



# An Embedding of ReLU Networks and an Analysis of their Identifiability

Pierre Stock, Rémi Gribonval

## ► To cite this version:

Pierre Stock, Rémi Gribonval. An Embedding of ReLU Networks and an Analysis of their Identifiability. 2021. hal-03292203v1

**HAL Id: hal-03292203**

**<https://hal.science/hal-03292203v1>**

Preprint submitted on 20 Jul 2021 (v1), last revised 31 Jan 2022 (v4)

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# AN EMBEDDING OF RELU NETWORKS AND AN ANALYSIS OF THEIR IDENTIFIABILITY

PIERRE STOCK AND RÉMI GRIBONVAL

**ABSTRACT.** Neural networks with the Rectified Linear Unit (ReLU) nonlinearity are described by a vector of parameters  $\theta$ , and realized as a piecewise linear continuous function  $\mathbf{R}_\theta : x \in \mathbb{R}^d \mapsto \mathbf{R}_\theta(x) \in \mathbb{R}^k$ . Natural scalings and permutations operations on the parameters  $\theta$  leave the realization unchanged, leading to equivalence classes of parameters that yield the same realization. These considerations in turn lead to the notion of identifiability – the ability to recover (the equivalence class of)  $\theta$  from the sole knowledge of its realization  $\mathbf{R}_\theta$ . The overall objective of this paper is to introduce an embedding for ReLU neural networks of any depth,  $\Phi(\theta)$ , that is invariant to scalings and that provides a locally linear parameterization of the realization of the network. Leveraging these two key properties, we derive some conditions under which a deep ReLU network is indeed locally identifiable from the knowledge of the realization on a finite set of samples  $x_i \in \mathbb{R}^d$ . We study the shallow case in more depth, establishing necessary and sufficient conditions for the network to be identifiable from a bounded subset  $\mathcal{X} \subseteq \mathbb{R}^d$ .

## CONTENTS

1. Introduction	2
Around Functional Identifiability	3
Embeddings Zoology	3
Applications	4
2. General setting and main results	4
2.1. Network architectures	5
2.2. Realization of a network	5
2.3. Invariance to permutation and scaling	5
2.4. An invariant embedding of ReLU networks	7
2.5. Some consequences of PS-identifiability	8
2.6. Identifiability conditions in the shallow case	11
2.7. A glimpse at the analysis of local identifiability	12
2.8. Non-degeneracy and irreducibility in shallow <i>vs</i> deeper architectures	13
2.9. Discussion	14
3. Rescaling invariance of the embedding	16
4. Analyzing local identifiability	19
4.1. Activation status of neurons and paths, and activation spaces	19
4.2. “Algebraic” expressions of the realization	20
4.3. Non-degeneracy and local S-identifiability	22

5. Identifiability for shallow neural networks	24
5.1. Activation spaces and twin neurons	24
5.2. Proof of Lemma 5: no twins implies non-degeneracy	24
5.3. Proof of Theorem 3: irreducibility and no twins implies PS-identifiability	25
5.4. Local S-identifiability despite the presence of twins	27
5.5. Discussion of the role of activation spaces	29
Acknowledgements	30
References	30
Appendix A. Proof of Lemma 7	32
Appendix B. Proof of Theorem 2	33
Appendix C. Proof of Lemma 3	34
Appendix D. Proof of Lemma 4	35
Appendix E. Proof of Lemma 9 and Lemma 10	36
Appendix F. Proof of Lemma 11	38
Appendix G. Proof of Lemma 13	39
Appendix H. Proof of Lemma 16	40
Appendix I. Details on Example 4	42
Appendix J. Details on Example 5	43

## 1. INTRODUCTION

The empirical success of Deep Neural Networks (DNNs) for traditional machine learning tasks such as image classification is a well-known fact for the research community [1]. While this empirical success percolates to areas ranging from protein folding to symbolic mathematics, a second well-known fact is that the theoretical tools to grasp DNNs and uncover the reasons of their success are still lagging behind the fast-paced experimental results. We argue that a deeper understanding of the expressivity and stability properties of such networks could lead to practical improvements [2, 3]. In this paper, we introduce an embedding,  $\Phi(\theta)$ , of the vector  $\theta$  of network parameters (weights and biases) that exhibits interesting properties for networks based on the popular Rectified Linear Unit (ReLU): in particular,  $\Phi(\theta)$  is invariant to natural rescalings of the parameters that leave unchanged the function implemented by the network. To showcase the potential of this tool, we leverage it to study the expressivity of DNNs from the perspective of their functional equivalence classes.

In the remainder of this section, we first list the papers tackling identifiability of neural networks with a given non-linearity function. We next present some related work that construct an embedding for ReLU networks. Finally, we list applications directly or indirectly derived from the two previous theoretical considerations.

**Around Functional Identifiability.** First, various results dating back from the 90’s identify conditions that allow to identify neural networks *with one hidden layer* equipped with various non-linearities<sup>1</sup> like the hyperbolic tangent [4, 5, 6, 7]. Such results *do not* encompass the ReLU case. Simultaneously, Fefferman derived identifiability conditions for deep networks equipped with the tanh nonlinearity using complex analysis [8]. More recently, the work of Rolnick and Kording [9] reflects a renewed interest for this subject and its application to ReLU networks. The authors propose to reverse-engineer deep ReLU networks and present a constructive algorithm that samples network realizations  $\mathbf{R}_\theta(x)$  for carefully chosen input points  $x$  to deduce the architecture of the network and its parameters, up to rescalings and permutations. The authors prove that their algorithm terminates, except for a measure-zero set of parameters<sup>2</sup>. Similarly, Fornasier *et al.* [10] propose to recover the parameters of a *two-hidden-layer* neural network with smooth nonlinearity by actively sampling finite difference approximations to Hessians of the network, and by combining the insights gained from the sampling with a heuristic for precise attribution of the parameters to the architecture. The authors demonstrate the empirical effectiveness of their approach and claim that the proposed method can be generalized to networks with any depth. Finally, Phuong and Lampert provide a result related to identifiability for ReLU networks under some more restrictive assumptions [11].

**Embeddings Zoology.** We list here the neural network embeddings in the literature that are the closest to our own embedding,  $\Phi(\theta)$ , which is introduced in Definition 6. Its main property is that it is invariant under the action of rescalings, as stated more formally in Theorem 1. Schematically,  $\Phi(\theta)$  lives in the linear space indexed by network *paths* and each coordinate is a product of weights and/or biases along a particular network path. In a similar fashion, Malgouyres and Landsberg [12] consider a particular class of *linear* structured networks called *Deep Structured Neural Networks*, without biases, and consider only layer-wise rescalings. In [12, Section 6], the authors provide sufficient and necessary conditions for local identifiability by studying complex algebraic varieties leveraging the *Segre embedding* of such networks. The Segre embedding bears a resemblance with  $\Phi(\theta)$  since it is also made of product of network parameters, but it does not encompass the biases. Moreover, we consider neuron-wise rescalings in this paper as opposed to less general layer-wise rescalings considered by the authors, and also encompass the ReLU non-linearity in our approach. Malgouyres later leverages the Segre embedding to study local stability properties of linear neural networks [13]. Finally, Neyshabur *et al.* introduce a family of *path regularizers* to derive an optimization procedure that takes the invariance of the realization  $\mathbf{R}_\theta$  under the action of the rescalings into account and called Path-SGD [14, 15]. Such path regularizers are scalars – as opposed to vectorial embeddings – that are obtained by summing, for any network path, the norm of the product of all weights along this paths. Moreover, this approach does not take the biases into account, as opposed to our embedding,  $\Phi(\theta)$ .

---

<sup>1</sup>It should be noted that the functional equivalence class generally depends on the considered non-linearity. For instance, with the hyperbolic tangent, the authors only consider permutations and sign flips.

<sup>2</sup>This measure-zero set of parameters is not explicitly described by the authors.

**Applications.** As illustrated with Path-SGD [15], several papers attempt to perform the optimization in the space of network parameters *quotiented* by the rescaling operation. Several work follow and perform the optimization by alternating between a standard SGD step and a projection step that modifies the rescaling coefficients without changing the function implemented by the network [16, 17, 18, 19]. The main difference between these papers is the projection step, that is performed either implicitly (with a regularizer) or explicitly (by computing the optimal<sup>3</sup> rescaling coefficients). In the latter case, the proposed empirical methods may not yield the optimal rescaling coefficients but rather more or less stable and good approximations. Another advantage of rescalings is to improve post-training scalar quantization of neural networks by carefully selecting the rescaling coefficients such that the dynamical range of the weights within a layer is relatively small, with as few outliers as possible [20, 21]. More related to the concept of (local) identifiability, Carlini *et al.* [22] design a differential attack to efficiently recover the parameters of remote model up to floating point precision, by sending carefully designed queries  $x$  to the remote network and receiving only its output.

After introducing the main notations, we define  $\Phi(\theta)$  and state the main results of the paper in Section 2. Then, we formally state and prove the main properties of the embedding  $\Phi(\theta)$  in Section 3. In particular, we prove that  $\Phi(\theta)$  is invariant under the action of the rescalings. Next, we leverage this embedding to derive partial and local identifiability results for ReLU neural networks of any depth in Section 4. To further demonstrate the validity of our approach, we fully study the shallow case in Section 5 and provide conditions under which a ReLU neural network with one hidden layer is identifiable. Finally, we argue that  $\Phi(\theta)$  may be leveraged to tackle other open problems in the Machine Learning community.

## 2. GENERAL SETTING AND MAIN RESULTS

We consider fully-connected feedforward ReLU neural networks with  $L \geq 2$  affine layers. Each network is supported on a graph  $G = (E, V)$  with vertex set  $V$  composed of neurons  $\nu$  and edge set  $E$  composed of connections. The set of neurons  $V$  is partitioned into the input layer  $N_0$ ,  $L - 1$  hidden layers  $N_\ell$ ,  $1 \leq \ell \leq L - 1$ , and the output layer  $N_L$ . Hidden neurons compose the set  $H = \cup_{\ell=1}^{L-1} N_\ell$ . Since we focus on fully-connected networks, the set of connections  $E$  is made of all oriented edges  $e = \nu \rightarrow \nu'$  between neurons belonging to consecutive layers,  $\nu \in N_{\ell-1}$ ,  $\nu' \in N_\ell$  for some  $1 \leq \ell \leq L$ . The subset of incoming edges of neuron  $\nu$  is denoted  $\bullet \rightarrow \nu$ , while  $\nu \rightarrow \bullet$  denotes its set of outgoing edges.

Each edge  $e \in E$  is equipped with a weight  $w_e$  and each hidden neuron  $\nu \in H$  with a bias  $b_\nu$ . Output neurons, i.e. neurons from the last layer  $\eta \in N_L$ , are also equipped with a bias  $b_\eta$ , which is sometimes constrained to be zero. The set of all neurons equipped with biases is  $\bar{H} := H \cup N_L$ . Parameters (weights and biases) are gathered in a parameter vector  $\theta \in \mathbb{R}^{E \cup \bar{H}}$  where  $E \cup \bar{H}$  indexes all possible weights and biases including biases on the output layer. For brevity we may denote  $\theta_e$  for weights and  $\theta_\nu$  for biases. When

<sup>3</sup>In the sense that such coefficients globally minimize a given objective function.

needed we also write  $\theta = (\theta_i)_{i \in E \cup \bar{H}}$ . Since we consider a fully connected architecture (this does not prevent some weights to possibly vanish on some edges),  $\theta$  can also be represented as a set of  $L$  matrices  $\mathbf{W}_\ell = (w_{\nu \rightarrow \nu'})_{\nu' \in N_\ell, \nu \in N_{\ell-1}} \in \mathbb{R}^{N_\ell \times N_{\ell-1}}$ ,  $1 \leq \ell \leq L$  and  $L$  vectors  $\mathbf{b}_\ell = (b_\nu)_{\nu \in N_\ell} \in \mathbb{R}^{N_\ell}$ ,  $1 \leq \ell \leq L$ .

**2.1. Network architectures.** Many of the notions of parameter identifiability or non-degeneracy that will be considered are relative to a choice of network “architecture”. This is represented both by the graph  $G$  (which determines how many layers there are, and how wide they are) but also by a possibly restricted set  $\Theta \subseteq \mathbb{R}^{E \cup \bar{H}}$  of network parameters, which may for example account for the following type of constraints:

- restricting to a convolutional structure;
- restricting to sparse networks, possibly with structured sparsity patterns;
- restricting to networks without output biases ( $b_\eta = 0$  for every  $\eta \in N_L$ );
- restricting to networks without biases ( $b_\nu = 0$  for every  $\nu \in H \cup N_L$ ).

**2.2. Realization of a network.** Given a parameter  $\theta$  and an input vector  $x \in \mathbb{R}^{N_0}$ , we sequentially define  $\mathbf{y}_0(\theta, x) = x$  and for each  $1 \leq \ell \leq L - 1$  the pre-activation  $\mathbf{z}_\ell(\theta, x) = \mathbf{W}_\ell \mathbf{y}_{\ell-1}(\theta, x) + \mathbf{b}_\ell \in \mathbb{R}^{N_\ell}$ , the post-activation  $\mathbf{y}_\ell(\theta, x) = \text{ReLU}(\mathbf{z}_\ell(\theta, x)) \in \mathbb{R}^{N_\ell}$  where the rectified linear unit (ReLU) activation function,  $\text{ReLU}(t) = \max(t, 0)$ , is applied entrywise. Finally we define the realization of the network as the function  $\mathbf{R}_\theta : x \mapsto \mathbf{R}_\theta(x) := \mathbf{z}_L(\theta, x) = \mathbf{W}_L \mathbf{y}_{L-1}(\theta, x) + \mathbf{b}_L \in \mathbb{R}^{N_L}$ . When needed we will use neuronwise versions of these notations, e.g.  $y_\nu(\theta, x) = (\mathbf{y}_\ell(\theta, x))_\nu$  where  $\nu \in N_\ell$ . Note the general convention to denote scalar-valued quantities in plain font to distinguish them from quantities that can be vector-valued, which are generally denoted in bold.

**2.3. Invariance to permutation and scaling.** A well known fact [23] is that the realization of any ReLU-network is invariant to permutations and scalings of the parameter  $\theta$ . The invariance to permutations is not specific to ReLU-networks, while the scaling-invariance is due to the homogeneity of the ReLU:  $\text{ReLU}(\lambda \cdot) = \lambda \text{ReLU}(\cdot)$  for every  $\lambda > 0$  and is also valid for other variants such as the leaky-ReLU. While various definitions co-exist in the literature [24, 25, 26], it is convenient to focus first on the practical *per-neuron* rescaling equivalence [23] as stated below.

**Rescaling equivalence.** Let  $\nu \in H$  and  $\lambda_\nu > 0$ . A neuron-wise scaling multiplies the incoming weights and the bias of  $\nu$  by  $\lambda_\nu$ , and divides the outgoing weights by  $\lambda_\nu$ . It is formally defined as  $s_{\nu, \lambda_\nu} : \theta = (w, b) \mapsto \theta' = (w', b')$  where for every connection  $e \in E$ ,

$$(1) \quad \forall e \in E, \quad w'_e = \begin{cases} w_e \lambda_\nu & \text{if } e \in \bullet \rightarrow \nu \\ \frac{1}{\lambda_\nu} w_e & \text{if } e \in \nu \rightarrow \bullet \\ w_e & \text{otherwise,} \end{cases} \quad \forall \nu \in H, \quad b'_\nu = b_\nu \lambda_\nu.$$

Let  $\mathcal{S}$  be the set of neuron-wise scalings. We observe that neuron-wise rescalings commute and are invertible, the inverse of  $s_{\nu, \lambda_\nu}$  being  $s_{\nu, 1/\lambda_\nu}$ . Let  $\langle \mathcal{S} \rangle$  be the commutative group

generated by  $\mathcal{S}$ . Every  $s \in \langle \mathcal{S} \rangle$  can be uniquely represented as the composition

$$s = \bigcirc_{\nu \in H} s_{\nu, \lambda_\nu}$$

where the  $\lambda_\nu$  are strictly positive. Note that in this representation, every hidden neuron  $\nu$  is associated to exactly one neuron-wise rescaling  $\lambda_\nu$ .

**Definition 1.**  $\theta$  and  $\theta'$  are *rescaling equivalent* if there exists  $s \in \langle \mathcal{S} \rangle$  such that  $\theta' = s(\theta)$ . We then denote  $\theta \sim_S \theta'$ .

Notice that if  $\theta' \sim_S \theta$ , then the output biases are equal:  $\theta'_\eta = \theta_\eta$  for all  $\eta \in N_L$ .

**Fact 1.**  $\theta' \sim_S \theta$  if, and only if, there exists diagonal matrices  $\mathbf{\Lambda}_\ell \in \mathbb{R}^{N_\ell \times N_\ell}$  with positive entries,  $0 \leq \ell \leq L$  such that  $\mathbf{\Lambda}_0 = \mathbf{I}_{N_0}$ ,  $\mathbf{\Lambda}_L = \mathbf{I}_{N_L}$ , and for every layer  $1 \leq \ell \leq L$

$$(2) \quad \mathbf{W}'_\ell = \mathbf{\Lambda}_\ell \mathbf{W}_\ell \mathbf{\Lambda}_{\ell-1}^{-1} \text{ and } \mathbf{b}'_\ell = \mathbf{\Lambda}_\ell \mathbf{b}_\ell.$$

**Permutation equivalence.** Consider  $\pi := (\pi_1, \dots, \pi_\ell)$  where  $\pi_\ell \in \mathfrak{S}_{N_\ell}$  is a permutation of the  $\ell$ -th hidden layer (input and output layers are never permuted),  $1 \leq \ell \leq L-1$ . Denote  $\mathfrak{S}_G = \mathfrak{S}_{N_1} \times \dots \times \mathfrak{S}_{N_{L-1}}$  the group of all such tuples of permutations. One can define a natural action of the group  $\mathfrak{S}_G$  on parameterizations via  $\theta \mapsto \pi \circ \theta := \theta'$  where each weight matrix  $\mathbf{W}'_\ell$  is obtained from  $\mathbf{W}_\ell$  by permuting rows according to  $\pi_\ell$  and columns according to  $\pi_{\ell-1}$ , while bias vector  $\mathbf{b}'_\ell$  is a permuted version of  $\mathbf{b}_\ell$  according to  $\pi_\ell$ .

**Definition 2.** Two parameters  $\theta, \theta'$  are *permutation-equivalent* if, and only if, there exists  $\pi \in \mathfrak{S}_G$  such that  $\theta' = \pi \circ \theta$ . This is denoted  $\theta \sim_P \theta'$ .

The parameters are *permutation-scaling equivalent* if, and only if, there exists  $\theta''$  such that  $\theta \sim_S \theta'' \sim_P \theta'$ . This is denoted  $\theta \sim_{PS} \theta'$ .

The parameters are *scaling-permutation equivalent* if, and only if, there exists  $\theta''$  such that  $\theta \sim_P \theta'' \sim_S \theta'$ . This is denoted  $\theta \sim_{SP} \theta'$ .

**Fact 2.**  $\theta' \sim_{PS} \theta$  if, and only if,  $\theta' \sim_{SP} \theta$ , if and only if there exists diagonal matrices  $\mathbf{\Lambda}_\ell \in \mathbb{R}^{N_\ell \times N_\ell}$  with positive entries and permutation matrices  $\mathbf{\Pi}_\ell \in \mathbb{R}^{N_\ell \times N_\ell}$ ,  $0 \leq \ell \leq L$ , such that  $\mathbf{\Pi}_0 = \mathbf{\Lambda}_0 = \mathbf{I}_{N_0}$ ,  $\mathbf{\Pi}_L = \mathbf{\Lambda}_L = \mathbf{I}_{N_L}$ , and for every layer  $1 \leq \ell \leq L$

$$(3) \quad \mathbf{W}'_\ell = \mathbf{\Pi}_\ell \mathbf{\Lambda}_\ell \mathbf{W}_\ell \mathbf{\Lambda}_{\ell-1}^{-1} \mathbf{\Pi}_{\ell-1}^{-1} \text{ and } \mathbf{b}'_\ell = \mathbf{\Pi}_\ell \mathbf{\Lambda}_\ell \mathbf{b}_\ell.$$

As widely documented [23, 24, 25, 26], PS-equivalent parameters share their realization as proven, e.g., in [9][Lemma 1].

**Lemma 1.** For any  $\theta, \theta' \in \mathbb{R}^{E \cup \bar{H}}$ , if  $\theta' \sim_{PS} \theta$  then  $\mathbf{R}_{\theta'} = \mathbf{R}_\theta$ .

A natural question is to determine conditions for the *identifiability* of (the equivalence class up to scaling and permutation of)  $\theta$  from  $\mathbf{R}_\theta$ . To be more specific, we consider identifiability with respect to a family of parameters  $\Theta$ , from a set  $\mathcal{X}$ . A case of particular interest will be when  $\mathcal{X}$  is finite, in order to characterize whether  $\theta$  can be recovered (up to scaling and permutations) from finitely many samples of the network realization  $\mathbf{R}_\theta$ .

**Definition 3** (PS-identifiability). A parameter  $\theta \in \Theta \subseteq \mathbb{R}^{E \cup \bar{H}}$  is PS-identifiable with respect to  $\Theta$  from  $\mathcal{X} \subseteq \mathbb{R}^{N_0}$  if for every  $\theta' \in \Theta$ , the equality  $\mathbf{R}_\theta = \mathbf{R}_{\theta'}$  on  $\mathcal{X}$  implies  $\theta' \sim_{PS} \theta$ . When considering  $\mathcal{X} = \mathbb{R}^{N_0}$ ,  $\theta$  is simply said to be PS-identifiable with respect to  $\Theta$ . When considering  $\Theta = \mathbb{R}^{E \cup \bar{H}}$ ,  $\theta$  is simply said to be PS-identifiable (from  $\mathcal{X}$ ).

A trivial observation is that if all outgoing weights of a hidden neuron are zero, then the realization of the network is unchanged under arbitrary modifications of the incoming weights and of the bias of his neuron, hence the corresponding parameter  $\theta$  cannot be PS-identifiable with respect to  $\Theta = \mathbb{R}^{E \cup \bar{H}}$ . A similar phenomenon occurs if all incoming weights to a hidden neuron are zero. This motivates the definition of admissible parameters and proves Lemma 2 below.

**Definition 4.**  $\theta$  is admissible if for each hidden neuron  $\nu \in H$  we have  $\mathbf{w}_{\bullet \rightarrow \nu} \neq 0$  and  $\mathbf{w}_{\nu \rightarrow \bullet} \neq 0$ . Equivalently, every hidden neuron belongs to a full path with nonzero weights.

**Lemma 2.** If  $\theta$  is PS-identifiable with respect to  $\Theta = \mathbb{R}^{E \cup \bar{H}}$  then it is admissible.

**2.4. An invariant embedding of ReLU networks.** The invariance with respect to (permutations and) scalings (Lemma 1) calls for an invariant representation of equivalence classes of network parameters. A central tool is a representation  $\Phi(\theta)$  mapping a network parameter  $\theta \in \mathbb{R}^{E \cup \bar{H}}$  to a vector  $\Phi(\theta)$  in a space indexed by paths of the network,  $\mathbb{R}^{\mathcal{P}}$ .

Before going further let us formally introduce paths, as illustrated in Figure 1

**Definition 5.** The set  $\mathcal{P}_\ell$ ,  $0 \leq \ell \leq L$  (resp.  $\mathcal{Q}_\ell$ ,  $1 \leq \ell \leq L-1$ ) consists of all partial paths from any neuron  $\nu_\ell \in N_\ell$  to a neuron of the last (resp. penultimate) layer  $\nu_L \in N_L$  (resp.  $\nu_{L-1} \in N_{L-1}$ ). Any path  $p \in \mathcal{P}_\ell$  is written as a tuple  $p = (p_\ell, \dots, p_L)$  where each  $p_i \in V$  is a neuron. We say that  $p$  is a full path if  $\ell = 0$ , that is, if  $p$  connects the input and the output layers. We may write  $p = p_\ell \rightarrow p_{\ell+1} \rightarrow \dots \rightarrow p_L$ , as well as  $p = \mu \rightarrow q \rightarrow \nu$  where  $\mu = p_\ell \in N_\ell$ ,  $\nu = p_L \in N_L$  and  $q = (p_{\ell+1}, \dots, p_{L-1}) \in \mathcal{Q}_{\ell+1}$ .

We next introduce the representation  $\Phi(\cdot)$ , which presents some connections with previous work [12, 13, 23] while being more generic as detailed in the introduction.

**Definition 6.** Given  $\theta \in \mathbb{R}^{E \cup \bar{H}}$ , the value of a path is

$$(4) \quad \Phi_p(\theta) = \Pi_{e \in p} \theta_e, \text{ for each full path } p \in \mathcal{P}_0,$$

$$(5) \quad \Phi_p(\theta) = \theta_{p_\ell} \Pi_{e \in p} \theta_e, \text{ for } p = (p_\ell, \dots, p_L) \in \mathcal{P}_\ell, 1 \leq \ell \leq L.$$

For  $p \in \mathcal{P}_L$ ,  $p = (\eta)$  with  $\eta \in N_L$ ,  $\Phi_p(\theta) = \theta_{p_L} = b_\eta$  is the corresponding output bias.

Define  $\mathcal{P} := \cup_{\ell=0}^L \mathcal{P}_\ell$ . For any  $\theta \in \mathbb{R}^{E \cup \bar{H}}$  we define

$$(6) \quad \Phi(\theta) := (\Phi_p(\theta))_{p \in \mathcal{P}} \in \mathbb{R}^{\mathcal{P}}$$

**Remark 1.** To streamline notations we say that an edge  $e = \mu \rightarrow \nu \in E$  belongs to  $p$  and also write  $e \in p$  if there exists  $\ell \leq i \leq L-1$  such that  $\mu = p_i$  and  $\nu = p_{i+1}$ . Similarly, we choose to denote  $\nu \in p$  if (and only if) the path  $p$  starts from neuron  $\nu$ , i.e., when  $p = (p_\ell, \dots, p_L) \in \mathcal{P}_\ell$ , if  $p_\ell = \nu$ . With these notations, we can write  $\Phi(\theta) = \Pi_{i \in p} \theta_i$ .

This representation characterizes the classes of scaling-equivalent admissible parameters.



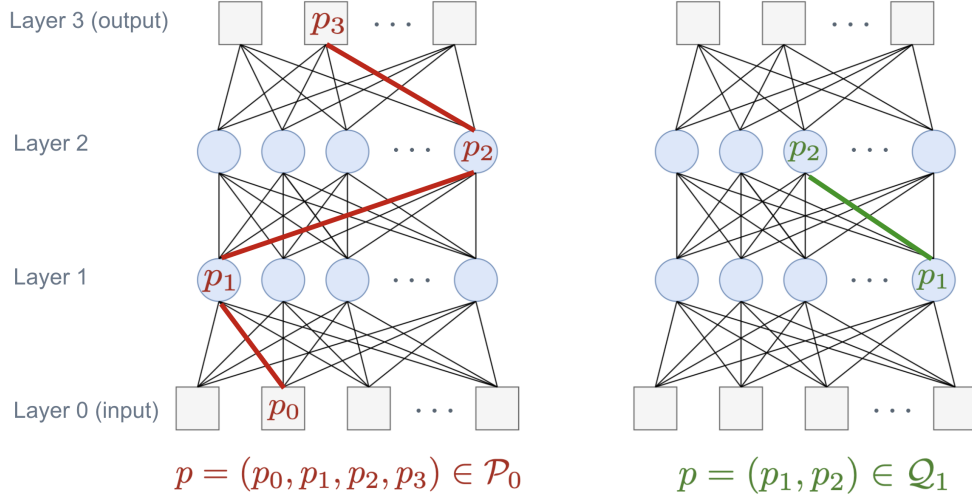


FIGURE 1. We consider a particular network architecture with  $L = 4$  layers (equivalently, with two hidden layers). Left: A particular path belonging to  $\mathcal{P}_0$ . Right: A particular path belonging to  $\mathcal{Q}_1$ .

