

A Proofs

Theorem 1 (Resource consumption of universal RLCs). *Let (x, y) be a binary classification task that admits a smooth separator as in Assumption 1. Then, there exists an RLC with neural network f_{θ^*} and absolutely continuous randomness source \mathbf{u} (Assumption 2) that is universal in the limit, i.e.,*

$$\mathcal{F}_{\theta^*}(x) = y(x), \forall x \in \mathcal{X},$$

and makes random predictions that are correct with probability

$$P(\text{maj}(\{\text{sgn}(\langle \mathbf{a}_{\theta^*}^{(j)}, x \rangle - \mathbf{b}_{\theta^*}^{(j)})\}_{j=1}^m) = y(x)) > 1 - \exp\{-2\epsilon^2 m^2\},$$

where ϵ is the minimum bias of \mathcal{F}_{θ^*} .

Further, if p^\dagger is the number of parameters used by a deterministic neural network with one hidden layer to achieve zero-error in the task, f_θ has at most

$$p \leq p^\dagger + \mathcal{O}(1) \text{ parameters.}$$

Proof of Theorem 1. Since Assumption 1 holds³, there exists a single hidden-layer neural network N that, like s , achieves zero-error in this task [8]. Further, since sgn is nonpolynomial, we can use it as the non-linearity of this network [21]. Putting it all together, there exists a number of hidden units M and parameters $b_j, o_j \in \mathbb{R}, w_j \in \mathbb{R}^d$ for $j = 1, \dots, M$ such that

$$N(x) := \sum_{j=1}^M o_j \text{sgn}(\langle w_j, x \rangle - b_j),$$

and

$$N(x) = \text{sgn}(s(x)), \forall x \in \mathcal{X}.$$

Note that this means we can achieve zero-error in classification, $N(x) = y(x), \forall x \in \mathcal{X}$. Having this in mind, we will now show that for any network N constructed as such, we can build a limiting classifier \mathcal{F}_θ equal to it. Let us first note that we can force positivity in the output weights of N by doing

$$\begin{aligned} o'_j &:= o_j \cdot \text{sgn}(o_j) \\ w'_j &:= w_j \cdot \text{sgn}(o_j) \end{aligned}$$

without changing the classification

$$\text{sgn}\left(\sum_{j=1}^M o_j \text{sgn}(\langle w_j, x \rangle - b_j)\right) = \text{sgn}\left(\sum_{j=1}^M o'_j \text{sgn}(\langle w'_j, x \rangle - b_j)\right), \forall x \in \mathcal{X}.$$

Further, note that we can also multiply the output weights by the constant $\sum_{j=1}^M o'_j$ without changing the sign and thus the classification. Overall, we can define a network N^* equivalent to N in classification with

$$N^*(x) := \sum_{j=1}^M \frac{1}{\sum_{j=1}^M o'_j} o'_j \text{sgn}(\langle w'_j, x \rangle - b_j).$$

Now, note that the limiting classifier can be written as $\mathcal{F}_\theta(x) = \text{sgn}(\mathbb{E}_{\mathbf{a}_\theta, \mathbf{b}_\theta} \text{sgn}(\langle \mathbf{a}_\theta, x \rangle - \mathbf{b}_\theta))$. We can then define its distribution as one that samples from the weights and bias w'_j, b_j with probabilities according to $\frac{1}{\sum_{j=1}^M o'_j} o'_j$. This implies that $\mathcal{F}_\theta(x) = N^*(x) = N(x)$, making our \mathcal{F}_θ achieve zero-error in the task. Note that if $p^\dagger := M$ is the number of parameters in N , we can generate the linear coefficients in $\mathcal{F}_\theta(x)$ with at most $p^\dagger + \mathcal{O}(1)$ —simply take the external noise from a uniform distribution and map it to the probabilities in the output layer to pick the coefficients.

