and orthogonal transformations (Appendix D.1), before formulating projected gradient descent in more detail (Appendix D.2), and introducing the so-called interpolating space (Appendix D.3). We restate Theorem 8 in more convenient notation (Appendix D.4) before proceeding to the proof (Appendix D.5).

Appendix E contains implementation details for the experiments summarized in Section 7. Our implementations use shifted radial rescaling activations, which we formulate in Appendix E.1.

Appendix F explains the connection between our constructions and radial basis functions networks. While radial neural networks turn out to be a specific type of radial basis functions network, our universality results are not implied by those for general radial basis functions networks.

B Universal approximation proofs and additional results

In this section, we provide full proofs of the universal approximation (UA) results for radial neural networks, as stated in Section 4. In order to do so, we first clarify our notational conventions (Appendix B.1), and collect basic topological results (Appendix B.2).

B.1 Notation

Recall that, for a point c in the Euclidean space \mathbb{R}^n and a positive real number r , we denote the r -ball around c by $B_r(c) = \{x \in \mathbb{R}^n \mid |x - c| < r\}$. All networks in this section have the Step-ReLU radial rescaling activation function, defined as:

$$\rho : \mathbb{R}^n \longrightarrow \mathbb{R}^n, \quad z \longmapsto \begin{cases} z & \text{if } |z| \geq 1 \\ 0 & \text{otherwise} \end{cases}$$

Throughout, \circ denotes the composition of functions. We identify a linear map with a corresponding matrix (in the standard bases). In the case of linear maps, the operation \circ

can be identified with matrix multiplication. Recall also that an affine map $L : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is one of the form $L(x) = Ax + b$ for a matrix $A \in \mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$.

B.2 Topology

Let K be a compact subset of \mathbb{R}^n and let $f : K \rightarrow \mathbb{R}^m$ be a continuous function.

Lemma 9. *For any $\epsilon > 0$, there exist $c_1, \dots, c_N \in K$ and $r_1, \dots, r_N \in (0, 1)$ such that, first, the union of the balls $B_{r_i}(c_i)$ covers K ; second, for all i , we have $f(B_{r_i}(c_i) \cap K) \subseteq B_\epsilon(f(c_i))$.*

Proof. The continuity of f implies that for each $c \in K$, there exists $r = r_c$ such that $f(B_{r_c}(c) \cap K) \subseteq B_\epsilon(f(c))$. The subsets $B_{r_c}(c) \cap K$ form an open cover of K . The compactness of K implies that there is a finite subcover. The result follows. \square

We also prove a variation of Lemma 9 that additionally guarantees that none of the balls in the cover of K contains the center point of another ball.

Lemma 10. *For any $\epsilon > 0$, there exist $c_1, \dots, c_M \in K$ and $r_1, \dots, r_M \in (0, 1)$ such that, first, the union of the balls $B_{r_i}(c_i)$ covers K ; second, for all i , we have $f(B_{r_i}(c_i)) \subseteq B_\epsilon(f(c_i))$; and, third, $|c_i - c_j| \geq r_i$.*

Proof. Because f is continuous on a compact domain, it is uniformly continuous. So, there exists $r > 0$ such that $f(B_r(c) \cap K) \subseteq B_\epsilon(f(c))$ for each $c \in K$. Because K is compact it has a finite volume, and so does $B_{r/2}(K) = \bigcup_{c \in K} B_{r/2}(c)$. Hence, there exists a finite maximal packing of $B_{r/2}(K)$ with balls of radius $r/2$. That is, a collection $c_1, \dots, c_M \in B_{r/2}(K)$ such that, for all i , $B_{r/2}(c_i) \subseteq B_{r/2}(K)$ and, for all $j \neq i$, $B_{r/2}(c_i) \cap B_{r/2}(c_j) = \emptyset$. The first condition implies that $c_i \in K$. The second condition implies that $|c_i - c_j| \geq r$. Finally, we argue that $K \subseteq \bigcup_{i=1}^M B_r(c_i)$. To see this, suppose, for a contradiction, that $x \in K$ does not belong to $\bigcup_{i=1}^M B_r(c_i)$. Then $B_{r/2}(c_i) \cap B_{r/2}(x) = \emptyset$, and x could be added to the packing, which contradicts the fact that the packing was chosen to be maximal. So the union of the balls $B_r(c_i)$ covers K . \square

We turn our attention to the minimal choices of N and M in Lemmas 9 and 10.

Definition 11. Given $f : K \rightarrow \mathbb{R}^m$ continuous and $\epsilon > 0$, let $N(f, K, \epsilon)$ be the minimal choice of N in Lemma 9, and let $M(f, K, \epsilon)$ be the minimal choice of M in Lemma 10.

Observe that $M(f, K, \epsilon) \geq N(f, K, \epsilon)$. In many cases, it is possible to give explicit bounds for the constants $N(f, K, \epsilon)$ and $M(f, K, \epsilon)$. As an illustration, we give the argument in the case that K is the closed unit cube in \mathbb{R}^n and $f : K \rightarrow \mathbb{R}^m$ is Lipschitz continuous.

Proposition 12. *Let $K = [0, 1]^n \subset \mathbb{R}^n$ be the (closed) unit cube and let $f : K \rightarrow \mathbb{R}^m$ be Lipschitz continuous with Lipschitz constant R . For any $\epsilon > 0$, we have:*

$$N(f, K, \epsilon) \leq \left\lceil \frac{R\sqrt{n}}{2\epsilon} \right\rceil^n \quad \text{and} \quad M(f, K, \epsilon) \leq \frac{\Gamma(n/2 + 1)}{\pi^{n/2}} \left(2 + \frac{2R}{\epsilon} \right)^n.$$

Proof. For the first inequality, observe that the unit cube can be covered with $\left\lceil \frac{R\sqrt{n}}{2\epsilon} \right\rceil^n$ cubes of side length $\frac{2\epsilon}{R\sqrt{n}}$. Each cube is contained in a ball of radius $\frac{\epsilon}{R}$ centered at the center of the cube. (In general, a cube of side length a in \mathbb{R}^n is contained in a ball of radius $\frac{a\sqrt{n}}{2}$.) Lipschitz continuity implies that, for all $x, x' \in K$, if $|x - x'| < \epsilon/R$ then $|f(x) - f(x')| \leq R|x - x'| < \epsilon$.

For the second inequality, let $r = \epsilon/R$. Lipschitz continuity implies that, for all $x, x' \in K$, if $|x - x'| < r$ then $|f(x) - f(x')| \leq R|x - x'| < \epsilon$. The n -dimensional volume of the set of points with distance at most $r/2$ to the unit cube is $\text{vol}(B_{r/2}(K)) \leq (1 + r)^n$. The volume

of a ball with radius $r/2$ is $\text{vol}(B_{r/2}(0)) = \frac{\pi^{n/2}}{\Gamma(n/2+1)}(r/2)^n$. Hence, any packing of $B_{r/2}(K)$ with balls of radius $r/2$ consists of at most

$$\frac{\text{vol}(B_{r/2}(K))}{\text{vol}(B_{r/2}(0))} \leq \frac{\Gamma(n/2+1)}{\pi^{n/2}} \left(2 + \frac{2R}{\epsilon}\right)^n$$

such balls. So there also exists a maximal packing with at most that many balls. This packing can be used in the proof of Lemma 10, which implies that it is a bound on $M(f, K, \epsilon)$. \square

We note in passing that any differentiable function $f : K \rightarrow \mathbb{R}^n$ on a compact subset K of \mathbb{R}^n is Lipschitz continuous. Indeed, the compactness of K implies that there exists R such that $|f'(x)| \leq R$ for all $x \in K$. Then one can take R to be the Lipschitz constant of f .

B.3 Proof of Theorem 3: UA for asymptotically affine functions

In this section, we restate and prove Theorem 3, which proves that radial neural networks are universal approximators of asymptotically affine functions. We recall the definition of such functions:

Definition 13. A function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is *asymptotically affine* if there exists an affine function $L : \mathbb{R}^n \rightarrow \mathbb{R}^m$ such that, for all $\epsilon > 0$, there exists a compact set $K \subset \mathbb{R}^n$ such that $|L(x) - f(x)| < \epsilon$ for all $x \in \mathbb{R}^n \setminus K$. We say that L is the limit of f .

Remark 14. An *asymptotically linear* function is defined in the same way, except L is taken to be linear (i.e., given just by applying matrix multiplication without translation). Hence any asymptotically linear function is in particular an asymptotically affine function, and Theorem 3 applies to asymptotically linear functions as well.

Given an asymptotically affine function $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and $\epsilon > 0$, let K be a compact set as in Definition 13. We apply Lemma 9 to the restriction $f|_K$ of f to K and produce a minimal constant $N = N(f|_K, K, \epsilon)$ as in Definition 11. We write simply $N(f, K, \epsilon)$ for this constant.

Theorem 3 (Universal approximation). *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be an asymptotically affine function. For any $\epsilon > 0$, there exists a compact set $K \subset \mathbb{R}^n$ and a function $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$ such that:*

1. *F is the feedforward function of a radial neural network with $N = N(f, K, \epsilon)$ layers whose hidden widths are $(n+1, n+2, \dots, n+N)$.*
2. *For any $x \in \mathbb{R}^n$, we have $|F(x) - f(x)| < \epsilon$.*

Proof. By the hypothesis on f , there exists an affine function $L : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and a compact set $K \subset \mathbb{R}^n$ such that $|L(x) - f(x)| < \epsilon$ for all $x \in \mathbb{R}^n \setminus K$. Abbreviate $N(f, K, \epsilon)$ by N . As in Lemma 9, fix $c_1, \dots, c_N \in K$ and $r_1, \dots, r_N \in (0, 1)$ such that, first, the union of the balls $B_{r_i}(c_i)$ covers K and, second, for all i , we have $f(B_{r_i}(c_i)) \subseteq B_\epsilon(f(c_i))$. Let $U = \bigcup_{i=1}^N B_{r_i}(c_i)$, so that $K \subset U$. Define $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$ as:

$$F(x) = \begin{cases} L(x) & \text{if } x \notin U \\ f(c_j) & \text{where } j \text{ is the smallest index with } x \in B_{r_j}(c_j) \end{cases}$$

If $x \notin U$, then $|F(x) - f(x)| = |L(x) - f(x)| < \epsilon$. Hence suppose $x \in U$. Let j be the smallest index such that $x \in B_{r_j}(c_j)$. Then $F(x) = f(c_j)$, and, by the choice of r_j , we have:

$$|F(x) - f(x)| = |f(c_j) - f(x)| < \epsilon.$$

We proceed to show that F is the feedforward function of a radial neural network. Let e_1, \dots, e_N be orthonormal basis vectors extending \mathbb{R}^n to \mathbb{R}^{n+N} . We regard each \mathbb{R}^{n+i-1} as a subspace of \mathbb{R}^{n+i} by embedding into the first $n+i-1$ coordinates. For $i = 1, \dots, N$, we set $h_i = \sqrt{1 - r_i^2}$ and define the following affine transformations:

$$\begin{aligned} T_i : \mathbb{R}^{n+i-1} &\rightarrow \mathbb{R}^{n+i} & S_i : \mathbb{R}^{n+i} &\rightarrow \mathbb{R}^{n+i} \\ z &\mapsto z - c_i + h_i e_i & z &\mapsto z - (1 + h_i^{-1}) \langle e_i, z \rangle e_i + c_i + e_i \end{aligned}$$

664 where $\langle e_i, z \rangle$ is the coefficient of e_i in z . Consider the radial neural network with widths
 665 $(n, n+1, \dots, n+N, m)$, whose affine transformations and activations are given by:

666 • For $i = 1, \dots, N$ the affine transformation from layer $i-1$ to layer i is given by
 667 $z \mapsto T_i \circ S_{i-1}(z)$, where $S_0 = \text{id}_{\mathbb{R}^n}$.

668 • The activation function at the i -th hidden layer is Step-ReLU on \mathbb{R}^{n+i} , that is:

$$\rho_i : \mathbb{R}^{n+i} \longrightarrow \mathbb{R}^{n+i}, \quad z \longmapsto \begin{cases} z & \text{if } |z| \geq 1 \\ 0 & \text{otherwise} \end{cases}$$

669 • The affine transformation from layer $i = N$ to the output layer is

$$z \mapsto \Phi_{L,f,\mathbf{c}} \circ S_N(z)$$

670 where $\Phi_{L,f,\mathbf{c}}$ is the affine transformation given by:

$$\Phi_{L,f,\mathbf{c}} : \mathbb{R}^{n+N} \rightarrow \mathbb{R}^m, \quad x + \sum_{i=1}^N a_i e_i \mapsto L(x) + \sum_{i=1}^N a_i (f(c_i) - L(c_i))$$

671 which can be shown to be affine when L is affine. Indeed, write $L(x) = Ax + b$
 672 where A is a matrix in $\mathbb{R}^{m \times n}$ and $b \in \mathbb{R}^m$ is a vector. Then $\Phi_{L,f,\mathbf{c}}$ is the composition
 673 of the linear map given by the matrix

$$\begin{bmatrix} A & f(c_1) - L(c_1) & f(c_2) - L(c_2) & \cdots & f(c_N) - L(c_N) \end{bmatrix} \in \mathbb{R}^{m \times (n+N)}$$

674 and translation by $b \in \mathbb{R}^m$. Note that we regard each $f(c_i) - L(c_i) \in \mathbb{R}^m$ as a
 675 column vector in the matrix above.

