Nonparametric Learning of Two-Layer ReLU Residual Units

Zhunxuan Wang Amazon^{*} London EC2A 2FA, United Kingdom

Linyun He Georgia Institute of Technology Atlanta, GA 30332, United States

Chunchuan Lyu Instituto de Telecomunicacões Torre Norte 1049-001 Lisbon, Portugal

Shay B. Cohen University of Edinburgh Edinburgh EH8 9AB, United Kingdom

Reviewed on OpenReview: https://openreview.net/forum?id=Yi0I0vqJ0n

Abstract

We describe an algorithm that learns two-layer residual units using rectified linear unit (ReLU) activation: suppose the input \mathbf{x} is from a distribution with support space \mathbb{R}^d and the ground-truth generative model is a residual unit of this type, given by $\mathbf{y} = \mathbf{B}^* \left[(\mathbf{A}^* \mathbf{x})^+ + \mathbf{x} \right]$, where ground-truth network parameters $\mathbf{A}^* \in \mathbb{R}^{d \times d}$ represent a full-rank matrix with nonnegative entries and $\mathbf{B}^* \in \mathbb{R}^{m \times d}$ is full-rank with $m \geq d$ and for $\mathbf{c} \in \mathbb{R}^d$, $[\mathbf{c}^+]_i = \max\{0, c_i\}$. We design layer-wise objectives as functionals whose analytic minimizers express the exact ground-truth network in terms of its parameters and nonlinearities. Following this objective landscape, learning residual units from finite samples can be formulated using convex optimization of a nonparametric function: for each layer, we first formulate the corresponding empirical risk minimization (ERM) as a positive semi-definite quadratic program (QP), then we show the solution space of the QP can be equivalently determined by a set of linear inequalities, which can then be efficiently solved by linear programming (LP). We further prove the strong statistical consistency of our algorithm, and demonstrate its robustness and sample efficiency through experimental results on synthetic data and a set of benchmark regression datasets.¹

1 Introduction

Neural networks have achieved remarkable success in various fields such as computer vision (LeCun et al., 1998; Krizhevsky et al., 2012; He et al., 2016a) and natural language processing (Kim, 2014; Sutskever et al., 2014). This success is largely due to the strong expressive power of neural networks (Bengio & Delalleau, 2011), where nonlinear activation units, such as rectified linear units (ReLU; Nair & Hinton 2010) and hyperbolic tangents (tanh) play a vital role to ensure the large learning capacity of the networks (Maas et al., 2013). However, the nonlinearity of neural networks makes them significantly more difficult to train than linear models (Livni et al., 2014). Therefore, with the development of neural network has become an important and a relatively new goal.

zhunxuan.wang@gmail.com

lhe85@gatech.edu

chunchuan.lv@gmail.com

scohen@inf.ed.ac.uk

^{*}Work done mostly at University of Edinburgh.

¹Our code is available at https://github.com/uuzeeex/relu-resunit-learning.

Residual networks, or ResNets (He et al., 2016a), are a class of deep neural networks that adopt skip connections to feed values between nonadjacent layers, where skipped layers may contain nonlinearities in between. Without loss of expressivity, ResNets avoid the vanishing gradient problem by directly passing gradient information from previous layers to current layers where otherwise gradients might vanish without skipping. In practice, ResNets have shown strong learning efficiency in several tasks, e.g. achieving at least 93% test accuracy on CIFAR-10 classification, lowering single-crop error to 20.1% on the 1000-class ImageNet dataset (Russakovsky et al., 2015; He et al., 2016b; Allen-Zhu & Li, 2019).



Common ResNets are often aggregated by many repeated shallow networks, where each network acts as a minimal unit with this kind of skip propagation, named a residual unit (He et al., 2016b).

Given the flexibility and simplicity of residual units, much theoretical work has been devoted to studying them and developing training algorithms for them in a way that sidesteps from the standard backpropagation regime and provides guarantees on the quality of estimation (see Section 1.1). In this paper, we propose algorithms that learn a general class of single-skip two-layer residual units with ReLU activation as shown on the right by the equation: $\mathbf{y} = \mathbf{B} \left[(\mathbf{A}\mathbf{x})^+ + \mathbf{x} \right]$, where for a scalar $c, c^+ = \max\{0, c\}$ (for a vector, this maximization is applied coordinate-wise), \mathbf{x} is a random vector as the network input with support space \mathbb{R}^d , and $\mathbf{A} \in \mathbb{R}^{d \times d}$ with nonnegative entries and $\mathbf{B} \in \mathbb{R}^{m \times d}$ are weight matrices of *layer 1* and *layer 2*, respectively.

Compared to previous work (Ge et al., 2019; 2018; Zhang et al., 2019; Wu et al., 2019; Tian, 2017; Du et al., 2018; Brutzkus & Globerson, 2017; Soltanolkotabi, 2017; Li & Yuan, 2017; Zhong et al., 2017), the introduction of residual connections simplifies the recovery of the network parameters by removing the permutation and scaling invariance. However, the naive mean square error minimization used in estimating the parameters remains nonconvex. Unlike most previous work, we do not assume a specific input distribution nor that the distribution is symmetric (Ge et al., 2019; Du & Goel, 2018).

We show that under regularity conditions on the weights of the residual unit, the problem of learning the ReLU unit can be formulated through quadratic programming (QP). We use nonparametric estimation (Guntuboyina & Sen, 2018) to estimate the ReLU function values in the networks. We further rewrite our constructed quadratic programs to linear programs (LPs). The LP formulation is simpler to optimize and has the same solution space as the QP for the network parameters.

1.1 Related Work

The study of two-layer ReLU networks (and two-layer networks in general, with one hidden layer) has received much attention in recent years because of the balance they present between their practical and theoretical properties. On one hand, two-layer networks represent non-convex learning problems and embody many of the same difficulties that we have with several layers. On the other hand, they are simple enough to study from a theoretical perspective. Arora et al. (2014) recover a multi-layer generative network with sparse connections and Livni et al. (2014) study the learning of multi-layer neural networks with polynomial activation. Goel et al. (2018) learns a one-layer convolution network with a perceptron-like rule. They prove the correctness of an iterative algorithm for exact recovery of the target network. Rabusseau et al. (2019) describe a spectral algorithm for two-layer linear networks. Recent work has connected optimization and two-layer network learning (Ergen & Pilanci, 2021a; Sahiner et al., 2021; Ergen & Pilanci, 2021b; Pilanci & Ergen, 2020; Mishkin et al., 2022) and has shown how to optimize networks layer by layer (Belilovsky et al., 2019) or their relationship to linear classification (Yang et al., 2021). Arjevani & Field (2021) study the symmetries that two-layer ReLU networks present.

For learning a one-layer ReLU network, Wu et al. (2019) optimize the norm and direction of neural network weight vectors separately, and Zhang et al. (2019) use gradient descent with a specific initialization. For learning a two-layer ReLU network, Ge et al. (2018) redesign the optimization landscape such that it is more amenable to theoretical analysis and Ge et al. (2019) use a moment-based method for estimating a neural network (Janzamin et al., 2015). Many others have studied ReLU networks in various settings (Tian, 2017; Du et al., 2018; Brutzkus & Globerson, 2017; Soltanolkotabi, 2017; Li & Yuan, 2017; Zhong et al., 2017; Goel & Klivans, 2019).

The study of ReLU networks with two hidden layers has also been gaining attention (Goel & Klivans, 2019; Allen-Zhu et al., 2019), with a focus on Probably Approximately Correct (PAC) learning (Valiant, 1984). In relation to our work, Allen-Zhu & Li (2019) examined the PAC-learnable function of a specific three-layer neural network with residual connections. Their work differs from ours in two aspects. First, their learnable functions include a smaller (in comparison to the student network) three-layer residual network. Second, the assumptions they make on their three-layer model are rather different than ours.

In relation to nonparametric estimation, Guntuboyina & Sen (2018) treat the final output of a shape-restricted regressor as a parameter, placing some restrictions on the type of function that can be estimated (such as convexity). They provide solutions for estimation with isotonic regression (Brunk, 1955; Ayer et al., 1955; van Eeden, 1956), convex regression (Seijo & Sen, 2011), shape-restricted additive models (Meyer, 2013; Chen & Samworth, 2016) and shape-restricted single index models (Kakade et al., 2011; Kuchibhotla et al., 2021).

1.2 Main Results

We design quadratic objective functionals with a linear bounded feasible domain that take network parameters and activation functions as variables to estimate the ground-truth network parameters and nonlinearities. The values of the objectives are moments over the input distribution.

In Theorems 4.1 and 5.1, we show that if a ground-truth residual unit has nondegenerate weights in both layers and nonnegative weights in layer 1, then there are quadratic functionals (see Eqs. (2) and (6)) with a domain defined by linear constraints whose minimizers a) are unique, and are the exact ground-truth, or b) are not unique, but can be adjusted to the exact ground-truth by a process of rescaling (see Theorem 5.2).

The minimizers above use moments of the input distribution. In practice, exact moments from the underlying distribution are not available. To address that, we use empirical risk minimization (ERM) of the moment-valued objectives by generated samples. With functions as variables in moment-valued objectives being optimized nonparametrically, the empirical objectives become quadratic functions with linear constraints (quadratic program; QP).

We further show the convexity of our QP which guarantees its solution in polynomial time w.r.t. sample size and dimension. With the available QP solution, our learning algorithm is strongly consistent (Theorem 6.3).² We show that if samples are generated by a ground-truth residual unit that has nondegenerate weights in both layers and nonnegative weights in layer 1, then our algorithm learns a network closer to the exact ground-truth (in Frobenius norm) as sample size grows.