Finally, note that we have to also characterize the probability that $\mathbf{y}_\theta^{(m)}$ outputs a different answer from \mathcal{F}_θ . The common tool for this is Hoeffding's inequality, which lets us upper-bound it for some $x \in \mathcal{X}$ with

$$P(\mathbf{y}_\theta^{(m)} \neq \mathcal{F}_\theta(x)) \leq \exp\{-2\epsilon_x^2 m^2\},$$

where $\epsilon_x = P(\mathbf{y}_\theta^{(m)} = \mathcal{F}_\theta(x)) - 0.5$. Note that ϵ_x depends on x and our definition of the minimum bias ϵ in the main text is precisely the lowest of all such ϵ_x . Thus, as stated we finally have that for all $x \in \mathcal{X}$

$$P(\mathbf{y}_\theta^{(m)} = \mathcal{F}_\theta(x)) > 1 - \exp\{-2\epsilon^2 m^2\}.$$

Note that the above is true for any \mathcal{F}_θ , including the universal and invariant ones later discussed in the main text. □

³Further taking the usual assumption that \mathcal{X} is compact.

Theorem 2 (*G*-invariant RLCs). *Let (\mathbf{x}, y) be a G -invariant task with a smooth separator as in Assumption 1. Then, the set of RLCs with a G -invariant distribution in the classifier weights, i.e.,*

$$\mathbf{a}_\theta \stackrel{d}{=} g \cdot \mathbf{a}_\theta, \forall g \in G,$$

and absolutely continuous randomness source (cf. Assumption 2) is both probabilistic G -invariant and universal in (\mathbf{x}, y) . That is,

$$\mathcal{F}_\theta(x) = \mathcal{F}_\theta(g \cdot x), \forall x \in \mathcal{X}, \forall g \in G, \forall \theta \in \mathbb{R}^p,$$

and

$$\exists \theta^* \in \mathbb{R}^p: \mathcal{F}_{\theta^*}(x) = y(x), \forall x \in \mathcal{X}, \forall g \in G.$$

Proof of Theorem 2. Let us start with Proposition 3, a central observation needed in Theorem 2. Put into words Proposition 3 says that we can transfer the action of G to the weights of the linear classifier. For the reader more familiar with group representation theory, the result follows immediately from noting that compact groups admit unitary representations, see [28] for a good resource on the matter.

Proposition 3. *If G is a compact group and $g \cdot x$ is an action of $g \in G$ on $x \in \mathbb{R}^d$, there exists a bijective mapping $\alpha: G \rightarrow G$ such that*

$$\langle g \cdot x, w \rangle = \langle x, \alpha(g) \cdot w \rangle, \forall x, w \in \mathbb{R}^d, \forall g \in G.$$

Proof. Since we are dealing with actions on a finite-dimensional vector space, the action \cdot of g can be associated with a linear representation ρ of G , i.e., $\rho(g) \in \mathbb{R}^{d \times d}$. That is, if $[x]$ is the column vector of x , we have that $\rho(g)[x] = g \cdot x$. Now, since G is compact and \mathbb{R}^d is finite-dimensional, G also admits a unitary linear representation ρ_u with associated action \cdot_u . Then, it is easy to see that

$$\langle g \cdot_u x, w \rangle = \langle \rho_u(g)[x], w \rangle = \langle x, \rho_u(g)^{-1}[w] \rangle = \langle x, g^{-1} \cdot_u [w] \rangle, \forall w \in \mathbb{R}^d.$$

Thus, since inversion defines a bijection on G , Proposition 3 holds for the unitary representation action. Now, note that—as any other group action— \cdot_u defines a homomorphism between G and the symmetric group of \mathbb{R}^d . Then, for any other action \cdot there exists a bijective function $\beta: G \rightarrow G$ such that $g \cdot x = \beta(g) \cdot_u x, \forall x \in \mathbb{R}^d$. Hence, we can write

$$\langle g \cdot x, w \rangle = \langle \beta(g) \cdot_u x, w \rangle = \langle \beta(g)^{-1} \cdot_u x, w \rangle = \langle \beta^{-1}(\beta(g)^{-1}) \cdot x, w \rangle, \forall x, w \in \mathbb{R}^d, \forall g \in G.$$

Finally, since β and group inversion are bijective, we finish the proof by defining the bijective mapping $\alpha(g) := \beta^{-1}(\beta(g)^{-1})$ for every $g \in G$. \square