676 We claim that the feedforward function of the above radial neural network is exactly F . To
 677 show this, we first state a lemma, whose (omitted) proof is an elementary computation.

678 **Lemma 3.1.** For $i = 1, \dots, N$, the composition $S_i \circ T_i$ is the embedding $\mathbb{R}^{n+i-1} \hookrightarrow \mathbb{R}^{n+i}$.

679 Next, recursively define $G_i : \mathbb{R}^n \rightarrow \mathbb{R}^{n+i}$ via

$$G_i = S_i \circ \rho_i \circ T_i \circ G_{i-1},$$

680 where $G_0 = \text{id}_{\mathbb{R}^n}$. The function G_i admits an direct formulation:

681 **Proposition 3.2.** For $i = 0, 1, \dots, N$, we have:

$$G_i(x) = \begin{cases} x & \text{if } x \notin \bigcup_{j=1}^i B_{r_j}(c_j) \\ c_j + e_j & \text{where } j \leq i \text{ is the smallest index with } x \in B_{r_j}(c_j) \end{cases}.$$

682 *Proof.* We proceed by induction. The base step $i = 0$ is immediate. For the induction step,
 683 assume the claim is true for $i-1$, where $0 \leq i-1 < N$. There are three cases to consider.

684 **Case 1.** Suppose $x \notin \bigcup_{j=1}^i B_{r_j}(c_j)$. Then in particular $x \notin \bigcup_{j=1}^{i-1} B_{r_j}(c_j)$, so the induction
 685 hypothesis implies that $G_{i-1}(x) = x$. Additionally, $x \notin B_{r_i}(c_i)$, so:

$$|T_i(x)| = |x - c_i + h_i e_i| = \sqrt{|x - c_i|^2 + h_i^2} \geq \sqrt{r_i^2 + 1 - r_i^2} = 1.$$

686 Using the definition of ρ_i and Lemma 3.1, we compute:

$$G_i(x) = S_i \circ \rho_i \circ T_i \circ G_{i-1}(x) = S_i \circ \rho_i \circ T_i(x) = S_i \circ T_i(x) = x.$$

687 **Case 2.** Suppose $x \in B_j \setminus \bigcup_{k=1}^{j-1} B_{r_k}(c_k)$ for some $j \leq i-1$. Then the induction hypothesis
 688 implies that $G_{i-1}(x) = c_j + e_j$. We compute:

$$|T_i(c_j + e_j)| = |c_j + e_j - c_i + h_i e_i| > |e_j| = 1.$$

Therefore,

$$G_i(x) = S_i \circ \rho_i \circ T_i(c_j + e_j) = S_i \circ T_i(c_j + e_j) = c_j + e_j.$$

Case 3. Finally, suppose $x \in B_i \setminus \bigcup_{j=1}^{i-1} B_{r_j}(c_j)$. The induction hypothesis implies that $G_{i-1}(x) = x$. Since $x \in B_{r_i}(c_i)$, we have:

$$|T_i(x)| = |x - c_i + h_i e_i| = \sqrt{|x - c_i|^2 + h_i^2} < \sqrt{r_i^2 + 1 - r_i^2} = 1.$$

Therefore:

$$G_i(x) = S_i \circ \rho_i \circ T_i(x) = S_i(0) = c_i + e_i.$$

This completes the proof of the proposition. \square

Finally, we show that the function F defined at the beginning of the proof is the feedforward function of the above radial neural network. The computation is elementary:

$$\begin{aligned} F_{\text{feedforward}} &= \Phi_{L,f,c} \circ S_N \circ \rho_N \circ T_N \circ S_{N-1} \circ \rho_{N-1} \circ T_{N-1} \circ \dots \circ S_1 \circ \rho_1 \circ T_1 \\ &= \Phi_{L,f,c} \circ G_N \\ &= F \end{aligned}$$

where the first equality follows from the definition of the feedforward function, the second from the definition of G_N , and the last from the case $i = N$ of Proposition 3.2 together with the definition of $\Phi_{L,f,c}$. This completes the proof of the theorem. \square

B.4 Proof of Theorem 5: bounded width UA for asymptotically affine functions

We restate and prove Theorem 5, which strengthens Theorem 3 by providing a bounded width radial neural network approximation of any asymptotically affine function.

Theorem 5. *Let $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ be an asymptotically affine function. For any $\epsilon > 0$, there exists a compact set $K \subset \mathbb{R}^n$ and a function $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$ such that:*

1. *F is the feedforward function of a radial neural network with $N = N(f, K, \epsilon)$ hidden layers whose widths are all $n + m + 1$.*
2. *For any $x \in \mathbb{R}^n$, we have $|F(x) - f(x)| < \epsilon$.*

Proof. By the hypothesis on f , there exists an affine function $L : \mathbb{R}^n \rightarrow \mathbb{R}^m$ and a compact set $K \subset \mathbb{R}^n$ such that $|L(x) - f(x)| < \epsilon$ for all $x \in \mathbb{R}^n \setminus K$. Given $\epsilon > 0$, let $N = N(f, K, \epsilon)$ and use Lemma 9 to choose $c_1, \dots, c_N \in K$ and $r_1, \dots, r_N \in (0, 1)$ such that the union of the balls $B_{r_i}(c_i)$ covers K , and, for all i , we have $f(B_{r_i}(c_i)) \subseteq B_\epsilon(f(c_i))$. Let s be the minimal non-zero value of $|f(c_i) - f(c_j)|$ for $i, j \in \{1, \dots, N\}$, that is, $s = \min_{i,j, f(c_i) \neq f(c_j)} |f(c_i) - f(c_j)|$.

Using the decomposition $\mathbb{R}^{n+m+1} \cong \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}$, we write elements of \mathbb{R}^{n+m+1} as (x, y, θ) , where $x \in \mathbb{R}^n$, $y \in \mathbb{R}^m$, and $\theta \in \mathbb{R}$. For $i = 1, \dots, N$, set:

$$T_i : \mathbb{R}^{n+m+1} \rightarrow \mathbb{R}^{n+m+1}, \quad (x, y, \theta) \mapsto \left(x - (1 - \theta)c_i, y - \theta \frac{f(c_i) - L(0)}{s}, (1 - \theta)h_i \right)$$

where $h_i = \sqrt{1 - r_i^2}$. Note that T_i is an invertible affine transformation, whose inverse is given by:

$$T_i^{-1}(x, y, \theta) = \left(x + \frac{\theta}{h_i}c_i, y + \left(1 - \frac{\theta}{h_i}\right) \frac{f(c_i) - L(0)}{s}, 1 - \frac{\theta}{h_i} \right)$$

For $i = 1, \dots, N$, define $G_i : \mathbb{R}^n \rightarrow \mathbb{R}^{n+m+1}$ via the following recursive definition:

$$G_i = T_i^{-1} \circ \rho \circ T_i \circ G_{i-1},$$

717 where $G_0(x) = (x, 0, 0) : \mathbb{R}^n \hookrightarrow \mathbb{R}^{n+m+1}$ is the inclusion, and $\rho : \mathbb{R}^{n+m+1} \rightarrow \mathbb{R}^{n+m+1}$ is
 718 Step-ReLU on \mathbb{R}^{n+m+1} . We claim that, for $x \in \mathbb{R}^n$, we have:

$$G_i(x) = \begin{cases} (x, 0, 0) & \text{if } x \notin \bigcup_{j=1}^i B_{r_j}(c_j) \\ \left(0, \frac{f(c_j) - L(0)}{s}, 1\right) & \text{where } j \leq i \text{ is the smallest index with } x \in B_{r_j}(c_j) \end{cases}$$

719 This claim can be verified by a straightforward induction argument, similar to the one
 720 given in the proof of Proposition 3.2, and using the following key facts:

- 721 • For $x \in \mathbb{R}^n$, $|T_i((x, 0, 0))| = |(x - c_i, 0, h_i)| < 1$ if and only if $|x - c_i| < r_i$.
- 722 • $T_i^{-1}(0) = \left(0, \frac{f(c_i) - L(0)}{s}, 1\right)$.
- 723 • $T_i\left(\left(0, \frac{f(c_j) - L(0)}{s}, 1\right)\right) = \left(0, \frac{f(c_j) - f(c_i)}{s}, 0\right)$, which, by the choice of s , has norm at
 724 least 1 if $f(c_j) \neq f(c_i)$, and is 0 if $f(c_j) = f(c_i)$.

725 Let $\Phi : \mathbb{R}^{n+m+1} \rightarrow \mathbb{R}^m$ denote the affine map sending (x, y, θ) to $L(x) + sy$. It follows that
 726 $F = \Phi \circ G_N$ satisfies

$$F(x) = \begin{cases} L(x) & \text{if } x \notin \bigcup_{j=1}^N B_{r_j}(c_j) \\ f(c_j) & \text{where } j \text{ is the smallest index with } x \in B_{r_j}(c_j) \end{cases}$$

727 By construction, F is the feedforward function of a radial neural network with N hidden
 728 layers whose widths are all $n + m + 1$. Let $x \in \mathbb{R}^n$. If $x \in K$, let j be the smallest index
 729 such that $x \in B_{r_j}(c_j)$. Then $F(x) = f(c_j)$, and, by the choice of r_j , we have $|F(x) - f(x)| =$
 730 $|f(c_j) - f(x)| < \epsilon$. Otherwise, $x \in \mathbb{R}^n \setminus K$, and $|F(x) - f(x)| = |L(x) - f(x)| < \epsilon$. \square

731 B.5 Additional result: bound of $\max(n, m) + 1$

732 We state and prove an additional bounded width result. In contrast to the results above, the
 733 theorem below only holds for functions defined on a compact domain, without assumptions
 734 about the asymptotic behavior. The proof is an adaptation of the proof of Theorem 5, so
 735 we give only a sketch.

736 **Theorem 15.** *Let $f : K \rightarrow \mathbb{R}^m$ be a continuous function, where K is a compact subset of \mathbb{R}^n . For
 737 any $\epsilon > 0$, there exists $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$ such that:*

- 738 1. *F is the feedforward function of a radial neural network with $N(f, K, \epsilon)$ hidden layers
 739 whose widths are all $\max(n, m) + 1$.*
- 740 2. *For any $x \in K$, we have $|F(x) - f(x)| < \epsilon$.*

741 *Sketch of proof.* The construction appearing in the proof of Theorem 5 with $L \equiv 0$ can
 742 be used to produce a radial neural network with $N(f, K, \epsilon)$ hidden layers with widths
 743 $n + m + 1$ that approximates f on K . (Note that the approximation works only on K , as f is
 744 not defined outside of K .) All values in the hidden layers are of the form $(x, 0, 0)$ or $(0, y, 1)$.
 745 We can therefore replace $(x, y, \theta) \in \mathbb{R}^{n+m+1}$ by $(x + y, \theta) \in \mathbb{R}^{\max(n, m)} \times \mathbb{R} \cong \mathbb{R}^{\max(n, m)+1}$
 746 everywhere, without affecting any statements about the hidden layers. In particular, the
 747 transformation T_i becomes

$$T_i : \mathbb{R}^{\max(n, m)+1} \rightarrow \mathbb{R}^{\max(n, m)+1}, \quad (x, \theta) \mapsto \left(x - (1 - \theta)c_i - \theta \frac{f(c_i)}{s}, (1 - \theta)h_i\right).$$

748 With this change the final affine map Φ sends (x, θ) to sx . From the rest of the proof
 749 of Theorem 5 it follows that the feedforward function F of the radial network satisfies
 750 $|F(x) - f(x)| < \epsilon$ for all $x \in K$. \square

751 **B.6 Additional result: bound of $\max(n, m)$**

752 In this section, we prove a different version of the result of the previous section. Specifically,
 753 we reduce the bound on the widths to $\max(n, m)$ at the cost of using more layers. Again,
 754 we focus on functions defined on a compact domain without assumptions about their
 755 asymptotic behavior. Recall the notation $M(f, K, \epsilon)$ from Lemma 10 and Definition 11.

756 **Theorem 16.** *Let $f : K \rightarrow \mathbb{R}^m$ be a continuous function, where K is a compact subset of \mathbb{R}^n for
 757 $n \geq 2$. For any $\epsilon > 0$, there exists $F : \mathbb{R}^n \rightarrow \mathbb{R}^m$ such that:*

- 758 1. *F is the feedforward function of a radial neural network with $2M(f, K, \epsilon/2)$ hidden layers
 759 whose widths are all $\max(n, m)$.*
 760 2. *For any $x \in K$, we have $|F(x) - f(x)| < \epsilon$.*

761 *Proof.* We first consider the proof in the case $n = m$. Set $M = M(f, K, \epsilon)$. As in Lemma 10,
 762 fix $c_1, \dots, c_M \in K$ and $r_1, \dots, r_M \in (0, 1)$ such that, first, the union of the balls $B_{r_i}(c_i)$ covers
 763 K ; second, for all i , we have $f(B_{r_i}(c_i)) \subseteq B_{\epsilon/2}(f(c_i))$; and third, $|c_i - c_j| \geq r_i$ for $i \neq j$. For
 764 $i = 1, \dots, M$, set

$$T_i : \mathbb{R}^n \rightarrow \mathbb{R}^n, \quad x \mapsto \frac{x - c_i}{r_i},$$

765 and recursively define $G_i : \mathbb{R}^n \rightarrow \mathbb{R}^n$ as $G_i = T_i^{-1} \circ \rho \circ T_i \circ G_{i-1}$, where $G_0 = \text{id}_{\mathbb{R}^n}$ is the
 766 identity on \mathbb{R}^n and $\rho : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is Step-ReLU.