Roadmap: Assume A and B are student network weights of layer 1 and 2. We first give a warm-up vanilla linear regression approach only knowing A is in $\mathbb{R}^{d \times d}$ and B is in $\mathbb{R}^{m \times d}$ (Section 3). We then move to the details of our nonparametric learning for layer 2 (Section 4), and similarly, layer 1 (Section 5). We formalize the quadratic functional minimization result above, showing how nonparametric learning allows us to select the values of A and B from reduced spaces that are seeded by A^* and B^* , under which we can use linear regression (LR) effectively compared to vanilla LR. In Section 6, we describe the strong consistency of our methods for respective layers, formalizing the consistency result using the continuous mapping theorem (Mann & Wald, 1943). In Appendix I, we describe a sample complexity analysis of our algorithm. We further discuss the multi-layer case in Appendix C.

2 Preliminaries

We describe the notations used in this paper, introduce the model and its underlying assumptions, and state conditions which simplify the problem but which can be removed without loss of learnability.

Notation The ReLU residual units we use take a vector $\boldsymbol{x} \in \mathbb{R}^d$ as input and return a vector $\boldsymbol{y} \in \mathbb{R}^m$. We use $\boldsymbol{A}^* \in \mathbb{R}^{d \times d}$ and $\boldsymbol{B}^* \in \mathbb{R}^{m \times d}$ to denote the ground-truth network parameters for layer 1 and 2, respectively. We use circumflex to denote predicted terms (e.g. an empirical objective function \hat{f} , estimated layer 1 weights

 $^{^{2}}$ By *strong* consistency, we are referring to almost sure convergence of the estimator to the true parameters, rather than convergence in probability, which is a weaker type of probabilistic convergence.

 \hat{A}). We use *n* to denote the number of independently and identically distributed (i.i.d.) samples available for the training algorithm, $\{\boldsymbol{x}^{(i)}, \boldsymbol{y}^{(i)}\}$ to denote the *i*-th sample drawn. For an integer *k*, We define [k] to be $\{1, 2, \ldots, k\}$, and $\boldsymbol{e}^{(j)}$ as the standard basis vector with a 1 at position *j*. All scalar-based operators are element wise in the case of vectors or matrices unless specified otherwise. For a matrix \boldsymbol{G} and an index *j*, $\boldsymbol{G}_{:,j}$ denotes its *j*th column and $\boldsymbol{G}_{j,:}$ denotes its *j*th row. We use **x** and **y** to refer to the input and output vectors, respectively, as random variables.

Linear Regression: Linear regression, or LR, in this paper refers to the problem of estimating L for unbiased model y = Lx. In this case, estimation is done by minimizing the empirical risk $\hat{R}(L) = \frac{1}{2n} \sum_{i \in [n]} \|Lx^{(i)} - y^{(i)}\|^2$ – using linear least squares (LLS). Its solution has a closed form which we use as given. We refer the reader to Hamilton (2020) for more information.

Models, Assumptions and Conditions Following previous work (Livni et al., 2014), we assume a given neural network structure specifies a hypothesis class that contains all the networks conforming to this structure. Learning such a class means using training samples to find a set of weights such that the neural networks predictions generalize well to unseen samples, where both the training and the unseen samples are drawn from an unobserved ground-truth distribution. In this paper, the hypothesis class is given by ReLU residual units and we assume it has sufficient expressive power to fit the ground-truth model. More specifically, we discuss the *realizable* case of learning, in which the ground-truth model is set to be a residual unit taken from the hypothesis class with the form:

$$\mathbf{y} = \boldsymbol{B}^* \left[\left(\boldsymbol{A}^* \mathbf{x} \right)^+ + \mathbf{x} \right], \tag{1}$$

and is used to draw samples for learning. Unlike other multi-layer ReLU-affine models which do not apply skip connections, we cannot permute the weight matrices of the residual unit and retain the same function because of skip-adding \mathbf{x} (it breaks symmetry). This helps us circumvent issues of identifiability,³ and allows us to precisely estimate the ground-truth weight matrices A^* and B^* .

Our general approach for residual unit layer 2 learns a scaled ground-truth weight matrix that also minimizes the layer 2 objective. The existence of such scaled equivalence of our layer 2 approach comes from what is defined below.

Definition 2.1 (component-wise scale transformation). A matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$ is said to be a scale transformation w.r.t. the *j*-th component if $(\mathbf{A}_{j,i})^{\top} = A_{j,j} \cdot \mathbf{e}^{(j)}$.

Additionally, estimation is more complex when the layer 2 weight matrix B^* is nonsquare. For simplicity of our algorithm presentation for layer 2, we use Condition 2.1 in the following sections.⁴

Condition 2.1 (layer 2 objective minimizer unique). A^* is not a scale transformation w.r.t. any components and B^* is a square matrix, i.e. m = d.

3 Warm-Up: Vanilla Linear Regression

Consider a ground-truth two-layer residual unit. If we assume that the inputs only contain vectors with negative entries, i.e. $\mathbf{x} < 0$, the effect of the ReLU function in the residual unit then disappears because of the nonnegativity of \mathbf{A}^* . The residual unit turns into linear model $\mathbf{y} = \mathbf{B}^* \mathbf{x}$. Thus, direct LR on samples with negative inputs can learn the exact ground-truth layer 2 parameter \mathbf{B}^* with at least d samples, when formulating the LR as a solvable full-rank linear equation system.

On the contrary, if the inputs only contain vectors with positive entries, i.e. $\mathbf{x} > 0$, all the neurons in the hidden layer are then activated by the ReLU and the nonlinearity is eliminated. The residual unit in this case turns into $\mathbf{y} = \mathbf{B}^* (\mathbf{A}^* + \mathbf{I}_d) \mathbf{x}$. Taking the value that left-multiplies \mathbf{x} as a single weight matrix \mathbf{D}^* , this is also a linear model. Direct LR on at least d samples with positive inputs by the residual unit can learn the exact \mathbf{D}^* . Since we have the access to the exact \mathbf{B}^* , solving for the exact \mathbf{A}^* can be accomplished through solving a full-rank linear equation system $\mathbf{B}^* \cdot \tilde{\mathbf{A}} = \mathbf{D}^*$, where the unique solution $\tilde{\mathbf{A}} = \mathbf{A}^* + \mathbf{I}_d$.

³Here, we are referring to the ability to identify the true model parameters using infinite samples.

⁴In Appendix E, we show that estimation of B^* remains possible without satisfying Condition 2.1.

While simple, this vanilla LR approach requires a large number of redundant samples, since sampled inputs usually have a small proportion of fully negative/positive vectors. Taking random input vectors i.i.d. with respect to each entry as an example, the probability of sampling a vector with all negative entries is p_{-}^{d} , where p_{-} is the probability of sampling a negative vector entry, then the expected number of samples to get one negative vector is $1/p_{-}^{d}$. Denoting p_{+} similarly, $1/p_{+}^{d}$ samples are expected for a positive vector. In addition, each LR step in this approach requires d such samples respectively to make the linear equation system full-rank, which implies the expected sample size to be d-exponential $d \cdot (1/p_{-}^{d} + 1/p_{+}^{d})$. For other common input distributions such as Gaussian, the proportions of fully negative/positive vectors in sampled inputs are also expected to decrease exponentially as d grows, as high-dimensional Gaussian random vectors are essentially concentrated uniformly in a sphere (Johnstone, 2007). Technical and experimental details about sample size expectations and the vanilla LR algorithm are further discussed in Appendix D.

4 Nonparametric Learning: Layer 2

We present how we learn a residual unit layer 2 under Condition 2.1 (estimating B^*): We first design an objective functional with the arguments being a matrix and a function. The objective uses expectation of a loss over the true distribution generating the data and is uniquely minimized by $[B^*]^{-1}$ and a rectifier function (ReLU). We then formulate its ERM using nonparametric estimation as a standard convex QP, further simplified as an LP that (in the noiseless case) has the same solution as the QP to learn layer 2.

4.1 Objective Design and Landscape

Consider the formulation of a residual unit as in Eq. (1). It is possible to rewrite the model as equation: $C^*\mathbf{y} = (A^*\mathbf{x})^+ + \mathbf{x}$, where the output of the hidden neuron with skip addition is on both sides of the equation, and $C^*B^* = I_d$. We aim to estimate the inverse of B^* by matrix variable C and the nonlinearity $\mathbf{x} \mapsto (A^*\mathbf{x})^+$ by a function variable h. The objective is formulated as risk functional by the L^2 error between values respectively computed by C and h,

$$G_{2}(\boldsymbol{C},h) = \frac{1}{2} \mathbb{E}_{\mathbf{x}} \left[\left\| h\left(\mathbf{x} \right) + \mathbf{x} - \boldsymbol{C} \mathbf{y} \right\|^{2} \right],$$
(2)

where the estimator $C \in \mathbb{R}^{d \times m}$, the domain of h is the nonnegative⁵ continuous⁶ $\mathbb{R}^d \to \mathbb{R}^d_{\geq 0}$ function space, written as $\mathbb{C}^0_{\geq 0}$ in shorthand. This objective is quadratic because the forward mapping $\mathbf{x} \mapsto h(\mathbf{x}) + \mathbf{x}$ and the backward mapping $\mathbf{y} \mapsto C\mathbf{y}$ are both linear w.r.t. C and h, and the two are linearly combined in an L^2 norm. The objective in Eq. (2) is minimized by the ground-truth, i.e. C^* and $\mathbf{x} \mapsto (\mathbf{A}^*\mathbf{x})^+$, is one of its minimizers. However, it is not simple to describe other variable values that minimize the objective if any exist. We give that detail under Condition 2.1, the minimizer of G_2 is unique, as the exact ground-truth in the given domain.