Now, we can proceed to prove the universality part of Theorem 2. Since the task admits a smooth separator, there exists a universal limiting classifier \mathcal{F}_θ for it. That is, for every $x \in \mathcal{X}$

$$\mathcal{F}_\theta(x) = y(x).$$

Since the task is G -invariant, we know that for every $g \in G$

$$\mathcal{F}_\theta(x) = y(g \cdot x).$$

Note that since G is compact, it admits a unique normalized Haar measure λ^4 . Then, the insight comes from considering a random group action $\mathbf{g} \sim \lambda$. By Fubini's theorem and Proposition 3, we have

$$\mathcal{F}_\theta(x) = \mathbb{E}_{\mathbf{g}}[\mathcal{F}_\theta(\mathbf{g} \cdot x)] = \text{sgn}(\mathbb{E}_{\mathbf{g}, \mathbf{a}_\theta, \mathbf{b}_\theta}[\text{sgn}(\langle x, \alpha(\mathbf{g}) \cdot \mathbf{a}_\theta \rangle - \mathbf{b}_\theta)]).$$

Now, we can define a limiting classifier \mathcal{F}_{θ^*} that samples coefficients according to $\mathbf{b}_{\theta^*} \stackrel{d}{=} \mathbf{b}_\theta$ and $\mathbf{a}_{\theta^*} \stackrel{d}{=} \mathbf{g} \cdot \mathbf{a}_\theta$. From above, we know that $\mathcal{F}_{\theta^*}(x) = \mathcal{F}_\theta(x), \forall x \in \mathcal{X}$. Since \mathcal{F}_θ is universal, \mathcal{F}_{θ^*} is also universal. Finally, since \mathbf{g} is sampled according to the Haar measure, we have that

$$\mathbf{a}_{\theta^*} \stackrel{d}{=} g \cdot \mathbf{a}_{\theta^*}, \forall g \in G.$$

At last, it is easy to see from Proposition 3 that the action on the weights can be transferred to the input as well, meaning that every model is probabilistic invariant, i.e., $\mathcal{F}_\theta(x) = \mathcal{F}_\theta(g \cdot x), \forall x \in \mathcal{X}, \forall g \in G, \forall \theta \in \mathbb{R}^p$, which finalizes the proof. \square

Proposition 1 (Infinitely G -invariant RLCs). *Let (\mathbf{x}, y) be a G -invariant task with infinitely G -invariant data (Assumption 3) with a smooth separator as in Assumption 1. Then, the set of RLCs with an infinitely G -invariant distribution in the linear classifier weights, i.e., as in Assumption 3 $(\mathbf{a}_{\theta_i})_{i=1}^\infty \stackrel{d}{=} g_\infty \cdot (\mathbf{a}_{\theta_i})_{i=1}^\infty, \forall g_\infty \in G_\infty$, where $\mathbf{a}_\theta \stackrel{d}{=} (\mathbf{a}_{\theta_i})_{i \in S}$, for $S \in \binom{\mathbb{N}}{d}$, and absolutely continuous randomness source (cf. Assumption 2) is probabilistic G -invariant and universal for (\mathbf{x}, y) as in Theorem 2.*

⁴The reader can think of λ as a uniform distribution over G .

Proof of Proposition 1. From Proposition 3 it also follows that the model is G -invariant. Let us proceed to prove it is universal. Since the task admits a smooth separator, from Theorem 1 there exists an RLC \mathcal{F}_θ such that

$$\mathbb{E}_{\mathbf{x}}[\mathcal{F}_\theta(\mathbf{x}) - y(\mathbf{x})] = 0.$$

Since the data is infinitely G -invariant, we have that for the infinite G -invariant sequence $(\mathbf{x}_i)_{i=1}^\infty$

$$\mathbb{E}_{\mathbf{x}}[\mathcal{F}_\theta(\mathbf{x}) - y(\mathbf{x})] = \mathbb{E}_{(\mathbf{x}_i)_{i \in S}}[\mathcal{F}_\theta((\mathbf{x}_i)_{i \in S}) - y((\mathbf{x}_i)_{i \in S})] = 0, \text{ for } S \in \binom{\mathbb{N}}{d},$$