**Theorem 1.** Consider any  $\theta', \theta \in \mathbb{R}^{E \cup \bar{H}}$ .

- a) Assume that  $\theta \sim_S \theta'$ . Then  $\Phi(\theta) = \Phi(\theta')$  and  $\text{sign}(\theta') = \text{sign}(\theta)$ .
- b) Assume that  $\theta$  is admissible, that  $\Phi(\theta') = \Phi(\theta)$ , and that  $\text{sign}(\theta'_E) = \text{sign}(\theta_E)$ . Then  $\theta \sim_S \theta'$  and  $\theta'$  is also admissible.

The proof is in Appendix A. A similar result is proven in [27, Theorem 3.3] without considering the biases and by replacing the condition on the signs by a condition on the activation statuses of all partial paths, which depend on the input variable  $x$  besides  $\theta$ .

**Remark 2.** The map  $\theta \mapsto \Phi(\theta)$  will be referred to as an embedding of network parameters. Stricto-sensu, as this map is not an injective function of network parameters, it does not match the definition of an embedding. However, since it characterizes equivalence classes of rescaling-equivalent admissible parameters, it can be used to define without ambiguity an embedding of these equivalence classes in  $\mathbb{R}^P$ .

**2.5. Some consequences of PS-identifiability.** Using the embedding  $\Phi(\cdot)$ , we show that if  $\theta$  is PS-identifiable then it is *locally identifiable up to scaling only*. Locality is measured in the sense of open balls  $B(\mathbf{c}, r) = \{\mathbf{c}' : \|\mathbf{c}' - \mathbf{c}\|_\infty < r\}$ , where the ambient linear space, equipped with the sup-norm, should always be clear from context.

**Definition 7** (local S-identifiability). Given  $\epsilon > 0$ , a parameter  $\theta \in \Theta \subseteq \mathbb{R}^{E \cup \bar{H}}$  is  $\epsilon$ -locally S-identifiable from  $\mathcal{X} \subset \mathbb{R}^{N_0}$  with respect to  $\Theta$ , if for every  $\theta' \in \Theta \cap B(\theta, \epsilon)$ , the identity  $\mathbf{R}_\theta = \mathbf{R}_{\theta'}$  on  $\mathcal{X}$  implies  $\theta' \sim_S \theta$ . If there exists  $\epsilon > 0$  such that  $\theta$  is  $\epsilon$ -locally S-identifiable

from  $\mathcal{X}$  then  $\theta$  is locally S-identifiable from  $\mathcal{X}$ . When  $\mathcal{X} = \mathbb{R}^{N_0}$  and/or  $\Theta = \mathbb{R}^{E \cup \bar{H}}$  we adopt the same simplified terminology as with the notion of PS-identifiability.

**Remark 3.** If  $\theta$  is PS-identifiable (resp. locally S-identifiable) from  $\mathcal{X} \subseteq \mathbb{R}^{N_0}$  with respect to  $\Theta \subseteq \mathbb{R}^{E \cup \bar{H}}$  then the same holds from any  $\mathcal{X}' \supseteq \mathcal{X}$  with respect to any  $\Theta' \subseteq \Theta$ .

Our first result is the following theorem.

**Theorem 2.** Consider  $\Theta \subseteq \mathbb{R}^{E \cup \bar{H}}$  and  $\mathcal{X} \subset \mathbb{R}^{N_0}$ . If  $\theta \in \Theta$  is admissible and PS-identifiable from  $\mathcal{X}$  with respect to  $\Theta$  then it is locally S-identifiable from  $\mathcal{X}$  with respect to  $\Theta$ .

The proof is in Appendix B and uses the embedding  $\Phi(\cdot)$ . By Lemma 2, PS-identifiability with respect to  $\Theta = \mathbb{R}^{E \cup \bar{H}}$  implies admissibility. Considering any  $\Theta$  with a similar property, a direct corollary of Theorem 2 is that PS-identifiability with respect to  $\Theta$  implies local S-identifiability with respect to  $\Theta$ .

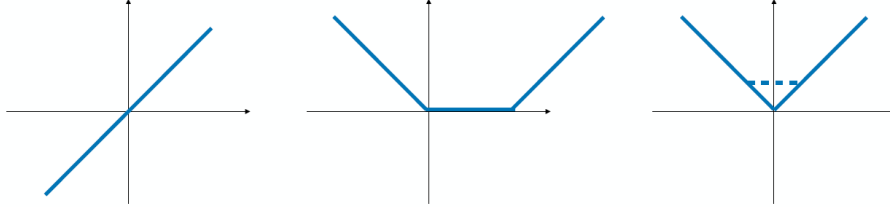


FIGURE 2. Realizations of networks from (a) Example 1; (b) Example 2 ; (c) Example 4

An example shows that indeed, local S-identifiability depends on the constraint set  $\Theta$ .

**Example 1** (see Figure 2-(a)). On a shallow network architecture with two hidden neurons  $\nu_1, \nu_2$ , the identity  $\text{id} : \mathbb{R} \rightarrow \mathbb{R}, x \mapsto x$  can be written as  $x = \text{ReLU}(x-t) - \text{ReLU}(-(x-t)) + t = \mathbf{R}_{\theta_t}$  with  $\theta_t = (w_{\mu \rightarrow \nu_1} = 1, w_{\mu \rightarrow \nu_2} = -1, b_{\nu_1} = -t, b_{\nu_2} = t, w_{\nu_1 \rightarrow \eta} = 1, w_{\nu_2 \rightarrow \eta} = -1, b_\eta = t)$  for every  $t \in \mathbb{R}$  ( $\mu$  is the input neuron,  $\eta$  the output neuron). Since  $\theta_t$  and  $\theta_{t'}$ ,  $t \neq t'$  have different output bias, they are not PS-equivalent. This shows that, e.g.,  $\theta_0$  is not locally S-identifiable with respect to  $\Theta = \mathbb{R}^{E \cup \bar{H}}$ . With respect to the set  $\Theta$  of networks without output bias ( $b_\eta = 0$ ), as detailed in Example 5,  $\theta_0$  becomes PS-identifiable from  $\mathcal{X} = \mathbb{R}$ .

The above example includes two neurons which are *twins* in the following sense.

**Definition 8** (Twin neurons). Consider a parameter  $\theta$  on a network architecture of any depth. Two hidden neurons  $\nu \neq \nu'$  from the same layer are said to be *twins* if there exists  $\lambda \in \mathbb{R}$  such that  $(\mathbf{w}_{\bullet \rightarrow \nu}, b_\nu) = \lambda(\mathbf{w}_{\bullet \rightarrow \nu'}, b_{\nu'})$ . If  $\theta$  is admissible then necessarily  $\lambda \neq 0$ , and  $\nu, \nu'$  are said to be *positive twins* if  $\lambda > 0$ , *negative twins* otherwise.

NB: Even though each hidden neuron  $\nu \in H$  is (positive) twin to itself, such a neuron is abusively said to have “no twin” if it is not twin with any  $\nu' \neq \nu$  from the same layer. We also say that  $\theta$  has no twins if none of its neurons have any twin.

Intuitively, if  $\nu, \nu'$  are twins then the corresponding pre-activation functions  $z_\nu(\theta, \cdot)$   $z_{\nu'}(\theta, \cdot)$  are collinear, and the resulting post-activation functions,  $y_\nu(\theta, \cdot)$   $y_{\nu'}(\theta, \cdot)$  are also collinear for positive twins. For negative twins, there exists linear combinations of the post-activations that are simply proportional to the pre-activations, somehow bypassing the effect of the ReLU nonlinearity. As proved in Appendix C, twins *always* prevent identifiability with respect to  $\Theta = \mathbb{R}^{E \cup \bar{H}}$ .

**Lemma 3.** *Consider  $\theta \in \Theta = \mathbb{R}^{E \cup \bar{H}}$ .*

- a) *Assume that  $\theta$  is locally S-identifiable with respect to  $\Theta$ .  
Then  $\theta$  has no positive twins.*
- b) *Assume that  $\theta$  is PS-identifiable from some bounded set  $\mathcal{X} \subseteq \mathbb{R}^{N_0}$  with respect to  $\Theta$ .  
Then  $\theta$  has no twins.*

We will see in Example 4 (in section 5) that the absolute value function (see Figure 2-(c)) is the realization of a shallow network with two hidden neurons that are negative twins, yet it is PS-identifiable (hence locally S-identifiable) with respect to  $\Theta = \mathbb{R}^{E \cup \bar{H}}$ . It is even locally S-identifiable from some finite set  $F \subseteq \mathbb{R}$ . Of course, by Lemma 3 such a network cannot be PS-identifiable from any bounded set with respect to  $\Theta = \mathbb{R}^{E \cup \bar{H}}$ .

Twins are a form of *local* degeneracy. For shallow networks, we will show that this is the only form of local degeneracy (see the upcoming Lemma 5 and Theorem 3), but we will see other forms for deeper networks (see Example 3). As illustrated next, there are also *non-local* degeneracies that can prevent identifiability.

**Example 2** (see Figure 2-(b)). *The function*

$$f(x) = \begin{cases} -x, & \text{if } x \leq 0 \\ 0, & \text{if } 0 \leq x \leq 1 \\ x - 1, & \text{if } x \geq 1 \end{cases}$$

*satisfies  $f(x) = \text{ReLU}(-x) + \text{ReLU}(x - 1) = \text{ReLU}(x) + \text{ReLU}(-(x - 1)) - 1$ . It is thus the realization of  $\theta = (w_{\mu \rightarrow \nu_1} = -1, w_{\mu \rightarrow \nu_2} = 1, b_{\nu_1} = 0, b_{\nu_2} = -1, w_{\nu_1 \rightarrow \eta} = w_{\nu_2 \rightarrow \eta} = 1, b_\eta = 0$ , but also of  $\theta' = (w'_{\mu \rightarrow \nu_1} = 1, w'_{\mu \rightarrow \nu_2} = -1, b'_{\nu_1} = 0, b'_{\nu_2} = 1, w'_{\nu_1 \rightarrow \eta} = w'_{\nu_2 \rightarrow \eta} = 1, b'_\eta = -1$ , which are not PS-equivalent since  $b_\eta \neq b'_\eta$ . Yet the theory we establish (see Lemma 5) shows that  $\theta$  and  $\theta'$  are both locally S-identifiable from some finite set  $F \subset \mathbb{R}$ .*

It turns out that the above example fails to be *irreducible* as we formalize next.

**Definition 9** (Irreducibility). *A parameter  $\theta$  is irreducible if for each hidden layer  $1 \leq \ell \leq L - 1$  and non-empty subset  $T \subset N_\ell$  we have*

$$(7) \quad \mathbf{W}_{\ell+1} \mathbf{I}_T \mathbf{W}_\ell \neq 0, \quad \text{with } \mathbf{I}_T = \text{diag}(\chi_T),$$

*with  $\chi_T \in \{0, 1\}^{N_\ell}$  the indicator function of  $T$ :  $(\chi_T)_\nu = 1$  if, and only if,  $\nu \in T$ . We denote  $\Theta_{\text{irr}} \subset \mathbb{R}^{E \cup \bar{H}}$  the set of all irreducible parameters.*

**Fact 3.** *Each irreducible parameter is also admissible.*

In fact, as established in Appendix D, any PS-identifiable parameter with no twins must be irreducible.

**Lemma 4.** *If  $\theta$  is PS-identifiable from  $\mathcal{X} \subseteq \mathbb{R}^{N_0}$  with respect to  $\Theta = \mathbb{R}^{E \cup \bar{H}}$  and has no twin, then it is irreducible.*

In particular, in light of Lemma 3, every PS-identifiable parameter is irreducible. In the shallow case, a direct consequence of irreducibility can be obtained using an “algebraic” expression of the realization  $\mathbf{R}_\theta$  (Lemma 9 in section 4): for every input vector  $x$  where  $\mathbf{R}_\theta$  is differentiable, the Jacobian of  $\mathbf{R}_\theta$  is given by  $\mathbf{W}_2 \mathbf{I}_1 \mathbf{W}_1$  with  $\mathbf{I}_1 = \text{diag}(\mathbf{a}_1(\theta, x))$ . Irreducibility thus implies that this Jacobian can only vanish if  $\mathbf{a}_1(\theta, x) = \mathbf{0}$ , i.e., if all neurons are inactive. As illustrated on Example 2 (see Figure 2-(b)) this however does not *characterize* irreducibility, and an intuitive characterization of irreducibility in terms of simple properties of  $\mathbf{R}_\theta$  is left to future work.

**2.6. Identifiability conditions in the shallow case.** For shallow neural networks, we prove that admissible parameters with no twins are locally S-identifiable from a *finite* set. Such results resonate with previous work on the identifiability of shallow networks equipped with various activation functions other than the ReLU [4, 5, 6, 7].

**Lemma 5.** *Consider a shallow architecture. If  $\theta$  is admissible with no twins, then there is a finite  $\mathcal{X} \subseteq \mathbb{R}^{N_0}$  with  $\text{card}(\mathcal{X}) \leq (|N_0| + 1)(|N_1| + 1)$  from which  $\theta$  is locally S-identifiable with respect to  $\Theta = \mathbb{R}^{E \cup \bar{H}}$ .*

The proof is in Section 5.2. Combined with irreducibility, the absence of twins is further shown to be equivalent to PS-identifiability from a *bounded* set. Whether this is also equivalent to PS-identifiability from a *finite* set is left to future work, as well as a possible explicit control of the cardinality of such a finite set.

**Theorem 3.** *Consider a shallow network architecture. The following are equivalent*

- a) *there is a bounded  $\mathcal{X} \subseteq \mathbb{R}^{N_0}$  from which  $\theta$  is PS-identifiable with respect to  $\Theta = \mathbb{R}^{E \cup \bar{H}}$ ;*
- b)  *$\theta$  has no twins and is irreducible.*

*Proof.* The implication a  $\Rightarrow$  b is a consequence of Lemma 3 and Lemma 4. The converse b  $\Rightarrow$  a follows by Theorem 6-b in Section 5.3.  $\square$

As established with Theorem 6-a in Section 5.3, the architecture itself is identifiable in the following sense for irreducible parameters with no twins.

**Theorem 4.** *Consider two shallow network architectures with the same input and output layers,  $N_0$  and  $N_2$ , and potentially distinct hidden layer  $H = N_1$ ,  $H' = N'_1$ . Let  $\theta, \theta'$  be parameters on each architecture. Assume that  $\theta$  is irreducible with no twins, and that  $\theta'$  is admissible with no twins. If  $\mathbf{R}_\theta = \mathbf{R}_{\theta'}$  on  $\mathbb{R}^{N_0}$  then  $\text{card}(N_1) = \text{card}(N'_1)$  and  $\theta' \sim_{PS} \theta$ .*

As illustrated by Example 4 in section 5, there are also shallow networks that are PS-identifiable from  $\mathcal{X} = \mathbb{R}^{N_0}$  but not from any bounded set. They are of course irreducible by Lemma 4, and have no positive twin by Lemma 3, but they have one or more pair of negative twins.

**2.7. A glimpse at the analysis of local identifiability.** Much of the local identifiability analysis, which is conducted in details in section 4, relies on an important property of the embedding  $\Phi$  (besides its ability to characterize scaling equivalence, see Theorem 1): it provides a *locally linear parameterization* of the realization of the network, in the sense that given  $\theta$  and for “most”  $x \in \mathbb{R}^{N_0}$  we have, for every  $\theta'$  in a (small enough) neighborhood  $\theta$

$$(8) \quad \mathbf{R}_{\theta'}(x) - \mathbf{R}_{\theta}(x) = \mathbf{C}_{\theta,x} \cdot (\Phi(\theta') - \Phi(\theta))$$

with  $\mathbf{C}_{\theta,x} \in \mathbb{R}^{N_L \times \mathcal{P}}$  some linear operator that is *independent of  $\theta'$* . This property holds provided  $x$  is a point where the gradient of  $\mathbf{R}_{\theta}$  (and of all pre-activations at intermediate hidden layers) is well-defined and continuous, which motivates the following definition.

**Definition 10.** *Consider any network architecture. Given a parameter  $\theta$  we define for each hidden neuron  $\nu \in H$  the set  $\Gamma_{\nu}(\theta)$  of input vectors where  $z_{\nu}(\theta, x) = 0$  and the gradient  $\nabla z_{\nu}(\theta, x)$  is well-defined and nonzero,*

$$(9) \quad \Gamma_{\nu}(\theta) := \{x \in \mathbb{R}^{N_0} : z_{\nu}(\theta, x) = 0 \text{ and } \nabla z_{\nu}(\theta, x) \neq 0\}.$$

We define  $\mathcal{X}_{\theta} \subseteq \mathbb{R}^{N_0}$  as the complement to  $\cup_{\nu \in H} \Gamma_{\nu}(\theta)$ .

Definition 10 is extremely close to the definition of *Bent Hyperplanes* [28] (except that we add the non-nullity condition on the gradient). Informally, and as previously stated [29, 30, 31], bent hyperplanes separate the input space into *linear regions* where the realization of the network  $x \mapsto \mathbf{R}_{\theta}(x)$  is affine, see Figure 3 for an illustration.

For our needs, we will provide in Lemma 11 an alternate characterization of  $\mathcal{X}_{\theta}$  which we have not found elsewhere in the literature. It will be used in Corollary 3 to formalize Property (8) for  $x \in \mathcal{X}_{\theta}$ , which motivates the definition of non-degenerate parameters.

**Definition 11** (Non-degeneracy). *Consider the finite dimensional linear space*

$$(10) \quad \mathbf{V}(\theta) := \cap_{x \in \mathcal{X}_{\theta}} \ker(\mathbf{C}_{\theta,x}) \subseteq \mathbb{R}^{\mathcal{P}}.$$

*A parameter  $\theta \in \Theta \subseteq \mathbb{R}^{E \cup \bar{H}}$  is  $\epsilon$ -non-degenerate with respect to  $\Theta$ , where  $\epsilon > 0$ , if it is admissible and for every  $\theta' \in \Theta \cap B(\theta, \epsilon)$  we have*

$$(11) \quad \Phi(\theta') - \Phi(\theta) \in \mathbf{V}(\theta) \Rightarrow \Phi(\theta') = \Phi(\theta).$$

*It is non-degenerate with respect to  $\Theta$  if there exists  $\epsilon > 0$  such that it is  $\epsilon$ -non-degenerate with respect to  $\Theta$ .*

Exploiting the fact that all considered spaces are finite dimensional, we characterize the space  $\mathbf{V}(\theta)$  in terms of certain *activation spaces* (Definition 14) and prove that non-degeneracy is equivalent (see Theorem 5, the main result of section 4) to the existence of some *finite set*  $F \subset \mathcal{X}_{\theta}$  such that  $\theta$  is locally S-identifiable from  $F$  (hence also locally S-identifiable from  $\mathcal{X} = \mathbb{R}^{N_0}$ ). The cardinality of  $F$  is bounded from above using the dimension of activation spaces.

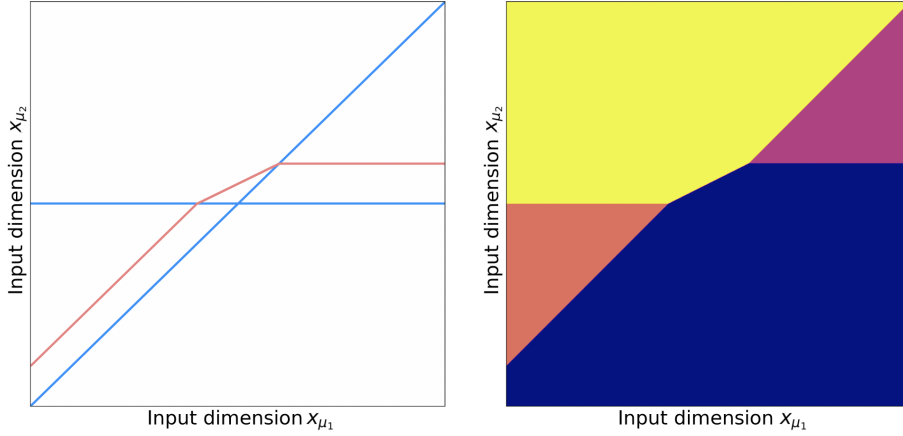


FIGURE 3. We consider a network architecture with  $|N_1| = 2$  neurons on the first hidden layer and  $|N_2| = 1$  neuron on the second hidden layer. The input  $x$  is two-dimensional:  $|N_0| = 2$  and the output is scalar:  $|N_4| = 1$ . Left: bent hyperplanes for the first hidden layer,  $\Gamma_\nu(\theta)$ ,  $\nu \in N_1$  (blue) and second hidden layer  $\Gamma_\nu(\theta)$ ,  $\nu \in N_2$  (red). Right: linear regions. All the weights and biases were initialized randomly. The figures are generated with a PyTorch script available at <https://github.com/pierrestock/linear-regions/blob/main/partition.ipynb>.

**2.8. Non-degeneracy and irreducibility in shallow *vs* deeper architectures.** An easy sufficient condition for non-degeneracy is to have a trivial space  $V(\theta) = \{0\}$ . For *scalar-valued* shallow networks ( $L = 2, |N_L| = 1$ ), we prove (cf Lemma 17 and Corollary 2 that non-degeneracy with respect to  $\Theta = \mathbb{R}^{E \cup H}$  is in fact *equivalent* to  $V(\theta) = \{0\}$ , and for shallow (possibly vector-valued) networks, the latter is proved to hold if, and only if, there are no twins (by Corollary 2 and Lemma 15). In light of Theorem 3, when combined with irreducibility, the fact that  $V(\theta) = \{0\}$  thus becomes equivalent (for shallow networks) to the PS-identifiability of  $\theta$  from some bounded set.

For networks of depth  $L \geq 3$ , any parameter such that  $V(\theta) = \{0\}$  is of course still non-degenerate (hence locally S-identifiable from a finite set, by Theorem 5), but this property is no longer equivalent to the absence of twins: further conditions between layers are required, as illustrated by the following example.

**Example 3.** In Figure 4, we exhibit a two-hidden-layer architecture valued with a parameter  $\theta$  that presents no twin neurons (see Definition 8) but such that  $\theta$  is not locally S-identifiable (see Definition 7).

Characterizing concrete conditions ensuring  $V(\theta) = \{0\}$  is left to future work. A particular challenge is to understand whether the condition  $V(\theta) = \{0\}$ , combined with (a possibly strengthened version of) irreducibility remains equivalent to PS-identifiability from a bounded set. We note that irreducibility in the shallow case is reminiscent of [32, Equation

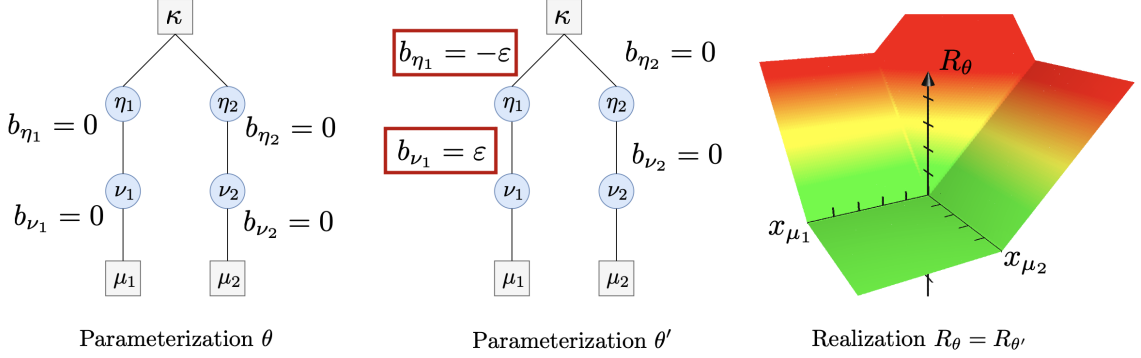


FIGURE 4. A network with two hidden layers that is not locally S-identifiable, while having no twin hidden neurons. Weights  $w_e$  and  $w'_e$ ,  $e \in E$  are set to one on the displayed edges and to zero on other edges, and are not depicted here for readability. Left: parameterization  $\theta$  valuing the architecture. Center: alternative parameterization  $\theta'$  such that  $\theta$  and  $\theta'$  are *not* rescaling equivalent. Right: the two realizations  $R_\theta$  and  $R_{\theta'}$  coincide: for every input point  $x = (x_{\mu_1}, x_{\mu_2}) \in \mathbb{R}^2$ ,  $R_\theta(x) = R_{\theta'}(x)$ . The construction is valid for arbitrary  $\varepsilon > 0$ .

(8)], a condition used to define so-called “general ReLU networks” to provide sufficient identifiability conditions in deeper settings. This may serve as a guide to identify stronger notions of irreducibility for deep networks. Preliminary investigations suggest that certain tensor products of activation vectors play a role when analyzing non-degeneracy. This is reminiscent of the tools studied by Fornasier *et al.* [10] with two hidden layers  $L = 3$  in a smooth context that cannot cover ReLU networks.

**2.9. Discussion.** Before diving into the technical contributions in the next Sections, we discuss some topics of interest for the reader that are mostly out of the scope of this work. We refer the reader to Figure 5 for a brief summary of the results proven in this paper.

*Local identifiability and optimization.* First, we argue that studying *local* (instead of global) S-identifiability is of practical interest, as discussed e.g. in [12, 13]. Indeed, neural networks are traditionally optimized with a variant of stochastic gradient descent, or SGD [33]. Hence, (1) during training, the optimization yields parameters that are close to the previous ones and (2) the parameters obtained after convergence can be expected to be locally optimal up to natural permutation and rescaling equivalences.

*Identifiability from a finite set.* Since we are mainly interested in the problem of recovering (the equivalence class of)  $\theta$  from the knowledge of its realization  $\mathbf{R}_\theta$ , we list below some questions calling for extensions of Theorem 3. Indeed, it is not always possible to recover  $\theta$  from its realization. Even when such a recovery is theoretically possible, it may

involve having full access to the function  $\mathbf{R}_\theta$ , which is not a concrete input to provide to any reconstruction algorithm. A more practical question is: when can we recover (the equivalence class of)  $\theta$  from the knowledge of *finitely many samples*  $\mathbf{R}_\theta(x_i), 1 \leq i \leq n$ ? When there exists a choice (that may depend on  $\theta$ ) of  $n$  and  $x_i, 1 \leq i \leq n$  such that this is feasible, we also get as a byproduct a reconstruction of  $\mathbf{R}_\theta$  from the sole knowledge of its samples at these points. Hence, another question of interest is: when can the function  $\mathbf{R}_\theta$  be identified from the knowledge of finitely many of its samples? This is possibly less demanding, as here it is not required to be able to reconstruct (the equivalence class of)  $\theta$  from its realization. In both cases, since  $\theta$  is not known beforehand, it is important to ensure that the choice of the sampling set is algorithmically feasible, for example if it is done iteratively at least the first sample must be chosen without any knowledge on  $\theta$  or  $\mathbf{R}_\theta$ . Of course, answers to these questions lead to further ones, that we do not touch upon: if  $\theta$  can be identified from finitely many samples, how many samples are sufficient<sup>4</sup> (resp. necessary)? Can we explicit an scheme (possibly randomized) to choose these samples? Can we explicit an algorithm to perform reconstruction? How stable is it to inaccuracies in the evaluation of  $\mathbf{R}_\theta(x_i)$  or to the knowledge of  $x_i$ ?