Theorem 4.1 (objective minimizer, layer 2). Define $G_2(\mathbf{C}, h)$ as Eq. (2), where $\mathbf{C} \in \mathbb{R}^{d \times m}$, $h \in \mathbb{C}^0_{\geq 0}$. Then under Condition 2.1, $G_2(\mathbf{C}, h)$ reaches its zero minimum iff $\mathbf{C} = [\mathbf{B}^*]^{-1}$ and $h : \mathbf{x} \mapsto (\mathbf{A}^* \mathbf{x})^+$.

Technical details are in Appendix F, where we use Lemma 4.2 (in Section 4.3) to prove a more general theorem which does not require Condition 2.1 and is sufficient for Theorem 4.1. In the next subsection, we construct the ERM of G_2 and present our convex QP formulation.

4.2 ERM with Nonparametric Estimation is Convex QP

Consider the second layer objective (Eq. (2)) with nonnegative continuous function space as the domain of h. We follow Vapnik (1991) and define its standard empirical risk functional:

$$\hat{G}_{2}(\boldsymbol{C},h) = \frac{1}{2n} \sum_{i \in [n]} \left\| h(\boldsymbol{x}^{(i)}) + \boldsymbol{x}^{(i)} - \boldsymbol{C} \boldsymbol{y}^{(i)} \right\|^{2}.$$
(3)

⁵Setting h as nonnegative ensures that a) only ReLU nonlinearity minimizes G_2 (see Theorem 4.1), and b) h's nonparametric estimator is linearly constrained (explained in Section 4.2).

 $^{^{6}}$ If h is not a continuous function, only a null set of discontinuities is possible to make G_{2} reach zero as its minimum. Setting h as continuous simplifies our theoretical results that still strictly support empirical discussion.

The function variable $h \in \mathbb{C}^{0}_{\geq 0}$ can be optimized either parametrically or nonparametrically. If we parameterize h, such nonlinearity w.r.t. its parameters would make \hat{G}_{2} lose its quadratic form. Instead, we estimate h nonparametrically: for each sample input $\boldsymbol{x}^{(i)}$, we introduce variables $\boldsymbol{\xi}^{(i)}$ that estimate mapped values by h. This avoids introducing nonlinearity to the objective and keeps \hat{G}_{2} quadratic. On the other hand, the domain of h, i.e. nonnegative continuous function space, turns into a set of linear inequalities as constraints when optimizing the learning objective nonparametrically. In this sense, learning the second layer of the residual unit can be formulated as the following QP:

$$\min_{\boldsymbol{C},\boldsymbol{\Xi}} \hat{G}_{2}^{\text{NPE}}(\boldsymbol{C},\boldsymbol{\Xi}) \coloneqq \frac{1}{2n} \sum_{i \in [n]} \left\| \boldsymbol{\xi}^{(i)} + \boldsymbol{x}^{(i)} - \boldsymbol{C} \boldsymbol{y}^{(i)} \right\|^{2}, \text{ s.t. } \boldsymbol{\xi}^{(i)} \ge \boldsymbol{0}, \forall i \in [n],$$
(4)

where $\boldsymbol{\Xi} = \{\boldsymbol{\xi}^{(i)}\}_{i=1}^{n}$ is the nonparametric estimator of $\boldsymbol{x} \mapsto (\boldsymbol{A}^* \boldsymbol{x})^+$.

Nonparametric Estimation Validation: A solution to the ERM with nonparametric estimation, i.e. the QP, is guaranteed to be sufficient for minimizing the standard empirical risk functional (Eq. (3)). More specifically, assuming C and $\Xi = \{\boldsymbol{\xi}^{(i)}\}_{i=1}^{n}$ are a solution to layer 2 QP (Eq. (4)), we see now that C and $h \in \mathbb{C}_{\geq 0}^{0}$ such that $h(\boldsymbol{x}^{(i)}) = \boldsymbol{\xi}^{(i)}$ minimize the empirical risk functional Eq. (3). Conversely, a minimizer of the standard empirical risk Eq. (3), C and h, corresponds to a solution to layer 2 QP as we set $\boldsymbol{\xi}^{(i)} = h(\boldsymbol{x}^{(i)})$. Therefore, minimizing \hat{G}_{2} and solving layer 2 QP are equivalent problems.

Convexity: The convexity of the QP: Eq. (4) is also guaranteed. First, its constraints are linear. Second, for each sample with index $i \in [n]$, the L^2 norm wraps linearity w.r.t. C and $\xi^{(i)}$. Such formulation ensures the quadratic coefficient matrix is positive semidefinite. Thus, with the sum of convex functions still being convex, the QP objective (Eq. (4)) becomes convex. Even without the knowledge of how samples are generated, this QP would be a convex program. Strict proofs are in Appendix G.

4.3 LP Simplification

Consider single-sample error written as $g_2(\boldsymbol{C}, h; \boldsymbol{x}, \boldsymbol{y}) = \frac{1}{2} \|h(\boldsymbol{x}) + \boldsymbol{x} - \boldsymbol{C}\boldsymbol{y}\|^2$. If there is a feasible \boldsymbol{C} such that $\boldsymbol{C}\boldsymbol{y} - \boldsymbol{x} \ge 0$ holds for all $\boldsymbol{x} \in \mathbb{R}^d$, then \boldsymbol{C} and $h: \boldsymbol{x} \mapsto \boldsymbol{C}\boldsymbol{y} - \boldsymbol{x}$ always minimize g_2 as zero, and thereby minimize the layer 2 objective (Eq. (2)). Thus, we obtain a condition that is equivalent to G_2 reaching the minimum in Theorem 4.1 and avoids randomness (see Lemma 4.2).

Lemma 4.2. $G_2(\mathbf{C}, h)$ reaches its zero minimum iff $\mathbf{C}\mathbf{y} - \mathbf{x} \ge 0$ holds for any $\mathbf{x} \in \mathbb{R}^d$ and its corresponding residual unit output \mathbf{y} , and $h : \mathbf{x} \mapsto \mathbf{C}\mathbf{y} - \mathbf{x}$.

The pointwise (for every $\boldsymbol{x} \in \mathbb{R}^d$) satisfaction of the inequality in Lemma 4.2 describes the solution space for the minimization of G_2 . The sufficiency of satisfying this inequality to minimize G_2 is directly obtained by assigning $h: \boldsymbol{x} \mapsto \boldsymbol{C}\boldsymbol{y} - \boldsymbol{x}$ where \boldsymbol{C} complies with $\boldsymbol{C}\boldsymbol{y} - \boldsymbol{x} \ge 0$ for all \boldsymbol{x} in the support space \mathbb{R}^d . The necessity of satisfying this inequality comes from its contraposition: If a \boldsymbol{C} violates the inequality, there must be a non-null set of \boldsymbol{x} that yield the violation due to the continuity of $\boldsymbol{C}\boldsymbol{y} - \boldsymbol{x}$. In this sense, the resulting G_2 value becomes nonzero.

Empirically speaking, we cannot solve an inequality that holds w.r.t. to every point in the support space if we only observe finite samples. We can only estimate C by solving the inequality that holds w.r.t. each sample. Following this, we formulate such estimation as to find a feasible point in the space defined by a set of linear inequalities, each of which corresponds to a sample: $Cy^{(i)} - x^{(i)} \ge 0$. Each point in the feasibility defined by the inequalities has a one-to-one correspondence in the solution space to layer 2 QP (Eq. (4)): $C \leftrightarrow (C, \{Cy^{(i)} - x^{(i)}\}_{i \in [n]})$. The set of inequalities can be solved by a standard LP with a constant objective and with constraints that are inequalities, i.e.

$$\min_{\boldsymbol{C}} \text{ const, s.t. } \boldsymbol{C}\boldsymbol{y}^{(i)} - \boldsymbol{x}^{(i)} \ge 0, \, \forall i \in [n].$$
(5)

With the one-to-one correspondences, LP: Eq. (5) and QP: Eq. (4) have identical solution spaces (feasible regions) for layer 2. Moreover, LP: Eq. (5) is also a convex program.

Algorithm 1 Learn a ReLU residual unit, layer 2.

- 1: Input: $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$, samples drawn by Eq. (1).
- 2: Output: \hat{B} , $\hat{\Xi}$, a layer 2 and $\boldsymbol{x} \mapsto (\boldsymbol{A}^*\boldsymbol{x})^+$ estimate.
- 3: Go to line 4 if QP, line 5 if LP.
- 4: Solve QP: Eq. (4) and obtain a \hat{G}_2^{NPE} minimizer, denoted by \hat{C} , $\hat{\Xi}$. Go to line 6.
- 5: Solve LP: Eq. (5) and obtain a minimizer \hat{C} , then assign $\hat{\xi}^{(i)} \leftarrow \hat{C} y^{(i)} x^{(i)}$.
- 6: return $\hat{C}^{-1}, \hat{\Xi}$.

Algorithm 1 summarizes the layer 2 estimator: Simply solve the QP/LP^7 and return the inverse of \hat{C} as the estimate of layer 2 weights and $\hat{\Xi}$ as $\boldsymbol{x} \mapsto (\boldsymbol{A}^*\boldsymbol{x})^+$ estimate. Regardless of time complexity, QP and LP in Algorithm 1 work equivalently since their solution spaces are equivalent to each other.

Our nonparametric learning directly finds a unique layer 2 estimate under Condition 2.1. In the general case without Condition 2.1 (discussed in Appendix E), nonparametric learning essentially reduces the possible values of \boldsymbol{B} from $\mathbb{R}^{m \times d}$ to a \boldsymbol{B}^* scale-equivalent matrix space, where LR uses sampled data much more efficiently than vanilla LR on layer 2 (Section 3).