767 **Lemma 16.1.** For $i = 0, 1, \dots, N$, we have:

$$G_i(x) = \begin{cases} x & \text{if } x \notin \bigcup_{j=1}^i B_{r_j}(c_j) \\ c_j & \text{where } j \leq i \text{ is the smallest index with } x \in B_{r_j}(c_j). \end{cases}$$

768 We omit the full proof of Lemma 16.1, as it is a standard induction argument similar
 769 to Proposition 3.2, relying on the following two facts. First, $|T_i(x)| < 1$ if and only if
 770 $x \in B_{r_i}(c_i)$. Second, by the choice of c_i , we have $|c_i - c_j| \geq r_i$ for all $i \neq j$. This implies that
 771 $|T_i(c_j)| \geq 1$ for $i \neq j$.

772 Next, perform the following loop over $i = 1, \dots, M$:

- 773 • Set $P_{i-1} = \{c_1, \dots, c_M\} \cup \{d_1, \dots, d_{i-1}\}$
 774 • Choose d_i in $B_{\epsilon/2}(f(c_i))$ that is not colinear with any pair of points in P_{i-1} . This is
 775 where we use the hypothesis that $n \geq 2$.
 776 • Let s_i be the minimum distance between any point on the line through c_i and d_i
 777 and any point in $P_{i-1} \setminus \{c_i\}$.
 778 • Let $U_i : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be the following affine transformation:

$$U_i : \mathbb{R}^n \rightarrow \mathbb{R}^n, \quad x \mapsto \frac{x - d_i}{s_i} + \left(\frac{1}{|c_i - d_i|} - \frac{1}{s_i} \right) \frac{\langle x - d_i, c_i - d_i \rangle}{|c_i - d_i|^2} (c_i - d_i)$$

- 779 • Define $H_i : \mathbb{R}^n \rightarrow \mathbb{R}^n$ recursively as $H_i = U_i^{-1} \circ \rho \circ U_i \circ H_{i-1}$, where $H_0 = \text{id}_{\mathbb{R}^n}$.

780 We note that the transformation U_i can also be written as $A_i(x - d_i)$ where A_i is the linear
 781 map given by $A_i = \frac{1}{s_i} \text{proj}_{\langle c_i - d_i \rangle^\perp} + \frac{1}{|c_i - d_i|} \text{proj}_{\langle c_i - d_i \rangle}$, which involves the projections onto
 782 the line spanned by $c_i - d_i$ and onto the orthogonal complement of this line.

783 **Lemma 16.2.** For $i, j = 1, \dots, M$, we have:

$$H_i(c_j) = \begin{cases} d_j & \text{if } j \leq i \\ c_j & \text{if } j > i \end{cases}$$

784 *Proof.* It is immediate that $U_i(d_i) = 0$ and $|U_i(c_i)| = 1/2$. It is also straightforward to show,
 785 using the choice of s_i , that $|U_i(p)| \geq 1$ for all $p \in P_{i-1} \setminus \{c_i\}$. It follows that $U_i^{-1} \circ \rho \circ U_i$
 786 sends c_i to d_i and fixes all other points in P_{i-1} . \square

787 **Lemma 16.3.** For $x \in K$, we have $H_M \circ G_M(x) = d_i$ where i is the smallest index with
 788 $x \in B_{r_i}(c_i)$

789 *Proof.* Let $x \in K$. By Lemma 16.1, we have that $G_M(x) = c_i$ where i is the smallest index
 790 with $x \in B_{r_i}(c_i)$. (We use the fact that the balls $\{B_{r_i}(c_i)\}$ cover K .) By Lemma 16.2, we have
 791 that $H_M(c_i) = d_i$ for all i . The result follows. \square

792 Set $F = H_M \circ G_M$. We see that, for $x \in K$:

$$|F(x) - f(x)| = |d_i - f(x)| \leq |d_i - f(c_i)| + |f(c_i) - f(x)| < \epsilon/2 + \epsilon/2 = \epsilon$$

793 where i is the smallest index with $x \in B_{r_i}(c_i)$. We show that F is the feedforward function of
 794 a radial neural network with $2M$ hidden layers, all of width equal to n . Indeed, take the
 795 affine transformations and activations as follows:

- 796 • For $i = 1, \dots, M$ the affine transformation from layer $i - 1$ to layer i is given by
 797 $x \mapsto T_i \circ T_{i-1}^{-1}(x)$, where $T_0 = \text{id}_{\mathbb{R}^n}$.
- 798 • For $i = 1, \dots, M$ the affine transformation from layer $M + i - 1$ to layer $M + i$ is
 799 given by $x \mapsto U_i \circ U_{i-1}^{-1}(x)$, where $U_0 = T_N^{-1}$.
- 800 • The activation at each hidden layer is Step-ReLU on \mathbb{R}^n that is $\rho(x) = x$ if $|x| \geq 1$
 801 and 0 otherwise.
- 802 • Layer $2M + 1$ has the affine transformation U_M^{-1} .

803 It is immediate from definitions that the feedforward function of this network is F .

804 To conclude the proof, we discuss the cases where $n \neq m$. Suppose $n < m$ so that
 805 $\max(n, m) = m$. Then we can regard K as a compact subset of \mathbb{R}^m and apply the above
 806 constructions. Suppose $n > m$ so that $\max(n, m) = n$. Let $\text{inc} : \mathbb{R}^m \hookrightarrow \mathbb{R}^n$. Apply the
 807 above constructions to the function $\tilde{f} = \text{inc} \circ f : K \rightarrow \mathbb{R}^n$. \square

808 C Model compression proofs

809 The aim of this appendix is to give a proof of Theorem 6. In order to do so, we first (1)
 810 provide background on a relevant version of the QR decomposition, and (2) establish basic
 811 properties of radial rescaling activations.

812 C.1 The QR decomposition

813 In this section, we recall the QR decomposition and note several relevant facts. For integers
 814 n and m , let $(\mathbb{R}^{n \times m})^{\text{upper}}$ denote the vector space of upper triangular n by m matrices.

815 **Theorem 17** (QR Decomposition). *The following map is surjective:*

$$\begin{aligned} O(n) \times (\mathbb{R}^{n \times m})^{\text{upper}} &\longrightarrow \mathbb{R}^{n \times m} \\ Q, R &\mapsto Q \circ R \end{aligned}$$

816 In other words, any matrix can be written as the product of an orthogonal matrix and an
 817 upper-triangular matrix. When $m \leq n$, the last $n - m$ rows of any matrix in $(\mathbb{R}^{n \times m})^{\text{upper}}$
 818 are zero, and the top m rows form an upper-triangular m by m matrix. These observations
 819 lead to the following “complete” version of the QR decomposition, which coincides with
 820 the above result when $m \geq n$:

821 **Corollary 18** (Complete QR Decomposition). *The following map is surjective:*

$$\begin{aligned} \mu : O(n) \times \left(\mathbb{R}^{k \times m} \right)^{\text{upper}} &\longrightarrow \mathbb{R}^{n \times m} \\ Q, R &\mapsto Q \circ \text{inc} \circ R \end{aligned}$$

822 where $k = \min(n, m)$ and $\text{inc} : \mathbb{R}^k \hookrightarrow \mathbb{R}^n$ is the standard inclusion into the first k coordinates.

823 We make some remarks:

- 824 1. There are several algorithms for computing the QR decomposition of a given
825 matrix. One is Gram–Schmidt orthogonalization, and another is the method of
826 Householder reflections. The latter has computational complexity $O(n^2m)$ in
827 the case of a $n \times m$ matrix with $n \geq m$. The package `numpy` includes a func-
828 tion `numpy.linalg.qr` that computes the QR decomposition of a matrix using
829 Householder reflections.
- 830 2. In each iteration of the loop in Algorithm 1, the method `QR-decomp` with mode
831 = ‘complete’ takes as input a matrix A_i of size $n_i \times (n_{i-1}^{\text{red}} + 1)$, and pro-
832 duces an orthogonal matrix $Q_i \in O(n_i)$ and an upper-triangular matrix R_i
833 of size $\min(n_i, n_{i-1}^{\text{red}} + 1) \times (n_{i-1}^{\text{red}} + 1)$ such that $A_i = Q_i \circ \text{inc}_i \circ R_i$. Note that
834 $n_i^{\text{red}} = \min(n_i, n_{i-1}^{\text{red}} + 1)$.
- 835 3. The QR decomposition is not unique in general, or, in other words, the map μ is
836 not injective in general. For example, if $n > m$, each fiber of μ contains a copy of
837 the orthogonal group $O(n - m)$.
- 838 4. The QR decomposition is unique (in a certain sense) for invertible square matrices.
839 To be precise, let B_n^+ be the subset of $(\mathbb{R}^{n \times n})^{\text{upper}}$ consisting of upper triangular
840 n by n matrices with positive entries along the diagonal. Both B_n^+ and $O(n)$
841 are subgroups of the general linear group $\text{GL}_n(\mathbb{R})$, and the multiplication map
842 $O(n) \times B_n^+ \rightarrow \text{GL}_n(\mathbb{R})$ is bijective. However, the QR decomposition is not unique
843 for non-invertible square matrices.

844 C.2 Radial rescaling functions

845 We now prove the following basic facts about radial rescaling functions:

846 **Lemma 19.** *Let $\rho = h^{(n)} : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be a radial rescaling function on \mathbb{R}^n .*

- 847 1. *The function ρ commutes with any orthogonal transformation of \mathbb{R}^n . That is, $\rho \circ Q = Q \circ \rho$
848 for any $Q \in O(n)$.*
- 849 2. *If $m \leq n$ and $\text{inc} : \mathbb{R}^m \hookrightarrow \mathbb{R}^n$ is the standard inclusion into the first m coordinates, then:
850 $h^{(n)} \circ \text{inc} = \text{inc} \circ h^{(m)}$.*

851 *Proof.* Suppose $Q \in O(n)$ is an orthogonal transformation of \mathbb{R}^n . Since Q is norm-
852 preserving, we have $|Qv| = |v|$ for any $v \in \mathbb{R}^n$. Since Q is linear, we have $Q(\lambda v) = \lambda Qv$
853 for any $\lambda \in \mathbb{R}$ and $v \in \mathbb{R}^n$. Using the definition of $a = h^{(n)}$ we compute:

$$\rho(Qv) = \frac{h(|Qv|)}{|Qv|} Qv = \frac{h(|v|)}{|v|} Qv = Q \left(\frac{h(|v|)}{|v|} v \right) = Q(\rho(v)).$$

854 The first claim follows. The second claim is an elementary verification. □

855 More generally, the restriction of the radial rescaling function ρ to a linear subspace of \mathbb{R}^n
856 is a radial rescaling function on that subspace. Given a tuple radial rescaling functions $\rho =$
857 $(\rho_i : \mathbb{R}^{n_i} \rightarrow \mathbb{R}^{n_i})_{i=1}^L$ suited to widths $\mathbf{n} = (n_i)_{i=1}^L$, we write $\rho^{\text{red}} = (\rho_i^{\text{red}} : \mathbb{R}^{n_i^{\text{red}}} \rightarrow \mathbb{R}^{n_i^{\text{red}}})$

858 for the tuple of restrictions suited to the reduced widths \mathbf{n}^{red} , so that $\rho_i^{\text{red}} = \rho_i \Big|_{\mathbb{R}^{n_i^{\text{red}}}}$.