*Reverse-engineering ReLU networks.* Here, we discuss the work of Rolnick and Kording [9] more extensively than what was done in the Introduction. The goal is to position our work with respect to this interesting work. The authors present a sampling algorithm to recover a ReLU network’s architecture and parameters, up to permutations and rescalings. The authors prove that their algorithm terminates except for a measure-zero set of networks and do not provide the complexity of their method in terms of number of the samples needed to recover  $\mathbf{R}_\theta$ , except for recovering the first layer’s parameters. They reason in terms of so-called *activation* and *linear* regions [28] and make the following assumptions. Recall that the sets  $\Gamma_\nu(\theta)$  are introduced in Definition 9 for every hidden neuron  $\nu$ .  $\Gamma_\nu(\theta)$  is often called the *separating* or *bent* hyperplane for neuron  $\nu$ .

- (1) Linear Regions assumption as stated by the authors: “Each [activation]<sup>5</sup>region represents a maximal connected component of input space on which the [realization  $\mathbf{R}_\theta$ ] is given by a single linear function”. In other words, the authors assume that activation regions and linear regions coincide (Section 3.2 in the original paper).
- (2) All the sets  $\Gamma_\nu(\theta)$  for  $\nu \in H$  have codimension 1<sup>6</sup> hence the name *separating hyperplane* (implicitly assumed, see in particular the first paragraph of Section 3.3).
- (3) For every hidden neuron  $\nu$  in layer  $1 \leq \ell \leq L - 1$ ,  $\Gamma_\nu(\theta)$  intersects *all* the sets  $\Gamma_{\nu'}(\theta)$  for all neurons  $\nu'$  in a previous layer  $1 \leq \ell' < \ell \leq L - 1$  (Section 5.2 in the original paper).
- (4) For  $\nu \neq \nu'$  such that  $\nu$  belong to layer  $\ell$  and  $\nu'$  belongs to layer  $\ell' < \ell$ , “ $\Gamma_\nu(\theta)$  bends on  $\Gamma_{\nu'}(\theta)$ , but  $\Gamma_\nu(\theta)$  and  $\Gamma_{\nu'}(\theta)$  cannot both bend at their intersection” (implicitly assumed, see in particular the first paragraph of Section 3.3).

<sup>4</sup>Lemma 5 partly answers this question regarding local S-identifiability for shallow networks.

<sup>5</sup>What the authors denote as linear regions are in fact known as activation regions, see [28].

<sup>6</sup>This prevents the case where  $\nabla R_\theta(x) = 0$  for  $x \in B(x_0, r)$



- (5) For every hidden neuron  $\nu \in H$ ,  $\Gamma_\nu(\theta)$  is *not bounded* and *not disconnected* (Section 5.2 in the original paper).

According to the authors, parameters  $\theta$  that do not satisfy at least one of these assumptions constitute a measure-zero set of networks, hence the authors discard these cases from their analysis. In the remainder of this paper, we aim at more precisely characterizing this measure-null zero set. This is fully done in the shallow case, and the developed tools should be instrumental when pursuing this mathematical study in deeper settings.

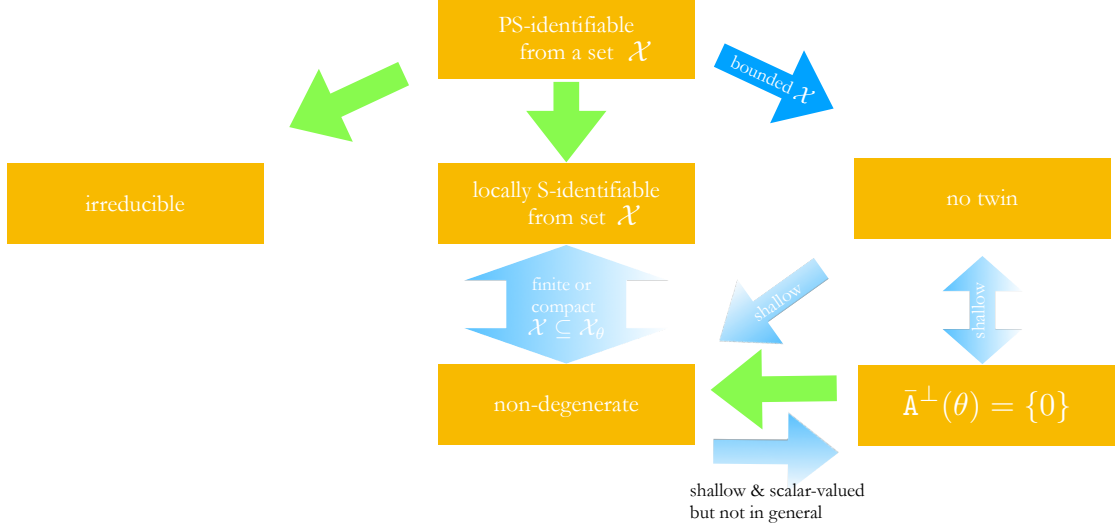


FIGURE 5. Summary of the various results proven in the paper. The space  $\bar{\mathbf{A}}^\perp(\theta)$  is defined in Definition 14. Theorem 3 further establishes that irreducibility and the absence of twins imply PS-identifiability from a bounded set in the shallow case.

### 3. RESCALING INVARIANCE OF THE EMBEDDING

The proof of the main property of the embedding  $\Phi(\cdot)$ , Theorem 1, exploits linear operators related to  $\Phi(\cdot)$ . The following definition is motivated by the obvious observation that, if  $\theta$  has positive entries  $\theta_i = e^{\alpha_i}$ ,  $i \in E \cup \bar{H}$ , then  $\Phi(\theta) = e^{\mathbf{P}\alpha}$  where the exponential is taken componentwise. This is related to the idea of updating weights *multiplicatively*, which is exploited in particular by Bernstein [34] to investigate learning stability.

**Definition 12.** Consider  $\mathbf{P} : \mathbb{R}^{E \cup \bar{H}} \rightarrow \mathbb{R}^{\mathcal{P}}$  the linear operator defined for  $\mathbf{u} \in \mathbb{R}^{E \cup \bar{H}}$  as

$$(12) \quad (\mathbf{P}\mathbf{u})_p := \begin{cases} \sum_{e \in p} u_e, & \text{for each full path } p \in \mathcal{P}_0; \\ u_{p_\ell} + \sum_{e \in p} u_e, & \text{for each partial path } p = (p_\ell, \dots, p_L) \in \mathcal{P}_\ell, \ 1 \leq \ell \leq L. \end{cases}$$

With the notations from Remark 1 we can also write  $(\mathbf{P}\mathbf{u})_p = \sum_{i \in p} u_i$ .

Before proving Theorem 1 we express a few technical lemmas.

**Lemma 6.** *For every  $\theta \in \mathbb{R}^{E \cup \bar{H}}$ , with  $\text{supp}(\Phi(\theta)) = \{p \in \mathcal{P} : \Phi_p(\theta) \neq 0\}$  we have*

$$(13) \quad \{i \in E \cup \bar{H} : \exists p \in \text{supp}(\Phi(\theta)), p \ni i\} \subseteq \{i \in E \cup \bar{H} : \theta_i \neq 0\} = \text{supp}(\theta).$$

*If  $\theta \in \mathbb{R}^{E \cup \bar{H}}$  is admissible then we further have*

$$(14) \quad \text{supp}(\theta) = \{i \in E \cup \bar{H} : \exists p \in \text{supp}(\Phi(\theta)), p \ni i\}.$$

*Proof.* For each path  $p \in \mathcal{P}$  denote  $I_p = \{i \in E \cup \bar{H} : i \in p\}$  and observe first that the left hand side in (13) is  $\cup_{p \in \text{supp}(\Phi(\theta))} I_p$ . Consider  $p \in \text{supp}(\Phi(\theta))$ . Since  $\Phi_p(\theta) = \prod_{i \in p} \theta_i$ , we have  $\theta_i \neq 0$  for each  $i \in I_p$ , i.e.,  $I_p \subseteq \text{supp}(\theta)$ . As this holds for every  $p \in \text{supp}(\Phi(\theta))$  we obtain  $\cup_{p \in \text{supp}(\Phi(\theta))} I_p \subseteq \text{supp}(\theta)$ . This establishes (13).

Assuming now that  $\theta$  is admissible, consider  $i \in \text{supp}(\theta)$  and distinguish three cases. If  $i = \eta \in N_L$  is an output neuron, then  $p = (\eta) \ni i$  yields  $\Phi_p(\theta) = \theta_\eta \neq 0$ . If  $i = \nu \in H$  is a hidden neuron, then since  $\theta$  is admissible there is a path  $p \ni i$  with nonzero weights connecting  $\nu$  to an output neuron. This path satisfies  $\Phi_p(\theta) \neq 0$ . Finally, if  $i = \nu \rightarrow \nu'$  is an edge, then since  $\theta$  is admissible there is a path connecting the input layer to  $\nu$  and a path connecting  $\nu'$  to the output layer, both with nonzero weights. Concatenating them yields a path  $p \ni i$  such that  $\Phi_p(\theta) \neq 0$ . In all cases, we obtain the existence of a path  $p \in \text{supp}(\Phi(\theta))$  such that  $p \ni i$ . This establishes (14).  $\square$

**Corollary 1.** *Consider  $\theta, \theta' \in \mathbb{R}^{E \cup \bar{H}}$  such that  $\Phi(\theta') = \Phi(\theta)$ . If  $\theta$  is admissible then  $\text{supp}(\theta') = \text{supp}(\theta)$  and  $\theta'$  is also admissible.*

*Proof.* By Lemma 6 and the equality  $\Phi(\theta') = \Phi(\theta)$  we have

$$\begin{aligned} \text{supp}(\theta) &= \{i \in E \cup \bar{H} : \exists p \in \mathcal{P}, \Phi_p(\theta) \neq 0, i \in p\} \\ &= \{i \in E \cup \bar{H} : \exists p \in \mathcal{P}, \Phi_p(\theta') \neq 0, i \in p\} \subseteq \text{supp}(\theta'). \end{aligned}$$

The fact that  $\theta$  is admissible is a property of its support, and the inclusion  $\text{supp}(\theta) \subseteq \text{supp}(\theta')$  implies that  $\theta'$  is also admissible. It follows using Lemma 6 again that the rightmost inclusion above is an equality.  $\square$

**Lemma 7.** *Given  $\theta \in \mathbb{R}^{E \cup \bar{H}}$  an admissible parameter, consider the spaces*

$$(15) \quad W_\theta := \{\alpha \in \mathbb{R}^{E \cup \bar{H}}, [\Phi(\theta) \odot \mathbf{P}\alpha]_{\mathcal{P}_0} = 0\}$$

$$(16) \quad V_\theta := \{\alpha \in W_\theta, \alpha_{\bar{H}} = 0, \text{supp}(\alpha) \subseteq \text{supp}(\theta)\}.$$

*Given  $\alpha \in W_\theta$ , define for each hidden neuron  $\nu \in H$*

$$(17) \quad (\mathbf{S}_\theta \alpha)_\nu \triangleq - \sum_{e \in p} \alpha_e$$

*with  $p$  any path with edges  $e \in E \cap \text{supp}(\theta)$  joining  $\nu$  to an output neuron  $\eta$ .*

*a) The linear map  $\mathbf{S}_\theta: W_\theta \rightarrow \mathbb{R}^H$  is well-defined and independent of the choice of  $p$  and  $\eta$ ;*

b) Its restriction  $\mathbf{S}_\theta : V_\theta \rightarrow \mathbb{R}^H$  is an isomorphism. Its inverse  $\mathbf{S}_\theta^{-1} : \mathbb{R}^H \rightarrow V_\theta$  is such that for any  $\beta \in \mathbb{R}^H$ ,  $\mathbf{S}_\theta^{-1}\beta = \alpha$  where  $\alpha_{\bar{H}} = 0$  and for each edge  $e = \mu \rightarrow \nu \in E \cap \text{supp}(\theta)$ ,

$$(18) \quad \alpha_e \triangleq \begin{cases} -\beta_\mu & \text{if } \mu \in N_{L-1} \text{ (and } \nu \in N_L) \\ \beta_\nu - \beta_\mu & \text{if } \mu \in N_\ell, 1 \leq \ell \leq L-2 \\ \beta_\nu & \text{if } \mu \in N_0. \end{cases}$$

while  $\alpha_e = 0$  for each  $e \in E \setminus \text{supp}(\theta)$ .

The proof is postponed to Appendix A to keep the reading flow.

*Proof of Theorem 1.* By Definition 1,  $\theta \sim_S \theta'$  if, and only if, there are  $\{\lambda_\nu\}_{\nu \in H \cup N_0 \cup N_L}$  such that

$$(19) \quad \lambda_\nu > 0, \quad \forall \nu \in H, \quad \text{and } \lambda_\nu = 1, \quad \forall \nu \in N_0 \cup N_L$$

$$(20) \quad \theta'_e = \lambda_\mu^{-1} \theta_e \lambda_\nu, \quad \forall e = \mu \rightarrow \nu \in E \quad \text{and } \theta'_\nu = \theta_\nu \lambda_\nu, \quad \forall \nu \in \bar{H}.$$

Thus, if  $\theta' \sim_S \theta$  then  $\text{sign}(\theta') = \text{sign}(\theta)$ , and for every path  $p = (p_0, \dots, p_L) \in \mathcal{P}_0$  we get

$$\Phi_p(\theta') = \prod_{k=0}^{L-1} \theta'_{p_k \rightarrow p_{k+1}} = \prod_{k=0}^{L-1} (\lambda_{p_k}^{-1} \theta_{p_k \rightarrow p_{k+1}} \lambda_{p_{k+1}}) = \prod_{k=0}^{L-1} \theta_{p_k \rightarrow p_{k+1}} = \Phi_p(\theta),$$

while for  $p = (p_\ell, \dots, p_L) \in \mathcal{P}_\ell$ ,  $1 \leq \ell \leq L$

$$\Phi_p(\theta') = \theta'_{p_\ell} \prod_{k=\ell}^{L-1} \theta'_{p_k \rightarrow p_{k+1}} = \theta_{p_\ell} \lambda_{p_\ell} \prod_{k=\ell}^{L-1} (\lambda_{p_k}^{-1} \theta_{p_k \rightarrow p_{k+1}} \lambda_{p_{k+1}}) = \theta_{p_\ell} \prod_{k=0}^{L-1} \theta_{p_k \rightarrow p_{k+1}} = \Phi_p(\theta).$$

This shows  $\Phi(\theta') = \Phi(\theta)$ .

Conversely, assume that  $\theta$  is admissible and that  $\Phi(\theta') = \Phi(\theta)$  and  $\text{sign}(\theta'_E) = \text{sign}(\theta_E)$ . By Corollary 1, since  $\theta$  is admissible and  $\Phi(\theta') = \Phi(\theta)$ , we have  $\text{supp}(\theta') = \text{supp}(\theta)$  hence there are  $\gamma_i \neq 0, i \in \text{supp}(\theta)$  such that  $\theta'_i = \gamma_i \theta_i$  for each  $i \in \text{supp}(\theta)$ . Since  $\text{sign}(\theta'_E) = \text{sign}(\theta_E)$ , we have  $\gamma_e > 0$  for every  $e \in E \cap \text{supp}(\theta)$ . Consider  $\alpha \in \mathbb{R}^{E \cup \bar{H}}$  such that  $\alpha_{\bar{H}} = 0$ ,  $e^{\alpha_e} = \gamma_e$  for  $e \in E \cap \text{supp}(\theta)$ , and  $\alpha_e = 0$  for  $e \in E \setminus \text{supp}(\theta)$ . For each  $p \in \mathcal{P}_0$

$$\Phi_p(\theta') = \prod_{e \in p} \theta'_e = \prod_{e \in p} (\theta_e e^{\alpha_e}) = \Phi_p(\theta) e^{\sum_{e \in p} \alpha_e} = \Phi_p(\theta) \odot e^{(\mathbf{P}\alpha)_p}.$$

Since  $\Phi(\theta') = \Phi(\theta)$ , it follows that for each  $p \in \mathcal{P}_0$  such that  $\Phi_p(\theta) \neq 0$  we have  $e^{(\mathbf{P}\alpha)_p} = 1$ , i.e.,  $(\mathbf{P}\alpha)_p = 0$ . Thus,  $\Phi_p(\theta)(\mathbf{P}\alpha)_p = 0$  for all  $p \in \mathcal{P}_0$ , i.e.,  $[\Phi(\theta) \odot \mathbf{P}\alpha]_{\mathcal{P}_0} = 0$ . Since  $\alpha_{\bar{H}} = 0$ , we get that  $\alpha$  belongs to the space  $V_\theta$  defined in (16) in Lemma 7. Since  $\theta$  is admissible, the linear operator  $\mathbf{S}_\theta$  defined in Lemma 7 is a well-defined bijection from  $V_\theta$  to  $\mathbb{R}^H$ , hence  $\alpha$  is related to  $\beta := \mathbf{S}_\theta \alpha \in \mathbb{R}^H$  through the relation (18). Considering  $\delta \in \mathbb{R}^{N_0 \cup H \cup N_L}$  with  $\delta_\nu := \beta_\nu$  for  $\nu \in H$ ,  $\delta_\nu = 0$  for  $\nu \in N_0 \cup N_L$ , relation (18) implies

$$\alpha_e = \delta_\nu - \delta_\mu, \quad \forall e = \mu \rightarrow \nu \in E \cap \text{supp}(\theta).$$

Setting  $\lambda_\nu := e^{\delta_\nu}$  for each  $\nu \in N_0 \cup H \cup N_L$ , it follows that for each  $e = \mu \rightarrow \nu \in E$  we have  $\theta'_e = \lambda_\mu^{-1} \theta_e \lambda_\nu$ . Since  $\text{supp}(\theta) = \text{supp}(\theta')$ , this also trivially holds for  $e \in E \setminus \text{supp}(\theta)$ .

To conclude, we show that  $\theta'_\nu = \theta_\nu \lambda_\nu$  for each  $\nu \in \bar{H}$ . As this holds trivially for  $\nu \in \bar{H} \cap \text{supp}(\theta)$ , we focus on  $\nu \in \bar{H} \setminus \text{supp}(\theta)$ . First, we treat the case of  $\eta \in N_L \cap \text{supp}(\theta)$  by observing that, with  $p = (\eta) \in \mathcal{P}_L$  we have  $\theta'_\eta = \Phi_p(\theta') = \Phi_p(\theta) = \theta_\eta = \theta_\eta \lambda_\eta$  since  $\lambda_\eta = e^{\delta_\eta} = 1$  by definition of  $\delta_\eta := 0$ . Now consider  $\nu \in H \cap \text{supp}(\theta)$ . Since  $\theta$  is admissible,

there is a partial path  $p$  connecting  $\nu$  to some output neuron  $\eta$  with edges in  $\text{supp}(\theta)$ . Since  $-\sum_{e \in p} \alpha_e = (\mathbf{S}_\theta \alpha)_\nu = \beta_\nu = \delta_\nu$  we have

$$\begin{aligned} \theta'_\nu \Pi_{e \in p} \theta'_e &= \Phi_p(\theta') = \Phi_p(\theta) = \theta_\nu \Pi_{e \in p} \theta_e = \theta_\nu \Pi_{e \in p} \theta'_e e^{-\alpha_e} = \theta_\nu (\Pi_{e \in p} \theta'_e) e^{-\sum_{e \in p} \alpha_e} \\ &= \theta_\nu (\Pi_{e \in p} \theta'_e) e^{\delta_\nu} = \theta_\nu (\Pi_{e \in p} \theta'_e) \lambda_\nu. \end{aligned}$$

We conclude using that  $\Pi_{e \in p} \theta'_e \neq 0$  since all edges  $e \in p$  belong to  $\text{supp}(\theta') = \text{supp}(\theta)$ .  $\square$

#### 4. ANALYZING LOCAL IDENTIFIABILITY

Equipped with the rescaling-invariant embedding  $\Phi(\cdot)$  we now establish the claimed local identifiability results. First, we need to introduce notations for the activation status of neurons and paths and use them to provide several expressions of the realization  $\mathbf{R}_\theta$  before providing the main result of the section, Theorem 5.

**4.1. Activation status of neurons and paths, and activation spaces.** The forthcoming analysis heavily involves the activation status of each hidden neuron  $\nu \in H$ ,  $a_\nu(\theta, x) = 1_{z_\nu(\theta, x) > 0} \in \{0, 1\}$ , which gives rise to the activation status of each hidden layer  $\mathbf{a}_\ell(\theta, x) = (a_\nu(\theta, x))_{\nu \in N_\ell}$ ,  $1 \leq \ell \leq L-1$ , and the global activation status  $\mathbf{a}(\theta, x) = (a_\nu(\theta, x))_{\nu \in H} = (\mathbf{a}_\ell(\theta, x))_{1 \leq \ell \leq L-1}$ .

**Definition 13.** *The activation of a path  $p$  (full or partial) is defined as*

$$\alpha_p(\theta, x) := \Pi_{\nu \in H \cap p} a_\nu(\theta, x) \in \{0, 1\}$$

where for  $p = (p_\ell, \dots, p_L) \in \mathcal{P}_\ell$  we denote the set of hidden neurons visited by the path  $p$  using the shorthand  $H \cap p := \{\nu \in H, \exists i \in [\max(\ell, 1), L-1], \nu = p_i\} \subset H$ .

**Remark 4.** *By convention, a product over an empty set is 1. If  $p$  contains no hidden neuron (e.g., if  $p = (\eta) \in \mathcal{P}_L$ ,  $L \geq 1$ ) its activation is  $\alpha_p(\theta, x) = 1$  for every  $x$ .*

With  $\mathcal{Q} := \cup_{\ell=1}^{L-1} \mathcal{Q}_\ell$  the set of all “partial” paths  $q \in (q_\ell, \dots, q_{L-1})$  from a hidden layer  $1 \leq \ell \leq L-1$  to the penultimate layer  $L-1$ , we define the binary-vector-valued function  $\boldsymbol{\alpha}(\theta, x) := (\alpha_q(\theta, x))_{q \in \mathcal{Q}} \in \{0, 1\}^{\mathcal{Q}}$ . We also define variants that are notably useful to account for output biases

$$\bar{\mathbf{a}}_\ell(\theta, x) = \begin{pmatrix} \mathbf{a}_\ell(\theta, x) \\ 1 \end{pmatrix} \in \{0, 1\}^{N_\ell+1} \text{ and } \bar{\boldsymbol{\alpha}}(\theta, x) := \begin{pmatrix} \boldsymbol{\alpha}(\theta, x) \\ 1 \end{pmatrix} \in \{0, 1\}^{\mathcal{Q}+1}$$

where for any set  $A, B$  we use the shorthand  $A^{B+1} = A^B \times A$ .

To state the connections between non-degeneracy and local S-identifiability from finite sets, it is convenient to observe that the linear space  $\mathbf{V}(\theta)$  from Definition 11 can be characterized using simpler linear spaces called *activation spaces*.

**Definition 14** (Activation spaces, activation dimension). *The activation spaces associated to  $\theta \in \mathbb{R}^{E \cup \tilde{H}}$  are*

$$(21) \quad \bar{\mathbf{A}}(\theta) := \text{span} \{ \bar{\alpha}(\theta, x), x \in \mathcal{X}_\theta \} \subseteq \mathbb{R}^{\mathcal{Q}+1}.$$

$$(22) \quad \mathbf{A}(\theta) = \text{span} \{ \mathbf{Q} \bar{\alpha}(\theta, x), x \in \mathcal{X}_\theta \} = \mathbf{Q} \bar{\mathbf{A}}(\theta) \subseteq \mathbb{R}^{\mathcal{Q}_1}$$

with  $\mathbf{Q}$  as in Lemma 10. We define its activation dimension as  $\text{actdim}(\theta) = \dim(\bar{\mathbf{A}}(\theta))$ .

**Remark 5.** Observe that if  $\theta$  and  $\tilde{\theta}$  share the same  $L - 1$  first affine layers  $(\mathbf{W}_\ell, \mathbf{b}_\ell) = (\tilde{\mathbf{W}}_\ell, \tilde{\mathbf{b}}_\ell)_{\ell=0}^{L-1}$  then their activation spaces are identical. This holds even if the dimension of the output layers of  $\tilde{\theta}$  and  $\theta$  differ.

**Lemma 8.** Viewing  $\mathbb{R}^{\mathcal{P}}$  as the product of  $N_L \times N_0$  copies of  $\mathbb{R}^{\mathcal{Q}_1}$  and  $N_L$  copies of  $\mathbb{R}^{\mathcal{Q}+1}$ ,  $\mathbf{V}(\theta) \subset \mathbb{R}^{\mathcal{P}}$  is the product of  $N_L \times N_0$  copies of  $\mathbf{A}^\perp(\theta) \subseteq \mathbb{R}^{\mathcal{Q}_1}$  and  $N_L$  copies of  $\bar{\mathbf{A}}^\perp(\theta)$ .

The proof is postponed to after Corollary 3 as it uses notations introduced there.

**Corollary 2.**  $\mathbf{V}(\theta) = \{0\}$  if, and only if,  $\bar{\mathbf{A}}(\theta) = \mathbb{R}^{\mathcal{Q}+1}$ .

*Proof.* If  $\mathbf{V}(\theta) = \{0\}$  then by Lemma 8 we have  $\bar{\mathbf{A}}^\perp(\theta) = \{0\}$  hence  $\bar{\mathbf{A}}(\theta) = \mathbb{R}^{\mathcal{Q}+1}$ . Vice-versa if  $\bar{\mathbf{A}}(\theta) = \mathbb{R}^{\mathcal{Q}+1}$  then  $\mathbf{A}(\theta) = \mathbf{Q} \bar{\mathbf{A}}(\theta) = \mathbb{R}^{\mathcal{Q}_1}$ . By Lemma 8 it follows that  $\mathbf{V}(\theta) = \{0\}$ ,  $\square$

**4.2. “Algebraic” expressions of the realization.** We can express the realization using weight matrices, bias vectors and layerwise binary activation vectors. A similar formula is stated without taking the biases into account in [16][Lemma A.2] whereas [35] performs analogous computations for gradient computations, still without biases.

**Lemma 9.** Consider  $\theta$  a network parameter of depth  $L \geq 1$ . Denote  $\mathbf{I}_0 = \mathbf{Id}_{\mathbb{R}^{N_0}}$  and for each  $x$  and  $1 \leq \ell \leq L - 1$ ,  $\mathbf{I}_\ell = \text{diag}(\mathbf{a}_\ell(\theta, x))$ . The realization of  $\theta$  satisfies

$$(23) \quad \mathbf{R}_\theta(x) = (\Pi_{\ell=1}^L \mathbf{W}_\ell \mathbf{I}_{\ell-1}) x + \sum_{\ell'=1}^L (\Pi_{\ell=\ell'+1}^L \mathbf{W}_\ell \mathbf{I}_{\ell-1}) \mathbf{b}_{\ell'}$$

with the convention that a product over an empty set is the identity matrix.