5 Nonparametric Learning: Layer 1

With layer 2 learned, the outputs of the hidden neurons become observable. The two-layer problem is thereby reduced to single-layer. Consider a ground-truth single-layer model: $\mathbf{h} = (\mathbf{A}^* \mathbf{x})^+$. To construct a learning objective for this model, we rewrite the model as a nonlinearity plus a linear mapping by \mathbf{A}^* : $\mathbf{h} = (-\mathbf{A}^* \mathbf{x})^+ + \mathbf{A}^* \mathbf{x}$, where on both sides of the equation is the output of layer 1, and the nonlinearity is $\mathbf{x} \mapsto (-\mathbf{A}^* \mathbf{x})^+$. The objective of layer 1 is formulated as:

$$G_1(\boldsymbol{A}, r) = \frac{1}{2} \mathbb{E}_{\mathbf{x}} \left[\| r(\mathbf{x}) + \boldsymbol{A}\mathbf{x} - \mathbf{h} \|^2 \right],$$
(6)

where $\mathbf{A} \in \mathbb{R}^{d \times d}$ is of the layer 1 weights estimator, the domain of r is also $\mathbb{C}^{0}_{\geq 0}$. The minimizer \mathbf{A} of the risk G_1 falls into a matrix space such that for any matrix \mathbf{A} in the space, each row of \mathbf{A} is a scaled-down version of the same row of \mathbf{A}^* without changing the direction (Theorem 5.1).

Theorem 5.1 (objective minimizer space, layer 1). Define $G_1(\mathbf{A}, r)$ as Eq. (6), where $\mathbf{A} \in \mathbb{R}^{d \times d}$, $r \in \mathbb{C}^0_{\geq 0}$. Then $G_1(\mathbf{A}, r)$ reaches its zero minimum iff for each $j \in [d]$, $\mathbf{A}_{j,:} = k_j \mathbf{A}^*_{j,:}$ where $0 \leq k_j \leq 1$, and $r : \mathbf{x} \mapsto (\mathbf{A}^* \mathbf{x})^+ - \operatorname{diag}(\mathbf{k}) \cdot \mathbf{A}^* \mathbf{x}$.

The scale equivalence in the solution space is derived from a ReLU inequality: $(x)^+ \ge kx$ where $0 \le k \le 1$. Let the *j*-th row of \boldsymbol{A} be a scaled-down version of \boldsymbol{A}^* , i.e. $\boldsymbol{A}_{j,:} = k_j \boldsymbol{A}_{j,:}^*$ where $0 \le k_j \le 1$. According to the inequality, we have $(\boldsymbol{A}_{j,:}^*\boldsymbol{x})^+ \ge k_j \boldsymbol{A}_{j,:}^*\boldsymbol{x}$, which indicates $(\boldsymbol{A}^*\boldsymbol{x})^+ \ge \text{diag}(\boldsymbol{k}) \cdot \boldsymbol{A}^*\boldsymbol{x}$. Thus, *r* minimizing G_1 in Theorem 5.1 lies in feasibility $\mathbb{C}_{\ge 0}^0$ when $\boldsymbol{A} = \text{diag}(\boldsymbol{k}) \cdot \boldsymbol{A}^*$.

Due to the existence of scale equivalence, we must compute the scale factor to obtain the ground-truth weights \mathbf{A}^* . The scale factor k_j is sufficiently obtainable with $(\mathbf{A}_{j,:}\mathbf{x})^+$ and $(\mathbf{A}_{j,:}^*\mathbf{x})^+$ observable: Conditioned on nonnegative ReLU input, we have a linear model $\mathbf{A}_{j,:}\mathbf{x} = k_j \mathbf{A}_{j,:}^*\mathbf{x}$ where $\mathbf{A}_{j,:}\mathbf{x}$ and $\mathbf{A}_{j,:}^*\mathbf{x}$ are observed. Theorem 5.2 summarizes layer 1 scale factor property, allowing us to correct a minimizer of G_1 to the ground-truth weights \mathbf{A}^* by computing a scalar for each row.

Theorem 5.2 (scale factor, layer 1). Assume A is a minimizer of G_1 . Then for any $j \in [d]$, $(A_{j,:}\mathbf{x})^+ / (A_j^* \cdot \mathbf{x})^+$ is always equal to the scale factor k_j given that $(A_{j,:}\mathbf{x})^+ > 0$.

The elimination of randomness for G_1 follows the same pattern as layer 2. If there is a feasible A such that $(A^*x)^+ - Ax \ge 0$ holds for all $x \in \mathbb{R}^d$, such A is a solution to our objective by Eq. (6). We can also have the following proposition that avoids randomness.

⁷We use CVX (Grant & Boyd, 2014; 2008) that calls SDPT3 (Toh et al., 1999) (a free solver under GPLv3 license) and solves our convex QP/LP in polynomial time. See Appendix B for technical details.

Algorithm 2 Learn a ReLU residual unit, layer 1.

- 1: Input: $\{(x^{(i)}, h^{(i)})\}_{i=1}^{n}$, layer 1 samples.
- 2: **Output:** \hat{A} , a layer 1 estimate.
- 3: Solve QP: Eq. (7) or LP: Eq. (8) and obtain a \hat{G}_1^{NPE} minimizer, denoted by \hat{A} , $\hat{\Phi}$. { $\hat{\Phi}$ is no longer needed.}
- 4: for all $j \in [d]$ do
- 5: $\hat{k}_j \leftarrow \operatorname{LR}\{h_j^{(i)}, \hat{A}_{j,:} \boldsymbol{x}^{(i)}\}_{h_j^{(i)} > 0}.$
- 6: $\hat{A}_{j,:} \leftarrow \hat{A}_{j,:}/\hat{k}_j$. {Rescale \hat{A} .}
- 7: end for
- 8: return \hat{A} .

Lemma 5.3. $G_1(\mathbf{A}, r)$ reaches its zero minimum iff $\mathbf{h} - \mathbf{A}\mathbf{x} \ge 0$ holds for any $\mathbf{x} \in \mathbb{R}^d$ and its corresponding hidden output \mathbf{h} , and $r : \mathbf{x} \mapsto \mathbf{h} - \mathbf{A}\mathbf{x}$.

We now turn into empirical discussion. Similar to layer 2, we formulate layer 1 QP by its empirical objective with nonparametric estimation and linear constraints representing the nonnegativity of r:

$$\min_{\boldsymbol{A}, \boldsymbol{\Phi}} \hat{G}_{1}^{\text{NPE}}(\boldsymbol{A}, \boldsymbol{\Phi}) \coloneqq \frac{1}{2n} \sum_{i \in [n]} \left\| \boldsymbol{\phi}^{(i)} + \boldsymbol{A} \boldsymbol{x}^{(i)} - \boldsymbol{h}^{(i)} \right\|^{2}, \text{ s.t. } \boldsymbol{\phi}^{(i)} \ge \boldsymbol{0}, \forall i \in [n].$$
(7)

where $\Phi = \{\phi^{(i)}\}_{i=1}^n$ is the function estimator of $\boldsymbol{x} \mapsto (\boldsymbol{A}^*\boldsymbol{x})^+ - \operatorname{diag}(\boldsymbol{k}) \cdot \boldsymbol{A}^*\boldsymbol{x}$ where \boldsymbol{k} refers to the scaling equivalence. Similarly, the solution space of QP: Eq. (7) can be represented by a set of linear inequalities as constraints of the following efficiently solvable LP

$$\min_{\boldsymbol{A}} \text{ const, s.t. } \boldsymbol{h}^{(i)} - \boldsymbol{A}\boldsymbol{x}^{(i)} \ge \boldsymbol{0}, \, \forall \, i \in [n].$$
(8)

Algorithm 2 describes how a layer 1 is learned: First, a G_1 minimizer estimate \hat{A} is obtained by solving the QP/LP. Then for the *j*-th row, the scale factor k_j is estimated by running LR on $h_j^{(i)}$ and $\hat{A}_{j,:}\boldsymbol{x}^{(i)}$ s.t. $h_j^{(i)} > 0$ to correct \hat{A} , as $A_{j,:}\boldsymbol{x} = k_j h_j$ is an unbiased linear model given that $h_j > 0$. By nonparametric learning, the value space of \boldsymbol{A} is reduced from $\mathbb{R}^{d \times d}$ to {diag(\boldsymbol{k}) $\cdot \boldsymbol{A}^* \mid 0 \leq k_j \leq 1$ }, where LR uses sampled data much more efficiently than vanilla LR layer 1 (Section 3).

6 Full Algorithm and Analysis

The full algorithm concatenates Algorithms 1 and 2 in a layerwise fashion, with observations of input/output samples from the ground-truth network and follows these steps: a) Estimates layer 2 and nonlinearity: $\mathbf{x} \mapsto (\mathbf{A}^* \mathbf{x})^+$ by Algorithm 1. b) Estimates layer 1 by running Algorithm 2 on input samples and the nonlinearity estimate. It has provable guarantees. For empirical analysis, we use $\hat{\mathbf{C}}_n$ to denote the estimation of \mathbf{C}^* from n random samples. Similar notation is applied to other estimations. First of all, our methods to solve respective layers, Algorithms 1 and 2, are strongly consistent if any convex QPs/LPs involved can be solved exactly.

Lemma 6.1 (layer 2 strong consistency). Under Condition 2.1, $\hat{\mathbf{C}}_n \xrightarrow{a.s.} \mathbf{C}^*$ and $\hat{\mathbf{B}}_n \xrightarrow{a.s.} \mathbf{B}^*$, $n \to \infty$.

For $\mathbf{\hat{C}}_n$: In Appendix H, we prove its more general a.s. (almost sure) convergence without satisfying Condition 2.1, where the solution space to layer 2 objective is a non-compact continuous set where all the elements are scaling equivalences. We use Hausdorff distance (Rockafellar & Wets, 2009) as metric and prove that the empirical solution space a.s. converges to the theoretical solution space. Then the more general a.s. convergence holds with the strong consistency of layer 2 scale factor estimator where we use LR to estimate the scale factors.

For $\hat{\mathbf{B}}_n$: Using the continuous mapping theorem (Mann & Wald, 1943), we directly propagate the strong consistency of C^* estimator to its inverse B^* 's estimator.