859 C.3 Proof of Theorem 6

860 Adopting notation from above and Section 5, we now restate and prove Theorem 6.

861 **Theorem 6.** *Let $(\mathbf{W}, \mathbf{b}, \rho)$ be a radial neural network with widths \mathbf{n} . Let \mathbf{W}^{red} and \mathbf{b}^{red} be the*
 862 *weights and biases of the compressed network produced by Algorithm 1. The feedforward function*
 863 *of the original network $(\mathbf{W}, \mathbf{b}, \rho)$ coincides with that of the compressed network $(\mathbf{W}^{\text{red}}, \mathbf{b}^{\text{red}}, \rho^{\text{red}})$.*

864 *Proof.* Let $(\mathbf{W}^{\text{red}}, \mathbf{b}^{\text{red}}, \mathbf{Q}) = \text{QR-Compress}(\mathbf{W}, \mathbf{b})$ be the output of Algorithm 1, so that
 865 $\mathbf{Q} \in O(\mathbf{n}^{\text{hid}})$ and $(\mathbf{W}^{\text{red}}, \mathbf{b}^{\text{red}}, \rho^{\text{red}})$ is a neural network with widths n^{red} and radial
 866 rescaling activations $\rho_i^{\text{red}} = \rho_i|_{\mathbb{R}^{n_i^{\text{red}}}}$. Let $F = F_{(\mathbf{W}, \mathbf{b}, \rho)}$ denote the feedforward function
 867 of the radial neural network with parameters (\mathbf{W}, \mathbf{b}) and activations ρ . Similarly, let
 868 $F^{\text{red}} = F_{(\mathbf{W}^{\text{red}}, \mathbf{b}^{\text{red}}, \rho^{\text{red}})}$ denote the feedforward function of the radial neural network with
 869 parameters $(\mathbf{W}^{\text{red}}, \mathbf{b}^{\text{red}})$ and activations ρ^{red} . Additionally, we have the partial feedforward
 870 functions F_i and F_i^{red} . We show by induction that

$$F_i = Q_i \circ \text{inc}_i \circ F_i^{\text{red}}$$

871 for any $i = 0, 1, \dots, N$. (Continuing conventions from Sections 5.1 and 5.2, we set $Q_0 =$
 872 $\text{id}_{\mathbb{R}^{n_0}}$, $Q_L = \text{id}_{\mathbb{R}^{n_L}}$, and $\text{inc}_i : \mathbb{R}^{n_i^{\text{red}}} \rightarrow \mathbb{R}^{n_i}$ to be the inclusion map.) The base step $i = 0$
 873 immediate. For the induction step, let $x \in \mathbb{R}^{n_0}$. Then:

$$\begin{aligned} F_i(x) &= \rho_i(W_i \circ F_{i-1}(x) + b_i) \\ &= \rho_i\left(W_i \circ Q_{i-1} \circ \text{inc}_{i-1} \circ F_{i-1}^{\text{red}}(x) + b_i\right) \\ &= \rho_i\left(\begin{bmatrix} b_i & W_i \circ Q_{i-1} \circ \text{inc}_{i-1} \end{bmatrix} \begin{bmatrix} 1 \\ F_{i-1}^{\text{red}}(x) \end{bmatrix}\right) \\ &= \rho_i\left(Q_i \circ \text{inc}_i \circ \begin{bmatrix} b_i^{\text{red}} & W_i^{\text{red}} \end{bmatrix} \begin{bmatrix} 1 \\ F_{i-1}^{\text{red}}(x) \end{bmatrix}\right) \\ &= Q_i \circ \text{inc}_i \circ \rho_i|_{\mathbb{R}^{n_i^{\text{red}}}}\left(W_i^{\text{red}} \circ F_{i-1}^{\text{red}}(x) + b_i^{\text{red}}\right) \\ &= Q_i \circ \text{inc}_i \circ F_i^{\text{red}} \end{aligned}$$

874 The first equality relies on the definition of the partial feedforward function F_i ; the second
 875 on the induction hypothesis; the fourth on an inspection of Algorithm 1, noting that
 876 $R_i = \begin{bmatrix} b_i^{\text{red}} & W_i^{\text{red}} \end{bmatrix}$; the fifth on the results of Lemma 19, observing that $\rho_i \circ \text{inc}_i = \rho_i|_{\mathbb{R}^{n_i^{\text{red}}}} =$
 877 $\text{inc}_i \circ \rho_i^{\text{red}}$; and the sixth on the definition of F_i^{red} . In the case $i = L$, we have:

$$F = F_L = Q_L \circ \text{inc}_L \circ F_L^{\text{red}} = F^{\text{red}}$$

878 since $Q_L = \text{inc}_L = \text{id}_{\mathbb{R}^{n_L}}$ and $F_L^{\text{red}} = F^{\text{red}}$. The theorem now follows. \square

879 The techniques of the above proof can be used to show that the action of the group $O(\mathbf{n}^{\text{hid}})$
 880 of orthogonal change-of-basis symmetries on the parameter space $\text{Param}(\mathbf{n})$ leaves the
 881 feedforward function unchanged. We do not use this result directly, but state it precisely it
 882 nonetheless:

883 **Proposition 20.** *Let $(\mathbf{W}, \mathbf{b}, \rho)$ be a radial neural network with widths vector \mathbf{n} . Suppose $\mathbf{g} \in$
 884 $O(\mathbf{n}^{\text{hid}})$. Then the original and transformed networks have the same feedforward function:*

$$F_{(\mathbf{g}\mathbf{W}, \mathbf{g}\mathbf{b}, \rho)} = F_{(\mathbf{W}, \mathbf{b}, \rho)}$$

885 In other words, fix parameters $(\mathbf{W}, \mathbf{b}) \in \text{Param}(\mathbf{n})$, radial rescaling activations ρ , and $\mathbf{g} \in$
 886 $O(\mathbf{n}^{\text{hid}})$. Then the radial neural network with parameters (\mathbf{W}, \mathbf{b}) has the same feedforward

function as the radial neural network with transformed parameters $(\mathbf{g} \cdot \mathbf{W}, \mathbf{g} \cdot \mathbf{b})$, where we take radial rescaling activations ρ in both cases.

We remark that Proposition 20 is analogous to the “non-negative homogeneity” (or “positive scaling invariance”) of the pointwise ReLU activation function³. In that setting, instead of considering the product of orthogonal groups $O(\mathbf{n}^{\text{hid}})$, one considers the rescaling action of the following subgroup of $\prod_{i=1}^{L-1} \text{GL}_{n_i}$:

$$G = \left\{ \mathbf{g} = (g_i) \in \prod_{i=1}^{L-1} \text{GL}_{n_i} \mid \text{each } g_i \text{ is diagonal with positive diagonal entries} \right\}$$

Note that G is isomorphic to the product $\prod_{i=1}^{L-1} \mathbb{R}_{>0}^{n_i}$, and the action on $\text{Param}(\mathbf{n})$ is given by the same formulas as those appearing near the end of Section 5.1. The feedforward function of a MLP with pointwise ReLU activations is invariant for the action of G on $\text{Param}(\mathbf{n})$.

D Projected gradient descent proofs

In this section, we give a proof of Theorem 8, which relates projected gradient descent for a representation with dimension \mathbf{n} to (usual) gradient descent for the corresponding reduced representation with dimension vector \mathbf{n}^{red} . This proof requires some set up and background results.

D.1 Gradient descent and orthogonal symmetries

We first prove a result that gradient descent commutes with invariant orthogonal transformations. This section is general and departs from the specific case of radial neural networks.

D.1.1 Setting

Let $\mathcal{L} : V = \mathbb{R}^p \rightarrow \mathbb{R}$ be a smooth function. Semantically, V is the parameter space of a neural network and \mathcal{L} the loss function with respect to a batch of training data. The differential $d\mathcal{L}_v$ of \mathcal{L} at $v \in V$ is row vector, while the gradient $\nabla_v \mathcal{L}$ of \mathcal{L} at v is a column vector⁴:

$$d\mathcal{L}_v = \left[\left. \frac{\partial \mathcal{L}}{\partial x_1} \right|_v \quad \cdots \quad \left. \frac{\partial \mathcal{L}}{\partial x_p} \right|_v \right] \quad \nabla_v \mathcal{L} = \begin{bmatrix} \left. \frac{\partial \mathcal{L}}{\partial x_1} \right|_v \\ \vdots \\ \left. \frac{\partial \mathcal{L}}{\partial x_p} \right|_v \end{bmatrix}$$

Hence $\nabla_v \mathcal{L}$ is the transpose of $d\mathcal{L}_v$, that is: $\nabla_v \mathcal{L} = (d\mathcal{L}_v)^T$. A step of gradient descent with respect to \mathcal{L} at learning rate $\eta > 0$ is defined as:

$$\begin{aligned} \gamma &= \gamma_\eta : V \longrightarrow V \\ v &\longmapsto v - \eta \nabla_v \mathcal{L} \end{aligned}$$

³See Armenta and Jodoin, *The Representation Theory of Neural Networks*, arXiv:2007.12213; Dinh, Pascanu, Bengio, and Bengio, *Sharp Minima Can Generalize For Deep Nets*, ICML 2017; Meng, Zheng, Zhang, Chen, Ye, Ma, Yu, and Liu, *G-SGD: Optimizing ReLU Neural Networks in its Positively Scale-Invariant Space*, 2019; and Neyshabur, Salakhutdinov, and Srebro. *Path-SGD: path-normalized optimization in deep neural networks*, NIPS’15.

⁴Following usual conventions, we regard column vectors as elements of V and row vectors as elements of the dual vector space V^* . The differential $d\mathcal{L}_v$ of \mathcal{L} at $v \in V$ is also known as the Jacobian of \mathcal{L} at $v \in V$.

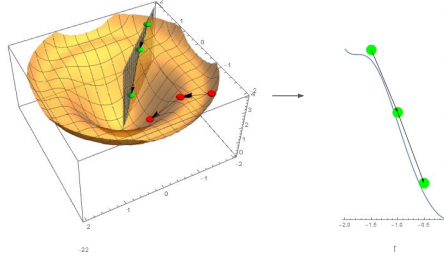


Figure 5: Illustration of Lemma 22. If the loss is invariant with respect to an orthogonal transformation Q of the parameter space, then optimization of the network by gradient descent is also invariant with respect to Q . (Note: in this example, projected and usual gradient descent match; this is not the case in higher dimensions, as explained in D.6.)

We drop η from the notation when it is clear from context. For any $k \geq 0$, we denote by γ^k the k -fold composition of the gradient descent map γ :

$$\gamma^k = \overbrace{\gamma \circ \gamma \circ \dots \circ \gamma}^k$$

D.1.2 Invariant group action

Now suppose $\rho : G \rightarrow \text{GL}(V)$ is an action of a Lie group G on V such that \mathcal{L} is G -invariant, i.e.:

$$\mathcal{L}(\rho(g)(v)) = \mathcal{L}(v)$$

for all $g \in G$ and $v \in V$. We write simply $g \cdot v$ for $\rho(g)(v)$, and g for $\rho(g)$.

Lemma 21. For any $v \in V$ and $g \in G$, we have:

$$\nabla_v \mathcal{L} = g^T \cdot (\nabla_{g \cdot v} \mathcal{L})$$

Proof. The proof is a computation:

$$\begin{aligned} \nabla_v \mathcal{L} &= (d_v \mathcal{L})^T = (d(\mathcal{L} \circ g)_v)^T = (d\mathcal{L}_{g \cdot v} \circ dg_v)^T = (d\mathcal{L}_{g \cdot v} \circ g)^T = g^T \cdot (d\mathcal{L}_{g \cdot v})^T \\ &= g^T \cdot (\nabla_{g \cdot v} \mathcal{L}) \end{aligned}$$

The second equality relies on the hypothesis that $\mathcal{L} \circ g = \mathcal{L}$, the third on the chain rule, and the fourth on the fact that $dg_v = g$ since g is a linear map. \square

One can perform the computation of the proof in coordinates, for $i = 1, \dots, p$:

$$\begin{aligned} (\nabla_v \mathcal{L})_i &= (d\mathcal{L}_v)^i = \frac{\partial \mathcal{L}}{\partial x_i} \Big|_v = \frac{\partial (\mathcal{L} \circ g)}{\partial x_i} \Big|_v = \frac{\partial \mathcal{L}}{\partial x_j} \Big|_{g \cdot v} \frac{\partial g_j}{\partial x_i} \Big|_v \\ &= (\nabla_{g \cdot v} \mathcal{L})_j g_j^i = (g^T)_i^j (\nabla_{g \cdot v} \mathcal{L})_j = (g^T \cdot \nabla_{g \cdot v} \mathcal{L})_i \end{aligned}$$

D.1.3 Orthogonal case

Furthermore, suppose the action of G is by orthogonal transformations, so that $\rho(g)^T = \rho(g)^{-1}$ for all $g \in G$. Then Lemma 21 implies that

$$\nabla_{g \cdot v} \mathcal{L} = g \cdot \nabla_v \mathcal{L} \tag{D.1}$$

for any $v \in V$ and $g \in G$. The proof of the following lemma is immediate from Equation D.1, together with the definition of γ . See Figure 5 for an illustration.

Lemma 22. Suppose the action of G on V is by orthogonal transformations, and that \mathcal{L} is G -invariant. Then the action of G commutes with gradient descent (for any learning rate). That is,

$$\gamma^k(g \cdot v) = g \cdot \gamma^k(v)$$

for any $v \in V$, $g \in G$, and $k \geq 0$.

928 D.2 Gradient descent notation and set-up

929 We now turn our attention back to radial neural networks. In this section, we recall notation
930 from above, and introduce new notation that will be relevant for the formulation and proof
931 of Theorem 8.