The proof is in Appendix E. To conduct an analysis of the local S-identifiability of a parameter, another expression of  $\mathbf{R}_\theta$  where the embedding  $\Phi(\theta)$  appears more explicitly will be useful. We rewrite (23) using  $\Phi(\theta)$  and the activation vector  $\bar{\alpha}(\theta, x)$ .

**Lemma 10.** Consider  $\theta$  a network parameter of depth  $L \geq 2$ . For each  $\eta \in N_L$ , denote<sup>7</sup>

$$(24) \quad \Phi_\eta^{\mathbf{i}}(\theta) := (\Phi_{\mu \rightarrow q \rightarrow \eta}(\theta))_{q \in \mathcal{Q}_1, \mu \in N_0} \in \mathbb{R}^{\mathcal{Q}_1 \times N_0}$$

$$(25) \quad \Phi_\eta^{\mathbf{h}}(\theta) := \left( \begin{array}{c} (\Phi_{q \rightarrow \eta}(\theta))_{q \in \mathcal{Q}} \\ \theta_\eta \end{array} \right) \in \mathbb{R}^{\mathcal{Q}+1}$$

<sup>7</sup>Superscripts **i** and **h** stand for “input” and “hidden”, as  $\Phi^{\mathbf{i}}$  is associated to full paths starting from the input layer, while  $\Phi^{\mathbf{h}}$  corresponds to partial paths starting from a hidden (or the output) layer.

Up to reshaping,  $\Phi(\theta) \in \mathbb{R}^{\mathcal{P}}$  is the concatenation of matrices  $\Phi_\eta^i(\theta) \in \mathbb{R}^{\mathcal{Q}_1 \times N_0}$  and vectors  $\Phi_\eta^h(\theta) \in \mathbb{R}^{\mathcal{Q}+1}$  over all output neurons  $\eta \in N_L$ . For each output neuron  $\eta \in N_L$  we have

$$(26) \quad \mathbf{R}_\theta(x)_\eta = \langle \mathbf{Q}\bar{\alpha}(\theta, x), \Phi_\eta^i(\theta)x \rangle + \langle \bar{\alpha}(\theta, x), \Phi_\eta^h(\theta) \rangle$$

where  $\mathbf{Q} : \mathbb{R}^{\mathcal{Q}+1} \rightarrow \mathbb{R}^{\mathcal{Q}_1}$  is the canonical restriction to  $\mathcal{Q}_1 \subset \mathcal{Q}$ .

The proof of Lemma 10 is in Appendix E. It yields an expression of  $\mathbf{R}_\theta$  that perfectly fits the upcoming analysis of local S-identifiability. A more abstract (but probably somewhat more digestible) version of the same result implies Property (8) as claimed.

**Corollary 3.** Consider  $\theta$  a network parameter of depth  $L \geq 2$ . For each  $x \in \mathbb{R}^{N_0}$  let  $\mathbf{L}_{\theta,x}$  be the linear form on  $\mathbb{R}^{\mathcal{Q}_1 \times N_0} \times \mathbb{R}^{\mathcal{Q}+1}$  defined as

$$\mathbf{L}_{\theta,x}\{(\mathbf{M}, \mathbf{v})\} := \langle \mathbf{Q}\bar{\alpha}(\theta, x), \mathbf{M}x \rangle + \langle \bar{\alpha}(\theta, x), \mathbf{v} \rangle, \quad \mathbf{M} \in \mathbb{R}^{\mathcal{Q}_1 \times N_0}, \quad \mathbf{v} \in \mathbb{R}^{\mathcal{Q}+1}$$

Define  $\mathbf{C}_{\theta,x} \in \mathbb{R}^{N_L \times \mathcal{P}}$  the matrix associated to the linear operator mapping each  $\phi \in \mathbb{R}^{\mathcal{P}}$ , seen as a reshaped concatenation of matrices  $\phi_\eta^i \in \mathbb{R}^{\mathcal{Q}_1 \times N_0}$  and vectors  $\phi_\eta^h \in \mathbb{R}^{\mathcal{Q}+1}$  as in Lemma 10, to  $\mathbf{r} := (r_\eta)_{\eta \in N_L}$ , with  $r_\eta = \mathbf{L}_{\theta,x}\{(\phi_\eta^i, \phi_\eta^h)\}$ . We have

$$(27) \quad \mathbf{R}_\theta(x) = \mathbf{C}_{\theta,x} \cdot \Phi(\theta)$$

We are now equipped with the notations needed to prove Lemma 8. The proof also relies on the following alternative characterization of the set  $\mathcal{X}_\theta$  from Definition 10 that we did not find elsewhere. It is proved in Appendix F.

**Lemma 11.** Given a parameter  $\theta$ , consider the open set of input variables  $x$  such that  $(\theta', z) \mapsto \mathbf{a}(\theta', z)$  is locally constant in some neighborhood of  $(\theta, x)$ .

$$(28) \quad \mathcal{X}'_\theta := \{x \in \mathbb{R}^{N_0} : \exists \epsilon, r > 0, \forall (\theta, z) \in B(\theta, \epsilon) \times B(x, r), \mathbf{a}(\theta', z) = \mathbf{a}(\theta, x)\}.$$

with the convention that  $\mathcal{X}'_\theta = \mathbb{R}^{N_0}$  if the network depth is  $L = 1$ . This set coincides exactly with the set  $\mathcal{X}_\theta$  from Definition 10.

*Proof of Lemma 8.* Consider a vector  $\phi \in \mathbb{R}^{\mathcal{P}}$  and its representation as  $\phi_\eta^i \in \mathbb{R}^{\mathcal{Q}_1 \times N_0}$ ,  $\phi_\eta^h \in \mathbb{R}^{\mathcal{Q}+1}$ ,  $\eta \in N_L$ . By definition  $\phi \in \mathbf{V}(\theta)$  if, and only if,  $\mathbf{C}_{\theta,x}\phi = 0$ ,  $\forall x \in \mathcal{X}_\theta$ , i.e., for each  $\eta \in N_L$  we have

$$(29) \quad \langle \mathbf{Q}\bar{\alpha}(\theta, x), \phi_\eta^i x \rangle + \langle \bar{\alpha}(\theta, x), \phi_\eta^h \rangle = 0, \quad \forall x \in \mathcal{X}_\theta.$$

By Lemma 11,  $x' \mapsto \bar{\alpha}(\theta, x')$  is locally constant in the neighborhood of each  $x \in \mathcal{X}_\theta$ , hence the left-hand-side in (29) is locally affine with respect to  $x$ , and (29) is thus equivalent to

$$(30) \quad \begin{cases} [\mathbf{Q}\bar{\alpha}(\theta, x)]^\top \phi_\eta^i &= \mathbf{0}_{1 \times N_0}, \quad \forall x \in \mathcal{X}_\theta, \\ \langle \bar{\alpha}(\theta, x), \phi_\eta^h \rangle &= 0 \end{cases}$$

that is to say each column of  $\phi_\eta^i$  is orthogonal to  $\mathbf{Q}\bar{\alpha}(\theta, x)$ , and  $\phi_\eta^h$  is orthogonal to  $\bar{\alpha}(\theta, x)$  for every  $x \in \mathcal{X}_\theta$ . We conclude using the definition of  $\mathbf{A}(\theta), \bar{\mathbf{A}}(\theta)$ .  $\square$

**4.3. Non-degeneracy and local S-identifiability.** We can now state the main result of this section.

**Theorem 5.** *Consider  $\theta \in \Theta \subseteq \mathbb{R}^{E \cup \bar{H}}$ . The following are equivalent:*

- i)  $\theta$  is non-degenerate with respect to  $\Theta$ ;
- ii) there is a finite  $F \subset \mathcal{X}_\theta$  such that  $\theta$  is locally S-identifiable from  $F$  with respect to  $\Theta$ .
- iii) there is a compact  $K \subset \mathcal{X}_\theta$  such that  $\theta$  is locally S-identifiable from  $K$  wrt  $\Theta$ .

When they hold, the finite set  $F$  can be chosen such that

$$(31) \quad \text{card}(F) \leq (N_0 + 1) \text{actdim}(\theta).$$

**Remark 6.** We exhibit in Example 4 a PS-identifiable (hence locally S-identifiable) parameter  $\theta$  that is degenerate, i.e., not locally S-identifiable from any compact  $K \subseteq \mathcal{X}_\theta$ .

*Proof.* **i)  $\Rightarrow$  ii)** Consider  $\epsilon > 0$  such that  $\theta$  is  $\epsilon$ -non-degenerate with respect to  $\Theta$ . To establish the existence of  $F$  such that  $\theta$  is locally S-identifiable from  $F$  with respect to  $\Theta$ , we use a Lemma whose proof is postponed.

**Lemma 12.** *Consider  $\theta \in \mathbb{R}^{E \cup \bar{H}}$ .*

- a) *There exists  $\epsilon > 0$  and a set  $F \subset \mathcal{X}_\theta$  of cardinality at most  $(N_0 + 1) \text{actdim}(\theta)$  such that: for each  $\theta' \in B(\theta, \epsilon)$ , if  $\mathbf{R}_{\theta'} = \mathbf{R}_\theta$  on  $F$ , then  $\Phi(\theta') - \Phi(\theta) \in \mathbf{V}(\theta)$ .*
- b) *For every compact set  $K \subset \mathcal{X}_\theta$ , there exists  $\epsilon' > 0$  such that: for each  $\theta' \in B(\theta, \epsilon')$ , if  $\Phi(\theta') - \Phi(\theta) \in \mathbf{V}(\theta)$ , then  $\mathbf{R}_{\theta'}(x) = \mathbf{R}_\theta(x)$  for all  $x \in K$ .*

Let  $\epsilon_0, F$  be given by Lemma 12-a and set  $\epsilon_1 := \min(\epsilon_0, \epsilon, \eta/2)$  where  $\eta := \min_{i \in \text{supp}(\theta)} |\theta_i|$ . We will show that  $\theta$  is  $\epsilon_1$ -locally S-identifiable from  $F$ . For this, consider  $\theta' \in \Theta \cap B(\theta, \epsilon_1)$  and assume that  $\mathbf{R}_{\theta'} = \mathbf{R}_\theta$  on  $K$ . By Lemma 12-a, since  $\theta' \in B(\theta, \epsilon_0)$ , we have  $\Phi(\theta') - \Phi(\theta) \in \mathbf{V}(\theta)$ . Since  $\theta' \in B(\theta, \epsilon)$  and  $\theta$  is  $\epsilon$ -non-degenerate, this implies  $\Phi(\theta') = \Phi(\theta)$  hence (recall that, since  $\theta$  is non-degenerate, it is admissible by definition) by Lemma 6 we have  $\text{supp}(\theta') = \text{supp}(\theta)$ . Since  $\theta' \in B(\theta, \eta/2)$  we further have  $\text{sign}(\theta'_i) = \text{sign}(\theta_i)$  for every  $i \in \text{supp}(\theta)$ , hence  $\text{sign}(\theta') = \text{sign}(\theta)$ . By Theorem 1 we obtain  $\theta' \sim_S \theta$ .

**ii)  $\Rightarrow$  iii)** Simply observe that a finite set is compact.

**iii)  $\Rightarrow$  i)** Consider  $\epsilon > 0$  such that  $\theta$  is  $\epsilon$ -locally identifiable from  $K$  with respect to  $\Theta$ . By Lemma 12-b for the compact set  $K$ , there is  $\epsilon_0 > 0$  such that: for each  $\theta' \in B(\theta, \epsilon_0)$ ,  $\Phi(\theta') - \Phi(\theta) \in \mathbf{V}(\theta) \Rightarrow (\mathbf{R}_{\theta'}(x) = \mathbf{R}_\theta(x), \forall x \in K)$ . Set  $\epsilon_1 := \min(\epsilon, \epsilon_0)$ . We will show that  $\theta$  is  $\epsilon_1$ -non-degenerate with respect to  $\Theta$ . Considering  $\theta' \in \Theta \cap B(\theta, \epsilon_1)$  such that  $\Phi(\theta') - \Phi(\theta) \in \mathbf{V}(\theta)$  we now show that  $\Phi(\theta') = \Phi(\theta)$ . Since  $\theta' \in B(\theta, \epsilon_0)$  and  $\Phi(\theta') - \Phi(\theta) \in \mathbf{V}(\theta)$ , we have  $\mathbf{R}_{\theta'}(x) = \mathbf{R}_\theta(x)$  for all  $x \in K$ . Since  $\theta' \in \Theta \cap B(\theta, \epsilon)$  and  $\theta$  is locally S-identifiable from  $K$  with respect to  $\Theta$  this implies  $\theta' \sim_S \theta$ , hence by Theorem 1 we have  $\Phi(\theta') = \Phi(\theta)$ .  $\square$

*Proof of Lemma 12.* We begin with some preliminaries. Since  $\bar{\mathbf{A}}(\theta) \subseteq \mathbb{R}^{\mathcal{Q}+1}$ , it is finite dimensional hence there is a finite set  $\mathcal{Z}_\theta \subset \mathcal{X}_\theta$  such that  $\text{card}(\mathcal{Z}_\theta) = \text{actdim}(\theta)$  and

$$(32) \quad \bar{\mathbf{A}}(\theta) = \text{span} \{ \bar{\alpha}(\theta, z), z \in \mathcal{Z}_\theta \}.$$

By definition of  $\mathcal{X}_\theta$ , for each  $z \in \mathcal{Z}_\theta$  there exists  $\epsilon(z), r(z) > 0$  such that, for every  $\theta' \in B(\theta, \epsilon(z))$  and  $x \in B(z, r(z))$ , we have  $\mathbf{a}(\theta', x) = \mathbf{a}(\theta, z)$ , hence  $\bar{\alpha}(\theta', x) = \bar{\alpha}(\theta, z)$ . Since  $\mathcal{Z}_\theta$  is finite,

$$\epsilon := \min_{z \in \mathcal{Z}_\theta} \epsilon(z) > 0.$$

Consider  $\theta' \in B(\theta, \epsilon)$ ,  $z \in \mathcal{Z}_\theta$ ,  $x \in B(z, r(z))$ . Since  $\bar{\alpha}(\theta', x) = \bar{\alpha}(\theta, z)$ , we have for each output neuron  $\eta \in N_L$

$$\begin{cases} [Q\bar{\alpha}(\theta', x)]^\top \Phi_\eta^i(\theta') &= [Q\bar{\alpha}(\theta', z)]^\top \Phi_\eta^i(\theta') \\ \bar{\alpha}(\theta', x)^\top \Phi_\eta^h(\theta') &= \bar{\alpha}(\theta', z)^\top \Phi_\eta^h(\theta') \end{cases}$$

hence using (26) we get

$$(33) \quad \mathbf{R}_{\theta'}(x)_\eta - \mathbf{R}_\theta(x)_\eta = [Q\bar{\alpha}(\theta, z)]^\top (\Phi_\eta^i(\theta') - \Phi_\eta^i(\theta)) x + \bar{\alpha}^\top(\theta, z) (\Phi_\eta^h(\theta') - \Phi_\eta^h(\theta)).$$

Considering  $z \in \mathcal{Z}_\theta$ , define  $F_z := \{x_i\}_{i=0}^{N_0} \subset B(z, r(z)) \subset \mathbb{R}^{N_0}$  where  $x_0 = z$  and for  $1 \leq i \leq N_0$ ,  $x_i = z + \frac{r(z)}{2} \delta_i$  with  $\delta_i$  the  $i$ -th vector of the canonical basis. Observe that if  $\mathbf{u} \in \mathbb{R}^{N_0}$ ,  $b \in \mathbb{R}$  are such that  $\mathbf{u}^\top x + b = 0$  for every  $x \in F_z$ , then  $\mathbf{u} = \mathbf{0}$  (since  $r(z)\mathbf{u}^\top \delta_i = \mathbf{u}^\top (x_i - x_0) = \mathbf{u}^\top x_i + b - (\mathbf{u}^\top x_0 + b) = 0$  for every  $i$ ), and therefore  $b = 0$  too.

**a)** The finite set  $F := \cup_{z \in \mathcal{Z}_\theta} F_z$  satisfies  $F \subset \cup_{z \in \mathcal{Z}_\theta} B(z, r(z)) \subset \mathcal{X}_\theta$ . Assume that  $\mathbf{R}_{\theta'} = \mathbf{R}_\theta$  on  $F$  where  $\theta' \in B(\theta, \epsilon)$ . By the preliminaries, this implies that the right hand side in (33) is zero for each  $\eta \in N_L$ ,  $z \in \mathcal{Z}_\theta$ ,  $x \in F_z$ , hence

$$\begin{cases} [Q\bar{\alpha}(\theta, z)]^\top (\Phi_\eta^i(\theta') - \Phi_\eta^i(\theta)) = \mathbf{0}_{1 \times N_0} \\ \bar{\alpha}^\top(\theta, z) (\Phi_\eta^h(\theta') - \Phi_\eta^h(\theta)) = 0. \end{cases}$$

Since this holds for every  $\eta \in N_L$ ,  $z \in \mathcal{Z}_\theta$ , in light of (32) this establishes that

$$(34) \quad \forall \eta \in N_L, \begin{cases} \Phi_\eta^i(\theta') - \Phi_\eta^i(\theta) & \in \underbrace{\mathbf{A}^\perp(\theta) \times \dots \times \mathbf{A}^\perp(\theta)}_{N_0 \text{ times}} \\ \Phi_\eta^h(\theta') - \Phi_\eta^h(\theta) & \in \bar{\mathbf{A}}^\perp(\theta) \end{cases}$$

and we conclude using Lemma 8 and the fact that  $\text{card}(F) \leq \text{card}(\mathcal{Z}_\theta) \times (N_0 + 1)$ .

**b)** Since  $K \subset \mathcal{X}_\theta$ , for each  $z \in K$  there are  $\epsilon(z), r(z) > 0$  such that: for each  $\theta' \in B(\theta, \epsilon(z))$ ,  $x \in B(z, r(z))$ ,  $\bar{\alpha}(\theta', x) = \bar{\alpha}(\theta, z)$ . Since  $K$  is compact and  $K \subset \cup_{z \in K} B(z, r(z))$ , there is a finite set  $\mathcal{Z} \subset K$  such that  $K \subset \cup_{z \in \mathcal{Z}} B(z, r(z))$ . Denote  $\epsilon' := \min_{z \in \mathcal{Z}} \epsilon(z) > 0$ . Considering  $\theta' \in B(\theta, \epsilon')$  such that  $\Phi(\theta') - \Phi(\theta) \in \mathbf{V}(\theta)$ , we now show that  $\mathbf{R}_{\theta'}(x) = \mathbf{R}_\theta(x)$  for each  $x \in K$ . Given  $x \in K$ , since there is  $z \in \mathcal{Z}_\theta$  such that  $x \in B(z, r(z))$ , we have

$$(35) \quad \bar{\alpha}(\theta', x) = \bar{\alpha}(\theta, z) = \bar{\alpha}(\theta, x).$$

For each  $\eta \in N_L$ , since by Lemma 8  $\Phi(\theta') - \Phi(\theta) \in \mathbf{V}(\theta)$  is equivalent to (34), we get

$$\begin{aligned} [Q\bar{\alpha}(\theta', x)]^\top \Phi_\eta^i(\theta') &\stackrel{(35)}{=} [Q\bar{\alpha}(\theta, x)]^\top \Phi_\eta^i(\theta') \stackrel{(34)}{=} [Q\bar{\alpha}(\theta, x)]^\top \Phi_\eta^i(\theta) \\ \bar{\alpha}^\top(\theta', x) \Phi_\eta^h(\theta') &\stackrel{(35)}{=} \bar{\alpha}^\top(\theta, x) \Phi_\eta^h(\theta') \stackrel{(34)}{=} \bar{\alpha}^\top(\theta, x) \Phi_\eta^h(\theta). \end{aligned}$$

Using (26) we conclude that  $\mathbf{R}_{\theta'}(x)_\eta = \mathbf{R}_\theta(x)_\eta$  for all  $\eta \in N_L$ , i.e.,  $\mathbf{R}_{\theta'}(x) = \mathbf{R}_\theta(x)$ .  $\square$



## 5. IDENTIFIABILITY FOR SHALLOW NEURAL NETWORKS

In this section we focus on shallow networks, for which the set  $\mathcal{Q} = \mathcal{Q}_1$  of paths is in bijection with the set  $H = N_1$  of hidden neurons. Identifying these sets the activation vectors also coincide  $\alpha(\theta, x) = \mathbf{a}_1(\theta, x) \in \mathbb{R}^{\mathcal{Q}} = \mathbb{R}^{N_1} = \mathbb{R}^H$ . After giving a complete characterization of the activation space  $\bar{\mathbf{A}}(\theta)$  using the notion of twin neurons, we show that the absence of twin neurons implies non-degeneracy (hence local S-identifiability), and that its combination with irreducibility implies PS-identifiability. Finally, we discuss what happens in the presence of twin neurons.

**5.1. Activation spaces and twin neurons.** Whenever  $\theta$  is admissible, each hidden neuron  $\nu \in H$  is not dead, i.e.  $\mathbf{w}_{\bullet \rightarrow \nu} \neq 0$  and  $\mathbf{w}_{\nu \rightarrow \bullet} \neq 0$ . According to Definition 8, neurons are twins if their extended vectors  $(\mathbf{w}_{\bullet \rightarrow \nu}, b_\nu)$  are colinear. This defines an equivalence relation, and the hidden layer  $H = N_1$  can be partitioned into *equivalence classes of twin neurons*, denoted

$$T_c \subset H, 1 \leq c \leq C.$$

Each equivalence class  $T_c$  is partitioned into  $I_c, J_c$ , where all neurons in  $I_c$  are positive twins, all neurons in  $J_c$  are positive twins, and  $\nu \in I_c, \nu' \in J_c$  are negative twins. By convention  $I_c$  is always non-empty, while  $J_c$  may be empty if there are no negative twins in  $T_c$ . For each class, we can define a *class signature* vector

$$\mathbf{s}_c = \mathbf{1}_{I_c} - \mathbf{1}_{J_c} \in \mathbb{R}^H,$$

which is zero out of  $T_c$ , with  $\pm 1$  entries on  $T_c$ , and has at least one  $+1$  entry. When  $T_c$  contains both positive and negative twins,  $\mathbf{s}_c$  is only defined up to a global sign. An equivalence class is said to be nontrivial if its cardinal is at least two. Equipped with these notions, we prove in Appendix G the following characterization of activation spaces.

**Lemma 13.** *Consider an admissible parameter  $\theta$  on a shallow network architecture. Using the notations introduced above, its activation spaces are*

$$(36) \quad \mathbf{A}(\theta) = \text{span} \{ \mathbf{1}_H, \mathbf{s}_c, 1 \leq c \leq C \} \subseteq \mathbb{R}^H$$

$$(37) \quad \bar{\mathbf{A}}(\theta) = \text{span} \{ (\mathbf{1}_H, 2), (\mathbf{s}_c, 0), 1 \leq c \leq C \} \subseteq \mathbb{R}^{H+1}.$$

**5.2. Proof of Lemma 5: no twins implies non-degeneracy.** Lemma 5 is a direct consequence of the combination of Theorem 5 with the following two results.

**Lemma 14.** *On any network architecture, if  $\theta \in \mathbb{R}^{E \cup \bar{H}}$  is admissible and  $\bar{\mathbf{A}}(\theta) = \mathbb{R}^{\mathcal{Q}+1}$  then<sup>8</sup>  $\theta$  is non-degenerate with respect to any  $\Theta \subset \mathbb{R}^{E \cup \bar{H}}$  that contains it.*

*Proof.* Since  $\bar{\mathbf{A}}(\theta) = \mathbb{R}^{\mathcal{Q}+1}$ , by Corollary 2  $\mathbf{V}(\theta) = \{0\}$ , hence  $\Phi(\theta') - \Phi(\theta) \in \mathbf{V}(\theta)$  is equivalent to  $\Phi(\theta') = \Phi(\theta)$ . Since  $\theta$  is admissible, this shows that  $\theta$  is non-degenerate.  $\square$

**Lemma 15.** *Consider a shallow architecture and  $\theta \in \mathbb{R}^{E \cup \bar{H}}$ . The equality  $\bar{\mathbf{A}}(\theta) = \mathbb{R}^{\mathcal{Q}+1}$  holds if, and only if, there is no twin. When this holds,  $\text{actdim}(\theta) = |H| + 1 = |N_1| + 1$ .*

<sup>8</sup>The converse does not hold: there are non-degenerate parameters with  $\bar{\mathbf{A}}(\theta) \neq \mathbb{R}^{\mathcal{Q}+1}$ , see Lemma 16.

*Proof.* Equivalence classes of twin neurons form a partition of  $H$ , hence  $|C| \leq |H| = |Q|$ . By Lemma 13,  $\bar{A}(\theta)$  is the span of  $|C| + 1$  vectors, hence its dimension is at most  $|C| + 1$ . In the presence of twins we get  $|C| < |H|$  hence  $\bar{A}(\theta) \neq \mathbb{R}^{Q+1}$ . In the absence of twins, each equivalence class  $T_c$  is trivial, i.e.  $|T_c| = 1$ . We obtain that  $|C| = |H|$ , that each signature vector  $\mathbf{s}_c$  is a distinct canonical vector  $\delta_c$ , and obtain  $\bar{A}(\theta) = \mathbb{R}^{Q+1}$  by Lemma 13.  $\square$

### 5.3. Proof of Theorem 3: irreducibility and no twins implies PS-identifiability.

By Lemma 3, PS-identifiability from a bounded set with respect to  $\Theta = \mathbb{R}^{E \cup H}$  implies that  $\theta$  has no twins, hence by Lemma 4), it is irreducible (hence admissible), and local S-identifiable, by Theorem 2. For shallow networks, we show that conversely, irreducibility and the absence of twins imply PS-identifiability from a bounded set.

**Theorem 6.** Consider  $N_1, N'_1$  two finite sets of indices, empty or not<sup>9</sup>, and integers  $d, k \geq 1$ . Consider  $\mathbf{c} \in \mathbb{R}^k$  and for each  $\nu \in N_1$ , let  $\mathbf{v}_\nu \in \mathbb{R}^k$ ,  $\mathbf{w}_\nu \in \mathbb{R}^d$  and  $b_\nu \in \mathbb{R}$ . Define

$$\varphi(x) = \sum_{\nu \in N_1} \mathbf{v}_\nu \text{ReLU}(\langle \mathbf{w}_\nu, x \rangle + b_\nu) + \mathbf{c}, \quad x \in \mathbb{R}^d.$$

Similarly define  $\psi(x)$  with  $\mathbf{v}'_\nu \in \mathbb{R}^k$ ,  $\mathbf{w}'_\nu \in \mathbb{R}^d$ ,  $b'_\nu \in \mathbb{R}$  for  $\nu \in N'_1$ , and  $\mathbf{c}' \in \mathbb{R}^k$ .

a) Assume that

- $\{(\mathbf{w}_\nu, b_\nu)\}_{\nu \in N_1}$  are pairwise not collinear, and  $\mathbf{v}_\nu, \mathbf{w}_\nu \neq 0$ ;
- $\{(\mathbf{w}'_\nu, b'_\nu)\}_{\nu \in N'_1}$  are pairwise not collinear, and  $\mathbf{v}'_\nu, \mathbf{w}'_\nu \neq 0$ .