According to full algorithm description, layer 1 estimation uses $\hat{\mathbf{C}}_n \mathbf{y}^{(i)} - \mathbf{x}^{(i)}$ as the outputs where $i \in [n]$. Thus, the strong consistency of the hidden neuron estimator is also guaranteed by the continuous mapping theorem. Following the proof sketch of the $\hat{\mathbf{C}}_n$ a.s. convergence, we obtain the strong consistency of layer 1 estimator (See Lemma 6.2).

Lemma 6.2 (layer 1 strong consistency). $\hat{\mathbf{A}}_n \xrightarrow{a.s.} \mathbf{A}^*, n \to \infty$.

By the continuous mapping theorem, the strong consistency of the full algorithm (Theorem 6.3), which is commonly defined by a loss function that is continuous on network weights, is implied by the strong consistency of network weights estimators (Lemmas 6.1 and 6.2). See Appendix H for proofs of strong consistency discussions in this section.

Theorem 6.3 (strong consistency). Define L as L^2 output loss. $L(\hat{\mathbf{A}}_n, \hat{\mathbf{B}}_n) \xrightarrow{a.s.} 0, n \to \infty$.

7 Experiments

In our experiments, we describe a first set of synthetic dataset experiments, that show that our algorithm identifies the true parameters from which the data were generated (Section 7.1), and then a second set of experiments that use our algorithm on standard real-world benchmark datasets (Section 7.2). In the appendices we provide further experiments. For example, in Appendix J, we describe what happens when the conditions required for the correctness of our algorithms are not satisfied.

7.1 Synthetic Data Experiments

We provide experimental analysis to demonstrate the effectiveness and robustness of our approach in comparison to stochastic gradient descent (SGD) on L^2 output loss: $L(\mathbf{A}, \mathbf{B}) = \frac{1}{2} \mathbb{E}_{\mathbf{x}} \|\hat{\mathbf{y}} - \mathbf{y}\|^2$, where we parameterize the output prediction by $\hat{\mathbf{y}} = \mathbf{B} \left[(\mathbf{A}\mathbf{x})^+ + \mathbf{x} \right]$. Our proposed methods outperform SGD in terms of sample efficiency and robustness to different network weights and noise strengths, which indicates a poor optimization landscape of L^2 output loss for ReLU residual units.

Setup: The ground-truth weights are generated through i.i.d. folded standard Gaussian⁸ and standard Gaussian for layer 1 and 2 respectively, i.e. $\mathbf{A}^* \stackrel{\text{i.i.d.}}{\sim} |\mathcal{N}|(0,1), \mathbf{B}^* \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0,1)$. The input distribution is set to be an i.i.d. zero mean Gaussian-uniform equal mixture $\mathcal{N}(-0.1, 1) - \mathcal{U}(-0.9, 1.1)$. SGD is conducted on mini-batch empirical losses of $L(\mathbf{A}, \mathbf{B})$ with batch size 32 for 256 epochs in each learning trial. We apply time-based learning rate decay $\eta = \eta_0/(1 + \gamma \cdot T)$ with initial rate $\eta_0 = 10^{-3}$ and decay rate $\gamma = 10^{-5}$, where T is the epoch number. The above hyperparameters are tuned to outperform other hyperparameters in learning ReLU residual units in terms of output errors.

Evaluation: We use relative errors to measure the accuracy of our vector/matrix estimates: For a network with weights \mathbf{A} and \mathbf{B} and its teacher network with weights \mathbf{A}^* and \mathbf{B}^* , a) layer 1 error refers to $\|\mathbf{A} - \mathbf{A}^*\| / \|\mathbf{A}^*\|$, similar to layer 2. b) output error refers to $\hat{\mathbb{E}}[\|\hat{\mathbf{y}} - \mathbf{y}\| / \|\mathbf{y}\|]$ by test data. Due to the equivalence between the solution spaces of our QP and LP without label noise, we choose LP in noiseless experiments, referred to as "ours". In addition, to reduce variance, the results of learning the same ground-truths are computed as means across 16 trials.

Sample Efficiency: Consider Figure 1. It depicts the variation in the prediction errors (warmer color, larger error) as a function of d (input dimension) and the number of samples the learning algorithm is using. We compare SGD against our algorithm. We observe that our approach to the estimating the neural network is more sample efficient. For SGD, the estimation is relatively easy with only up to 10 dimensions. As expected, once the dimension grows, the sample size required for the same level of error as our method is larger. Still, overall, our method is capable of learning robustly with small sample sizes and more efficiently than SGD even for larger sample sizes.

Network Weight Robustness: This experiment aims to verify whether our method can learn a broader class of residual units. In Table 1, we see our method shows a light-tailed distribution with nearly-zero means

⁸A folded Gaussian is the absolute value of a Gaussian, with p.d.f. $p(|\mathbf{x}|)$ where $\mathbf{x} \sim \mathcal{N}$, denoted as $|\mathcal{N}|$. We use folded Gaussian to ensure layer 1 weights are nonnegative.



Figure 1: Output errors by SGD and our method for different dimensions and sample sizes. For the same d, we fix the ground-truth network as sample size grows.

Table 1: Means / Standard deviations of network estimate errors for different network weights. Values are computed from the process of learning 128 different ground-truth networks with d = 16. 512 training samples are drawn for each learning trial.

| | Lay | er 1 | Lay | er 2 | Output | | |
|------|----------|-------|-------------|-------------|--------|-------|--|
| | Mean Std | | Mean | Mean Std | | Std | |
| SGD | 0.715 | 0.090 | 1.203 | 0.134 | 0.431 | 0.038 | |
| Ours | 0.039 | 0.008 | ≈ 0 | ≈ 0 | 0.055 | 0.008 | |



Figure 2: Respective errors of layers 1 and 2 and outputs of different label noise strengths for SGD, LP and QP. We fix the ground-truth weights with d = 10 and only the noise strength varies. 512 training samples are drawn for each learning trial.

and standard deviations for layers 1 and 2, and output errors across various ground-truth networks, whereas SGD is less robust in the same context. Our method shows strong robustness to network weight changes, indicating its applicability across the whole hypothesis class.

Noise Robustness: Figure 2 confirms the robustness of our methods when output noise exists. Samples are generated by a ground-truth residual unit with output noise being i.i.d. zero-mean Gaussian in different strengths (i.e. standard deviations). We try both QP and LP because in noisy setting the two approaches are not equivalent w.r.t. the solution space.⁹ First, SGD always gives larger errors than our methods, even though it is hardly affected by tuning the noise strength. For QP/LP, all the errors for layer 1, 2 and output grow almost linearly as noise strength increases, indicating that both QP/LP learn the optima robustly when output noise is present, where QP slightly outperforms LP.

Running Time Efficiency versus SGD In Figure 3, we compare the running time of our algorithm against the running time of SGD for ground-truth networks sampled at random. The running time of our algorithm is significantly lower, with an error level which is also significantly lower. Furthermore, this is demonstrated across input dimensions, where the difference between the running time of SGD until convergence and our algorithm's running time increases as d increases. In addition, the running time of our algorithm is fixed, and does not depend on the ground-truth network, while the running time of SGD varies based on the network sampled.

Layer 2 Weights Condition Number Robustness Our algorithm depends on a matrix-inversion step to obtain the estimated parameters. This affects the estimation of B^* . Due to this, we further investigate the effect of the condition number of B^* in the ground-truth parameters on the accuracy of estimation. For

⁹See Appendix A for noisy model discussion.



Figure 3: Output errors against running time by SGD and our LP algorithm for different dimensions of networks. For each dimension size of d = 8, d = 16 and d = 32, we report the learning curves of SGD and the final output error of our LP algorithm against time used on 5 different ground-truth networks. We used 512 training samples, drawn for learning each ground-truth network. SGD has different epochs, while the LP algorithm runs once. (CPU specification: 2.8 GHz Quad-Core Intel Core i7.)



Figure 4: Layer 2 errors against layer 2 ground-truth weight condition numbers for different dimensions of networks by SGD and our LP algorithm. Each data point is the mean across learning 32 different ground-truth networks with the same pair of d and a condition number for B^* . The condition number is denoted by $\kappa(B^*)$. We used 512 training samples, drawn for learning each ground-truth network.

this experiment, to obtain ground-truth network parameters with different condition number, we follow the procedure of Ge et al. (2019), and multiply a diagonal matrix with exponentially-dropping values on the diagonal (λ^{-i} for the *i*-th element) by two random orthonormal matrices, **U** and **V**^T, on the left and on the right respectively: $\mathbf{B} = \mathbf{U} \operatorname{diag} (\lambda^{-1}, \ldots, \lambda^{-d}) \mathbf{V}^{\mathsf{T}}$.

In Figure 4, we compare the effect of the condition number of B^* on SGD and on our algorithm. SGD turns out to be quite sensitive to the condition number, especially for lower dimensions. Our algorithm, on the other hand, is almost not affected by the condition number, with a final estimation error consistently close to 0, recovering the matrix B^* .