with $\binom{\mathbb{N}}{d}$ being the set of all d -size subsets of \mathbb{N} . Let λ be the unique normalized Haar measure of G . As in the proof of Theorem 2, we can take a random group action $\mathbf{g} \sim \lambda$ and have from Proposition 3 and Fubini's theorem that

$$\begin{aligned} \mathbb{E}_{(\mathbf{x}_i)_{i \in S}}[\mathcal{F}_\theta((\mathbf{x}_i)_{i \in S}) - y((\mathbf{x}_i)_{i \in S})] &= \mathbb{E}_{((\mathbf{x}_i)_{i \in S})_{i \in S}, \mathbf{g}}[\mathcal{F}_\theta(\mathbf{g} \cdot ((\mathbf{x}_i)_{i \in S}))] \\ &= 0, \text{ for } S \in \binom{\mathbb{N}}{d}. \end{aligned}$$

Now, let us define a task (\mathbf{x}', y) on a higher dimension $d' \geq d$, $\text{supp}(\mathbf{x}) \subset \text{supp}(\mathbf{x}') \subseteq \mathbb{R}^{d'}$ where $y(\mathbf{x}') = y(\mathbf{x}'_{[1:d]})$. Further, \mathbf{g}' is a random group action (sampled from the Haar measure) from G' , the homomorphism of G into the dimension d' . From above, together with Proposition 3 we have that

$$\begin{aligned} \mathbb{E}_{(\mathbf{x}_i)_{i \in S'}}[\mathcal{F}_\theta((\mathbf{x}_i)_{i \in S'}) - y((\mathbf{x}_i)_{i \in S'})] &= \mathbb{E}_{((\mathbf{x}_i)_{i \in S'})_{i \in S'}, \mathbf{g}'}[\mathcal{F}_\theta(\mathbf{g}' \cdot ((\mathbf{x}_i)_{i \in S'}))] \\ &= \mathbb{E}_{((\mathbf{x}_i)_{i \in S'})_{i \in S'}, \mathbf{g}'}[\text{sgn}(\mathbb{E}_{\mathbf{a}_\theta, \mathbf{b}_\theta, \mathbf{g}'}[\text{sgn}(\langle x, \alpha(\mathbf{g}') \cdot \mathbf{a}'_\theta \rangle - \mathbf{b}'_\theta)])] \\ &= 0, \text{ for } S' \in \binom{\mathbb{N} \setminus \{1, \dots, d\}}{d'}. \end{aligned}$$

Now, from above, since the lower dimensions $\{1, \dots, n\}$ are G -invariant if all higher dimensions are G -invariant, it is easy to see that there is an infinitely G -invariant sequence of weights $(\mathbf{a}_{\theta_i})_{i=1}^\infty$ that achieves zero error. \square

Proposition 2 (Universality and resource consumption of RSetCs). *Let (\mathbf{x}, y) be a permutation-invariant task with a smooth separator as in Assumption 1 and infinitely permutation-invariant data as in Assumption 3. Then, RSetCs as in Definition 1 with absolutely continuous randomness source (cf. Assumption 2) are probabilistic G -invariant and universal for (\mathbf{x}, y) (as in Theorem 2). Further, the number of parameters needed by RSetCs in this task will depend only on the smallest finite absolute moments of the weight and bias distributions in the zero-error solution, i.e., the ones given by $f_{\theta_f^*}(\mathbf{u}, \mathbf{u}_i)$ and $g_{\theta_g^*}(\mathbf{u}, \mathbf{u}_b)$.*

Proof of Proposition 2. The result follows directly from the combination of de Finetti's theorem [9], Kallenberg's noise transfer theorem [16] and the recent results on the capacity of neural networks to generate distributions from noise [34]. Let us outline it in more detail. From de Finetti's theorem [9] we know that any infinite sequence of exchangeable random variables can be expressed by i.i.d. random variables conditioned on a common latent measure. Combining this with Kallenberg's noise transfer theorem we have that the weights and biases can be written as $f(\mathbf{u}, \mathbf{u}_i)$ and $g(\mathbf{u}, \mathbf{u}_b)$ where f, g are measurable maps and the noises are sampled from a uniform distribution. Finally, the work of [34] showed that we can replace the uniform noise with absolutely continuous noise and f, g with a sufficiently expressive ReLU multi-layer perceptron. Moreover, Theorem 2.1 in [34] says that the number of parameters depends only on the output dimension and the smallest finite absolute moments of the distributions. Since the output dimension k is constant for us, it depends solely on the latter and we finalize the proof. \square