If  $\varphi(x) = \psi(x)$  for every  $x \in \mathbb{R}^d$  then  $\text{card}(N_1) = \text{card}(N'_1)$ .

b) Assume that  $\{(\mathbf{w}_\nu, b_\nu)\}_{\nu \in N_1}$  are pairwise not collinear, and

$$(38) \quad \sum_{\nu \in T} \mathbf{v}_\nu \mathbf{w}_\nu^\top \neq 0, \quad \text{for all non-empty } T \subset N_1.$$

There exists a bounded set  $\mathcal{X} \subseteq \mathbb{R}^{N_0}$  (which depends on  $\theta$ ) such that: if  $N'_1 = N_1$  and  $\varphi(x) = \psi(x)$  for every  $x \in \mathcal{X}$ , then<sup>10</sup>  $\mathbf{c} = \mathbf{c}'$  and there exists a permutation  $\pi$  of  $N_1$  and  $\lambda_\nu > 0$ ,  $\nu \in N_1$  such that

$$(39) \quad \forall \nu \in N_1 : \mathbf{v}'_{\pi(\nu)} = \lambda_\nu^{-1} \mathbf{v}_\nu; \quad \mathbf{w}'_{\pi(\nu)} = \lambda_\nu \mathbf{w}_\nu \quad \text{and} \quad b'_{\pi(\nu)} = \lambda_\nu b_\nu.$$

*Proof.* As a preliminary, consider  $\nu \in N_1$  and denote  $\mathcal{V}_\nu := \{x \in \mathbb{R}^{N_0} : \langle \mathbf{w}_\nu, x \rangle + b_\nu = 0\}$ . Since  $\mathbf{w}_\nu \neq 0$ , the set  $\mathcal{V}_\nu$  is a hyperplane which matches the set  $\Gamma_\nu(\theta)$  from Definition 10 when considering  $\theta$  such that  $\varphi = \mathbf{R}_\theta$ . As none of the  $(\mathbf{w}_\nu, b_\nu)$  is collinear to another, the hyperplanes associated to  $\nu \neq \nu' \in N_1$  are distinct. As  $\mathbf{v}_\nu \neq 0$  for every  $\nu \in N_1$  and  $\varphi$  is continuous and piecewise affine, this function is differentiable exactly on the complement of  $\mathcal{T} := \cup_{\nu \in N_1} \mathcal{V}_\nu$ , which is a union of  $\text{card}(N_1)$  distinct hyperplanes.

a) Similarly, since none of the  $(\mathbf{w}'_\nu, b'_\nu)$  is collinear to another and  $\mathbf{v}'_\nu \neq 0, \mathbf{w}'_\nu \neq 0$  for each  $\nu \in N'_1$ , the function  $\psi$  is differentiable exactly on the complement of a union of  $\text{card}(N'_1)$  distinct hyperplanes,  $\mathcal{T}' = \cup_{\nu \in N'_1} \mathcal{V}'_\nu$ , where  $\mathcal{V}'_\nu := \{x \in \mathbb{R}^{N_0} : \langle \mathbf{w}'_\nu, x \rangle + b'_\nu = 0\}$ . Note that  $\mathcal{T}$  may be empty since  $N_1$  may be empty, and similarly for  $\mathcal{T}'$ . Since  $\varphi = \psi$ , we

<sup>9</sup>We use the convention:  $\sum_\emptyset = 0$ .

<sup>10</sup>Let us emphasize that here no further assumption is made on  $\mathbf{w}'_\nu, b'_\nu, \mathbf{v}'_\nu, \nu \in N_1$ .

have  $\mathcal{T} = \mathcal{T}'$  hence  $\text{card}(N_1) = \text{card}(N'_1)$ , otherwise there would exist one point  $x \in \mathbb{R}^d$  where one function would be differentiable and the other not.

**b)** We now assume  $N'_1 = N_1$ , but make no specific assumption on  $\mathbf{v}'_\nu \in \mathbb{R}^k$ ,  $\mathbf{w}'_\nu \in \mathbb{R}^d$ ,  $b'_\nu \in \mathbb{R}$  for  $\nu \in N_1$  or on  $\mathbf{c}' \in \mathbb{R}^k$ . By (38) with  $T = \{\nu\}$  we have  $\mathbf{v}_\nu \mathbf{w}_\nu^\top \neq 0$  hence, as in the preliminary,  $\mathcal{V}_\nu$ ,  $\nu \in N_1$  are pairwise distinct hyperplanes. Consider an arbitrary hidden neuron  $\nu \in N_1$ . As the hyperplanes  $\{\mathcal{V}_\mu\}_{\mu \in N_1}$  are pairwise distinct, there exist  $x_\nu \in \mathcal{V}_\nu$  and  $\epsilon_\nu > 0$  such that  $\Omega_\nu := B(x_\nu, \epsilon_\nu)$  satisfies  $\Omega_\nu \cap \mathcal{T} = \Omega_\nu \cap \mathcal{V}_\nu$ . We will show that the result holds with  $\mathcal{X} := \cup_{\nu \in N_1} \Omega_\nu$ , which is easily seen to be bounded.

From now, assume that  $\psi(x) = \varphi(x)$  for every  $x \in \mathcal{X}$ .

For each  $\nu \in \hat{N}_1 := \{\nu \in N_1 : \mathbf{v}'_\nu \neq 0, \mathbf{w}'_\nu \neq 0\}$ , since  $\mathbf{w}'_\nu \neq 0$ , the set  $\mathcal{V}'_\nu$  is a hyperplane. Consider the equivalence relation on  $\hat{N}_1$  defined by:  $\nu \sim \mu \Leftrightarrow \mathcal{V}'_\nu = \mathcal{V}'_\mu$ , and the resulting quotient set  $\bar{N}_1 = \hat{N}_1 / \sim$ . For each equivalence class  $\bar{\nu} \in \bar{N}_1$ , denote  $\mathcal{V}'_{\bar{\nu}}$  the common hyperplane associated to every  $\nu \in \bar{\nu}$ , and set  $\bar{\mathcal{T}} = \cup_{\bar{\nu} \in \bar{N}_1} \mathcal{V}'_{\bar{\nu}}$ . We will prove below that there exists an injective map  $\pi : N_1 \rightarrow \bar{N}_1$  such that  $\mathcal{V}_\nu = \mathcal{V}'_{\pi(\nu)}$  for every  $\nu \in N_1$ . This will imply that  $\text{card}(\bar{N}_1) \geq \text{card}(N_1)$ , and since  $\text{card}(\bar{N}_1) \leq \text{card}(\hat{N}_1) \leq \text{card}(N_1)$ , it will follow that  $\hat{N}_1 = N_1$  (hence  $\mathbf{v}'_\nu \neq 0, \mathbf{w}'_\nu \neq 0$  for every  $\nu \in N_1$ ) and that each equivalence class  $\bar{\nu}$  is a singleton. In other words,  $\pi$  is indeed a permutation of  $N_1$ , and the hyperplanes  $\mathcal{V}'_{\{\nu\}}$ ,  $\nu \in N_1$  are pairwise distinct.

To build  $\pi$ , consider a hidden neuron  $\nu \in N_1$ . For the sake of contradiction, assume that  $\mathcal{V}'_{\bar{\mu}} \neq \mathcal{V}_\nu$  for every  $\bar{\mu} \in \bar{N}_1$ . This implies the existence of  $x'_\nu \in \Omega_\nu \cap \mathcal{V}_\nu$  and of  $\epsilon'_\nu > 0$  such that  $\Omega'_\nu := B(x'_\nu, \epsilon'_\nu) \subseteq \Omega_\nu$  and  $\Omega'_\nu \cap \bar{\mathcal{T}} = \emptyset$  and  $\Omega'_\nu \cap \mathcal{V}_\nu = \mathcal{V}_\nu$ . Since  $\Omega'_\nu \cap \bar{\mathcal{T}} = \emptyset$ , the function  $\psi$  is affine linear on  $\Omega'_\nu$ , hence it has constant Jacobian on  $\Omega'_\nu$ . Denote  $\Omega_\nu^+ := \{x \in \Omega'_\nu : \langle \mathbf{w}_\nu, x \rangle + b_\nu > 0\}$ ,  $\Omega_\nu^- := \{x \in \Omega'_\nu : \langle \mathbf{w}_\nu, x \rangle + b_\nu < 0\}$ , and observe that both sets are non-empty. For any  $x \in \Omega'_\nu \setminus \mathcal{V}_\nu = \Omega_\nu^+ \cup \Omega_\nu^-$ , the function  $\varphi$  is differentiable and its Jacobian is  $\varphi'(x) = \mathbf{v}_\nu \mathbf{w}_\nu^\top H(\langle \mathbf{w}_\nu, x \rangle + b_\nu) + \mathbf{d}$  where  $\mathbf{d} \in \mathbb{R}^k$  and

$$H(t) := \begin{cases} 1, & \text{if } t > 0 \\ 0, & \text{otherwise.} \end{cases}$$

For each  $x_\nu^+ \in \Omega_\nu^+, x_\nu^- \in \Omega_\nu^-$  we have  $H(\langle \mathbf{w}_\nu, x_\nu^+ \rangle + b_\nu) - H(\langle \mathbf{w}_\nu, x_\nu^- \rangle + b_\nu) = 1$ , hence  $\varphi'(x_\nu^+) - \varphi'(x_\nu^-) = \mathbf{v}_\nu \mathbf{w}_\nu^\top$ . As  $\psi = \varphi$  on  $\mathcal{X} \supseteq \Omega_\nu \supseteq \Omega'_\nu$  and  $\psi$  has constant Jacobian on  $\Omega'_\nu$ , it follows that  $\mathbf{v}_\nu \mathbf{w}_\nu^\top = 0$ , which contradicts our assumptions. Hence, there is  $\bar{\mu} \in \bar{N}_1$  such that  $\mathcal{V}'_{\bar{\mu}} = \mathcal{V}_\nu$ . Since the hyperplanes  $\{\mathcal{V}'_{\bar{\nu}}\}_{\bar{\nu} \in \bar{N}_1}$  are pairwise disjoint by construction, such a  $\bar{\mu}$  is unique and we define  $\pi(\nu) := \bar{\mu}$ . Since this holds for every  $\nu \in N_1$ , we can define the map  $\pi : N_1 \rightarrow \bar{N}_1$  with  $\pi(\nu) := \bar{\mu}$ . For  $\nu \neq \nu'$  we have  $\mathcal{V}'_{\pi(\nu')} = \mathcal{V}_{\nu'} \neq \mathcal{V}_\nu = \mathcal{V}'_{\pi(\nu)}$  since the hyperplanes  $\{\mathcal{V}_\nu\}_{\nu \in N_1}$  are pairwise distinct. This proves the injectivity of  $\pi$ . As we have seen, this means that indeed  $\pi$  is a permutation of  $N_1$ . Without loss of generality, to simplify notations, we assume from now on that  $\pi$  is the identity.

For each  $\nu \in N_1$ , since  $\mathcal{V}'_{\{\nu\}} = \mathcal{V}'_{\pi(\nu)} = \mathcal{V}_\nu$  there is a nonzero  $\lambda_\nu \in \mathbb{R}$  such that

$$(\mathbf{w}'_\nu, b'_\nu) = \lambda_\nu (\mathbf{w}_\nu, b_\nu).$$

Reasoning as above, with  $\Omega_\nu^\pm$  defined using  $\Omega'_\nu := \Omega_\nu$ , we obtain that

$$\varphi'(x_\nu^+) - \varphi'(x_\nu^-) = \mathbf{v}_\nu \mathbf{w}_\nu^\top.$$

for each  $x_\nu^+ \in \Omega_\nu^+$ ,  $x_\nu^- \in \Omega_\nu^-$ , and that for each  $x \in \Omega_\nu \setminus \mathcal{V}_\nu$ , the Jacobian of  $\psi$  satisfies  $\psi'(x) = \mathbf{v}'_\nu(\mathbf{w}'_\nu)^\top H(\langle \mathbf{w}'_\nu, x \rangle + b'_\nu) + \mathbf{d}'$  with some  $\mathbf{d}' \in \mathbb{R}^k$ , hence for each  $x_\nu^+ \in \Omega_\nu^+$ ,  $x_\nu^- \in \Omega_\nu^-$

$$\psi'(x_\nu^+) - \psi'(x_\nu^-) = \mathbf{v}'_\nu(\mathbf{w}'_\nu)^\top (H(\langle \mathbf{w}'_\nu, x_\nu^+ \rangle + b'_\nu) - H(\langle \mathbf{w}'_\nu, x_\nu^- \rangle + b'_\nu)).$$

Since  $(\mathbf{w}'_\nu, b'_\nu) = \lambda_\nu(\mathbf{w}_\nu, b_\nu)$  and  $\text{sign}(\langle \mathbf{w}_\nu, x_\nu^\pm \rangle + b_\nu) = \pm 1$ , we have

$$H(\langle \mathbf{w}'_\nu, x_\nu^+ \rangle + b'_\nu) - H(\langle \mathbf{w}'_\nu, x_\nu^- \rangle + b'_\nu) = \text{sign}(\lambda_\nu).$$

Moreover, as  $\psi = \varphi$  on  $\mathcal{X}$ , we have  $\varphi'(x_\nu^+) - \varphi'(x_\nu^-) = \psi'(x_\nu^+) - \psi'(x_\nu^-)$ , hence

$$\mathbf{v}_\nu \mathbf{w}_\nu^\top = \mathbf{v}'_\nu(\mathbf{w}'_\nu)^\top \text{sign}(\lambda_\nu).$$

Since  $\mathbf{w}'_\nu = \lambda_\nu \mathbf{w}_\nu$ , this simplifies to

$$(40) \quad \mathbf{v}_\nu \mathbf{w}_\nu^\top = \mathbf{v}'_\nu(\mathbf{w}'_\nu)^\top \text{sign}(\lambda_\nu) = \mathbf{v}'_\nu \mathbf{w}_\nu^\top \lambda_\nu \text{sign}(\lambda_\nu) = |\lambda_\nu| \mathbf{v}'_\nu \mathbf{w}_\nu^\top$$

Hence,  $\mathbf{v}'_\nu = \mathbf{v}_\nu / |\lambda_\nu|$  for each  $\nu \in N_1$ .

To conclude, it is enough to prove that  $\lambda_\nu > 0$  for every  $\nu \in N_\nu$ . Using (40), we can re-write the equality  $\varphi(x) = \psi(x)$  for every  $x$  as follows:

$$(41) \quad \sum_{\nu \in N_1} \mathbf{v}_\nu \left[ \text{ReLU}(\langle \mathbf{w}_\nu, x \rangle + b_\nu) - |\lambda_\nu|^{-1} \text{ReLU}(\underbrace{\langle \mathbf{w}'_\nu, x \rangle + b'_\nu}_{\lambda_\nu(\langle \mathbf{w}_\nu, x \rangle + b_\nu)}) \right] + \mathbf{c} - \mathbf{c}' = 0.$$

Now, we observe that

$$(42) \quad \text{ReLU}(\langle \mathbf{w}_\nu, x \rangle + b_\nu) - |\lambda_\nu|^{-1} \text{ReLU}(\lambda_\nu(\langle \mathbf{w}_\nu, x \rangle + b_\nu)) = \begin{cases} 0 & \text{if } \text{sign}(\lambda_\nu) = 1 \\ \langle \mathbf{w}_\nu, x \rangle + b_\nu & \text{if } \text{sign}(\lambda_\nu) = -1 \end{cases}$$

We now show that  $T := \{\nu \in N_1 \mid \text{sign}(\lambda_\nu) = -1\} = \emptyset$ . Using (42), we re-write (41) as:

$$(43) \quad \sum_{\nu \in T} \mathbf{v}_\nu (\langle \mathbf{w}_\nu, x \rangle + b_\nu) + \mathbf{c} - \mathbf{c}' = 0.$$

Since this is valid for all  $x \in \mathbb{R}^{N_0}$  we get  $\mathbf{c} = \mathbf{c}'$  and  $\sum_{\nu \in T} \mathbf{v}_\nu \mathbf{w}_\nu^\top = 0$ . In light of (38) the latter implies  $T = \emptyset$ , hence  $\text{sign}(\lambda_\nu) = 1$  for all  $\nu \in N_1$ .  $\square$

**5.4. Local S-identifiability despite the presence of twins.** It is natural to wonder if there exists shallow networks with twins that are nevertheless either non-degenerate, or locally S-identifiable, or PS-identifiable. Positive twins are excluded (for any network depth) by Lemma 3, hence we can focus on the case where there are  $K \geq 1$  nontrivial classes of twins, each made of a single pair of (distinct) negative twins (as any equivalence class with at least three twins necessarily contains two positive ones). We detail here the case  $K = 1$  and leave to future work a more detailed analysis of what happens for  $K \geq 2$ .

**Lemma 16 (Single pair of negative twins).** *Consider a shallow network architecture. If  $\theta \in \Theta \subseteq \mathbb{R}^{E \cup \bar{H}}$  is admissible with a single pair of negative twins,  $\{\nu_1, \nu_2\} \subseteq H$ , then*

$$(44) \quad \mathbf{A}^\perp(\theta) = \{0\} \text{ and } \bar{\mathbf{A}}^\perp(\theta) = \text{span}\{\delta_{\nu_1} + \delta_{\nu_2} - \delta_\star\} \neq \{0\}$$

where  $\delta_\nu \in \mathbb{R}^{H+1}$ ,  $\nu \in H$  is the  $\nu$ -th canonical eigenvector, and  $\delta_\star = (\mathbf{0}_H, 1)$ .

Moreover, if at least one of the following conditions holds:

- i)  $\mathbf{w}_{\nu_1 \rightarrow \bullet}, \mathbf{w}_{\nu_2 \rightarrow \bullet}$  are linearly independent (which is only possible if  $|N_2| \geq 2$ ); or
- ii)  $\Theta$  is contained in the set of parameters with zero output bias;

then  $\theta$  is non-degenerate with respect to  $\Theta$ . Conversely, if

- iii)  $\mathbf{w}_{\nu_1 \rightarrow \bullet}, \mathbf{w}_{\nu_2 \rightarrow \bullet}$  are linearly dependent and  $\theta$  belongs to the interior of  $\Theta$ ,

then  $\theta$  is degenerate with respect to  $\Theta$ .

**Remark 7.** *Inspecting the proof shows that the assumption in iii) that  $\theta$  is in the interior of  $\Theta$  can be relaxed to: for small enough  $\epsilon$ , each parameter  $\theta' \in B(\theta, \epsilon)$  differing from  $\theta$  only in terms of biases belongs to  $\Theta \cap B(\theta, \epsilon)$ .*

The proof is in Appendix H. We are now equipped to show with an example that non-degeneracy and local-identifiability are distinct concepts.

**Example 4 (Absolute value).** *Consider a shallow architecture with scalar input and output and two hidden neurons. The absolute value can be written as  $|x| = \text{ReLU}(x) + \text{ReLU}(-x) = \mathbf{R}_\theta$  where  $\theta = (w_{\mu \rightarrow \nu_1} = 1, w_{\mu \rightarrow \nu_2} = -1, b_{\nu_1} = b_{\nu_2} = 0, w_{\nu_1 \rightarrow \eta} = w_{\nu_2 \rightarrow \eta} = 1, b_\eta = 0)$  has a single pair of negative twins. This parameter  $\theta$  satisfies the following properties*

- i) *it is not PS-identifiable from any bounded set  $\mathcal{X} \subset \mathbb{R}$  (by Lemma 3);*
- ii) *it is not locally S-identifiable from any finite  $F \subset \mathcal{X}_\theta$  (i.e., it is degenerate, see below);*
- iii) *it is PS-identifiable (hence locally S-identifiable) from  $\mathcal{X} = \mathbb{R}$ ;*
- iv) *it is locally S-identifiable from  $F \cup \{0\}$  for some finite set  $F \subset \mathcal{X}_\theta$ ;*

The last two points are detailed in Appendix I. Let us detail ii) here. Since  $|N_2| = 1$ , by Lemma 16-iii) we get that  $\theta$  is degenerate with respect to  $\Theta = \mathbb{R}^{E \cup \bar{H}}$ , i.e. not locally S-identifiable from any finite  $F \subseteq \mathcal{X}_\theta$ . Indeed, if  $F \subseteq \mathcal{X}_\theta = \mathbb{R} \setminus \{0\}$  is finite then  $F \subset (-\infty, -t] \cup [t, +\infty)$  for some  $t > 0$ , and  $\text{abs}$  coincides on  $F$  with (see Figure 2-(c))

$$\text{ReLU}(x - t) + \text{ReLU}(-(x + t)) + t = \begin{cases} -x, & x \leq -t \\ t, & |x| \leq t = \mathbf{R}_{\theta'}(x) \\ x, & x \geq t. \end{cases}$$

where  $\theta'$  has nonzero biases, so that  $\theta' \not\sim_{PS} \theta$ .

**Example 5 (Revisiting the identity function from Example 1).** *The identity function from Example 1 is another example with a single pair of twin neurons. With  $\Theta = \mathbb{R}^{E \cup \bar{H}}$  the parameter  $\theta_0$  is not locally S-identifiable (from  $\mathcal{X} = \mathbb{R}$ ) as already explained in Example 1. With  $\Theta = \Theta_0 \subsetneq \mathbb{R}^{E \cup \bar{H}}$  the set of parameters with zero output bias,  $\theta_0$  is on the contrary PS-identifiable from  $\mathbb{R}$  (see details in Appendix J). It can also be shown that  $\theta_0$  is non-degenerate with respect to  $\Theta_0$ , using arguments similar to those used in Appendix I to prove*

item iv) of Example 4. This illustrates the fact that, in the presence of a pair of negative twins, many things can happen: the parameter can be PS-identifiable and non-degenerate, or not even locally S-identifiable.

**5.5. Discussion of the role of activation spaces.** For shallow irreducible networks, PS-identifiability from a bounded set is equivalent (Theorem 3) to the absence of twin neurons, which corresponds (Lemma 15) to a completeness property of the activation space that reads  $\bar{\mathbf{A}}^\perp(\theta) = \{0\}$ . The property  $\bar{\mathbf{A}}^\perp(\theta) = \{0\}$  also implies non-degeneracy (Lemma 14), yet a consequence of Lemma 16 is that the converse does not generally hold (and that the weaker assumption  $\mathbf{A}^\perp(\theta) = \{0\}$  is no longer sufficient to imply non-degeneracy). An exception occurs for scalar-valued shallow networks.

**Lemma 17.** *Consider a scalar-valued shallow architecture ( $|N_2| = 1$ ). If  $\theta$  belongs to the interior of  $\Theta \subseteq \mathbb{R}^{E \cup H}$  and is non-degenerate with respect to  $\Theta$  then  $\bar{\mathbf{A}}^\perp(\theta) = \{0\}$ .*

This exception is a consequence of the following result.

**Lemma 18.** *Consider a scalar-valued shallow network architecture. If  $\theta$  is admissible then there is  $0 < C < \infty$  such that: for each  $\mathbf{z} \in \mathbb{R}^{\mathcal{Q}+1}$ , there exists  $\theta' \in B(\theta, C\|\mathbf{z}\|_\infty)$*

$$(45) \quad \Phi_\eta^i(\theta') - \Phi_\eta^i(\theta) = \mathbf{0}_{\mathcal{Q}_1 \times N_0},$$

$$(46) \quad \Phi_\eta^h(\theta') - \Phi_\eta^h(\theta) = \mathbf{z}$$

where  $\eta$  is the single output neuron constituting the output layer  $N_L$ . The parameters  $\theta$  and  $\theta'$  differ only in terms of biases.

*Proof.* Write  $\mathbf{z} = (\mathbf{y}, \gamma)$  with  $\mathbf{y} \in \mathbb{R}^{\mathcal{Q}} = \mathbb{R}^H$  and  $\gamma \in \mathbb{R}$ . For each hidden neuron  $\nu \in H = N_1$ , denote  $v_\nu = v_{\nu \rightarrow \eta}$  the unique weight from neuron  $\nu$  to the single output neuron. For each input neuron  $\mu \in N_0$ , the  $\mu$ -th column of  $\Phi_\eta^i(\theta)$  is  $\Phi_{\mu \rightarrow \eta}^i(\theta) := (w_{\mu \rightarrow \nu} v_\nu)_{\nu \in H}$ , and  $\Phi_\eta^h(\theta) = ((b_\nu v_\nu)_{\nu \in H}, b_\eta)^\top$ . To prove the result we define  $\theta'$  with identical weights as  $\theta$ ,  $v'_\nu := v_\nu$ ,  $w'_{\mu \rightarrow \nu} := w_{\mu \rightarrow \nu}$ , and set the output bias to  $b'_\eta := b_\eta + \gamma$ . This implies  $w'_{\mu \rightarrow \nu} v'_\nu = w_{\mu \rightarrow \nu} v_\nu$  for every  $\nu \in H, \mu \in N_0$ , hence  $\Phi_\eta^i(\theta') = \Phi_\eta^i(\theta)$ . We now seek  $b'_\nu$  such that  $b'_\nu v'_\nu = b_\nu v_\nu + y_\nu$ , for each  $\nu \in H$ . Since  $\theta$  is admissible,  $v_\nu \neq 0$  for all  $\nu \in H$ , hence we can choose  $b'_\nu := b_\nu + y_\nu / v_\nu$ . We conclude with  $C_\theta$  driven by  $\min_\nu 1/|v_\nu|$ .  $\square$

*Proof of Lemma 17.* We prove the contraposition. Assume that  $\theta$  is admissible, that it belongs to the interior of  $\Theta \subset \mathbb{R}^{E \cup H}$ , and that  $\bar{\mathbf{A}}^\perp(\theta) \neq \{0\}$ . Since  $\theta$  is in the interior of  $\Theta$ , there is  $\eta > 0$  such that  $B(\theta, \eta) \subseteq \Theta$ . For each  $0 < \epsilon < \eta$  there exists  $\mathbf{z} \in \bar{\mathbf{A}}^\perp(\theta)$  with norm  $\|\mathbf{z}\|_\infty = \epsilon/C$  where  $C$  is the constant from Lemma 18. Since  $\theta$  is admissible, by Lemma 18 and the characterization of  $\mathbf{V}(\theta)$  (Lemma 8) there exists  $\theta' \in B(\theta, C\|\mathbf{z}\|_\infty) = B(\theta, \epsilon) = \Theta \cap B(\theta, \epsilon)$  such that  $\mathbf{0} \neq \Phi(\theta') - \Phi(\theta) \in \mathbf{V}(\theta)$ . This shows that  $\theta$  is degenerate with respect to  $\Theta$ .  $\square$

**Remark 8.** *The assumption that  $\theta$  is in the interior of  $\Theta$  can be relaxed to: each parameter  $\theta' \in B(\theta, \epsilon)$  differing from  $\theta$  only in terms of biases belongs to  $\Theta \cap B(\theta, \epsilon)$ .*

## ACKNOWLEDGEMENTS

The authors are thankful to François Malgouyres for the interesting discussions on invariant embeddings of linear and ReLU networks we had at different stages of advancement of this work and Joachim Bona-Pellissier for his technical comments on early versions of this article. The authors thank Elisa Riccietti for her feedback that helped a lot improve the readability of this paper. The authors are thankful to Hervé Jégou and Benjamin Graham for their continuous support since the genesis of this project a few years ago.