7.2 Experiments with Regression Benchmark Datasets

In this set of experiments, we test our algorithm and compare it against stochastic gradient descent for four datasets: HOUSING (506 examples, 13 features), DELTAELEVATORS (9,517 examples, 6 features), DELTAAILERONS (7,129 examples, 5 features), AILERONS (7,154 examples, 41 features), REDWINE (1,599 examples, 11 features), WHITEWINE (4,898 examples, 11 features) and JIGSAW (47,991 examples, 100 features). The first six datasets are standard benchmark datasets taken from Delve¹⁰ and the UCI Machine

¹⁰https://www.cs.toronto.edu/~delve/data/datasets.html

Table 2: **Results on regression benchmarks.** We give RMSE on several datasets with our algorithm (QP), backpropagation with SGD and backpropagation initialized with our algorithm estimates (SGD+QP). Numbers after slash denote the (average) number of epochs required for backpropagation to converge.

| Detect | Root Mean Squared Error | | | | | | |
|----------------|-------------------------|-----------------------|----------------------|--|--|--|--|
| Dataset | QP | SGD | SGD+QP | | | | |
| Housing | 19.46 | 18.08 / 459 | 12.43 / 1379 | | | | |
| DeltaElevators | 0.00240 | $0.01360 \ / \ 8638$ | 0.00209 / 843 | | | | |
| DeltaAilerons | 0.00030 | $0.00216 \ / \ 16555$ | $0.00030 \ / \ 3$ | | | | |
| AILERONS | 0.00070 | $0.00193 \ / \ 8039$ | $0.00023 \ / \ 1047$ | | | | |
| RedWine | 2.73 | 2.49 / 5270 | 1.48 / 5610 | | | | |
| WHITEWINE | 2.99 | 2.59 / 8200 | 1.61 / 3129 | | | | |
| JIGSAW | 1.17 | $2.45 \ / \ 6731$ | 0.92 / 2413 | | | | |

Learning Repository.¹¹ The last dataset, JIGSAW, is bigger with its goal to use word FastText¹² embeddings of tweets to predict their level of toxicity. For this set of experiments, with all datasets, we use the MOSEK solver with an academic license (ApS, 2019).

In all of our experiments, we report five-fold cross-validation results, where the first fold is used to tune the hyperparameters, and the last fold is used as both a validation set for early stopping (first half) and to report the results (second half). The models we use for backpropagation and our algorithm are identical, in the form of Eq. (1). We report (the fold-average) Root Mean Square Error (RMSE), defined as the square-root of the average of the squared deviations between the predicted value and the true value. Each fold is run with five different random seeds to initialize the backpropagation algorithm (results are averaged).

To satisfy the constraint on the length of \mathbf{y} , we duplicate the regression target multiple times, each time adding Gaussian noise as large as 0.1 of the standard deviation of the regression target. More specifically, let $\mathbf{y}^{(i)}$ be the *i*th example in a dataset. Then, we create a new $\mathbf{y}^{*,(i)} \in \mathbb{R}^d$ such that $\mathbf{y}_j^{*,(i)} = (1+0.1\zeta_{i,j})\mathbf{y}_{1+(j-1 \mod d)}^{(i)}$ for $j \geq 2$ (for i = 1, we add no noise, and just copy $\mathbf{y}^{(i)}$), where $\zeta_{i,j}$ are Gaussians with the mean being the standard deviation of \mathbf{y} over the dataset.

With backpropagation, we use SGD with a learning rate of 0.000001 and batch size of 500. We run backpropagation until the mean squared error does not change between epochs within a fraction of 1/10000. We experiment with two types of initializations for the backpropagation algorithm: one in which we initialize all the weights randomly with a standard Gaussian distribution, and one in which we initialize it with the result of our quadratic programming algorithm. It has recently been shown that initializing a ResNet with positive values (or even just zero-one values) yields an improvement in estimating the network with backpropagation Zhao et al. (2022), supporting our use of the nonnegativity constraint for the first layer. An additional advantage of such an approach is that it removes the randomness that characterizes neural network learning.

Table 2 gives the RMSE values, comparing our QP algorithm and the backpropagation algorithm. We note that our algorithm is especially potent when used to initialize the backpropagation algorithm: not only then it achieves lower RMSE on the test set, but it also does so in fewer iterations.

8 Conclusion

In this paper, we address the problem of learning a general class of two-layer residual units and propose an algorithm based on landscape design and convex optimization: We demonstrate firstly that minimizers of our objective functionals can express the *exact* ground-truth network. Then, we show that the corresponding ERM with nonparametric function estimation can be solved using *convex* QP/LP, which indicates *polynomial-time*

¹¹https://archive.ics.uci.edu/ml/index.php

¹²https://fasttext.cc/

solvability w.r.t. sample size and dimension. Moreover, our algorithms that are used to estimate both layers as well as the whole networks are *strongly consistent*, with very weak conditions on input distributions.

Limitations and Future Work The main limitation of our work is that the use of L_2 loss with the QP objective is not readily adaptable to other loss functions. Our preliminary experiments with classification datasets demonstrate that while our algorithm behaves better than L_2 backpropagation on such datasets, using the log-loss is more effective for these datasets. We leave it as future work to generalize our QP to a convex program with an arbitrary loss.

Our algorithm is not yet scalable for large datasets. We believe the learning algorithm could still be added to the ML toolkit for problems at a smaller scale for which we need a predictor, problems where linear regression is a good fit. Finally, we note that while we have nonnegativity constraints on the parameters (\mathbf{A}) , the class of networks we learn is expressive. Previous work that has similar nonnegativity constraints includes binary neural networks (Courbariaux et al., 2015) and others, without a 0/1 constraint (Chorowski & Zurada, 2014; Hosseini-Asl et al., 2015). We leave it for future work to alleviate this constraint. Recently, Zhao et al. (2022) showed that initializing ResNets with 0/1 weights is as effective for training as random initialization.

Acknowledgments

We thank Yftah Ziser, the anonymous reviewers and the action editor for their useful feedback and comments. Computational resources and software were partially provided by Edinburgh Parallel Computing Centre and through the Amazon Enterprise License for MATLAB.

References

- Zeyuan Allen-Zhu and Yuanzhi Li. What can ResNet learn efficiently, going beyond kernels? In Advances in Neural Information Processing Systems, volume 32, 2019.
- Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in overparameterized neural networks, going beyond two layers. In *Advances in Neural Information Processing Systems*, volume 32, 2019.
- Mosek ApS. MOSEK optimization toolbox for MATLAB. User's Guide and Reference Manual, Version, 4, 2019.
- Yossi Arjevani and Michael Field. Analytic study of families of spurious minima in two-layer ReLU neural networks: a tale of symmetry ii. In Advances in Neural Information Processing Systems, volume 34, 2021.
- Sanjeev Arora, Aditya Bhaskara, Rong Ge, and Tengyu Ma. Provable bounds for learning some deep representations. In *Proceedings of the 31st International Conference on Machine Learning*, volume 32, pp. 584–592. PMLR, 2014.
- Miriam Ayer, H. Daniel Brunk, George M. Ewing, William T. Reid, and Edward Silverman. An empirical distribution function for sampling with incomplete information. *The Annals of Mathematical Statistics*, pp. 641–647, 1955.
- Eugene Belilovsky, Michael Eickenberg, and Edouard Oyallon. Greedy layerwise learning can scale to ImageNet. In Proceedings of the 36th International Conference on Machine Learning, volume 97, pp. 583–593. PMLR, 2019.
- Yoshua Bengio and Olivier Delalleau. On the expressive power of deep architectures. In Algorithmic Learning Theory, pp. 18–36. Springer, 2011.
- M. Émile Borel. Les probabilités dénombrables et leurs applications arithmétiques. Rendiconti del Circolo Matematico di Palermo (1884-1940), 27(1):247–271, 1909.
- Hugh D. Brunk. Maximum likelihood estimates of monotone parameters. The Annals of Mathematical Statistics, pp. 607–616, 1955.

- Alon Brutzkus and Amir Globerson. Globally optimal gradient descent for a ConvNet with Gaussian inputs. In Proceedings of the 34th International Conference on Machine Learning, volume 70, pp. 605–614. PMLR, 2017.
- Yining Chen and Richard J. Samworth. Generalized additive and index models with shape constraints. Journal of the Royal Statistical Society: Series B (Statistical Methodology), 78(4):729–754, 2016.
- Jan Chorowski and Jacek M. Zurada. Learning understandable neural networks with nonnegative weight constraints. *IEEE Transactions on Neural Networks and Learning Systems*, 26(1):62–69, 2014.
- Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. Binaryconnect: Training deep neural networks with binary weights during propagations. In *Advances in Neural Information Processing Systems*, volume 28, 2015.
- Simon S. Du and Surbhi Goel. Improved learning of one-hidden-layer convolutional neural networks with overlaps. arXiv preprint arXiv:1805.07798, 2018.
- Simon S. Du, Jason D. Lee, Yuandong Tian, Aarti Singh, and Barnabas Poczos. Gradient descent learns one-hidden-layer CNN: Don't be afraid of spurious local minima. In *Proceedings of the 35th International Conference on Machine Learning*, volume 80, pp. 1339–1348. PMLR, 2018.
- Tolga Ergen and Mert Pilanci. Implicit convex regularizers of CNN architectures: Convex optimization of twoand three-layer networks in polynomial time. In *International Conference on Learning Representations*, 2021a.
- Tolga Ergen and Mert Pilanci. Revealing the structure of deep neural networks via convex duality. In *Proceedings of the 38th International Conference on Machine Learning*, volume 139, pp. 3004–3014. PMLR, 2021b.
- Julius Farkas. Theorie der einfachen ungleichungen. Journal für die reine und angewandte Mathematik (Crelles Journal), 1902(124):1–27, 1902.
- Rong Ge, Jason D. Lee, and Tengyu Ma. Learning one-hidden-layer neural networks with landscape design. In *International Conference on Learning Representations*, 2018.
- Rong Ge, Rohith Kuditipudi, Zhize Li, and Xiang Wang. Learning two-layer neural networks with symmetric inputs. In *International Conference on Learning Representations*, 2019.
- Surbhi Goel and Adam R. Klivans. Learning neural networks with two nonlinear layers in polynomial time. In Proceedings of the 32nd Conference on Learning Theory, volume 99, pp. 1470–1499. PMLR, 2019.
- Surbhi Goel, Adam R. Klivans, and Raghu Meka. Learning one convolutional layer with overlapping patches. In Proceedings of the 35th International Conference on Machine Learning, volume 80, pp. 1783–1791. PMLR, 2018.
- Michael Grant and Stephen Boyd. Graph implementations for nonsmooth convex programs. In V. Blondel, S. Boyd, and H. Kimura (eds.), *Recent Advances in Learning and Control*, Lecture Notes in Control and Information Sciences, pp. 95–110. Springer-Verlag Limited, 2008. http://stanford.edu/~boyd/graph_ dcp.html.
- Michael Grant and Stephen Boyd. CVX: MATLAB software for disciplined convex programming, version 2.1. http://cvxr.com/cvx, March 2014.
- Adityanand Guntuboyina and Bodhisattva Sen. Nonparametric shape-restricted regression. *Statistical Science*, 33(4):568–594, 2018.
- James Douglas Hamilton. Time series analysis. Princeton University Press, 2020.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pp. 770–778, 2016a.

- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In European Conference on Computer Vision, pp. 630–645. Springer, 2016b.
- Ehsan Hosseini-Asl, Jacek M. Zurada, and Olfa Nasraoui. Deep learning of part-based representation of data using sparse autoencoders with nonnegativity constraints. *IEEE Transactions on Neural Networks and Learning Systems*, 27(12):2486–2498, 2015.
- Majid Janzamin, Hanie Sedghi, and Anima Anandkumar. Beating the perils of non-convexity: Guaranteed training of neural networks using tensor methods. arXiv preprint arXiv:1506.08473, 2015.
- Florian Jarre. On the convergence of the method of analytic centers when applied to convex quadratic programs. *Mathematical Programming*, 49(1-3):341–358, 1990.
- Robert I. Jennrich. Asymptotic properties of non-linear least squares estimators. The Annals of Mathematical Statistics, 40(2):633–643, 1969.
- Iain M. Johnstone. High dimensional statistical inference and random matrices. In Proceedings of the International Congress of Mathematicians Madrid, August 22–30, 2006, pp. 307–333, 2007.
- Sham M. Kakade, Varun Kanade, Ohad Shamir, and Adam Kalai. Efficient learning of generalized linear and single index models with isotonic regression. In Advances in Neural Information Processing Systems, volume 24, 2011.
- Narendra Karmarkar. A new polynomial-time algorithm for linear programming. In Proceedings of the 16th Annual ACM Symposium on Theory of Computing, pp. 302–311, 1984.
- Yoon Kim. Convolutional neural networks for sentence classification. In Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP), pp. 1746–1751. Association for Computational Linguistics, 2014. doi: 10.3115/v1/D14-1181. URL https://aclanthology.org/D14-1181.
- Mikhail K. Kozlov, Sergei P. Tarasov, and Leonid G. Khachiyan. The polynomial solvability of convex quadratic programming. USSR Computational Mathematics and Mathematical Physics, 20(5):223–228, 1980.
- Alex Krizhevsky, Ilya Sutskever, and Geoffrey E. Hinton. ImageNet classification with deep convolutional neural networks. In Advances in Neural Information Processing Systems, volume 25, 2012.
- Arun K. Kuchibhotla, Rohit K. Patra, and Bodhisattva Sen. Semiparametric efficiency in convexity constrained single-index model. Journal of the American Statistical Association, pp. 1–15, 2021.
- Yann LeCun, Léon Bottou, Yoshua Bengio, Patrick Haffner, et al. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Yuanzhi Li and Yang Yuan. Convergence analysis of two-layer neural networks with ReLU activation. In Advances in Neural Information Processing Systems, volume 30, 2017.
- Roi Livni, Shai Shalev-Shwartz, and Ohad Shamir. On the computational efficiency of training neural networks. In Advances in Neural Information Processing Systems, volume 27, 2014.
- David G. Luenberger. Optimization by vector space methods. John Wiley & Sons, 1997.
- Andrew L. Maas, Awni Y. Hannun, and Andrew Y. Ng. Rectifier nonlinearities improve neural network acoustic models. In Proceedings of the 30th International Conference on Machine Learning Workshop on Deep Learning for Audio, Speech and Language Processing, 2013.
- Henry B. Mann and Abraham Wald. On stochastic limit and order relationships. The Annals of Mathematical Statistics, 14(3):217–226, 1943.
- Mary C. Meyer. A simple new algorithm for quadratic programming with applications in statistics. Communications in Statistics - Simulation and Computation, 42:1126–1139, 2013.

- Aaron Mishkin, Arda Sahiner, and Mert Pilanci. Fast convex optimization for two-layer ReLU networks: Equivalent model classes and cone decompositions. In *Proceedings of the 39th International Conference on Machine Learning*, volume 162, pp. 15770–15816. PMLR, 2022.
- Renato D. C. Monteiro and Ilan Adler. Interior path following primal-dual algorithms. part ii: Convex quadratic programming. *Mathematical Programming*, 44(1):43–66, 1989.
- Katta G. Murty. Computational complexity of parametric linear programming. *Mathematical Programming*, 19(1):213–219, 1980.
- Vinod Nair and Geoffrey E. Hinton. Rectified linear units improve restricted Boltzmann machines. In Proceedings of the 27th International Conference on Machine Learning, pp. 807–814, 2010.
- Mert Pilanci and Tolga Ergen. Neural networks are convex regularizers: Exact polynomial-time convex optimization formulations for two-layer networks. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119, pp. 7695–7705. PMLR, 2020.
- Guillaume Rabusseau, Tianyu Li, and Doina Precup. Connecting weighted automata and recurrent neural networks through spectral learning. In *Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics*, volume 89, pp. 1630–1639. PMLR, 2019.
- R. Tyrrell Rockafellar and Roger J-B Wets. Variational analysis, volume 317. Springer Science & Business Media, 2009.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. ImageNet large scale visual recognition challenge. International Journal of Computer Vision, 115(3):211–252, 2015.
- Arda Sahiner, Tolga Ergen, John M. Pauly, and Mert Pilanci. Vector-output ReLU neural network problems are copositive programs: Convex analysis of two layer networks and polynomial-time algorithms. In International Conference on Learning Representations, 2021.
- Brian K. Schmidt and T. H. Mattheiss. The probability that a random polytope is bounded. *Mathematics of Operations Research*, 2(3):292–296, 1977.
- Emilio Seijo and Bodhisattva Sen. Nonparametric least squares estimation of a multivariate convex regression function. *The Annals of Statistics*, 39(3):1633–1657, 2011.
- Pranab K. Sen and Julio M. Singer. Large sample methods in statistics: an introduction with applications, volume 25. CRC Press, 1994.
- Alexander Shapiro, Darinka Dentcheva, and Andrzej Ruszczyński. Lectures on stochastic programming: modeling and theory. SIAM, 2014.
- Mahdi Soltanolkotabi. Learning ReLUs via gradient descent. In Advances in Neural Information Processing Systems, volume 30, 2017.
- Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. Sequence to sequence learning with neural networks. In Advances in Neural Information Processing Systems, volume 27, 2014.
- Yuandong Tian. An analytical formula of population gradient for two-layered ReLU network and its applications in convergence and critical point analysis. In *Proceedings of the 34th International Conference* on Machine Learning, volume 70, pp. 3404–3413, 2017.
- Kim-Chuan Toh, Michael J. Todd, and Reha H. Tütüncü. SDPT3 a MATLAB software package for semidefinite programming, version 1.3. Optimization Methods and Software, 11(1-4):545–581, 1999.
- Csaba D. Toth, Joseph O'Rourke, and Jacob E. Goodman. *Handbook of discrete and computational geometry*. CRC Press, 2017.
- Leslie G. Valiant. A theory of the learnable. Communications of the ACM, 27(11):1134–1142, 1984.

- Constance van Eeden. Maximum likelihood estimation of ordered probabilities, 2. Stichting Mathematisch Centrum. Statistische Afdeling, (S 196/56), 1956.
- Vladimir Vapnik. Principles of risk minimization for learning theory. In Advances in Neural Information Processing Systems, volume 4, 1991.
- Vladimir Vapnik. The nature of statistical learning theory. Springer Science & Business Media, 1999.
- Shanshan Wu, Alexandros G. Dimakis, and Sujay Sanghavi. Learning distributions generated by one-layer ReLU networks. In Advances in Neural Information Processing Systems, volume 32, 2019.
- Qiuling Yang, Alireza Sadeghi, Gang Wang, and Jian Sun. Learning two-layer ReLU networks is nearly as easy as learning linear classifiers on separable data. *IEEE Transactions on Signal Processing*, 69:4416–4427, 2021.
- Yinyu Ye and Edison Tse. An extension of Karmarkar's projective algorithm for convex quadratic programming. Mathematical Programming, 44(1-3):157–179, 1989.
- Xiao Zhang, Yaodong Yu, Lingxiao Wang, and Quanquan Gu. Learning one-hidden-layer ReLU networks via gradient descent. In *Proceedings of the 22nd International Conference on Artificial Intelligence and Statistics*, volume 89, pp. 1524–1534. PMLR, 2019.
- Yin Zhang, Richard A. Tapia, and John E. Dennis, Jr. On the superlinear and quadratic convergence of primal-dual interior point linear programming algorithms. SIAM Journal on Optimization, 2(2):304–324, 1992.
- Jiawei Zhao, Florian Tobias Schaefer, and Anima Anandkumar. ZerO initialization: Initializing neural networks with only zeros and ones. *Transactions on Machine Learning Research*, 2022. URL https://openreview.net/forum?id=1AxQpKmiTc.
- Kai Zhong, Zhao Song, Prateek Jain, Peter L. Bartlett, and Inderjit S. Dhillon. Recovery guarantees for one-hidden-layer neural networks. In *Proceedings of the 34th International Conference on Machine Learning*, volume 70, pp. 4140–4149, 2017.