932 D.2.1 Merging widths and biases

933 Let $\mathbf{n} = (n_0, n_1, n_2, \dots, n_{L-1}, n_L)$ be the widths vector of an MLP. Recall the definition of
934 $\text{Param}(\mathbf{n})$ as the parameter space of all possible choices of trainable parameters:

$$\text{Param}(\mathbf{n}) = (\mathbb{R}^{n_1 \times n_0} \times \mathbb{R}^{n_2 \times n_1} \times \dots \times \mathbb{R}^{n_L \times n_{L-1}}) \times (\mathbb{R}^{n_1} \times \mathbb{R}^{n_2} \times \dots \times \mathbb{R}^{n_L})$$

935 We have been denoting an element therein as a pair of tuples (\mathbf{W}, \mathbf{b}) where $\mathbf{W} = (W_i \in$
936 $\mathbb{R}^{n_i \times n_{i-1}})_{i=1}^L$ are the weights and $\mathbf{b} = (b_i \in \mathbb{R}^{n_i})_{i=1}^L$ are the biases. However, in this
937 appendix we adopt different notation. Observe that, placing each bias vector as a extra
938 column on the left of the weight matrix, we obtain matrices:

$$A_i = [b_i \ W_i] \in \mathbb{R}^{n_i \times (1+n_{i-1})}.$$

939 Thus, there is an isomorphism:

$$\text{Param}(\mathbf{n}) \simeq \bigoplus_{i=1}^L \mathbb{R}^{n_i \times (n_{i-1}+1)} = \mathbb{R}^{n_1 \times (n_0+1)} \times \mathbb{R}^{n_2 \times (n_1+1)} \times \dots \times \mathbb{R}^{n_L \times (n_{L-1}+1)}$$

940 In this appendix, we regard an element of $\text{Param}(\mathbf{n})$ as a tuple of ‘merged’ matrices
941 $\mathbf{A} = (A_i \in \mathbb{R}^{n_i \times (1+n_{i-1})})_{i=1}^L$. We now define convenient maps to translate between the
942 merged notation and the split notation. For each i , define the extension-by-one map from
943 \mathbb{R}^{n_i} to $\mathbb{R} \times \mathbb{R}^{n_i} \simeq \mathbb{R}^{n_i+1}$ as follows:

$$\text{ext}_i : \mathbb{R}^{n_i} \rightarrow \mathbb{R}^{n_i+1} \quad v = (v_1, v_2, \dots, v_{n_i}) \mapsto (1, v_1, v_2, \dots, v_{n_i}) \quad (\text{D.2})$$

Observe that, for any i and $x \in \mathbb{R}^{n_{i-1}}$, we have

$$A_i \circ \text{ext}_{i-1}(x) = W_i x + b_i.$$

944 Consequently, the i -th partial feedforward function can be defined recursively as:

$$F_i = \rho_i \circ A_i \circ \text{ext}_{i-1} \circ F_{i-1} \quad (\text{D.3})$$

945 where $\rho_i : \mathbb{R}^{n_i} \rightarrow \mathbb{R}^{n_i}$ is the activation⁵ at the i -th layer, and F_0 is the identity on \mathbb{R}^{n_0} .

946 D.2.2 Orthogonal change-of-basis action

947 To describe the orthogonal change-of-basis symmetries of the parameter space in the
948 merged notation, recall the following product of orthogonal groups, with sizes correspond-
949 ing to the widths of the hidden layers:

$$O(\mathbf{n}^{\text{hid}}) = O(n_1) \times O(n_2) \times \dots \times O(n_{L-1})$$

950 In the merged notation, the element $\mathbf{Q} = (Q_i)_{i=1}^{L-1} \in O(\mathbf{n}^{\text{hid}})$ transforms $\mathbf{A} \in \text{Param}(\mathbf{n})$ as:

$$\mathbf{A} \mapsto \mathbf{Q} \cdot \mathbf{A} := \left(Q_i \circ A_i \circ \begin{bmatrix} 1 & 0 \\ 0 & Q_{i-1}^{-1} \end{bmatrix} \right)_{i=1}^L \quad (\text{D.4})$$

951 where $Q_0 = \text{id}_{n_0}$ and $Q_L = \text{id}_{n_L}$.

⁵In this general formulation, ρ_i can be any piece-wise differentiable function; for most of the rest of the paper we will be interested in the case where ρ_i is a radial rescaling function.

952 D.2.3 Model compression algorithm

953 We now restate Algorithm 1 in the merged notation. We emphasize that Algorithms 1 and
954 2 are mathematically equivalent; the later simply uses more compact notation.

Algorithm 2: QR Model Compression (QR-compress)

```

input   :  $\mathbf{A} \in \text{Param}(\mathbf{n})$ 
output  :  $\mathbf{Q} \in O(\mathbf{n}^{\text{hidden}})$  and  $\mathbf{V} \in \text{Param}(\mathbf{n}^{\text{red}})$ 

 $\mathbf{Q}, \mathbf{V} \leftarrow [], []$  // initialize output matrix lists
 $M_1 \leftarrow A_1$ 
for  $i \leftarrow 1$  to  $L - 1$  do // iterate through layers
     $Q_i, R_i \leftarrow \text{QR-decomp}(M_i, \text{mode} = \text{'complete'})$  //  $M_i = Q_i \circ \text{inc}_i \circ R_i$ 
    Append  $Q_i$  to  $\mathbf{Q}$ 
    Append  $R_i$  to  $\mathbf{V}$  // reduced merged weights for layer  $i$ 
    Set  $M_{i+1} \leftarrow A_{i+1} \circ \begin{bmatrix} 1 & 0 \\ 0 & Q_i \circ \text{inc}_i \end{bmatrix}$  // transform next layer
end
Append  $M_L$  to  $\mathbf{V}$ 
return  $\mathbf{Q}, \mathbf{V}$ 

```

956 We explain the notation. As noted in Appendix B.1, the symbol ‘ \circ ’ denotes composition
957 of maps, or matrix multiplication in the case of linear maps. The standard inclusion
958 $\text{inc}_i : \mathbb{R}^{n_i^{\text{red}}} \hookrightarrow \mathbb{R}^{n_i}$ maps into the first n_i^{red} coordinates. As a matrix, $\text{Inc}_i \in \mathbb{R}^{n_i \times n_i^{\text{red}}}$ has
959 ones along the main diagonal and zeros elsewhere. The method QR-decomp with mode =
960 ‘complete’ computes the complete QR decomposition of the $n_i \times (1 + n_{i-1}^{\text{red}})$ matrix M_i as
961 $Q_i \circ \text{inc}_i \circ R_i$ where $Q_i \in O(n_i)$ and R_i is upper-triangular of size $n_i^{\text{red}} \times (1 + n_{i-1}^{\text{red}})$. The
962 definition of n_i^{red} implies that either $n_i^{\text{red}} = n_{i-1}^{\text{red}} + 1$ or $n_i^{\text{red}} = n_i$. The matrix R_i is of size
963 $n_i^{\text{red}} \times n_i^{\text{red}}$ in the former case and of size $n_i \times (1 + n_{i-1}^{\text{red}})$ in the latter case.

964 D.2.4 Gradient descent definitions

965 As in Section 6, we fix:

- 966 • a widths vector $\mathbf{n} = (n_0, n_1, \dots, n_L)$.
- 967 • a tuple $\boldsymbol{\rho} = (\rho_1, \dots, \rho_L)$ of radial rescaling activations, where $\rho_i : \mathbb{R}^{n_i} \rightarrow \mathbb{R}^{n_i}$ for
968 $i = 1, \dots, L$.
- 969 • a batch of training data $\{(x_j, y_j)\} \subseteq \mathbb{R}^{n_0} \times \mathbb{R}^{n_L} = \mathbb{R}^{n_0^{\text{red}}} \times \mathbb{R}^{n_L^{\text{red}}}$.
- 970 • a cost function $\mathcal{C} : \mathbb{R}^{n_L} \times \mathbb{R}^{n_L} \rightarrow \mathbb{R}$

971 As a result, we have a loss function on $\text{Param}(\mathbf{n})$:

$$\mathcal{L} : \text{Param}(\mathbf{n}) \rightarrow \mathbb{R} \quad \mathcal{L}(\mathbf{A}) = \sum \mathcal{C}(F_{(\mathbf{A}, \boldsymbol{\rho})}(x_j), y_j)$$

972 where $F_{(\mathbf{A}, \boldsymbol{\rho})}$ is the feedforward of the radial neural network with (merged) parameters \mathbf{A}
973 and activations $\boldsymbol{\rho}$. We emphasize that the loss function \mathcal{L} depends on the batch of training
974 data chosen above; however, for clarity, we omit extra notation indicating this dependency
975 since the batch of training data is fixed throughout this discussion. Similarly, we have:

- 976 • the reduced widths vector $\mathbf{n}^{\text{red}} = (n_0^{\text{red}}, n_1^{\text{red}}, \dots, n_L^{\text{red}})$.
- 977 • the restrictions $\boldsymbol{\rho}^{\text{red}} = (\rho_1^{\text{red}}, \dots, \rho_L^{\text{red}})$, where $\rho_i^{\text{red}} : \mathbb{R}^{n_i^{\text{red}}} \rightarrow \mathbb{R}^{n_i^{\text{red}}}$ for $i = 1, \dots, L$.

978 Using the fact that $n_0^{\text{red}} = n_0$ and $n_L^{\text{red}} = n_L$, there is a loss function on $\text{Param}(\mathbf{n}^{\text{red}})$:

$$\mathcal{L}_{\text{red}} : \text{Param}(\mathbf{n}^{\text{red}}) \rightarrow \mathbb{R} \quad \mathcal{L}_{\text{red}}(\mathbf{B}) = \sum \mathcal{C}(F_{(\mathbf{B}, \boldsymbol{\rho}^{\text{red}})}(x_j), y_j)$$

where $F_{(\mathbf{B}, \rho^{\text{red}})}$ is the feedforward of the radial neural network with parameters $\mathbf{B} \in \text{Param}(\mathbf{n}^{\text{red}})$ and activations ρ^{red} . (Again, technically speaking, the loss function \mathcal{L}_{red} depends on the batch of training data fixed above.) For any learning rate $\eta > 0$, we obtain a gradient descent maps:

$$\begin{aligned} \gamma : \text{Param}(\mathbf{n}) &\rightarrow \text{Param}(\mathbf{n}) & \gamma_{\text{red}} : \text{Param}(\mathbf{n}^{\text{red}}) &\rightarrow \text{Param}(\mathbf{n}^{\text{red}}) \\ \mathbf{A} &\mapsto \mathbf{A} - \eta \nabla_{\mathbf{A}} \mathcal{L} & \mathbf{B} &\mapsto \mathbf{B} - \eta \nabla_{\mathbf{B}} \mathcal{L}_{\text{red}} \end{aligned}$$

D.3 The interpolating space

In this section, we introduce a subspace $\text{Param}^{\text{int}}(\mathbf{n})$ of $\text{Param}(\mathbf{n})$, that, as we will later see, interpolates between $\text{Param}(\mathbf{n})$ and $\text{Param}(\mathbf{n}^{\text{red}})$.

Let $\text{Param}^{\text{int}}(\mathbf{n})$ denote the subspace of $\text{Param}(\mathbf{n})$ consisting of those $\mathbf{T} = (T_1, \dots, T_L) \in \text{Param}(\mathbf{n})$ for which the bottom left $(n_i - n_i^{\text{red}}) \times (1 + n_{i-1}^{\text{red}})$ block of T_i is zero for each i . Schematically:

$$T_i = \begin{bmatrix} * & * \\ 0 & * \end{bmatrix}$$

where the rows are divided as n_i^{red} on top and $n_i - n_i^{\text{red}}$ on the bottom, while the columns are divided as $(1 + n_{i-1}^{\text{red}})$ on the left and $n_{i-1} - n_{i-1}^{\text{red}}$ on the right. Let

$$\iota_1 : \text{Param}^{\text{int}}(\mathbf{n}) \hookrightarrow \text{Param}(\mathbf{n})$$

be the inclusion. The following proposition follows from an elementary analysis of the workings of Algorithm 2 (or, equivalently, Algorithm 1).

Proposition 23. *Let $\mathbf{A} \in \text{Param}(\mathbf{n})$ and let $\mathbf{Q} \in O(\mathbf{n}^{\text{hid}})$ be the tuple of orthogonal matrices produced by Algorithm 2. Then $\mathbf{Q}^{-1} \cdot \mathbf{A}$ belongs to $\text{Param}^{\text{int}}(\mathbf{n})$.*

Define a map

$$q_1 : \text{Param}(\mathbf{n}) \rightarrow \text{Param}^{\text{int}}(\mathbf{n})$$

by taking $\mathbf{A} \in \text{Param}(\mathbf{n})$ and zeroing out the bottom left $(n_i - n_i^{\text{red}}) \times (1 + n_{i-1}^{\text{red}})$ block of A_i for each i . Schematically:

$$\mathbf{A} = \left(A_i = \begin{bmatrix} * & * \\ * & * \end{bmatrix} \right)_{i=1}^L \mapsto q_1(\mathbf{A}) = \left(\begin{bmatrix} * & * \\ 0 & * \end{bmatrix} \right)_{i=1}^L$$

It is straightforward to check that q_1 is a well-defined, surjective linear map. The transpose of q_1 is the inclusion ι_1 . We summarize the situation in the following diagram:

$$\begin{array}{ccc} \text{Param}^{\text{int}}(\mathbf{n}) & \xrightleftharpoons[q_1]{\iota_1} & \text{Param}(\mathbf{n}) \end{array} \quad (\text{D.5})$$

We observe that the composition $q_1 \circ \iota_1$ is the identity on $\text{Param}^{\text{int}}(\mathbf{n})$.