## REFERENCES

- [1] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, 2012.
- [2] I. Daubechies, R. DeVore, S. Foucart, B. Hanin, and G. Petrova. Nonlinear approximation and (deep) relu networks, 2019.
- [3] Ronald DeVore, Boris Hanin, and Guergana Petrova. Neural network approximation, 2020.
- [4] Héctor J. Sussmann. Uniqueness of the weights for minimal feedforward nets with a given input-output map. *Neural Networks*, 1992.
- [5] Paul Kainen, Vera Kurková, Vladik Kreinovich, and Ongard Sirisengtaksin. Uniqueness of network parameterizations and faster learning. *Preprint*, 1994.
- [6] Věra Kůrková and Paul C. Kainen. Functionally equivalent feedforward neural networks. *Neural Comput.*, 1993.
- [7] Francesca Albertini, Eduardo D. Sontag, and Vincent Maillot. Uniqueness of weights for neural networks. In *Artificial Neural Networks with Applications in Speech and Vision*, 1993.
- [8] Charles Fefferman. Reconstructing a neural net from its output. *Revista Matemática Iberoamericana*, 1994.
- [9] David Rolnick and Konrad P. Kording. Reverse-engineering deep relu networks, 2019.
- [10] Massimo Fornasier, Timo Klock, and Michael Rauchensteiner. Robust and resource efficient identification of two hidden layer neural networks, 2019.
- [11] Mary Phuong and Christoph H Lampert. Functional vs. parametric equivalence of ReLU networks. *ICLR*, 2020.
- [12] Francois Malgouyres and Joseph Landsberg. Multilinear compressive sensing and an application to convolutional linear networks. *SIAM*, 2018.
- [13] Francois Malgouyres. On the stable recovery of deep structured linear networks under sparsity constraints. *Proceedings of Machine Learning Research*, 2020.
- [14] Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. Norm-Based Capacity Control in Neural Networks. *Journal of Machine Learning Research*, 2015.
- [15] Behnam Neyshabur, Ruslan Salakhutdinov, and Nathan Srebro. Path-SGD - Path-Normalized Optimization in Deep Neural Networks. *NIPS*, 2015.

- [16] Qi Meng, Shuxin Zheng, Huishuai Zhang, Wei Chen 0034, Qiwei Ye, Zhi-Ming Ma, Nenghai Yu, and Tie-Yan Liu. G-SGD - Optimizing ReLU Neural Networks in its Positively Scale-Invariant Space. *ICLR*, 2019.
- [17] Mingyang Yi, Qi Meng, Wei Chen, Zhi-ming Ma, and Tie-Yan Liu. Positively Scale-Invariant Flatness of ReLU Neural Networks. *arXiv:1903.02237 [cs, stat]*, March 2019. arXiv: 1903.02237.
- [18] Pierre Stock, Benjamin Graham, Remi Gribonval, and Hervé Jégou. Equi-normalization of Neural Networks. In *ICLR 2019 - Seventh International Conference on Learning Representations*, pages 1–20, New Orleans, United States, May 2019.
- [19] Qunyang Yuan and Nanfeng Xiao. Scaling-Based Weight Normalization for Deep Neural Networks. *IEEE Access*, 7:7286–7295, January 2019. Publisher: IEEE.
- [20] Eldad Meller, Alexander Finkelstein, Uri Almog, and Mark Grobman. Same, Same But Different - Recovering Neural Network Quantization Error Through Weight Factorization. *arXiv:1902.01917 [cs, stat]*, February 2019. arXiv: 1902.01917.
- [21] Markus Nagel, Mart van Baalen, Tijmen Blankevoort, and Max Welling. Data-Free Quantization Through Weight Equalization and Bias Correction. *arXiv:1906.04721 [cs, stat]*, November 2019. arXiv: 1906.04721.
- [22] Nicholas Carlini, Matthew Jagielski, and Ilya Mironov. Cryptanalytic extraction of neural network models, 2020.
- [23] Behnam Neyshabur, Ruslan Salakhutdinov, and Nathan Srebro. Path-sgd: Path-normalized optimization in deep neural networks. *arXiv preprint arXiv:1506.02617*, 2015.
- [24] Markus Nagel, Mart van Baalen, Tijmen Blankevoort, and Max Welling. Data-free quantization through weight equalization and bias correction, 2019.
- [25] Eldad Meller, Alexander Finkelstein, Uri Almog, and Mark Grobman. Same, same but different - recovering neural network quantization error through weight factorization, 2019.
- [26] Mingyang Yi, Qi Meng, Wei Chen, Zhi ming Ma, and Tie-Yan Liu. Positively scale-invariant flatness of ReLU neural networks, 2019.
- [27] Qi Meng, Shuxin Zheng, Huishuai Zhang, Wei Chen, Zhi-Ming Ma, and Tie-Yan Liu. G-sgd: Optimizing relu neural networks in its positively scale-invariant space, 2018.
- [28] Boris Hanin and David Rolnick. Deep relu networks have surprisingly few activation patterns, 2019.
- [29] Razvan Pascanu, Guido Montufar, and Yoshua Bengio. On the number of response regions of deep feed forward networks with piece-wise linear activations, 2013.
- [30] Guido Montúfar, Razvan Pascanu, Kyunghyun Cho, and Yoshua Bengio. On the number of linear regions of deep neural networks. In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2*, 2014.
- [31] Maithra Raghu, Ben Poole, Jon Kleinberg, Surya Ganguli, and Jascha Sohl-Dickstein. On the expressive power of deep neural networks. In *Proceedings of the 34th International Conference on Machine Learning*, 2017.
- [32] Mary Phuong and Christoph H Lampert. Functional vs. parametric equivalence of relu networks. In *International Conference on Learning Representations*, 2019.



- [33] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. In *Proceedings of the IEEE*, volume 86, pages 2278–2324, 1998.
- [34] Jeremy Bernstein, Jiawei Zhao, Markus Meister, Ming-Yu Liu, Anima Anandkumar, and Yisong Yue. Learning compositional functions via multiplicative weight updates. *Advances in neural information processing systems*, 33, 2020.
- [35] David Balduzzi, Marcus Frean, Lennox Leary, JP Lewis, Kurt Wan-Duo Ma, and Brian McWilliams. The shattered gradients problem: If resnets are the answer, then what is the question?, 2018.

#### APPENDIX A. PROOF OF LEMMA 7

To lighten notations we omit the dependence of  $W, V$  and  $\mathbf{S}$  on  $\theta$ . The proof follows three steps.

1)  $\mathbf{S}$  is well-defined. Consider  $\alpha \in W$  and  $\nu \in H$ . First, since  $\theta$  is admissible, there exists at least one path  $p$  connecting  $\nu$  to some output neuron  $\eta$  through edges  $e \in E \cap \text{supp}(\theta)$ . We wish to show that if  $p$  and  $p'$  are two such partial paths then  $\sum_{e \in p} \alpha_e = \sum_{e \in p'} \alpha_e$ . Since  $\theta$  is admissible, there exists a partial path  $\bar{p}$  going from some input neuron  $\mu$  to  $\nu$  through edges  $e \in E \cap \text{supp}(\theta)$ . Since  $\bar{p}_\ell = p_\ell = p'_\ell$ , define by concatenation the full paths  $q = (\bar{p}_0, \dots, \bar{p}_\ell, p_{\ell+1}, \dots, p_L)$  and  $q' = (\bar{p}_0, \dots, \bar{p}_\ell, p'_{\ell+1}, \dots, p'_L)$ . As  $q$  and  $q'$  have all their edges in  $\text{supp}(\theta)$ , we have  $\Phi_q(\theta) \neq 0 \neq \Phi_{q'}(\theta)$ . Since  $\alpha \in W$  and  $q, q' \in \mathcal{P}_0$ , we have  $\Phi_q(\theta) \cdot (\mathbf{P}\alpha)_q = \Phi_{q'}(\theta) \cdot (\mathbf{P}\alpha)_{q'} = 0$ , hence  $(\mathbf{P}\alpha)_q = (\mathbf{P}\alpha)_{q'} = 0$  and

$$\sum_{e \in \bar{p}} \alpha_e + \sum_{e \in p} \alpha_e = \sum_{e \in q} \alpha_e = (\mathbf{P}\alpha)_q = 0 = (\mathbf{P}\alpha)_{q'} = \sum_{e \in \bar{p}} \alpha_e + \sum_{e \in p'} \alpha_e.$$

Thus,  $\sum_{e \in q} \alpha_e = \sum_{e \in q'} \alpha_e$  and  $\mathbf{S}$  is well-defined.

2)  $\mathbf{S}$  is injective on  $V$ . Note that  $\mathbf{S}$  is linear hence it is sufficient to show that its kernel is reduced to zero. Let  $\alpha \in V$  such that  $\mathbf{S}\alpha = 0$ . Since  $\alpha_{\bar{H}} = 0$  and  $\text{supp}(\alpha) \subseteq \text{supp}(\theta)$ , we have  $\alpha_e = 0$  for each  $e \in E \setminus \text{supp}(\theta)$ , hence it is sufficient to show  $\alpha_e = 0$  for any edge  $e = \mu \rightarrow \nu \in E \cap \text{supp}(\theta)$ . Since  $\theta$  is admissible, there is a partial path  $\bar{p}$  going from  $\nu$  to some output neuron  $\eta$  through edges  $e' \in E \cap \text{supp}(\theta)$ . Since  $e \in E \cap \text{supp}(\theta)$ , the extended path  $p := \mu \rightarrow \bar{p}$  also has all its edges in  $E \cap \text{supp}(\theta)$  and joins  $\mu$  to an output neuron. We distinguish three cases: in the first case,  $\mu, \nu$  are two hidden neurons, and

$$\alpha_e = \sum_{e' \in p} \alpha_{e'} - \sum_{e' \in \bar{p}} \alpha_{e'} = -(\mathbf{S}\alpha)_\mu + (\mathbf{S}\alpha)_\nu = 0.$$

In the second case,  $\mu \in N_0$  is an input neuron, hence  $p$  is a full path with edges  $e' \in E \cap \text{supp}(\theta)$ , so that  $\Phi_p(\theta) \neq 0$ . Since  $\alpha \in W$  it follows that  $\sum_{e' \in p} \alpha_{e'} = (\mathbf{P}\alpha)_p = 0$  and we also obtain  $\alpha_e = 0$ . Finally, in the third case,  $\nu \in N_L$  is an output neuron hence  $\bar{p} = (\nu)$  contains no edge, so that  $\sum_{e' \in \bar{p}} \alpha_{e'} = 0$  and we get  $\alpha_e = 0$  as well.

c)  $\mathbf{S}$  is surjective. Consider  $\beta \in \mathbb{R}^H$ , and  $\alpha \in \mathbb{R}^{E \cup \bar{H}}$  as defined around (18). Consider  $p = (p_\ell, \dots, p_L) \in \mathcal{P}$  (with  $0 \leq \ell \leq L-1$ ) a full or partial path going from some neuron  $\mu = p_\ell \in N_0 \cup H$  to an arbitrary output neuron  $\eta = p_L \in N_L$  through edges  $e \in E \cap \text{supp}(\theta)$ . In the case of a full path,  $\mu \in N_0$  and

$$\sum_{e \in p} \alpha_e = -\beta_{p_{L-1}} + \sum_{j=1}^{L-1} (\beta_{p_j} - \beta_{p_{j-1}}) + \beta_{p_0} = 0.$$

As this holds for any full path with edges  $e \in \text{supp}(\theta)$ , and since  $\alpha_{\bar{H}} = 0$ , we get  $\alpha \in W$ . Since  $\alpha_e = 0$  for  $e \in E \setminus \text{supp}(\theta)$  we have indeed  $\text{supp}(\alpha) \subseteq \text{supp}(\theta)$  hence  $\alpha \in V$ .

In the case of a partial path ( $\ell \geq 1$ ), we have  $\mu \in H$  and similarly

$$\sum_{e \in p} \alpha_e = -\beta_{p_{L-1}} + \sum_{j=\ell+1}^{L-1} (\beta_{p_j} - \beta_{p_{j-1}}) = -\beta_{p_\ell} = -\beta_\mu$$

hence  $(\mathbf{S}\alpha)_\mu = \beta_\mu$ . As this holds for any  $\mu \in H$ , this shows that  $(\mathbf{S}\alpha) = \beta$ .

## APPENDIX B. PROOF OF THEOREM 2

Denote  $\eta = \min_{i \in \text{supp}(\theta)} |\theta_i|$ . Assume by contradiction that  $\theta$  is not locally S-identifiable from  $\mathcal{X}$  with respect to  $\Theta$ . This implies that for each  $n \geq 1$  there is  $\theta_n \in \Theta \cap B(\theta, \min(\eta, 1/n))$  which is not scaling-equivalent to  $\theta$  such that  $\mathbf{R}_{\theta_n} = \mathbf{R}_\theta$  on  $\mathcal{X}$ . For  $n \geq 1/\epsilon$ , since  $\theta$  is PS-identifiable from  $\mathcal{X}$  with respect to  $\Theta$  and since  $\theta_n \in \Theta$  satisfies  $\mathbf{R}_{\theta_n} = \mathbf{R}_\theta$  on  $\mathcal{X}$ , we have  $\theta_n \sim_{PS} \theta$ , hence there is a permutation  $\pi_n \in \mathfrak{S}_G$  such that  $\pi_n \circ \theta_n \sim_S \theta$ , hence by Theorem 1

$$(47) \quad \text{sign}(\pi_n \circ \theta_n) = \text{sign}(\theta), \quad \text{and} \quad \Phi(\pi_n \circ \theta_n) = \Phi(\theta).$$

Since the set of permutations is finite, there exists  $\pi \in \mathfrak{S}_G$ , and an increasing subsequence  $n_k$  such that  $\pi_{n_k} = \pi$  for each  $k \geq 1$ . By (47), for every  $k$  we have

$$(48) \quad \text{sign}(\pi \circ \theta_{n_k}) = \text{sign}(\theta) \quad \text{and} \quad \Phi(\pi \circ \theta_{n_k}) = \Phi(\theta).$$

There is a permutation matrix  $\Pi \in \mathbb{R}^{P \times P}$  such that  $\Phi(\pi \circ \theta') = \Pi \Phi(\theta')$  for all  $\theta' \in \mathbb{R}^{E \cup \bar{H}}$ . Since  $\lim_{k \rightarrow \infty} \theta_{n_k} = \theta$ , by continuity of  $\Phi$  we obtain  $\Pi \Phi(\theta) = \Phi(\pi \circ \theta) = \Phi(\theta)$  hence for every  $k \geq 1$

$$(49) \quad \Phi(\theta) = \Pi^{-1} \Phi(\theta) \stackrel{(48)}{=} \Pi^{-1} \Pi \Phi(\theta'_{n_k}) = \Phi(\theta_{n_k}).$$

Since  $\theta$  is admissible, by Corollary 1, the equality  $\Phi(\theta_{n_k}) = \Phi(\theta)$  implies  $\text{supp}(\theta_{n_k}) = \text{supp}(\theta)$ . This implies that  $\text{sign}((\theta_{n_k})_i) = 0 = \text{sign}(\theta_i)$  for each  $i \notin \text{supp}(\theta)$ . Since  $\theta_{n_k} \in B(\theta, \eta)$  for each  $k$ , we also have  $\text{sign}((\theta'_{n_k})_i) = \text{sign}(\theta_i) \in \{-1, 1\}$  for every  $i \in \text{supp}(\theta)$ , hence  $\text{sign}(\theta_{n_k}) = \text{sign}(\theta)$  for every  $k$ . Since  $\theta$  is admissible, by Theorem 1, the fact that  $\text{sign}(\theta_{n_k}) = \text{sign}(\theta)$  and  $\Phi(\theta_{n_k}) = \Phi(\theta)$  implies  $\theta'_n \sim_S \theta$ . This contradicts our assumption that  $\theta'_n \not\sim_S \theta$  for every  $n$ .

## APPENDIX C. PROOF OF LEMMA 3

We will prove the contraposition using the following observation.

**Fact 4.** Consider  $\alpha, \beta \in \mathbb{R}$  and  $M > 0$ . For any  $t \in [-M, \infty)$  we have

$$\alpha \text{ReLU}(t) + \beta \text{ReLU}(-t) = \begin{cases} \alpha t, & t \geq 0 \\ -\beta t, & t \leq 0 \end{cases} = (\alpha + \beta) \text{ReLU}(t) - \beta \text{ReLU}(t + M) + M.$$

Assuming that  $\theta$  has twins, consider a hidden layer  $1 \leq \ell \leq L - 1$  and  $T \subseteq N_\ell$  a pair of twin neurons  $T = \{\nu_1, \nu_2\}$ . Denote  $\mathbf{w}_i = \mathbf{w}_{\bullet \rightarrow \nu_i}$ ,  $b_i = b_{\nu_i}$ ,  $\mathbf{v}_i = \mathbf{w}_{\nu_i \rightarrow \bullet}$ . As these neurons are twins, there is  $\lambda \in \mathbb{R}$  such that for every  $x \in \mathbb{R}^{N_0}$ ,

$$z_{\nu_2}(\theta, x) = \langle \mathbf{w}_2, y_{\ell-1}(\theta, x) \rangle + b_2 = \lambda (\langle \mathbf{w}_1, y_{\ell-1}(\theta, x) \rangle + b_1) = \lambda z_{\nu_1}(\theta, x)$$

In the case of positive twins we have  $\lambda > 0$ , hence  $y_{\nu_1}(\theta, x) = \lambda y_{\nu_2}(\theta, x)$  for every  $x$ . Given  $\epsilon > 0$  consider  $\theta'(\epsilon)$  obtained by keeping unchanged all weights and biases in  $\theta$  except the weights outgoing from neurons  $\nu_i$ ,  $i = 1, 2$ :  $\mathbf{v}'_1 = \mathbf{v}_1 + \lambda \epsilon \mathbf{1}_{N_{\ell+1}}$ ,  $\mathbf{v}'_2 = \mathbf{v}_2 - \epsilon \mathbf{1}_{N_{\ell+1}}$ . Since the linear layers and biases of hidden neurons up to layer  $\ell$  are unchanged, we have  $\mathbf{y}_\ell(\theta, x) = \mathbf{y}_\ell(\theta', x)$  for all  $x$ . For every neuron  $\nu \in N_\ell \setminus T$ , since the outgoing weights are unchanged, we get for every  $x$

$$y_\nu(\theta', x) \mathbf{w}'_{\nu \rightarrow \bullet} = y_\nu(\theta, x) \mathbf{w}_{\nu \rightarrow \bullet}.$$

Moreover, since  $y_{\nu_2}(\theta, x) = \lambda y_{\nu_1}(\theta, x)$  and  $\mathbf{y}_\ell(\theta, x) = \mathbf{y}_\ell(\theta', x)$ , we obtain for every  $x$

$$\begin{aligned} y_{\nu_1}(\theta', x) \mathbf{w}'_{\nu_1 \rightarrow \bullet} + y_{\nu_2}(\theta', x) \mathbf{w}'_{\nu_2 \rightarrow \bullet} &= y_{\nu_1}(\theta, x) \mathbf{v}'_1 + y_{\nu_2}(\theta, x) \mathbf{v}'_2 \\ &= y_{\nu_1}(\theta, x) (\mathbf{v}_1 + \lambda \epsilon \mathbf{1}_{N_{\ell+1}}) + \lambda y_{\nu_1}(\theta, x) (\mathbf{v}_2 - \epsilon \mathbf{1}_{N_{\ell+1}}) \\ &= y_{\nu_1}(\theta, x) \mathbf{v}_1 + \lambda y_{\nu_1}(\theta, x) \mathbf{v}_2 \\ &= y_{\nu_1}(\theta, x) \mathbf{w}_{\nu_1 \rightarrow \bullet} + y_{\nu_2}(\theta, x) \mathbf{w}_{\nu_2 \rightarrow \bullet}. \end{aligned}$$

Summing over all hidden neurons we obtain  $\mathbf{z}_{\ell+1}(\theta, x) = \mathbf{z}_{\ell+1}(\theta', x)$  for every  $x$ , and since all the next affine layers are unchanged, we obtain  $\mathbf{R}_{\theta'} = \mathbf{R}_\theta$ . It is not difficult to check that  $\theta' = \theta'(\epsilon)$  is not scaling equivalent to  $\theta$  and can be made arbitrarily close to it. This shows that  $\theta$  is not locally S-identifiable from  $\mathbb{R}^{N_0}$ . By contraposition, if  $\theta$  is locally S-identifiable from  $\mathbb{R}^{N_0}$  then it has no positive twins.

In the case of negative twins we have  $y_{\nu_1}(\theta, x) = \text{ReLU}(t)$  and  $y_{\nu_2}(\theta, x) = |\lambda| \text{ReLU}(-t)$  with  $t = t(x) := z_{\nu_1}(\theta, x) = \langle \mathbf{w}_1, y_{\ell-1}(\theta, x) \rangle + b_1$ . Since  $\mathcal{X}$  is bounded there is some finite  $M > 0$  such that  $|z_{\nu_1}(\theta, x)| \leq M$  for every  $x \in \mathcal{X}$ . Consider  $\theta'$  obtained by keeping all weights and biases unchanged from  $\theta$  except the incoming and outgoing weights of  $\nu_1, \nu_2$ , their biases, and the biases of the neurons of the next layer,  $\eta \in N_{\ell+1}$ , which are set as:

- $\mathbf{w}'_{\bullet \rightarrow \nu_1} = \mathbf{w}'_{\bullet \rightarrow \nu_2} = \mathbf{w}_{\bullet \rightarrow \nu_1}$ ;
- $b'_{\nu_1} = b_{\nu_1}$ ;  $b'_{\nu_2} = b_{\nu_1} + M$ ;
- $\mathbf{w}'_{\nu_1 \rightarrow \bullet} = \mathbf{w}_{\nu_1 \rightarrow \bullet} + |\lambda| \mathbf{w}_{\nu_2 \rightarrow \bullet}$ ;  $\mathbf{w}'_{\nu_2 \rightarrow \bullet} = -|\lambda| \mathbf{w}_{\nu_2 \rightarrow \bullet}$ ;
- $b'_\eta = b_\eta + M$

For each  $\eta \in N_{\ell+1}$ , using Fact 4 for  $\alpha = w_{\nu_1 \rightarrow \eta}$ ,  $\beta = |\lambda|w_{\nu_2 \rightarrow \eta}$  we obtain

$$\begin{aligned}
y_{\nu_1}(\theta, x)w_{\nu_1 \rightarrow \eta} + y_{\nu_2}(\theta, x)w_{\nu_2 \rightarrow \eta} &= \alpha \text{ReLU}(t) + \beta \text{ReLU}(-t) \\
&= (\alpha + \beta) \text{ReLU}(t) - \beta \text{ReLU}(t + M) + M \\
&= (w_{\nu_1 \rightarrow \eta} + |\lambda|w_{\nu_2 \rightarrow \eta}) \text{ReLU}(t) - |\lambda|w_{\nu_2 \rightarrow \eta} \text{ReLU}(t + M) + M \\
&= w'_{\nu_1 \rightarrow \eta} \text{ReLU}(t) + w'_{\nu_2 \rightarrow \eta} \text{ReLU}(t + M) + M \\
&= w'_{\nu_1 \rightarrow \eta} \text{ReLU}(z_{\nu_1}(\theta', x)) + w'_{\nu_2 \rightarrow \eta} \text{ReLU}(z_{\nu_2}(\theta', x)) + M \\
&= w'_{\nu_1 \rightarrow \eta} y_{\nu_1}(\theta', x) + w'_{\nu_2 \rightarrow \eta} y_{\nu_2}(\theta', x) + M.
\end{aligned}$$

Reasoning as in the case of positive twins we obtain  $\mathbf{z}_{\ell+1}(\theta', x) = \mathbf{z}_{\ell+1}(\theta, x)$  and eventually  $\mathbf{R}_{\theta'}(x) = \mathbf{R}_{\theta}(x)$  for every  $x$  in the bounded set  $\mathcal{X}$ . Since the sign of  $w_{\nu_2 \rightarrow \bullet}$  has changed,  $\theta'$  is not PS-equivalent to  $\theta$ . This shows that  $\theta$  is not PS-identifiable from  $\mathcal{X}$  with respect to  $\Theta = \mathbb{R}^{E \cup \bar{H}}$ .

By contraposition, assuming that  $\theta$  is PS-identifiable from a bounded set  $\mathcal{X}$  with respect to  $\Theta = \mathbb{R}^{E \cup \bar{H}}$ , there is no negative twin. Besides, by Theorem 2, such a  $\theta$  is also locally S-identifiable from  $\mathcal{X}$  with respect to  $\Theta = \mathbb{R}^{E \cup \bar{H}}$ , hence it is locally S-identifiable from  $\mathbb{R}^{N_0}$  with respect to  $\Theta = \mathbb{R}^{E \cup \bar{H}}$ . By the first part of the lemma, we conclude that  $\theta$  has no positive twins either. Hence, it has no twins.

#### APPENDIX D. PROOF OF LEMMA 4

We will use the following observation.

**Fact 5.** for  $\chi \in \{0, 1\}$  and  $e = (-1)^\chi$  we have  $\text{ReLU}(t) = \chi t + e \text{ReLU}(et)$  for every  $t \in \mathbb{R}$ .

Assume for the sake of contradiction that  $\theta$  is not irreducible:  $\mathbf{W}_{\ell+1} \mathbf{I}_T \mathbf{W}_{\ell} = 0$  for some non-empty  $T \subset N_{\ell}$  with some  $1 \leq \ell \leq L - 1$ . Denote  $\theta'$  a network with the same weights and biases as  $\theta$  except on layers  $\ell$  and  $\ell + 1$ , where  $\mathbf{W}'_{\ell}, \mathbf{W}'_{\ell+1}$  and  $\mathbf{b}'_{\ell}, \mathbf{b}'_{\ell+1}$  will soon be described. By an easy induction we have  $\mathbf{y}_{\ell'}(\theta', x) = \mathbf{y}_{\ell'}(\theta, x)$  for  $0 \leq \ell' \leq \ell - 1$ .