Theorem 3. *Let (\mathbf{x}, y) be a graph isomorphism-invariant task that either i) has a smooth separator as in Assumption 1 or ii) is an inner-product decision graph problem as in Definition 3. Further, the task has infinitely graph isomorphism-invariant data as in Assumption 3. Then, there exists an RGraphC as in Definition 2 with absolutely continuous randomness source (cf. Assumption 2) that is probabilistic G -invariant and universal for (\mathbf{x}, y) (as in Theorem 2). Further, the number of parameters of this RGraphC will depend only on the smallest finite absolute moments of its weight and bias distributions, i.e., the ones given by $f_{\theta_f}(\mathbf{u}, \mathbf{u}_i)$ and $g_{\theta_g}(\mathbf{u}, \mathbf{u}_b)$.*

Proof of Theorem 3. Note that it suffices to show that there exists an RLC that takes as input graphs of size d and decides correctly tasks in Definition 3 with probability > 0.5 —this implies a universal limiting classifier

\mathcal{F}_θ . Then, we can replace the smoothness assumption in Theorem 2 and Proposition 1 with Definition 3 and we will have that there exists an infinitely jointly exchangeable sequence that is universal for tasks as in Definition 3. Finally, we follow Proposition 2’s proof by simply replacing de Finetti’s with Aldous-Hoover’s theorem. Thus, let us now proceed to prove the main part of this theorem, *i.e.*, there exists an RLC that takes as input graphs of size d and decides correctly tasks in Definition 3 with probability > 0.5 .

Define an RLC that samples the linear coefficients as follows.

- Let \mathbf{a}'_θ be a random variable such that $\text{supp}(\mathbf{a}'_\theta) = S$. It is important to note that we do not require that this random vector to have the same distribution as the input, or any other distribution on S . We simply need it to be supported on S .
- Now, let $t \sim \text{Bernoulli}(0.5 + \gamma)$ for some $\gamma > 0$. We define $\mathbf{a}_\theta := t \cdot \mathbf{a}'_\theta$ and $\mathbf{b}_\theta := t \cdot b$. Our random prediction can then be rewritten as $\text{sgn}(t \cdot (\langle \mathbf{a}'_\theta, x \rangle - b))$.

Let us now calculate our probability of success. For convenience, let us define $\text{sgn}(0) = 1^5$. Then, the probability of a positive graph x^+ , *i.e.*, $y(x^+) = 1$, being classified correctly is

$$\begin{aligned} P(\text{sgn}(\langle \mathbf{a}_\theta, x^+ \rangle - \mathbf{b}_\theta) = 1) &= P(\text{sgn}(t \cdot (\langle \mathbf{a}'_\theta, x^+ \rangle - b)) = 1) \\ &= 0.5 - \gamma + P(\langle \mathbf{a}'_\theta, x^+ \rangle \geq b) \cdot (0.5 + \gamma), \end{aligned}$$

while the probability of a negative graph x^- , *i.e.*, $y(x^-) = -1$, being classified correctly is

$$P(\text{sgn}(\langle \mathbf{a}_\theta, x^- \rangle - \mathbf{b}_\theta) = 1) = P(\text{sgn}(t \cdot (\langle \mathbf{a}'_\theta, x^- \rangle - b)) = 1) = 0.5 + \gamma.$$

Then, since \mathcal{X} is a finite set for any input distribution \mathbf{x} there exists a constant γ where $0 < \gamma/(0.5 + \gamma) < P(\langle \mathbf{a}'_\theta, x \rangle \geq b)$ for all $x \in \mathcal{X}$. Given such gamma, every input has a probability of success greater than 0.5, which implies that $\mathcal{F}_\theta(x) = y(x), \forall x \in \mathcal{X}$ as we wanted to show. \square

B RLCs for spherical data

Here we want to show how the same ideas presented in Sections 3.1 and 3.2 can be applied to spherical data using Freedman’s theorem [13]. More specifically, we let our input be d -dimensional vectors, *i.e.*, $\text{supp}(\mathbf{x}) \subseteq \mathbb{R}^d$ and our task invariant to the orthogonal group, *i.e.*, $y(x) = y(g \cdot x), \forall x \in \text{supp}(\mathbf{x}), \forall g \in G$, where $G := O(d)$. We refer to tasks like this as tasks over spherical data since it is invariant to rotations and reflections.