D.4 Projected gradient descent and model compression

Recall from Section 6 that the *projected gradient descent* map on $\text{Param}(\mathbf{n})$ is given by:

$$\gamma_{\text{proj}} : \text{Param}(\mathbf{n}) \rightarrow \text{Param}(\mathbf{n}), \quad \mathbf{A} \mapsto \text{Proj}(\mathbf{A} - \eta \nabla_{\mathbf{A}} \mathcal{L})$$

where $\mathbf{A} = (\mathbf{W}, \mathbf{b})$ are the merged parameters (Appendix D.2), and, in the notation of the previous section, the map Proj is $\iota_1 \circ q_1$. To reiterate, while all entries of each weight matrix and each bias vector contribute to the computation of the gradient $\nabla_{\mathbf{A}} \mathcal{L} = \nabla_{(\mathbf{W}, \mathbf{b})} \mathcal{L}$, only those not in the bottom left submatrix get updated under the projected gradient descent map γ_{proj} .

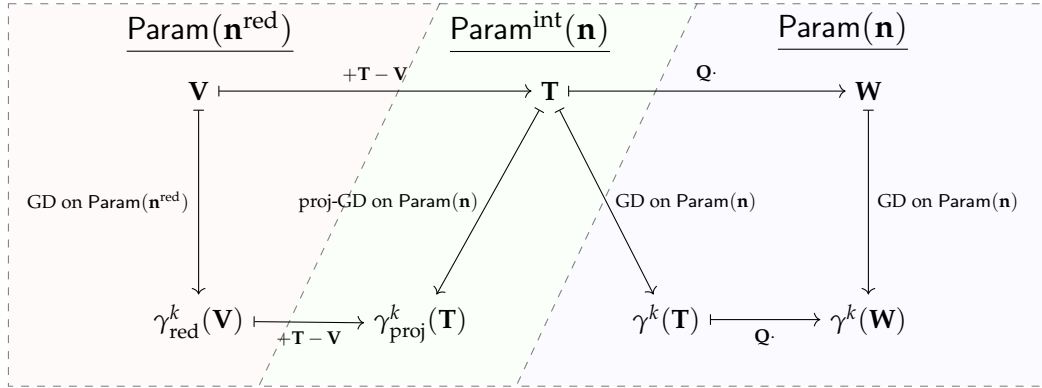
Let $\mathbf{V}, \mathbf{Q} = \text{QR-Compress}(\mathbf{A})$ be the outputs of Algorithm 2 (which is equivalent to Algorithm 1), so that $\mathbf{V} = (\mathbf{W}^{\text{red}}, \mathbf{b}^{\text{red}}) \in \text{Param}(\mathbf{n}^{\text{red}})$ are the parameters of the compressed model corresponding to the full model with merged parameters $\mathbf{A} = (\mathbf{W}, \mathbf{b})$, and $\mathbf{Q} \in O(\mathbf{n}^{\text{hid}})$ is an orthogonal change-of-basis symmetry of the parameter space. Moreover, set $\mathbf{T} = \mathbf{Q}^{-1} \cdot \mathbf{A} \in \text{Param}^{\text{int}}(\mathbf{n})$, where we use the change-of-basis action from Appendix D.2 and Proposition 23. We have the following rephrasing of Theorem 8.

Theorem 24 (Theorem 8). *Let $\mathbf{A} \in \text{Param}(\mathbf{n})$, and let $\mathbf{V}, \mathbf{Q}, \mathbf{T}$ be as above. For any $k \geq 0$:*

1. $\gamma^k(\mathbf{A}) = \mathbf{Q} \cdot \gamma^k(\mathbf{T})$
2. $\gamma_{\text{proj}}^k(\mathbf{T}) = \gamma_{\text{red}}^k(\mathbf{V}) + \mathbf{T} - \mathbf{V}$.

More precisely, the second equality is $\gamma_{\text{proj}}^k(\mathbf{T}) = \iota(\gamma_{\text{red}}^k(\mathbf{V})) + \mathbf{T} - \iota(\mathbf{V})$ where $\iota : \text{Param}(\mathbf{n}^{\text{red}}) \hookrightarrow \text{Param}(\mathbf{n})$ is the inclusion into the top left corner in each coordinate. Also, in the statement of Theorem 8, we have $\mathbf{U} = \mathbf{T} - \mathbf{V}$.

We summarize this result in the following diagram. The left horizontal maps indicate the addition of $\mathbf{U} = \mathbf{T} - \mathbf{V}$, the right horizontal arrows indicate the action of \mathbf{Q} , and the vertical maps are various versions of gradient descent. The shaded regions indicate the (smallest) vector space to which the various representations naturally belong.



D.5 Proof of Theorem 8

We begin by explaining the sense in which $\text{Param}^{\text{int}}(\mathbf{n})$ interpolates between $\text{Param}(\mathbf{n})$ and $\text{Param}(\mathbf{n}^{\text{red}})$. One extends Diagram D.5 as follows:

$$\text{Param}(\mathbf{n}^{\text{red}}) \xrightleftharpoons[q_2]{\iota_2} \text{Param}^{\text{int}}(\mathbf{n}) \xrightleftharpoons[q_1]{\iota_1} \text{Param}(\mathbf{n})$$

- The map

$$\iota_2 : \text{Param}(\mathbf{n}^{\text{red}}) \hookrightarrow \text{Param}^{\text{int}}(\mathbf{n})$$

takes $\mathbf{B} = (B_i) \in \text{Param}(\mathbf{n}^{\text{red}})$ and pad each matrix with $n_i - n_i^{\text{red}}$ rows of zeros on the bottom and $n_{i-1} - n_{i-1}^{\text{red}}$ columns of zeros on the right:

$$\mathbf{B} = (B_i)_{i=1}^L \mapsto \iota_2(\mathbf{B}) = \left(\begin{bmatrix} B_i & 0 \\ 0 & 0 \end{bmatrix} \right)_{i=1}^L$$

It is straightforward to check that ι_2 is a well-defined injective linear map.

- The map

$$q_2 : \text{Param}^{\text{int}}(\mathbf{n}) \twoheadrightarrow \text{Param}(\mathbf{n}^{\text{red}})$$

1032

extracts from \mathbf{T} the top left $n_i^{\text{red}} \times (1 + n_{i-1}^{\text{red}})$ matrix:

$$\mathbf{T} = \left(T_i = \begin{bmatrix} T_i^{(1)} & T_i^{(2)} \\ 0 & T_i^{(4)} \end{bmatrix} \right)_{i=1}^L \mapsto q_2(\mathbf{T}) = \left(T_i^{(1)} \right)_{i=1}^L$$

1033

It is straightforward to check that q_2 is a surjective linear map. The transpose of q_2 is the inclusion ι_2 .

1034

1035

1036 **Lemma 25.** *We have the following:*

1037

1. The inclusion $\iota : \text{Param}(\mathbf{n}^{\text{red}}) \hookrightarrow \text{Param}(\mathbf{n})$ coincides with the composition $\iota_1 \circ \iota_2$, and commutes with the loss functions:

1038

$$\begin{array}{ccc} \text{Param}(\mathbf{n}^{\text{red}}) & \xrightarrow{\iota_1 \circ \iota_2 = \iota} & \text{Param}(\mathbf{n}) \\ & \searrow \mathcal{L}_{\text{red}} & \swarrow \mathcal{L} \\ & \mathbb{R} & \end{array}$$

1039

2. The following diagram commutes:

$$\begin{array}{ccc} \text{Param}^{\text{int}}(\mathbf{n}) & \xrightarrow{q_2} & \text{Param}(\mathbf{n}^{\text{red}}) \\ \downarrow \iota_1 & & \downarrow \mathcal{L}_{\text{red}} \\ \text{Param}(\mathbf{n}) & \xrightarrow{\mathcal{L}} & \mathbb{R} \end{array}$$

1040

3. For any $\mathbf{T} \in \text{Param}^{\text{int}}(\mathbf{n})$, we have: $q_1 \left(\nabla_{\iota_1(\mathbf{T})} \mathcal{L} \right) = \iota_2 \left(\nabla_{q_2(\mathbf{T})} \mathcal{L}_{\text{red}} \right)$.

1041

Proof. We have the following standard inclusions into the first coordinates and projections onto the first coordinates, for $i = 0, 1, \dots, L$:

1042

$$\text{inc}_i = \text{inc}_{n_i^{\text{red}}, n_i} : \mathbb{R}^{n_i^{\text{red}}} \hookrightarrow \mathbb{R}^{n_i}, \quad \widetilde{\text{inc}}_i = \text{inc}_{1+n_i^{\text{red}}, 1+n_i} : \mathbb{R}^{1+n_i^{\text{red}}} \hookrightarrow \mathbb{R}^{1+n_i},$$

1043

$$\pi_i : \mathbb{R}^{n_i} \twoheadrightarrow \mathbb{R}^{n_i^{\text{red}}}, \quad \widetilde{\pi}_i : \mathbb{R}^{1+n_i} \twoheadrightarrow \mathbb{R}^{1+n_i^{\text{red}}}.$$

1044

Observe that $\text{Param}^{\text{int}}(\mathbf{n})$ is the subspace of $\text{Param}(\mathbf{n})$ consisting of those $\mathbf{T} = (T_1, \dots, T_L) \in \text{Param}(\mathbf{n})$ such that:

1045

$$(\text{id}_{n_i} - \text{inc}_i \circ \pi_i) \circ T_i \circ \widetilde{\text{inc}}_{i-1} \circ \widetilde{\pi}_{i-1} = 0$$

1046

for $i = 1, \dots, L$.

1047

By the definition of radial rescaling functions, for each $i = 1, \dots, L$, there is a piece-wise

1048

differentiable function $h_i : \mathbb{R} \rightarrow \mathbb{R}$ such that $\rho_i = h_i^{(n_i)}$. Note that $\rho_i^{\text{red}} = h_i^{(n_i^{\text{red}})}$, and

1049

$$h^{(n_i)} \circ \text{inc}_i = \text{inc}_i \circ h^{(n_i^{\text{red}})}.$$

1050

The identity $\iota = \iota_1 \circ \iota_2$ follows directly from definitions. To prove the commutativity of the first diagram, it is enough to show that, for any \mathbf{X} in $\text{Param}(\mathbf{n}^{\text{red}})$, the feedforward functions of \mathbf{X} and $\iota(\mathbf{X})$ coincide. This follows easily from the fact that, for $i = 1, \dots, L$, we have:

1053

$$\pi_i \circ h^{(n_i)} \circ \text{inc}_i = \pi_i \circ \text{inc}_i \circ h^{(n_i^{\text{red}})} = h^{(n_i^{\text{red}})}.$$

1054

For the second claim, let $\mathbf{T} \in \text{Param}^{\text{int}}(\mathbf{n})$. It suffices to show that $\iota_1(\mathbf{T})$ and $q_2(\mathbf{T})$ have the same feedforward function. Recall the ext_i maps and the formulation of the feedforward function in the merged notation given in Equation D.3. Using this set-up, the key computation is:

1055

1056

1057

$$\begin{aligned} \text{inc}_i \circ h^{(n_i^{\text{red}})} \circ \pi_i \circ T_i \circ \text{ext}_{n_{i-1}} \circ \text{inc}_{i-1} &= h^{(n_i)} \circ \text{inc}_i \circ \pi_i \circ T_i \circ \widetilde{\text{inc}}_{i-1} \circ \text{ext}_{n_{i-1}} \\ &= h^{(n_i)} \circ T_i \circ \widetilde{\text{inc}}_{i-1} \circ \text{ext}_{n_{i-1}} \\ &= h^{(n_i)} \circ T_i \circ \text{ext}_{n_{i-1}} \circ \text{inc}_{i-1} \end{aligned}$$

1058 which uses the fact that $(\text{id}_{n_i} - \text{inc}_i \circ \pi_i) \circ T_i \circ \widetilde{\text{inc}}_{i-1} = 0$, or, equivalently, $\text{inc}_i \circ \pi_i \circ T_i \circ$
 1059 $\widetilde{\text{inc}}_{i-1} = T_i \circ \widetilde{\text{inc}}_{i-1}$, as well as the fact that $\text{ext}_i \circ \text{inc}_i = \widetilde{\text{inc}}_i \circ \text{ext}_i$. Applying this relation
 1060 successively starting with the second-to-last layer ($i = L - 1$) and ending in the first ($i = 1$),
 1061 one obtains the result. For the last claim, one computes $\nabla_{\mathbf{T}}(\mathcal{L} \circ \iota_1)$ in two different ways.
 1062 The first way is:

$$\begin{aligned}\nabla_{\mathbf{T}}(\mathcal{L} \circ \iota_1) &= (d(\mathcal{L}_{\mathbf{T}} \circ \iota_1))^T = \left(d\mathcal{L}_{\iota_1(\mathbf{T})} \circ d\mathbf{T}\iota_1\right)^T = \left(d\mathcal{L}_{\iota_1(\mathbf{T})} \circ \iota_1\right)^T \\ &= \iota_1^T \left(d\mathcal{L}_{\iota_1(\mathbf{T})}^T\right) = q_1 \left(\nabla_{\iota_1(\mathbf{T})}\mathcal{L}\right)\end{aligned}$$

1063 where we use the fact that ι_1 is a linear map whose transpose is q_1 . The second way uses
 1064 the commutative diagram of the second part of the Lemma:

$$\begin{aligned}\nabla_{\mathbf{T}}(\mathcal{L} \circ \iota_1) &= \nabla_{\mathbf{T}}(\mathcal{L}_{\text{red}} \circ q_2) = (d(\mathcal{L}_{\text{red}})_{\mathbf{T}} \circ q_2)^T = \left(d(\mathcal{L}_{\text{red}})_{q_2(\mathbf{T})} \circ d(q_2)_{\mathbf{Z}}\right)^T \\ &= \left(d(\mathcal{L}_{\text{red}})_{q_2(\mathbf{T})} \circ q_2\right)^T = q_2^T \left(d(\mathcal{L}_{\text{red}})_{q_2(\mathbf{T})}^T\right) = \iota_2 \left(\nabla_{q_2(\mathbf{T})}\mathcal{L}_{\text{red}}\right).\end{aligned}$$

1065 We also use the fact that q_2 is a linear map whose transpose is ι_2 . □

1066 *Proof of Theorem 8.* As above, let $\mathbf{R}, \mathbf{Q} = \text{QR-compress}(\mathbf{A})$ be the outputs of Algorithm
 1067 1, so that $\mathbf{V} = (\mathbf{W}^{\text{red}}, \mathbf{b}^{\text{red}}) \in \text{Param}(\mathbf{n}^{\text{red}})$ is the dimensional reduction of the merged
 1068 parameters $\mathbf{A} = (\mathbf{W}, \mathbf{b})$, and $\mathbf{Q} \in O(\mathbf{n}^{\text{hid}})$. Set $\mathbf{T} = \mathbf{Q}^{-1} \cdot \mathbf{A} \in \text{Param}^{\text{int}}(\mathbf{n})$.