Defining  $\mathbf{J}_T = \text{diag}(e_{\nu})_{\nu \in N_{\ell}}$  with  $e_{\nu} = -1$  if  $\nu \in T$  and  $e_{\nu} = 1$  otherwise, we obtain from Fact 5 that for every vector  $\mathbf{z}_{\ell} \in \mathbb{R}^{N_{\ell}}$ ,  $\text{ReLU}(\mathbf{z}_{\ell}) = \mathbf{I}_T \mathbf{z}_{\ell} + \mathbf{J}_T \text{ReLU}(\mathbf{J}_T \mathbf{z}_{\ell})$ . Define  $\mathbf{W}'_{\ell} = \mathbf{J}_T \mathbf{W}_{\ell}$ ,  $\mathbf{W}'_{\ell+1} = \mathbf{W}_{\ell+1} \mathbf{J}_T$ ,  $\mathbf{b}'_{\ell} = \mathbf{J}_T \mathbf{b}_{\ell}$ . For each  $x \in \mathbb{R}^{N_0}$ , since  $\mathbf{y}_{\ell-1}(\theta, x) = \mathbf{y}_{\ell-1}(\theta', x)$  and  $\mathbf{J}_T^2 = \mathbf{I}_{\mathbb{R}^{N_{\ell}}}$  we get using the shorthands  $\mathbf{z}_i = \mathbf{z}_i(\theta, x)$ ,  $\mathbf{z}'_i = \mathbf{z}_i(\theta', x)$ ,  $i \in \{\ell, \ell + 1\}$

$$\begin{aligned}
\mathbf{z}_{\ell} &= \mathbf{W}_{\ell} \mathbf{y}_{\ell-1}(\theta, x) + \mathbf{b}_{\ell} = \mathbf{J}_T (\mathbf{J}_T \mathbf{W}_{\ell} \mathbf{y}_{\ell-1}(\theta, x) + \mathbf{J}_T \mathbf{b}_{\ell}) = \mathbf{J}_T \mathbf{z}'_{\ell} \\
\mathbf{z}_{\ell+1} &= \mathbf{W}_{\ell+1} \text{ReLU}(\mathbf{z}_{\ell}) + \mathbf{b}_{\ell+1} = \mathbf{W}_{\ell+1} (\mathbf{I}_T \mathbf{z}_{\ell} + \mathbf{J}_T \text{ReLU}(\mathbf{J}_T \mathbf{z}_{\ell})) + \mathbf{b}_{\ell+1} \\
&= \underbrace{\mathbf{W}_{\ell+1} \mathbf{I}_T \mathbf{W}_{\ell}}_{=0} \mathbf{y}_{\ell-1}(\theta, x) + \mathbf{W}_{\ell+1} \mathbf{I}_T \mathbf{b}_{\ell} + \mathbf{W}'_{\ell+1} \text{ReLU}(\mathbf{z}'_{\ell}) + \mathbf{b}_{\ell+1} \\
&= \mathbf{W}'_{\ell+1} \text{ReLU}(\mathbf{z}'_{\ell}) + (\mathbf{W}_{\ell+1} \mathbf{I}_T \mathbf{b}_{\ell} + \mathbf{b}_{\ell+1}).
\end{aligned}$$

Defining  $\mathbf{b}'_{\ell+1} := \mathbf{W}_{\ell+1} \mathbf{I}_T \mathbf{b}_{\ell} + \mathbf{b}_{\ell+1}$ , we get  $\mathbf{z}_{\ell+1}(\theta, x) = \mathbf{z}'_{\ell+1}(\theta', x)$  for all  $x$ . Since all other layers of  $\theta$  and  $\theta'$  are identical, an easy induction yields  $\mathbf{R}_{\theta} = \mathbf{R}_{\theta'}$ , where  $\theta' \in \mathbb{R}^{E \cup \bar{H}} = \Theta$ . To conclude, we prove below that  $\theta'$  is not PS-equivalent to  $\theta$ : this contradicts the assumption that  $\theta$  is PS-identifiable and concludes the proof.

For the sake of (yet another) contradiction, assume that  $\theta' \sim_{PS}$ , so that there exists diagonal matrices  $\Lambda_{\ell'} \in \mathbb{R}^{N_{\ell'} \times N_{\ell'}}$  with positive entries and permutation matrices  $\Pi_{\ell'} \in \mathbb{R}^{N_{\ell'} \times N_{\ell'}}$ ,  $0 \leq \ell' \leq L$ , such that  $\Lambda_0 = \Pi_0 = \mathbf{I}_{N_0}$ ,  $\Lambda_L = \Pi_L = \mathbf{I}_{N_L}$ ,  $\mathbf{W}'_{\ell'} = \Pi_{\ell'} \Lambda_{\ell'} \mathbf{W}_{\ell'} \Lambda_{\ell'-1}^{-1} \Pi_{\ell'-1}^{-1}$ , and  $\mathbf{b}'_{\ell'} = \Pi_{\ell'} \Lambda_{\ell'} \mathbf{b}_{\ell'}$  for every  $1 \leq \ell' \leq L$ . We show by induction that  $\Lambda_{\ell'} = \Pi_{\ell'} = \mathbf{I}_{N_{\ell'}}$  for every  $0 \leq \ell' < \ell$ . This trivially holds for  $\ell' = 0$ . If it holds for some  $\ell' < \ell - 1$  then, as  $(\mathbf{W}'_{\ell'+1}, \mathbf{b}'_{\ell'+1}) = (\mathbf{W}_{\ell'+1}, \mathbf{b}_{\ell'+1})$  by construction of  $\theta'$ , we have

$$\begin{aligned} (\mathbf{W}_{\ell'+1}, \mathbf{b}_{\ell'+1}) &= (\mathbf{W}'_{\ell'+1}, \mathbf{b}'_{\ell'+1}) = (\Pi_{\ell'+1} \Lambda_{\ell'+1} \mathbf{W}_{\ell'+1} \Lambda_{\ell'}^{-1} \Pi_{\ell'}^{-1}, \Pi_{\ell'+1} \Lambda_{\ell'+1} \mathbf{b}_{\ell'+1}) \\ &= \Pi_{\ell'+1} \Lambda_{\ell'+1} (\mathbf{W}_{\ell'+1}, \mathbf{b}_{\ell'+1}), \end{aligned}$$

i.e.,  $(\mathbf{w}_{\bullet \rightarrow \nu}, b_{\nu}) = \lambda_{\pi(\nu)}(\mathbf{w}_{\bullet \rightarrow \pi(\nu)}, b_{\pi(\nu)})$  for every  $\nu \in N_{\ell'+1}$ , with  $\pi$  the permutation of  $N_{\ell'+1}$  associated to  $\Pi_{\ell'+1}$  and  $\Lambda_{\ell'+1} = \text{diag}(\lambda_{\nu})_{\nu \in N_{\ell'+1}}$ . Since  $\theta$  has no twin, it follows that  $\pi$  is the identity and  $\lambda_{\nu} = 1$  for every  $\nu \in N_{\ell'+1}$ , which concludes the induction. Now, since  $(\mathbf{W}'_{\ell}, \mathbf{b}'_{\ell}) = \mathbf{J}_T(\mathbf{W}_{\ell}, \mathbf{b}_{\ell})$  by construction of  $\theta'$ , we have

$$\mathbf{J}_T(\mathbf{W}_{\ell}, \mathbf{b}_{\ell}) = (\mathbf{W}'_{\ell}, \mathbf{b}'_{\ell}) = (\Pi_{\ell} \Lambda_{\ell} \mathbf{W}_{\ell} \Lambda_{\ell-1}^{-1} \Pi_{\ell-1}^{-1}, \Pi_{\ell} \Lambda_{\ell} \mathbf{b}_{\ell}) = \Pi_{\ell} \Lambda_{\ell} (\mathbf{W}_{\ell}, \mathbf{b}_{\ell}).$$

As a result, for each  $\nu \in T \neq \emptyset$  we have  $-(\mathbf{w}_{\bullet \rightarrow \nu}, b_{\nu}) = \lambda_{\pi(\nu)}(\mathbf{w}_{\bullet \rightarrow \pi(\nu)}, b_{\pi(\nu)})$  where  $\pi$  is the permutation of  $N_{\ell}$  associated to  $\Pi_{\ell}$  and  $\Lambda_{\ell} = \text{diag}(\lambda_{\nu})_{\nu \in N_{\ell}}$ . However, since  $\theta$  has no twin,  $(\mathbf{w}_{\bullet \rightarrow \nu}, b_{\nu})$  is not collinear to any  $(\mathbf{w}_{\bullet \rightarrow \nu'}, b_{\nu'})$ ,  $\nu' \in N_{\ell}$ ,  $\nu' \neq \nu$ , hence  $\pi(\nu) = \nu$ . It follows the  $-(\mathbf{w}_{\bullet \rightarrow \nu}, b_{\nu}) = \lambda_{\nu}(\mathbf{w}_{\bullet \rightarrow \nu}, b_{\nu})$ , and as  $\lambda_{\nu} > 0$  we obtain  $(\mathbf{w}_{\bullet \rightarrow \nu}, b_{\nu}) = 0$ , therefore  $\theta$  is not admissible. However, by Lemma 2, since  $\theta$  is PS-identifiable with respect to  $\Theta = \mathbb{R}^{E \cup \bar{H}}$ , it is admissible. Hence the desired contradiction.

#### APPENDIX E. PROOF OF LEMMA 9 AND LEMMA 10

*Proof of Lemma 9.* The proof is by induction on  $L$ . For  $L = 1$ , since  $\mathbf{I}_0$  is the identity

$$\mathbf{R}_{\theta}(x) = \mathbf{z}_1(\theta, x) = \mathbf{W}_1 x + \mathbf{b}_1 = \mathbf{W}_1 \mathbf{I}_0 x + \mathbf{b}_1.$$

With the convention that a product of matrices over an empty index set is the identity, this establishes (23) for  $L = 1$ . Now, assuming that (23) holds for every network of depth  $L$ , let us prove it for  $\theta$  of depth  $L + 1$ . For this, observe that  $\mathbf{z}_L(\theta, x)$  is the realization of a network  $\underline{\theta}$  of depth  $L$  made of the first  $L$  affine layers of  $\theta$ , hence by the induction hypothesis we can use (23) to get

$$\mathbf{z}_L(\theta, x) = \mathbf{R}_{\underline{\theta}}(x) = (\Pi_{\ell=1}^L \mathbf{W}_{\ell} \mathbf{I}_{\ell-1}) x + \sum_{\ell'=1}^{L-1} (\Pi_{\ell=\ell'+1}^L \mathbf{W}_{\ell} \mathbf{I}_{\ell}) \mathbf{b}_{\ell'}$$

Since  $\mathbf{y}_L(\theta, x) = \mathbf{a}_L(\theta, x) \odot \mathbf{z}_L(\theta, x) = \mathbf{I}_L \mathbf{z}_L(\theta, x)$  we get

$$\begin{aligned} \mathbf{R}_{\theta}(x) &= \mathbf{z}_{L+1}(\theta, x) = \mathbf{W}_{L+1} \mathbf{y}_L(\theta, x) + \mathbf{b}_{L+1} \\ &= \mathbf{W}_{L+1} \mathbf{I}_L \left( (\Pi_{\ell=1}^L \mathbf{W}_{\ell} \mathbf{I}_{\ell-1}) x + \sum_{\ell'=1}^L (\Pi_{\ell=\ell'+1}^L \mathbf{W}_{\ell} \mathbf{I}_{\ell-1}) \mathbf{b}_{\ell'} \right) + \mathbf{b}_{L+1}. \end{aligned}$$

To conclude simply observe that

$$\begin{aligned} \mathbf{W}_{L+1} \mathbf{I}_L (\Pi_{\ell=1}^L \mathbf{W}_\ell \mathbf{I}_{\ell-1}) &= \left( \Pi_{\ell=1}^{L+1} \mathbf{W}_\ell \mathbf{I}_{\ell-1} \right), \\ \mathbf{W}_{L+1} \mathbf{I}_L \left( \sum_{\ell'=1}^L (\Pi_{\ell=\ell'+1}^L \mathbf{W}_\ell \mathbf{I}_{\ell-1}) \mathbf{b}_{\ell'} \right) &= \sum_{\ell'=1}^L \left( \Pi_{\ell=\ell'+1}^{L+1} \mathbf{W}_\ell \mathbf{I}_{\ell-1} \right) \mathbf{b}_{\ell'} \\ \text{and that with } \ell' = L+1, \quad \mathbf{b}_{L+1} &= \mathbf{b}_{\ell'} = \left( \Pi_{\ell=\ell'+1}^{L+1} \mathbf{W}_\ell \mathbf{I}_{\ell-1} \right) \mathbf{b}_{\ell'}. \end{aligned} \quad \square$$

*Proof of Lemma 10.* With  $\mathcal{P}_H := \cup_{\ell=1}^{L-1} \mathcal{P}_\ell$  the set of all paths from a hidden neuron to an output neuron, for each output neuron  $\eta \in N_L$  we prove at the end of this section that

$$(50) \quad \mathbf{R}_\theta(x)_\eta = \sum_{p \in \mathcal{P}_0, p_L = \eta} \alpha_p(\theta, x) \Phi_p(\theta) x_{p_0} + \sum_{p \in \mathcal{P}_H, p_L = \eta} \alpha_p(\theta, x) \Phi_p(\theta) + \theta_\eta$$

Any  $p \in \mathcal{P}_0$  is uniquely written  $p = \mu \rightarrow q \rightarrow \eta$  with  $\mu = p_0$  its input neuron,  $\eta \in N_L$  its output neuron, and  $q = (p_1, \dots, p_{L-1}) \in \mathcal{Q}_1$  a partial path from the first layer to the penultimate layer, and  $\alpha_{\mu \rightarrow q \rightarrow \eta}(\theta, x) = \alpha_q(\theta, x)$  for every  $q \in \mathcal{Q}_1$  and any  $\mu \in N_0, \eta \in N_L$ . Similarly, for  $1 \leq \ell \leq L-1$ , every partial path  $p \in \mathcal{P}_\ell$  starting from the  $\ell$ -th hidden layer and ending at the output layer can be written as  $p = q \rightarrow \eta$  where  $\eta = p_L \in N_L$  and  $q \in \mathcal{Q}_\ell$  starts from the  $\ell$ -th hidden layer and ends at the penultimate layer, and we have  $\alpha_{q \rightarrow \eta}(\theta, x) = \alpha_q(\theta, x)$  for all  $\theta, x$ . Therefore, (50) can be rewritten as

$$\begin{aligned} \mathbf{R}_\theta(x)_\eta &= \sum_{q \in \mathcal{Q}_1} \alpha_q(\theta, x) \sum_{\mu \in N_0} \Phi_{\mu \rightarrow q \rightarrow \eta}(\theta) x_\mu + \sum_{q \in \mathcal{Q}} \alpha_q(\theta, x) \Phi_{q \rightarrow \eta}(\theta) + \theta_\eta \\ &= \sum_{q \in \mathcal{Q}_1} \alpha_q(\theta, x) [\Phi_\eta^i(\theta) x]_q + \sum_{q \in \mathcal{Q}+1} [\bar{\alpha}(\theta, x)]_q [\Phi_\eta^h(\theta)]_q \\ &= \langle \mathbf{Q} \bar{\alpha}(\theta, x), \Phi_\eta^i(\theta) x \rangle + \langle \bar{\alpha}(\theta, x), \Phi_\eta^h(\theta) \rangle. \end{aligned}$$

where we used that  $\bar{\alpha}(\theta, x) := (\alpha(\theta, x), 1)$  with  $\alpha(\theta, x) := (\alpha_q(\theta, x))_{q \in \mathcal{Q}}$ , and  $\mathbf{Q}$  is the canonical restriction from  $\mathcal{Q}+1$  to  $\mathcal{Q}_1$ .  $\square$

*Proof of Equation (50).* We prove the result by induction on the number of layers  $L$ .

For  $L = 1$  we have  $\mathcal{P}_0 = \{(\mu, \eta)\}_{\mu \in N_0, \eta \in N_1}$  and  $\mathcal{P}_1 = \{(\eta)\}_{\eta \in N_1}$ . Since  $H = \emptyset$ ,  $\alpha_p(\theta, x) = 1$  for all  $p \in \mathcal{P}$  and  $\mathcal{P}_H = \emptyset$ . We have  $\Phi_p(\theta) = w_{\mu \rightarrow \eta}$  for all  $p = (\mu, \eta) \in \mathcal{P}_0$  and  $\Phi_p(\theta) = b_\eta$  for each  $p = (\eta) \in \mathcal{P}_1$ . It follows that

$$\sum_{\substack{p \in \mathcal{P}_0 \\ p_L = \eta}} \alpha_p(\theta, x) \Phi_p(\theta) x_{p_0} + \sum_{\substack{p \in \mathcal{P}_H \\ p_L = \eta}} \alpha_p(\theta, x) \Phi_p(\theta) + \theta_\eta = \sum_{\mu \in N_0} w_{\mu \rightarrow \eta} x_\mu + b_\eta = (\mathbf{W}_1 x + \mathbf{b}_1)_\eta = (\mathbf{R}_\theta(x))_\eta.$$

This establishes (50) for  $L = 1$ . Assume now that (50) holds for networks of depth  $L \geq 1$ . With  $\theta$  a network of depth  $L+1$ , observe that  $\mathbf{z}_L(\theta, x) = \mathbf{R}_{\tilde{\theta}}(x)$  with  $\tilde{\theta}$  the network made

of the first  $L$  affine layers of  $\theta$ . Using the induction hypothesis, we get, for  $\nu \in N_{L-1}$ ,

$$(51) \quad z_\nu(\theta, x) = \sum_{\substack{\tilde{p} \in \tilde{\mathcal{P}}_0 \\ \tilde{p}_{L-1} = \nu}} \alpha_{\tilde{p}}(\tilde{\theta}, x) \Phi_{\tilde{p}}(\tilde{\theta}) x_{\tilde{p}_0} + \sum_{\substack{\tilde{p} \in \tilde{\mathcal{P}}_H \\ \tilde{p}_{L-1} = \nu}} \alpha_{\tilde{p}}(\tilde{\theta}, x) \Phi_{\tilde{p}}(\tilde{\theta})$$

with  $\tilde{\mathcal{P}}_0 = \{(p_0, \dots, p_{L-1}) \mid p \in \mathcal{P}_0\}$ ,  $\tilde{\mathcal{P}}_H = \{(p_\ell, \dots, p_{L-1}) \mid p = (p_\ell, \dots, p_L) \in \mathcal{P}_H\}$ . Since  $\text{ReLU}(z_\nu(\theta, x)) = a_\nu(\theta, x) z_\nu(\theta, x)$  we get

$$\begin{aligned} (\mathbf{R}_\theta(x))_\eta &= \sum_{\nu \in N_{L-1}} y_\nu(\theta, x) w_{\nu \rightarrow \eta} + b_\eta = \sum_{\nu \in N_{L-1}} \text{ReLU}(z_\nu(\theta, x)) w_{\nu \rightarrow \eta} + b_\eta \\ &= \sum_{\nu \in N_{L-1}} a_\nu(\theta, x) z_\nu(\theta, x) w_{\nu \rightarrow \eta} + b_\eta \\ &= \sum_{\nu \in N_{L-1}} \sum_{\substack{\tilde{p} \in \tilde{\mathcal{P}}_0 \\ \tilde{p}_{L-1} = \nu}} a_\nu(\theta, x) \alpha_{\tilde{p}}(\tilde{\theta}, x) \Phi_{\tilde{p}}(\tilde{\theta}) w_{\nu \rightarrow \eta} x_{\tilde{p}_0} \\ &\quad + \sum_{\nu \in N_{L-1}} \sum_{\substack{\tilde{p} \in \tilde{\mathcal{P}}_H \\ \tilde{p}_{L-1} = \nu}} a_\nu(\theta, x) \alpha_{\tilde{p}}(\tilde{\theta}, x) \Phi_{\tilde{p}}(\tilde{\theta}) w_{\nu \rightarrow \eta} \\ &\quad + b_\eta \end{aligned}$$

For each path such that  $\tilde{p}_{L-1} = \nu \in N_{L-1}$  we have  $a_\nu(\theta, x) \alpha_{\tilde{p}}(\tilde{\theta}, x) \Phi_{\tilde{p}}(\tilde{\theta}) w_{\nu \rightarrow \eta} = \alpha_{\tilde{p} \rightarrow \eta}(\theta, x) \Phi_{\tilde{p} \rightarrow \eta}(\theta)$ , and  $p := \tilde{p} \rightarrow \eta$  belongs to  $\mathcal{P}_0$  (resp. to  $\mathcal{P}_H$ ) if, and only if,  $\tilde{p} \in \tilde{\mathcal{P}}_0$  (resp.  $\tilde{p} \in \tilde{\mathcal{P}}_H$ ). Thus,

$$\begin{aligned} (\mathbf{R}_\theta(x))_\eta &= \sum_{\nu \in N_{L-1}} \sum_{\substack{\tilde{p} \in \tilde{\mathcal{P}}_0 \\ \tilde{p}_{L-1} = \nu}} \alpha_{\tilde{p} \rightarrow \eta}(\theta, x) \Phi_{\tilde{p} \rightarrow \eta}(\theta) x_{\tilde{p}_0} + \sum_{\nu \in N_{L-1}} \sum_{\substack{\tilde{p} \in \tilde{\mathcal{P}}_H \\ \tilde{p}_{L-1} = \nu}} \alpha_{\tilde{p} \rightarrow \eta}(\theta, x) \Phi_{\tilde{p} \rightarrow \eta}(\theta) + b_\eta \\ &= \sum_{\substack{p \in \mathcal{P}_0 \\ p_L = \eta}} \alpha_p(\theta, x) \Phi_p(\theta) x_{p_0} + \sum_{\substack{p \in \mathcal{P}_H \\ p_L = \eta}} \alpha_p(\theta, x) \Phi_p(\theta) + b_\eta. \end{aligned} \quad \square$$

## APPENDIX F. PROOF OF LEMMA 11

The result is proved by induction on the network's depth. The case  $L = 1$  is trivial with the convention that a union over an empty family is empty. For any depth, since

$$(\cup_{\nu \in H} \Gamma_\nu(\theta))^c = \cap_{\nu \in H} \Gamma_\nu^c(\theta) = \cap_{\ell=1}^{L-1} (\cap_{\nu \in N_\ell} \Gamma_\nu^c(\theta)) = \cap_{\ell=1}^{L-1} (\cup_{\nu \in N_\ell} \Gamma_\nu(\theta))^c$$

the result is equivalent to  $\mathcal{X}'_\theta = \cap_{\ell=1}^{L-1} (\cup_{\nu \in N_\ell} \Gamma_\nu(\theta))^c$ , which is the quantity manipulated in the induction. Assume that the result is valid for all parameters of depth  $L$  and consider  $\theta$  a parameter of depth  $L+1 \geq 2$ . Denoting  $\underline{\theta} = g(\theta)$  its restriction to its first  $L$  layers, we will show that  $\mathcal{X}'_\theta = \mathcal{X}'_{\underline{\theta}} \cap (\cup_{\nu \in N_L} \Gamma_\nu(\theta))^c$ . First we prove  $(\mathcal{X}'_\theta)^c \cup (\cup_{\nu \in N_L} \Gamma_\nu(\theta))^c \subset (\mathcal{X}'_{\underline{\theta}})^c$ .

- if  $x \notin \mathcal{X}'_\theta$  then (by definition of  $\mathcal{X}'_\theta$ ) the function  $(\underline{\theta}', x') \mapsto \mathbf{a}(\underline{\theta}', x')$  is not locally constant around  $(\underline{\theta}, x)$  hence there exists  $1 \leq \ell \leq L-1$  and  $\nu \in N_\ell$  such that  $a_\nu(\underline{\theta}', x')$  is not locally constant around  $(\underline{\theta}, x)$ . Since  $\ell \leq L-1$ , for every  $\theta', x'$  we

have  $a_\nu(\theta', x') = a_\nu(\underline{\theta}', x')$  with  $\underline{\theta}' = g(\theta')$  the restriction of  $\theta'$  to its first  $L$  layers. We obtain that  $a_\nu(\theta', x')$  is not locally constant around  $(\theta, x)$ , showing that  $x \notin \mathcal{X}'_\theta$ .

- If  $x \in \cup_{\nu \in N_L} \Gamma_\nu(\theta)$ , there exists  $\nu \in N_L$  such that  $x \in \Gamma_\nu(\theta)$  hence  $z_\nu(\theta, x) = 0$ , the gradient is well-defined, and  $\nabla z_\nu(\theta, x) \neq 0$ . This implies that the sign of  $z_\nu(\theta, x')$  is not locally constant around  $x$ , hence  $x' \mapsto a_\nu(\theta, x')$  is not locally constant around  $x$ , therefore  $(\theta', x') \mapsto \mathbf{a}(\theta', x')$  is not locally constant around  $(\theta, x)$ . Thus,  $x \notin \mathcal{X}'_\theta$ .

This establishes equivalently that  $\mathcal{X}'_\theta \subset \mathcal{X}'_{\underline{\theta}} \cap (\cup_{\nu \in N_L} \Gamma_\nu)^c$ .

Vice-versa, consider  $x \in \mathcal{X}'_{\underline{\theta}} \cap (\cup_{\nu \in N_L} \Gamma_\nu(\theta))^c$ . Since  $x \in \mathcal{X}'_\theta$ ,  $\mathbf{a}(\underline{\theta}', x')$  is locally constant around  $(\underline{\theta}, x)$ , hence  $(\theta', x') \mapsto \mathbf{a}_\ell(g(\theta'), x') = \mathbf{a}_\ell(\theta', x')$  is locally constant around  $(\theta, x)$  for each  $1 \leq \ell \leq L-1$ . There remains to show that  $\mathbf{a}_L(\theta', x')$  is locally constant around  $(\theta, x)$ . Indeed, since  $x \notin \cup_{\nu \in N_L} \Gamma_\nu$ , we have  $z_\nu(\theta, x) \neq 0$  for every  $\nu \in N_L$ . By continuity of  $(\theta', x') \mapsto \mathbf{z}_L(\theta', x')$ , there exists a neighborhood of  $(\theta, x)$  on which  $\text{sign}(z_\nu(\theta', x'))$  is constant for every  $\nu \in N_L$ , hence  $\mathbf{a}_L(\theta', x')$  is locally constant around  $(\theta, x)$ . Overall, we get that  $(\theta, x) \mapsto \mathbf{a}(\theta', x')$  is locally constant around  $(\theta, x)$ , i.e.  $x \in \mathcal{X}'_\theta$ . This concludes the proof that  $\mathcal{X}'_{\underline{\theta}} \cap (\cup_{\nu \in N_L} \Gamma_\nu(\theta))^c \subset \mathcal{X}'_\theta$ , hence the equality  $\mathcal{X}'_\theta = \mathcal{X}'_{\underline{\theta}} \cap (\cup_{\nu \in N_L} \Gamma_\nu(\theta))^c$ .

#### APPENDIX G. PROOF OF LEMMA 13

First we prove that  $\mathbf{s}_c \in \mathbf{A}(\theta)$  and  $(\mathbf{s}_c, 0) \in \bar{\mathbf{A}}(\theta)$  for each  $c$ . Since  $\theta$  is admissible, one can check (cf Definition 8) that two hidden neurons  $\nu, \nu' \in H$  of a shallow network are:

- positive twins if, and only if,  $a_\nu(\theta, x) = a_{\nu'}(\theta, x)$  for all  $x \in \mathcal{X}_\theta$ ;
- negative twins if, and only if,  $a_\nu(\theta, x) = 1 - a_{\nu'}(\theta, x)$  for all  $x \in \mathcal{X}_\theta$ ;

Since we are on a shallow architecture, we identify  $\mathcal{Q} = \mathcal{Q}_1$  with  $H$  and  $\mathbf{a}(\theta, x)$  with  $\boldsymbol{\alpha}(\theta, x)$ . Considering the  $c$ -th equivalence class  $T_c$  of twins, it follows that for every  $x$  there is  $\epsilon_c(x) \in \{-1, +1\}$  such that

$$(52) \quad 2\boldsymbol{\alpha}_{T_c}(\theta, x) = 2\mathbf{a}_{T_c}(\theta, x) = \mathbf{1}_{T_c} + \epsilon_c(x) \cdot \mathbf{s}_c$$

where for any  $\mathbf{u} \in \mathbb{R}^H$ ,  $\mathbf{u}_T \in \mathbb{R}^H$  is its restriction to  $T$  (which matches  $\mathbf{u}$  on its coordinates indexed by  $T$  and is zero elsewhere), and  $\mathbf{1}_H \in \mathbb{R}^H$  is the vector with all entries equal to one, while  $\mathbf{1}_T$  is its restriction to  $T$ . To continue we use the following result.