Just as de Finetti, Aldous and Hoover informed us about sets and graphs, Freedman did it for spherical data [13]. It follows from his work that an infinite sequence invariant to the action of the orthogonal group can be represented as a random scalar multiplying a sequence of i.i.d. standard Gaussian distributions. Leveraging this, we can define the following very simple invariant model for spherical data (assuming Assumption 3).

Definition 4 (Randomized Sphere Classifiers). *A Randomized Sphere Classifier (RSphereC) uses 2 neural networks $f_{\theta_f} : \mathbb{R}^2 \rightarrow \mathbb{R}$ and $g_{\theta_g} : \mathbb{R}^2 \rightarrow \mathbb{R}$ together with 3 sources of randomness: \mathbf{u}, \mathbf{u}_a and \mathbf{u}_b . The random linear classifier coefficients are generated with*

$$\mathbf{a}_\theta \stackrel{d}{=} f_{\theta_f}(\mathbf{u}, \mathbf{u}_a) \cdot [\mathcal{N}(0, 1), \dots, \mathcal{N}(0, 1)], \text{ and } \mathbf{b}_\theta \stackrel{d}{=} g_{\theta_g}(\mathbf{u}, \mathbf{u}_b),$$

where $[\mathcal{N}(0, 1), \dots, \mathcal{N}(0, 1)]$ is a vector of d i.i.d. standard Gaussians.

Proposition 4 (Universality and resource consumption of RSphereCs). *Let (x, y) be a $O(d)$ -invariant task with a smooth separator as in Assumption 1 and infinitely spherical data as in Assumption 3. Then, RSphereCs as in Definition 4 with absolutely continuous randomness source (cf. Assumption 2) are probabilistic G -invariant and universal for (x, y) (as in Theorem 2). Further, the number of parameters needed by RSphereCs in this task will depend only on the smallest finite absolute moments of the weight and bias distributions in the zero-error solution, *i.e.*, the ones given by $f_{\theta_f^*}(\mathbf{u}, \mathbf{u}_i)$ and $g_{\theta_g^*}(\mathbf{u}, \mathbf{u}_b)$.*

Proof. The proof is exactly as the one from Proposition 2 but replacing de Finetti’s theorem with Freedman’s [13]. For a clean and more accessible material on Freedman’s theorem, we refer the reader to the lectures by Kallenberg in <https://mysite.science.uottawa.ca/givanoff/wskallenberg.pdf>. \square

Finally, we note that tasks invariant to $O(n)$ are not that common. However, tasks invariant to $\text{SO}(n)$ are extremely relevant in computer vision, see [6]. Therefore, we think this model is a first step towards probabilistic invariance for $\text{SO}(n)$ ⁶.

⁵This can be done by always adding a sufficiently small constant to the input.

⁶Note that $\text{SO}(n)$ is the connected component of $O(n)$.

C Implementation Details

As mentioned in the main text, all models were trained for a maximum of 1000 epochs using a cosine annealing scheduler for the learning. We tuned all the hyperparameters on the validation set using a patience of 30. The Deep Sets models found a better learning rate of 0.001 and batch size of 250. The GNN model found a better learning rate of 0.01 and batch size 100. The RSetC model used a batch size of 250 and learning rate 0.5. The RGraphC model used a batch size of 100 and a learning rate of 0.5. Finally, we also note that RLCs can output different answers given the same weights (due to the random input). Thus, we amplify each mini-batch with 1000 samples for each model—we noted that this reduces the variance and helps convergence. The amplification size used in test was $m = 10$. Finally, for all RLCs we used a single dimensional standard normal distribution in each of their noise sources.

Source code is available at: <https://github.com/cottascience/invariant-rlcs>