1069 The action of $\mathbf{Q} \in O(\mathbf{n}^{\text{hid}})$ on $\text{Param}(\mathbf{n})$ is an orthogonal transformation, so the first claim
 1070 follows from Lemma 22.

1071 For the second claim, it suffices to consider the case $\eta = 1$. The general case follows
 1072 similarly. We proceed by induction. The base case $k = 0$ amounts to Theorem 6. For the
 1073 induction step, we set

$$\mathbf{Z}^{(k)} = \iota(\gamma_{\text{red}}^k(\mathbf{V})) + \mathbf{T} - \iota(\mathbf{V}).$$

1074 Each $\mathbf{Z}^{(k)}$ belongs to $\text{Param}^{\text{int}}(\mathbf{n})$, so $i_1(\mathbf{Z}^{(k)}) = \mathbf{Z}^{(k)}$. Moreover, $q_2(\mathbf{Z}^{(k)}) = \gamma_{\text{red}}^k(\mathbf{V})$. We
 1075 compute:

$$\begin{aligned}\gamma_{\text{proj}}^{k+1}(\mathbf{Q}^{-1} \cdot \mathbf{A}) &= \gamma_{\text{proj}}\left(\gamma_{\text{proj}}^k(\mathbf{Q}^{-1} \cdot \mathbf{A})\right) \\ &= \gamma_{\text{proj}}\left(\iota(\gamma_{\text{red}}^k(\mathbf{V})) + \mathbf{T} - \iota(\mathbf{V})\right) \\ &= \iota_1 \circ q_1 \left(\iota(\gamma_{\text{red}}^k(\mathbf{V})) + \mathbf{T} - \iota(\mathbf{V}) - \nabla_{\iota(\gamma_{\text{red}}^k(\mathbf{V})) + \mathbf{T} - \iota(\mathbf{V})}\mathcal{L}\right) \\ &= \iota(\gamma_{\text{red}}^k(\mathbf{V})) - \iota_1 \circ q_1 \left(\nabla_{\iota_1(\mathbf{Z}^{(k)})}\mathcal{L}\right) + \mathbf{T} - \iota(\mathbf{V}) \\ &= \iota(\gamma_{\text{red}}^k(\mathbf{V})) - \iota_1 \circ \iota_2 \left(\nabla_{q_2(\mathbf{Z}^{(k)})}\mathcal{L}_{\text{red}}\right) + \mathbf{T} - \iota(\mathbf{V}) \\ &= \iota \left(\gamma_{\text{red}}^k(\mathbf{V}) - \nabla_{\gamma_{\text{red}}^k(\mathbf{V})}\mathcal{L}_{\text{red}}\right) + \mathbf{T} - \iota(\mathbf{V}) \\ &= \iota \left(\gamma_{\text{red}}^{k+1}(\mathbf{V})\right) + \mathbf{T} - \iota(\mathbf{V})\end{aligned}$$

1076 where the second equality uses the induction hypothesis; the third invokes the definition
 1077 of γ_{proj} ; the fourth uses the fact that $\mathbf{Z}^{(k)} = \iota(\gamma_{\text{red}}^k(\mathbf{V})) + \mathbf{T} - \iota(\mathbf{V})$ belongs to $\text{Param}^{\text{int}}(\mathbf{n})$;
 1078 the fifth and sixth use Lemma 25 above; and the last uses the definition of γ_{red} . □

1079 D.6 Example

1080 We now discuss an example where projected gradient descent does not match usual
 1081 gradient descent.

1082 Let $\mathbf{n} = (1, 3, 1)$ be a widths vector. The space of parameters with this widths vector is
 1083 10-dimensional:

$$\text{Param}(\mathbf{n}) = \text{Hom}(\mathbb{R}^2, \mathbb{R}^3) \oplus \text{Hom}(\mathbb{R}^4, \mathbb{R}) \simeq \mathbb{R}^{10}.$$

1084 We identify a choice of parameters (in the merged notation)

$$\mathbf{A} = \left(A_1 = \begin{bmatrix} a & b \\ c & d \\ e & f \end{bmatrix}, A_2 = [g \quad h \quad i \quad j] \right) \in \text{Param}((1, 3, 1)) \quad (\text{D.6})$$

1085 with the point $p = (a, b, c, d, e, f, g, h, i, j)$ in \mathbb{R}^{10} . To be even more explicit, the weights for

1086 the first layer are $W_1 = \begin{bmatrix} b \\ d \\ f \end{bmatrix}$, the bias in the first hidden hidden layer is $b_1 = (a, c, e)$, the

1087 weights for the second layer are $W_2 = [h \quad i \quad j]$, and the bias for the output layer is $b_2 = g$.

1088 The action of the orthogonal group $O(\mathbf{n}) = O(3)$ on $\text{Param}(\mathbf{n}) \simeq \mathbb{R}^{10}$ can be expressed as:

$$Q \mapsto \begin{bmatrix} Q & 0 & 0 & 0 \\ 0 & Q & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & Q \end{bmatrix},$$

1089 where the rows and columns are divided according to the partition $3 + 3 + 1 + 3 = 10$.

1090 Consider the function⁶:

$$\begin{aligned} \mathcal{L} : \text{Param}(\mathbf{n}) &\rightarrow \mathbb{R} \\ p = (a, b, c, d, e, f, g, h, i, j) &\mapsto h(a + b) + i(c + d) + j(e + f) + g \end{aligned}$$

1091 By the product rule, we have:

$$\nabla_p \mathcal{L} = (h, h, i, i, j, j, 1, a + b, c + d, e + f)$$

1092 One easily checks that $\mathcal{L}(Q \cdot p) = \mathcal{L}(p)$ and that $\nabla_{Q \cdot p} \mathcal{L} = Q \cdot \nabla_p \mathcal{L}$ for any $Q \in O(3)$.

1093 The interpolating space is the eight-dimensional subspace of $\text{Param}(\mathbf{n}) \simeq \mathbb{R}^{10}$ with $e =$
 1094 $f = 0$ (using the notation of Equation D.6). Suppose $p' = (a, b, c, d, 0, 0, g, h, i, j)$ belongs to
 1095 the interpolating space. Then the gradient is

$$\nabla_{p'} \mathcal{L} = (h, h, i, i, j, j, 1, a + b, c + d, 0)$$

1096 which does not belong to the interpolating space. So one step of usual gradient descent,
 1097 with learning rate $\eta > 0$ yields:

$$\begin{aligned} \gamma : p' &= (a, b, c, d, 0, 0, g, h, i, j) \mapsto \\ &(a - \eta h, b - \eta h, c - \eta i, d - \eta i, -\eta j, -\eta j, g - \eta, h - \eta(a + b), i - \eta(c + d), j) \end{aligned}$$

1098 On the other hand, one step of projected gradient descent yields:

$$\begin{aligned} \gamma_{\text{proj}} : p' &= (a, b, c, d, 0, 0, g, h, i, j) \mapsto \\ &(a - \eta h, b - \eta h, c - \eta i, d - \eta i, 0, 0, g - \eta, h - \eta(a + b), i - \eta(c + d), j) \end{aligned}$$

1099 Direct computation shows that the difference between the evaluation of \mathcal{L} after one step of
 1100 gradient descent and the evaluation of \mathcal{L} after one step of projected gradient descent is:

$$\mathcal{L}(\gamma(p')) - \mathcal{L}(\gamma_{\text{proj}}(p')) = 2\eta j^2.$$

⁶For $\mathbf{A} \in \text{Param}(\mathbf{n})$, the neural function of the neural network with affine maps determined by \mathbf{A} and identity activation functions is $\mathbb{R} \rightarrow \mathbb{R}; x \mapsto \mathcal{L}(\mathbf{W})x$. The function \mathcal{L} can appear as a loss function for certain batches of training data and cost function on \mathbb{R} .

E Experiments

As mentioned in Section 7, we provide an implementation of Algorithm 1 in order to (1) empirically validate that our implementation satisfies the claims of Theorems 6 and Theorem 8 and (2) quantify real-world performance. Our implementation uses a generalization of radial neural networks, which we explain presently.

E.1 Radial neural networks with shifts

In this section, we consider radial neural networks with an extra trainable parameter in each layer that shifts the radial rescaling activation. Adding such parameters allows for more flexibility in the model, and (as shown in Theorem 26) the model compression of Theorem 6 holds for such networks. It is this generalization that we use in our experiments.

Let $h : \mathbb{R} \rightarrow \mathbb{R}$ be a function. For any $n \geq 1$ and any $t \in \mathbb{R}$, the corresponding *shifted radial rescaling function* on \mathbb{R}^n is given by:

$$\rho = h^{(n,t)} : v \mapsto \frac{h(|v| - t)}{|v|}v$$

if $v \neq 0$ and $\rho(0) = 0$. A *radial neural network with shifts* consists of the following data:

1. Hyperparameters: A positive integer L and a widths vector $\mathbf{n} = (n_0, n_1, n_2, \dots, n_L)$.
2. Trainable parameters:
 - (a) A choice of weights and biases $(\mathbf{W}, \mathbf{b}) \in \text{Param}(\mathbf{n})$.
 - (b) A vector of shifts $\mathbf{t} = (t_1, t_2, \dots, t_L) \in \mathbb{R}^L$.
3. Activations: A tuple $\mathbf{h} = (h_1, \dots, h_L)$ of piecewise differentiable functions $\mathbb{R} \rightarrow \mathbb{R}$. Together with the shifts, we have the shifted radial rescaling activation $\rho_i = h_i^{(n_i, t_i)} : \mathbb{R}^{n_i} \rightarrow \mathbb{R}^{n_i}$ in each layer.

The *feedforward function* of a radial neural network with shifts is defined in the usual recursive way, as in Section 3. The trainable parameters form the vector space $\text{Param}(\mathbf{n}) \times \mathbb{R}^L$, and the loss function of a batch of training data $\{(x_i, y_i)\} \subset \mathbb{R}^{n_0} \times \mathbb{R}^{n_L}$ is defined as

$$\mathcal{L} : \text{Param}(\mathbf{n}) \times \mathbb{R}^L \longrightarrow \mathbb{R}; \quad (\mathbf{W}, \mathbf{t}) \mapsto \sum_j \mathcal{C}(F_{(\mathbf{W}, \mathbf{b}, \mathbf{t}, \mathbf{h})}(x_j), y_j)$$

where $F_{(\mathbf{W}, \mathbf{b}, \mathbf{t}, \mathbf{h})}$ is the feedforward function of a radial neural network with weights \mathbf{W} , biases \mathbf{b} , shifts \mathbf{t} , and radial rescaling activations produced from \mathbf{h} . We have the gradient descent map:

$$\gamma : \text{Param}(\mathbf{n}) \times \mathbb{R}^L \longrightarrow \text{Param}(\mathbf{n}) \times \mathbb{R}^L$$

which updates the entries of \mathbf{W} , \mathbf{b} , and \mathbf{t} . The group $O(\mathbf{n}^{\text{hid}}) = O(n_1) \times \dots \times O(n_{L-1})$ acts on $\text{Param}(\mathbf{n})$ as usual (see Section 5.1), and on \mathbb{R}^L trivially. The neural function is unchanged by this action. We conclude that the $O(\mathbf{n}^{\text{hid}})$ action on $\text{Param}(\mathbf{n}) \times \mathbb{R}^L$ commutes with gradient descent γ . We now state a generalization of Theorem 6 for the case of radial neural networks with shifts. We omit a proof, as it uses the same techniques as the proof of Theorem 6.

Theorem 26. *Let $(\mathbf{W}, \mathbf{b}, \mathbf{t}, \mathbf{h})$ be a radial neural network with shifts and widths vector \mathbf{n} . Let \mathbf{W}^{red} and \mathbf{b}^{red} be the weights and biases of the compressed network produced by Algorithm 1. The feedforward function of the original network $(\mathbf{W}, \mathbf{b}, \mathbf{t}, \mathbf{h})$ coincides with that of the compressed network $(\mathbf{W}^{\text{red}}, \mathbf{b}^{\text{red}}, \mathbf{t}, \mathbf{h})$.*

Theorem 8 also generalizes to the setting of radial neural networks with shifts, using projected gradient descent with respect to the subspace $\text{Param}^{\text{int}}(\mathbf{n}) \times \mathbb{R}^L$ of $\text{Param}(\mathbf{n}) \times \mathbb{R}^L$.