**Lemma 19.** *Consider a shallow network with parameter  $\theta$ , and  $T \subset H$  an equivalence class of twin neurons. There are  $x_T^+, x_T^- \in \mathcal{X}_\theta$  such that*

$$(53) \quad |a_\nu(\theta, x_T^+) - a_\nu(\theta, x_T^-)| = \begin{cases} 1, & \text{if } \nu \in T \\ 0, & \text{otherwise} \end{cases}$$

*Proof.* For each  $\nu \in H$  denote  $\mathcal{V}_\nu = \{x \in \mathbb{R}^{N_0} : \langle \mathbf{w}_{\bullet \rightarrow \nu}, x \rangle + b_\nu = 0\}$ . Since  $\theta$  is admissible,  $\mathbf{w}_{\bullet \rightarrow \nu} \neq 0$  for each  $\nu \in H$ , hence these linear spaces are hyperplanes. The hyperplanes associated to two neurons coincide if, and only if, these neurons are twins. Choose an arbitrary  $\nu \in T$ . Since  $\mathcal{V}_\nu$  is distinct from each of the (finitely many)  $\mathcal{V}_{\nu'}, \nu' \notin T$ , there exists  $x_0 \in \mathcal{V}_\nu$  that belongs to the complement of  $\cup_{\nu' \notin T} \mathcal{V}_{\nu'}$ . As this complement is open, there exists  $\epsilon > 0$  such that  $B(x_0, \epsilon \|\mathbf{w}_{\bullet \rightarrow \nu}\|_2)$  does not intersect any of the hyperplanes  $\mathcal{V}_{\nu'}, \nu' \notin T$ . Since  $x_T^\pm := x_0 \pm \mathbf{w}_{\bullet \rightarrow \nu} \epsilon / 2 \in B(x_0, \epsilon \|\mathbf{w}_{\bullet \rightarrow \nu}\|_2)$  we obtain:  $a_{\nu'}(\theta, x_T^+) = a_{\nu'}(\theta, x_T^-)$



for every  $\nu' \notin T$ , and  $\text{sign}(\langle \mathbf{w}_{\bullet \rightarrow \nu}, x_T^\pm \rangle + b_\nu) = \pm 1$  hence  $a_\nu(\theta, x_T^+) = 1 - a_\nu(\theta, x_T^-)$ . The latter extends to each  $\nu' \in T$  by the twin property, and yields the conclusion.  $\square$

By Lemma 19 there are  $x_c^+, x_c^- \in \mathcal{X}_\theta$  such that

$$(54) \quad |a_\nu(\theta, x_c^+) - a_\nu(\theta, x_c^-)| = \begin{cases} 1, & \text{if } \nu \in T_c \\ 0, & \text{if } \nu \in H \setminus T_c \end{cases}$$

It follows that  $\alpha(\theta, x_c^+) - \alpha(\theta, x_c^-) = \alpha_{T_c}(\theta, x_c^+) - \alpha_{T_c}(\theta, x_c^-) = \pm \mathbf{s}_c$ . As a result

$$\begin{aligned} \mathbf{s}_c &= \pm (\alpha(\theta, x_c^+) - \alpha(\theta, x_c^-)) \in \mathbf{A}(\theta) \\ (\mathbf{s}_c, 0) &= \pm (\bar{\alpha}(\theta, x_c^+) - \bar{\alpha}(\theta, x_c^-)) \in \bar{\mathbf{A}}(\theta) \end{aligned}$$

as claimed. Using (52) and the partition of  $H$  into  $T_1, \dots, T_C$  we have for any  $x \in \mathcal{X}_\theta$

$$(55) \quad 2\alpha(\theta, x) = \sum_c 2\alpha_{T_c}(\theta, x) = \sum_c (1_{T_c} + \epsilon_c(x) \cdot \mathbf{s}_c) = 1_H + \sum_c \epsilon_c(x) \mathbf{s}_c$$

We obtain

$$(56) \quad 1_H = 2\alpha(\theta, x) - \sum_c \epsilon_c(x) \mathbf{s}_c,$$

and since  $\alpha(\theta, x) \in \mathbf{A}(\theta)$  and  $\mathbf{s}_c \in \mathbf{A}(\theta)$  for all  $c$ , it follows that  $1_H \in \mathbf{A}(\theta)$ . This proves  $\text{span}\{1_H, \mathbf{s}_c, 1 \leq c \leq C\} \subseteq \mathbf{A}(\theta)$ . Vice-versa, (55) shows  $\alpha(\theta, x) \in \text{span}\{1_H, \mathbf{s}_c, 1 \leq c \leq C\}$  for every  $x \in \mathcal{X}_\theta$ , hence  $\mathbf{A}(\theta) \subseteq \text{span}\{1_H, \mathbf{s}_c, 1 \leq c \leq C\}$ . By (56) we also get

$$(1_H, 2) = 2(\alpha(\theta, x), 1) - \sum_c \epsilon_c(x) (\mathbf{s}_c, 0),$$

and since  $(\alpha(\theta, x), 1) = \bar{\alpha}(\theta, x) \in \bar{\mathbf{A}}(\theta)$  and  $(\mathbf{s}_c, 0) \in \bar{\mathbf{A}}(\theta)$ , we get  $(1_H, 2) \in \bar{\mathbf{A}}(\theta)$ . This proves  $\text{span}\{(1_H, 2), (\mathbf{s}_c, 0), 1 \leq c \leq C\} \subseteq \bar{\mathbf{A}}(\theta)$ , and also implies

$$2\bar{\alpha}(\theta, x) = (1_H, 2) + \sum_c \epsilon_c(x) (\mathbf{s}_c, 0)$$

hence  $\bar{\mathbf{A}}(\theta) \subseteq \text{span}\{(1_H, 2), (\mathbf{s}_c, 0), 1 \leq c \leq C\}$ .

## APPENDIX H. PROOF OF LEMMA 16

We use the shorthands  $\mathbf{w}_\nu = \mathbf{w}_{\bullet \rightarrow \nu}$ ,  $\mathbf{v}_\nu = \mathbf{w}_{\nu \rightarrow \bullet}$ .

Given the assumption there are  $C = |H| - 1$  classes of twin neurons, all being trivial except one made of a pair of negative twins  $\{\nu, \nu'\}$ . Without loss of generality we enumerate the neurons and their classes such that  $T_1 = \{\nu_1, \nu_2\} = \{\nu, \nu'\}$  and  $T_c = \{\nu_{c+1}\}$ ,  $2 \leq c \leq C = |H| - 1$ . First we establish that, with this numbering,

$$(57) \quad \bar{\mathbf{A}}^\perp(\theta) = \text{span}\{(1, 1, 0, \dots, 0, -1)\} \text{ and } \mathbf{A}^\perp(\theta) = \{0\}.$$

The signatures of the classes are  $\mathbf{s}_1 = \delta_1 - \delta_2$  and  $\mathbf{s}_c = \delta_{c+1}$ ,  $2 \leq c \leq C$ . By Lemma 13 we have  $\mathbf{A}(\theta) = \text{span}\{1_H, \mathbf{s}_c, 1 \leq c \leq C\}$ . It is not difficult to check <sup>11</sup> that the  $C + 1 = |H|$

<sup>11</sup>If, instead of a single pair of negative twins, we consider a single pair of *positive* twins, then  $\mathbf{s}_1 = \delta_1 + \delta_2$  and the spanning vectors of  $\mathbf{A}(\theta)$  become linearly *dependent*, with  $\mathbf{A}^\perp(\theta) = \text{span}\{(1, -1, 0, \dots, 0)\} \neq \{0\}$ .

spanning vectors are linearly independent, hence  $\mathbf{A}(\theta) = \mathbb{R}^H$  and  $\mathbf{A}^\perp(\theta) = \{0\}$ . Now, consider  $\mathbf{v} = (v_1, \dots, v_{|H|+1}) \in \bar{\mathbf{A}}^\perp(\theta)$ . By Lemma 13, this vector is orthogonal to each  $(\mathbf{s}_c, 0)$ ,  $1 \leq c \leq C$ , and to  $(\mathbf{1}_H, 2)$ . For  $2 \leq c \leq C$ , orthogonality to  $(\mathbf{s}_c, 0) = (\delta_{c+1}, 0)$  implies  $v_{c+1} = 0$ , hence  $\mathbf{v} = (\alpha, \beta, 0, \dots, 0, \gamma)$  for some  $\alpha, \beta, \gamma \in \mathbb{R}$ . Orthogonality to  $(\mathbf{s}_1, 0) = (1, -1, 0, \dots, 0)$  implies  $\beta = \alpha$ , and orthogonality to  $(\mathbf{1}_H, 2)$  implies  $\gamma = -\alpha$ , hence  $\mathbf{v}$  is proportional to  $(1, 1, 0, \dots, 0, -1)$  as claimed. Since  $\bar{\mathbf{A}}(\theta)$  is spanned by  $C + 1 = |H|$  vectors, its dimension is at most  $|H|$ , hence the dimension of  $\bar{\mathbf{A}}^\perp(\theta)$  is at least one. This concludes the proof that  $\bar{\mathbf{A}}^\perp(\theta) = \text{span}\{(1, 1, 0, \dots, 0, -1)\}$ .

Since  $\theta$  is admissible, there is an input neuron  $\mu \in N_0$  such that  $w_{\mu \rightarrow \nu_1} \neq 0$ . Since  $\nu_1, \nu_2$  are twins, we also have  $w_{\mu \rightarrow \nu_2} \neq 0$ . Let  $\epsilon_0 := \min_{1 \leq j \leq 2} |w_{\mu \rightarrow \nu_j}|/2$ . Consider  $\theta' \in B(\theta, \epsilon_0)$  such that  $\Phi(\theta') - \Phi(\theta) \in \mathbf{V}(\theta)$ . First, observe that  $w'_{\mu \rightarrow \nu_j} \neq 0$  for  $j = 1, 2$ . Then, in light of Lemma 8 and (57), for every  $\eta \in N_2$ , we have  $\Phi_\eta^i(\theta') = \Phi_\eta^i(\theta)$  and

$$(58) \quad \Phi_\eta^h(\theta') - \Phi_\eta^h(\theta) \in \text{span}\{(1, 1, 0, \dots, 0, -1)\},$$

hence  $b'_{\nu_{c+1}} w'_{\nu_{c+1} \rightarrow \eta} = b_{\nu_{c+1}} w_{\nu_{c+1} \rightarrow \eta}$  for  $2 \leq c \leq C$ , and there are scalars  $\lambda_\eta \in \mathbb{R}$  such that

$$(59) \quad b'_{\nu_j} w'_{\nu_j \rightarrow \eta} - b_{\nu_j} w_{\nu_j \rightarrow \eta} = \lambda_\eta, \quad \forall j \in \{1, 2\} \quad \text{and} \quad b'_\eta - b_\eta = -\lambda_\eta.$$

When  $\Theta$  is the set of parameters with zero output biases, the fact that  $\theta, \theta' \in \Theta$  implies  $b'_\eta = b_\eta = 0$ , hence  $\lambda_\eta = 0$  and  $\Phi_\eta^h(\theta') = \Phi_\eta^h(\theta)$  for every  $\eta \in N_2$ . We show below that the same holds for arbitrary  $\Theta$  when  $\mathbf{w}_{\nu_1 \rightarrow \bullet}$  and  $\mathbf{w}_{\nu_2 \rightarrow \bullet}$  are linearly independent. This implies  $\Phi(\theta') = \Phi(\theta)$ , hence  $\theta$  is then  $\epsilon$ -non-degenerate with respect to  $\Theta$ .

Indeed, the equality  $\Phi_\eta^i(\theta') = \Phi_\eta^i(\theta)$  for all  $\eta \in N_2$  implies that for  $1 \leq j \leq 2$ ,

$$(60) \quad w'_{\mu \rightarrow \nu_j} w'_{\nu_j \rightarrow \eta} = w_{\mu \rightarrow \nu_j} w_{\nu_j \rightarrow \eta}, \quad \forall \eta \in N_2,$$

and since  $w'_{\mu \rightarrow \nu_j} \neq 0$  for  $j = 1, 2$ , we obtain from (59) and (60) that for each  $\eta \in N_2$ ,

$$\begin{aligned} \lambda_\eta &= b'_{\nu_j} w'_{\nu_j \rightarrow \eta} - b_{\nu_j} w_{\nu_j \rightarrow \eta} = b'_{\nu_j} \frac{w'_{\mu \rightarrow \nu_j} w'_{\nu_j \rightarrow \eta}}{w'_{\mu \rightarrow \nu_j}} - b_{\nu_j} w_{\nu_j \rightarrow \eta} \\ &= b'_{\nu_j} \frac{w_{\mu \rightarrow \nu_j} w_{\nu_j \rightarrow \eta}}{w'_{\mu \rightarrow \nu_j}} - b_{\nu_j} w_{\nu_j \rightarrow \eta} = \left( b'_{\nu_j} \frac{w_{\mu \rightarrow \nu_j}}{w'_{\mu \rightarrow \nu_j}} - b_{\nu_j} \right) w_{\nu_j \rightarrow \eta} \end{aligned}$$

We obtain  $\lambda = x_j \mathbf{w}_{\nu_j \rightarrow \bullet}$ ,  $j = 1, 2$  where  $\lambda := (\lambda_\eta)_{\eta \in N_2}$  and  $x_j := b'_{\nu_j} \frac{w_{\mu \rightarrow \nu_j}}{w'_{\mu \rightarrow \nu_j}} - b_{\nu_j}$ . Since  $\mathbf{w}_{\nu_1 \rightarrow \bullet}$  and  $\mathbf{w}_{\nu_2 \rightarrow \bullet}$  are linearly independent, it follows that  $x_1 = x_2 = 0$ , hence  $\lambda = \mathbf{0}$ .

Assume now that  $\mathbf{w}_{\nu_1 \rightarrow \bullet}$  and  $\mathbf{w}_{\nu_2 \rightarrow \bullet}$  are linearly *dependent*, and recall that since  $\theta$  is admissible they are both nonzero vectors, hence  $\mathbf{w}_{\nu_2 \rightarrow \bullet} = \alpha \mathbf{w}_{\nu_1 \rightarrow \bullet}$  for some  $\alpha \neq 0$ . Consider

$0 < \epsilon < \epsilon_0$  and set  $\theta'$  as follows:

$$\begin{aligned} \mathbf{W}'_\ell &= \mathbf{W}_\ell, \quad 1 \leq \ell \leq 2; \\ b'_\nu &= b_\nu, \quad \nu \in H \setminus \{\nu_1, \nu_2\}; \\ b'_{\nu_1} &= b_{\nu_1} + \gamma\epsilon; \\ b'_{\nu_2} &= b_{\nu_2} + \gamma\epsilon/\alpha, j = 1, 2; \\ b'_\eta &= b_\eta - w_{\nu_1 \rightarrow \eta} \gamma\epsilon, \eta \in N_2, \end{aligned}$$

with  $0 < \gamma < \min(1, |\alpha|, 1/\|\mathbf{w}_{\nu_1 \rightarrow \bullet}\|_\infty)$  so that  $\theta' \in B(\theta, \epsilon)$ . Since the weights of  $\theta'$  and  $\theta$  coincide we have  $\Phi_\eta^i(\theta') = \Phi_\eta^i(\theta)$  for every  $\eta \in N_2$ . It is not difficult to check that, with  $\lambda_\eta := w_{\nu_1 \rightarrow \eta}\epsilon$ , we also have  $\Phi_\eta^h(\theta') - \Phi_\eta^h(\theta) = \lambda_\eta(1, 1, 0, \dots, 0, -1)$ , hence  $\Phi(\theta') - \Phi(\theta) \in \mathbf{V}(\theta)$ . Yet,  $\Phi(\theta') \neq \Phi(\theta)$  since  $\boldsymbol{\lambda} := (\lambda_\eta)_{\eta \in N_2} = \epsilon \mathbf{w}_{\nu_1 \rightarrow \bullet} \neq \mathbf{0}$ . Assuming that  $\theta$  belongs to the interior of  $\Theta$ , we have  $\theta' \in \Theta \cap B(\theta, \epsilon)$  for small enough  $\epsilon$ . It follows that  $\theta$  is degenerate.

#### APPENDIX I. DETAILS ON EXAMPLE 4

**$\theta$  is PS-identifiable from  $\mathcal{X} = \mathbb{R}$ .** Consider an arbitrary  $\theta' \in \Theta = \mathbb{R}^{E \cup \bar{H}}$ . If  $\mathbf{R}_{\theta'}(x) = \mathbf{R}_\theta(x) = |x|$  on  $\mathbb{R}$  then  $\theta'$  is admissible (otherwise its realization would be, up to an additive constant, proportional to a single shifted version of the ReLU, which would prevent it from being equal to  $\mathbf{R}_\theta = \mathbf{abs}$ ) hence  $w'_{\nu \rightarrow \nu_i} \neq 0$ ,  $i = 1, 2$ . Writing  $\alpha_i = |w'_{\mu \rightarrow \nu_i}| w'_{\nu_i \rightarrow \eta}$  and  $\beta_i = -b'_{\nu_i}/|w'_{\mu \rightarrow \nu_i}|$ , and  $s_i = \mathbf{sign}(w'_{\mu \rightarrow \nu_i}) \in \{-1, +1\}$  for  $i = 1, 2$ , we have  $\alpha_i \neq 0$  and

$$\mathbf{R}_{\theta'}(x) = \alpha_1 \text{ReLU}(s_1(x - s_1\beta_1)) + \alpha_2 \text{ReLU}(s_2(x - s_2\beta_2)) + b'_\eta, \quad \forall x \in \mathbb{R}.$$

If we had  $s_1\beta_1 \neq s_2\beta_2$ ,  $\mathbf{R}_{\theta'}$  would be non-differentiable at two distinct points  $s_1\beta_1, s_2\beta_2$ . However  $\mathbf{R}_{\theta'} = \mathbf{R}_\theta = \mathbf{abs}$  is differentiable on  $\mathbb{R} \setminus \{0\}$ , hence  $s_1\beta_1 = s_2\beta_2$ , and a similar reasoning yields  $s_1\beta_1 = s_2\beta_2 = 0$ . Since  $|s_i| = 1$ , we get  $\beta_1 = \beta_2 = 0$  and

$$\mathbf{R}_{\theta'}(x) = \alpha_1 \text{ReLU}(s_1x) + \alpha_2 \text{ReLU}(s_2x) + b'_\eta, \quad \forall x \in \mathbb{R}.$$

If we had  $s_1 = s_2$ , the realization would be  $(\alpha_1 + \alpha_2)\text{ReLU}(s_1x) + b'_\eta$ , which cannot match  $\mathbf{abs}$ , hence  $s_2 = -s_1$ . Without loss of generality (up to a permutation of indices of the hidden layer)  $s_1 = 1, s_2 = -1$ . Now, for  $x < 0$  we have  $-x = |x| = \mathbf{R}_{\theta'}(x) = -\alpha_2x + b'_\eta$  while for  $x > 0$  we get  $x = |x| = \mathbf{R}_{\theta'}(x) = \alpha_1x + b'_\eta$ , hence  $\alpha_1 = 1, \alpha_2 = 1, b'_\eta = 0$ . Overall, up to the possible permutation of the hidden layer, we obtain  $\mathbf{sign}(\theta') = \mathbf{sign}(\theta)$  and  $\alpha_1s_1 = 1, \alpha_2s_2 = -1, \alpha_1\beta_1 = \alpha_2\beta_2 = 0, b'_\eta = 0$ , hence  $\Phi(\theta') = \Phi(\theta)$ . Since  $\theta$  is admissible, it follows by Theorem 1 that  $\theta' \sim_{PS} \theta$ . Since this holds for any  $\theta'$  such that  $\mathbf{R}_{\theta'} = \mathbf{R}_\theta$ , this shows that  $\theta$  is PS-identifiable from  $\mathcal{X} = \mathbb{R}$  with respect to  $\Theta = \mathbb{R}^{E \cup \bar{H}}$ .

**$\theta$  is locally S-identifiable from some finite set  $F \subset \mathbb{R}$  (with  $0 \in F$ )**

With the same notations as above, observe that there is  $\epsilon > 0$  such that for every  $\theta' \in B(\theta, \epsilon)$  we have  $s_i := \mathbf{sign}(w'_{\mu \rightarrow \nu_i}) = \mathbf{sign}(w_{\mu \rightarrow \nu_i})$ ,  $i = 1, 2$ , and  $\max(|\alpha_1 - 1|, |\alpha_2 - 1|, |\beta_1|, |\beta_2|, |b'_\eta|) \leq 1/2$ . Consider  $\theta' \in B(\theta, \epsilon)$  such that  $\mathbf{R}_{\theta'} = \mathbf{R}_\theta$  on  $F = \{-3, -2, -1, 0, 1, 2, 3\}$ . We have  $s_1 = +1, s_2 = -1$  hence

$$\mathbf{R}_{\theta'}(x) = \alpha_1 \text{ReLU}(x - \beta_1) + \alpha_2 \text{ReLU}(-x - \beta_2) + b'_\eta.$$

Since  $|\beta_i| \leq 1/2$  for  $i = 1, 2$ , we have  $\mathbf{R}_{\theta'}(x) = \alpha_1(x - \beta_1) + b'_\eta$  for every  $x \geq 1/2$ , hence  $\alpha_1(y - x) = \mathbf{R}_{\theta'}(y) - \mathbf{R}_{\theta'}(x) = \mathbf{R}_\theta(y) - \mathbf{R}_\theta(x) = |y| - |x| = y - x$  for  $(x, y) = (1, 2)$ . Therefore  $\alpha_1 = 1$ . A similar reasoning with  $(x, y) = (-2, -1)$  shows that  $\alpha_2 = 1$ , hence

$$\mathbf{R}_{\theta'}(x) = \text{ReLU}(x - \beta_1) + \text{ReLU}(-x - \beta_2) + b'_\eta$$

Specializing to  $x = 1$ , since  $x - \beta_1 > 0$  and  $-x - \beta_2 < 0$  we get  $1 = |x| = \mathbf{R}_{\theta'}(x) = x - \beta_1 + b'_\eta = 1 - \beta_1 + b'_\eta$  hence  $b'_\eta = \beta_1$ . Similarly, with  $x = -1$ , we get  $b'_\eta = \beta_2$  hence

$$\mathbf{R}_{\theta'}(x) = \text{ReLU}(x - b'_\eta) + \text{ReLU}(-x - b'_\eta) + b'_\eta.$$

Specializing to  $x = 0 \in F$  yields

$$0 = |x| = \mathbf{R}_{\theta'}(x) = b'_\eta + 2\text{ReLU}(-b'_\eta) = \begin{cases} b'_\eta & \text{if } b'_\eta \geq 0 \\ -b'_\eta & \text{if } b'_\eta \leq 0 \end{cases} = |b'_\eta|$$

hence  $b'_\eta = 0$ . Overall we have shown that for every  $\theta' \in B(\theta, \epsilon)$  such that  $\mathbf{R}_{\theta'} = \mathbf{R}_\theta$  on  $F = \{-2, -1, 0, 1, 2\}$  we have  $\alpha_1 = \alpha_2 = 1$ ,  $s_1 = 1, s_2 = -1$ ,  $\beta_1 = \beta_2 = b'_\eta = 0$ . These imply  $\Phi(\theta') = \Phi(\theta)$  and  $\text{sign}(\theta') = \text{sign}(\theta)$  hence  $\theta' \sim_S \theta$ . In other words,  $\theta$  is locally S-identifiable from  $F$ .

#### APPENDIX J. DETAILS ON EXAMPLE 5

Here we show, as claimed in Example 5 that the parameter  $\theta_0 \in \mathbb{R}^{E \cup \bar{H}}$  from Example 1 is PS-identifiable from  $\mathcal{X} = \mathbb{R}$  with respect to the set  $\Theta_0$  of parameters with zero output biases. Consider an arbitrary  $\theta' \in \Theta_0$ . If  $\mathbf{R}_{\theta'}(x) = \mathbf{R}_{\theta_0}(x) = x$  on  $\mathcal{X}$  then  $\theta'$  is admissible (otherwise its realization would be, up to an additive constant, proportional to a single shifted version of the ReLU, which would prevent it from being equal to  $\mathbf{R}_\theta = \text{id}$ ) hence  $w'_{\nu \rightarrow \nu_i} \neq 0$ ,  $i = 1, 2$ . Writing  $\alpha_i = |w'_{\mu \rightarrow \nu_i}| w'_{\nu_i \rightarrow \eta}$  and  $\beta_i = -b'_{\nu_i} / |w'_{\mu \rightarrow \nu_i}|$ , and  $s_i = \text{sign}(w'_{\mu \rightarrow \nu_i}) \in \{-1, +1\}$  for  $i = 1, 2$ , we have  $\alpha_i \neq 0$  and since the output bias is zero

$$\mathbf{R}_{\theta'}(x) = \alpha_1 \text{ReLU}(s_1(x - s_1\beta_1)) + \alpha_2 \text{ReLU}(s_2(x - s_2\beta_2)), \quad \forall x \in \mathcal{X}.$$

If we had  $s_1\beta_1 \neq s_2\beta_2$ ,  $\mathbf{R}_{\theta'}$  would be non-differentiable at two distinct points  $s_1\beta_1, s_2\beta_2$ . However  $\mathbf{R}_{\theta'} = \mathbf{R}_{\theta_0} = \text{id}$  is differentiable on  $\mathbb{R}$ , hence  $s_1\beta_1 = s_2\beta_2$ . It follows that  $s_2\beta_1 = s_1\beta_1 = \text{id}(s_1\beta_1) = \mathbf{R}_{\theta'}(s_1\beta_1) = 0$ . Since  $|s_i| = 1$ , we get  $\beta_1 = \beta_2 = 0$  and

$$\mathbf{R}_{\theta'}(x) = \alpha_1 \text{ReLU}(s_1x) + \alpha_2 \text{ReLU}(s_2x), \quad \forall x \in \mathcal{X}.$$

If we had  $s_1 = s_2$ , the realization would be  $(\alpha_1 + \alpha_2)\text{ReLU}(s_1x)$ , which cannot match  $\text{id}$  on  $\mathcal{X}$ , hence  $s_2 = -s_1$ . Without loss of generality (up to a permutation of indices of the hidden layer)  $s_1 = 1, s_2 = -1$ . Now, for  $x < 0$  we have  $x = \text{id}(x) = \mathbf{R}_{\theta'}(x) = -\alpha_2x$  while for  $x > 0$  we get  $x = \text{id}(x) = \mathbf{R}_{\theta'}(x) = \alpha_1x$ , hence  $\alpha_1 = 1, \alpha_2 = -1$  (and  $b'_\eta = 0$  because  $\theta' \in \Theta_0$ ). Overall, up to the possible permutation of the hidden layer, we obtain  $\text{sign}(\theta') = \text{sign}(\theta_0)$  and  $\alpha_1s_1 = 1, \alpha_2s_2 = -1, \alpha_1\beta_1 = \alpha_2\beta_2 = 0, b'_\eta = 0$ , hence  $\Phi(\theta') = \Phi(\theta_0)$ . Since  $\theta_0$  is admissible, it follows by Theorem 1 that  $\theta' \sim_{PS} \theta_0$ . Since this holds for any  $\theta' \in \Theta_0$  such that  $\mathbf{R}_{\theta'} = \mathbf{R}_{\theta_0}$ , this shows that  $\theta$  is PS-identifiable from  $\mathcal{X} = \mathbb{R}$  with respect to  $\Theta_0$ .