1139 E.2 Implementation details

1140 Our implementation is written in Python and uses the QR decomposition routine in
 1141 NumPy [21]. We also implement a general class RadNet for radial neural networks using
 1142 PyTorch [41]. For brevity, we write $\hat{\mathbf{W}}$ for (\mathbf{W}, \mathbf{b}) and $\hat{\mathbf{W}}^{\text{red}}$ for $(\mathbf{W}^{\text{red}}, \mathbf{b}^{\text{red}})$.

1143 **(1) Empirical verification of Theorem 6.** We use synthetic data to learn the function
 1144 $f(x) = e^{-x^2}$ with $N = 121$ samples $x_j = -3 + j/20$ for $0 \leq j < 121$. We model $f_{\hat{\mathbf{W}}}$
 1145 as a radial neural network with widths $\mathbf{n} = (1, 6, 7, 1)$ and activation the radial shifted
 1146 sigmoid $h(x) = 1/(1 + e^{-x+s})$. Applying QR-compress gives a radial neural network
 1147 $f_{\hat{\mathbf{W}}^{\text{red}}}$ with widths $\mathbf{n}^{\text{red}} = (1, 2, 3, 1)$. Theorem 6 implies that the neural functions of
 1148 $f_{\hat{\mathbf{W}}}$ and $f_{\hat{\mathbf{W}}^{\text{red}}}$ are equal. Over 10 random initializations of $\hat{\mathbf{W}}$, the mean absolute error
 1149 $(1/N) \sum_j |f_{\hat{\mathbf{W}}}(x_j) - f_{\hat{\mathbf{W}}^{\text{red}}}(x_j)| = 1.31 \cdot 10^{-8} \pm 4.45 \cdot 10^{-9}$. Thus $f_{\hat{\mathbf{W}}}$ and $f_{\hat{\mathbf{W}}^{\text{red}}}$ agree up to
 1150 machine precision.

1151 **(2) Empirical verification of Theorem 8.** Adopting the notation from above, the claim is
 1152 that training $f_{\mathbf{Q}^{-1}, \hat{\mathbf{W}}}$ with objective \mathcal{L} by projected gradient descent coincides with training
 1153 $f_{\hat{\mathbf{W}}^{\text{red}}}$ with objective \mathcal{L}_{red} by usual gradient descent. We verified this on synthetic data
 1154 using 3000 epochs at learning rate 0.01. Over 10 random initializations of $\hat{\mathbf{W}}$, the loss
 1155 functions match up to machine precision with $|\mathcal{L} - \mathcal{L}_{\text{red}}| = 4.02 \cdot 10^{-9} \pm 7.01 \cdot 10^{-9}$.

1156 **(3) Reduced model trains faster.** Due to the relation between projected gradient descent
 1157 of the full network $\hat{\mathbf{W}}$ and gradient descent of the reduced network $\hat{\mathbf{W}}^{\text{red}}$, our method may
 1158 be applied before training to produce a smaller model class which *trains* faster without
 1159 sacrificing accuracy. We test this hypothesis in learning the function $f : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ sending
 1160 $x = (t_1, t_2)$ to $(e^{-t_1^2}, e^{-t_2^2})$ using $N = 121^2$ samples $(-3 + j/20, -3 + k/20)$ for $0 \leq j, k < 121$.
 1161 We model $f_{\hat{\mathbf{W}}}$ as a radial neural network with layer widths $\mathbf{n} = (2, 16, 64, 128, 16, 2)$
 1162 and activation the radial sigmoid $h(r) = 1/(1 + e^{-r})$. Applying QR-compress gives a radial
 1163 neural network $f_{\hat{\mathbf{W}}^{\text{red}}}$ with widths $\mathbf{n}^{\text{red}} = (2, 3, 4, 5, 6, 2)$. We trained both models until
 1164 the training loss was ≤ 0.01 . Running on a system with an Intel i5-8257U@1.40GHz and
 1165 8GB of RAM and averaged over 10 random initializations, the reduced network trained in
 1166 15.32 ± 2.53 seconds and the original network trained in 31.24 ± 4.55 seconds.

1167 F Relation to radial basis function networks

1168 In this appendix, we show that radial neural networks are equivalent to a particular class of
 1169 multilayer radial basis functions networks. This class is obtained by imposing the condition
 1170 that the so-called ‘hidden dimension’ at each layer is equal to one; the total number of
 1171 layers, however, is unconstrained. To our knowledge, the literature contains no universal
 1172 approximation result for this class of radial basis functions networks.

1173 F.1 Single layer case

1174 We first recall the definition of a radial basis function network. A *local linear model extension*
 1175 *of a radial basis function network* (henceforth abbreviated simply by *RBFN*) consists of:

- 1176 • An input dimension n , an output dimension m , and a ‘hidden’ dimension N .
- 1177 • For $i = 1, \dots, N$, a matrix $W_i \in \mathbb{R}^{m \times n}$, a vector $b_i \in \mathbb{R}^n$, and a weight $a_i \in \mathbb{R}^m$.
- 1178 • A nonlinear function⁷ $\lambda : \mathbb{R} \rightarrow \mathbb{R}$.

⁷A more general version allows for a different nonlinear function for every $i = 1, \dots, N$.

The feedforward function of a RBFN is defined as:

$$F : \mathbb{R}^n \rightarrow \mathbb{R}^m \quad x \mapsto \sum_{i=1}^N (a_i + W_i(x + b_i)) \lambda(|x + b_i|).$$

1179 The integer N is commonly referred to as ‘the hidden number of neurons’. This is a bit of
 1180 a misnomer. Really there is only one layer with input dimension n and output dimension
 1181 m ; the integer N is part of the specification of the activation function.

We observe that if $N = 1$ and $a_1 = 0$, then the feedforward function is given by:

$$F : \mathbb{R}^n \rightarrow \mathbb{R}^m \quad x \mapsto W\rho(x + b)$$

1182 where ρ is the radial rescaling function determined by λ . In words, one adds $b_1 = b \in \mathbb{R}^n$
 1183 to the input vector x , applies the activation ρ to obtain new vector in \mathbb{R}^n , and then applies
 1184 the linear transformation determined by the matrix $W_1 = W$ to obtain the output vector in
 1185 \mathbb{R}^m . Motivated by this observation, we say that a RBFN is *constrained* if $N = 1$ and $a_1 = 0$.

1186 F.2 Constrained multilayer case

1187 Next, we consider the constrained multilayer case of a radial basis functions network.
 1188 Specifically, a *constrained multilayer* RBFN consists of:

- 1189 • A widths vector (n_0, \dots, n_L) where L is the number of layers.
- 1190 • A matrix $W_\ell \in \mathbb{R}^{n_\ell \times n_{\ell-1}}$ for $\ell = 1, \dots, L$.
- 1191 • A vector $b_\ell \in \mathbb{R}^{n_\ell}$ for $\ell = 0, 1, \dots, L - 1$.
- 1192 • A nonlinear function $\lambda_\ell : \mathbb{R} \rightarrow \mathbb{R}$ for $\ell = 0, 1, \dots, L - 1$. (Equivalently, the
 1193 corresponding radial rescaling function $\rho_\ell : \mathbb{R}^{n_\ell} \rightarrow \mathbb{R}^{n_\ell}$ for $\ell = 0, \dots, L - 1$.)

The feedforward function is defined as follows. For $\ell = 0, \dots, L$, we recursively define
 $F_\ell : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_\ell}$ by setting $F_0(x) = x$ and

$$F_\ell(x) = W_\ell \rho_{\ell-1}(F_{\ell-1}(x) + b_{\ell-1})$$

1194 for $\ell = 1, \dots, L$. The feedforward function is F_L .

1195 F.3 Relation to radial neural networks

1196 We now demonstrate that radial neural networks are equivalent to multilayer RBFNs.

1197 **Proposition 27.** *For any radial neural network, there is a constrained multilayer RBFN with the*
 1198 *same feedforward function. Conversely, for any constrained multilayer RBFN, there is a radial*
 1199 *neural network with the same feedforward function.*

Proof. For the first statement, let $(\mathbf{W}, \mathbf{b}, \rho)$ be a radial neural network with L layers and
 widths vector (n_0, \dots, n_L) . Recall the partial feedforward functions $G_\ell : \mathbb{R}^{n_0} \rightarrow \mathbb{R}^{n_\ell}$ defined
 recursively by setting $G_0(x) = x$ and

$$G_\ell(x) = \rho_\ell(W_\ell G_{\ell-1}(x) + b_\ell)$$

1200 The feedforward function is G_L . Consider the constrained multilayer RBFN with $L + 1$
 1201 layers and the following:

- 1202 • Widths vector $(n_0, n_1, \dots, n_{L-1}, n_L, n_L)$. The last two layers have the same dimen-
 1203 sion.
- 1204 • Weight matrices $W_\ell \in \mathbb{R}^{n_\ell \times n_{\ell-1}}$ for $\ell = 1, \dots, L$ and $W_{L+1} = \text{id}_{n_L} \in \mathbb{R}^{n_L \times n_L}$.
- 1205 • A vector $b_\ell \in \mathbb{R}^{n_\ell}$ for $\ell = 1, \dots, L$, and $b_0 = 0 \in \mathbb{R}^{n_0}$.
- 1206 • A radial rescaling activation $\rho_\ell : \mathbb{R}^{n_\ell} \rightarrow \mathbb{R}^{n_\ell}$ for $\ell = 1, \dots, L$, and $\rho_0 = \text{id}_{n_0}$.

Let F_ℓ be the partial feedforward functions for this RBFN, defined recursively as above. We claim that

$$F_\ell(x) = W_\ell \circ G_{\ell-1}(x)$$

for any $x \in \mathbb{R}^{n_0}$ and $\ell = 1, \dots, L$. We prove this by induction. The base case is $\ell = 1$:

$$F_1(x) = W_1 \circ \rho_0(F_0(x) + b_0) = W_1 x = W_1 \circ G_0(x)$$

For the induction step, take $\ell > 1$ and compute:

$$F_\ell(x) = W_\ell \circ \rho_{\ell-1}(F_{\ell-1}(x) + b_{\ell-1}) = W_\ell \circ \rho_{\ell-1}(W_{\ell-1}G_{\ell-2}(x) + b_{\ell-1}) = W_\ell \circ G_{\ell-1}(x)$$

1207 The first claim now follows from the case $\ell = L$, using the fact that W_{L+1} is the identity.

1208 For the second statement, let $(\mathbf{W}, \mathbf{b}, \rho)$ be a constrained multilayer RBFN with L layers and
1209 widths vector (n_0, \dots, n_L) . Consider the radial neural network with $L + 1$ layers and the
1210 following:

- 1211 • Widths vector $(n_0, n_0, n_1, \dots, n_{L-1}, n_L)$. The first two layers have the same dimen-
1212 sion.
- 1213 • Weight matrices given by $\tilde{W}_1 = \text{id}_{n_0}$ and $\tilde{W}_\ell = W_{\ell-1}$ for $\ell = 2, \dots, L + 1$.
- 1214 • Bias vectors given by $\tilde{b}_\ell = b_{\ell-1}$ for $\ell = 1, 2, \dots, L$, and $\tilde{b}_{L+1} = 0$.
- 1215 • Radial rescaling activations given by $\tilde{\rho}_\ell = \rho_{\ell-1}$ for $\ell = 1, \dots, L$, and $\tilde{\rho}_{L+1} = \text{id}_{n_L}$.

One uses the recursive definition of the partial feedforward functions to show that, for $\ell = 1, \dots, L$, we have $F_\ell(x) = W_\ell \circ G_\ell(x)$, where F_ℓ and G_ℓ are the partial feedforward functions of the RBFN and radial neural network, respectively. Then:

$$G_{L+1}(x) = \tilde{\rho}_{L+1}(\tilde{W}_{L+1} \circ G_L(x) + \tilde{b}_{L+1}) = W_L \circ G_L(x) = F_L(x),$$

1216 so the two feedforward functions coincide. □

1217 F.4 Conclusions

1218 While radial neural networks are equivalent to a certain class of radial basis function
1219 network, we point out differences between our results and the standard theory of radial
1220 basis functions network. First, RBFNs generally only have two layers; we consider ones
1221 with unbounded depth. Second, to our knowledge, ours is the first universal approximation
1222 result such that:

- 1223 • it uses networks in the subclass of multilayer RBFNs satisfying the constraint that
1224 all the number of ‘hidden neurons’ in each layer is equal to 1.
- 1225 • it approximates functions with networks of bounded width.
- 1226 • it can be used to approximate asymptotically affine functions, rather than functions
1227 defined on a compact domain.

1228 Our compressibility result may apply to multilayer RBFNs where the number of ‘hidden
1229 neurons’ N_ℓ at each layer is not equal to 1, but we expect the compression to be weaker,
1230 and that constrained multilayer RBFNs are in some sense the most compressible type of
1231 RBFN.