Appendices

We include an outline of the appendices:

- Appendix A: A brief discussion of our learning algorithms in a noisy context, referring to the LP slack variable technique that is used in our experiments with noise in the main paper. Also, this section provides an insight to potential future direction of this work.
- Appendix B: A detailed explanation of the computational complexity and the methods through which the QPs/LPs in this paper are solved.
- Appendix C: A preliminarily proposed general optimization formula for the extension of learning multi-layer ResNets, to show the scalability of our algorithm.
- Appendix D: A formal result and an empirical validation of the example given in the main paper in the context of the vanilla linear regression method (Section 3).
- Appendix E: A generalization of layer 2 learning to the case where Condition 2.1 is not satisfied.
- Appendix F: Proofs regarding the minimizers of the objective functions we use.
- Appendix G: Proofs that justify the convexity of our QPs.
- Appendix H: Proofs that show our estimation algorithm is strongly consistent.

- Appendix I: A sample complexity justification of how our QP/LP approaches improves the exponential convergence rate that the vanilla LR approach (Section 3) has.
- Appendix J: Experiments with negative entries in **A** and a demonstration of the practicalities of the nonnegative entries assumption.

A Discussion: Noisy Model

Here, we discuss our learning methods in the noisy case. We introduce output noise to our model, namely

$$\mathbf{y} = \boldsymbol{B}^* \left[\left(\boldsymbol{A}^* \mathbf{x} \right)^+ + \mathbf{x} \right] + \mathbf{z},$$

where the label noise $\mathbf{z} \in \mathbb{R}^d$ is an i.i.d. random vector with respect to each component, satisfying $\mathbb{E}[\mathbf{z}] = 0$. In addition \mathbf{z} and \mathbf{x} are statistically independent.

Taking layer 2 as an example, our original objective functional does not reach zero by substituting the ground-truth: $G_2(\mathbf{C}^*, \mathbf{x} \mapsto (\mathbf{A}^* \mathbf{x})^+) = \mathbb{E}_{\mathbf{z}}[\|\mathbf{C}^* \mathbf{z}\|^2] = \sigma^2 \operatorname{Tr} \mathbf{C}^* \mathbf{C}^*^\top$ where σ is the noise strength. However, if layer 2 is well conditioned, the ground-truth will still assign a value close to zero to the objective. In this sense, with G_2 's continuity, the ground-truth can approximately minimize G_2 , which validates our QP approach in terms of learning noisy residual units.

Our original LP (Eq. (5)) fails to give feasible solutions due to possible violation of the inequality $C^*y - x \ge 0$, since $C^*y - x = (A^*x)^+ + C^*z$ is not necessarily an entrywise nonnegative vector because the term C^*z might have negative entries. So we introduce slack variables $\zeta^{(i)}$ to soften the constraints:

$$\begin{split} \min_{\boldsymbol{C},\boldsymbol{Z}} & \frac{1}{n} \sum_{i \in [n]} \mathbf{1}^\top \cdot \boldsymbol{\zeta}^{(i)}, \\ \text{s.t.} & \boldsymbol{C} \boldsymbol{y}^{(i)} - \boldsymbol{x}^{(i)} \geq -\boldsymbol{\zeta}^{(i)}, \, \boldsymbol{\zeta}^{(i)} \geq \boldsymbol{0}. \end{split}$$

For a sample $(\boldsymbol{x}^{(i)}, \boldsymbol{y}^{(i)})$ with noise $\boldsymbol{z}^{(i)}$, if $(\boldsymbol{A}^* \boldsymbol{x}^{(i)})^+ + \boldsymbol{C}^* \boldsymbol{z}^{(i)} < 0$, then its L^1 norm would be added to the objective. This method remedies violations to the inequality $C^* \mathbf{y} - \mathbf{x} \ge 0$. With access to a sufficiently large sample and with the "stability" assumption, our solution \hat{C} would be close to C^* since large deviations seldom rise and their penalties are diluted in the objective.

B Discussion: Solving Convex Programs

The theoretical foundation of solving convex QP and LP has been driven to maturity in terms of computational complexity (Kozlov et al., 1980; Murty, 1980) and convergence analysis (Jarre, 1990; Zhang et al., 1992). The time complexity for a convex QP/LP is analyzed in terms of the number of scalar variables N and the number of bits L in the input (Ye & Tse, 1989). For example, under Condition 2.1, $N = d^2 + nd$ for the QP because the matrix variables have d^2 scalars and the function estimators have nd scalars. Similarly we have $N = d^2$ for the noiseless LP and $N = d^2 + nd$ for the noisy LP. It is guaranteed that primal-dual interior point methods can solve the convex QP/LP in a polynomial number of iterations $\mathcal{O}(\sqrt{NL})$, where each iteration costs at worst $\mathcal{O}(N^{2.5})$ arithmetic operations from the Cholesky decomposition for the needed matrix inversion (Monteiro & Adler, 1989; Karmarkar, 1984). This indicates that our QPs/LPs are guaranteed to be solvable in $\mathcal{O}(N^3L)$ arithmetic operations, which is generally an $\mathcal{O}(\text{poly}(n, m, d, L))$ complexity.

To solve convex QPs/LPs in this paper, we mostly use CVX, a commonly used package for specifying and solving convex programs (Grant & Boyd, 2014; 2008). For both QPs and LPs, CVX calls a solver, SDPT3 (Toh et al., 1999), which is specified for semidefinite-quadratic-linear programming and applies interior-point methods with the computational complexity mentioned above. Experimentally, SDPT3 indeed specifies and solves our programs fast and makes our numerical results robust and stable.

Algorithm 3 Learn a ReLU residual unit by LR.

- 1: Input: $\{(\boldsymbol{x}^{(i)}, \boldsymbol{y}^{(i)})\}_{i=1}^{n}$, *n* samples drawn by Eq. (1).
- 2: Output: \hat{A} , \hat{B} , estimated weight matrices.
- 3: $\hat{\boldsymbol{B}} \leftarrow \operatorname{LR}\{(\boldsymbol{x}, \boldsymbol{y}) \in \{(\boldsymbol{x}^{(i)}, \boldsymbol{y}^{(i)})\}_{i=1}^{\lfloor n/2 \rfloor} \mid \boldsymbol{x} < 0\}.$
- 4: $\hat{\boldsymbol{D}} \leftarrow \operatorname{LR}\{(\boldsymbol{x}, \boldsymbol{y}) \in \{(\boldsymbol{x}^{(i)}, \boldsymbol{y}^{(i)})\}_{i=\lfloor n/2 \rfloor + 1}^{n} \mid \boldsymbol{x} > 0\}.$ {the estimation of $\boldsymbol{B}^* (\boldsymbol{A}^* + \boldsymbol{I}_d)$ }
- 5: Solve full-rank linear equation system $\hat{B} \cdot \tilde{A} = \hat{D}$
- 6: return $A I_d$, B.

C Discussion: Extension to Multi-Layer Networks

In this appendix, we propose a generalized optimization formula of nonparametric learning of multi-layer residual networks to show that our algorithm is scalable. Assume we have an *L*-layer network of the form

$$egin{aligned} m{r}_0 &= m{x} \ m{r}_\ell &= ig(m{W}^*_{\ell-1} m{r}_{\ell-1} ig)^+ + m{r}_{\ell-1}, & \ell = 1, \dots, L-1 \ m{y} &= m{W}^*_{L-1} m{r}_{L-1}. \end{aligned}$$

We can scale the objective to use nonparametric variables $\{\boldsymbol{\xi}_{\ell}^{(i)}\}_{l \in [L-1]}$ to estimate every ReLU activation $(\boldsymbol{W}_{\ell}^* \boldsymbol{r}_{\ell})^+$ in this network. Generally, for the ℓ -th ReLU we can have the outer layer objective

$$\min_{\boldsymbol{\xi}_{[\ell]}^{(i)}, \boldsymbol{V}_{\ell}} \left\| \boldsymbol{\xi}_{\ell}^{(i)} + \dots + \boldsymbol{\xi}_{1}^{(i)} + \boldsymbol{x}^{(i)} - \boldsymbol{V}_{\ell} \hat{\boldsymbol{h}}_{\ell}^{(i)} \right\|, \text{ s.t. } \boldsymbol{\xi}_{[\ell]}^{(i)} \ge 0$$

where hat denotes estimated values and $\hat{h}_{L-1}^{(i)} = y^{(i)}$, and its corresponding inner layer objective

$$\min_{\boldsymbol{\phi}_{\ell}^{(i)}, \boldsymbol{W}_{\ell-1}, \boldsymbol{\xi}_{\ell}^{(i)}} \left\| \boldsymbol{\phi}_{\ell}^{(i)} + \boldsymbol{W}_{\ell-1} \left(\hat{\boldsymbol{\xi}}_{\ell-1}^{(i)} + \dots + \hat{\boldsymbol{\xi}}_{1}^{(i)} + \boldsymbol{x}^{(i)} \right) - \boldsymbol{\xi}_{\ell}^{(i)} \right\|, \text{ s.t. } \boldsymbol{\phi}_{\ell}^{(i)} \ge 0, \, \boldsymbol{\xi}_{\ell}^{(i)} \ge 0.$$

We can apply techniques, for example, tuning the combinations of those objectives with respect to the layers, using a stage-wise optimization of the ξ auxiliary variables, or weighting these variables, to keep the convex quadratic form and at the same time improve its optimization landscape. Specifically when L = 2, this optimization problem becomes the two-layer ResNet learning discussed in the main paper.

D Vanilla Linear Regression: More Details

Algorithm 3 gives the full list of steps of learning two-layer residual units by LR, where we first split the drawn samples into two halves for the respective two key steps, then for both halves we filter negative and positive vectors, respectively. By running LR on both filtered sets we obtain an estimation of the ground-truth network.

In the following sections, we formalize this intuition and describe the exponential sample complexity of this approach under entry-wise i.i.d. setting as an example of the inefficiency of this method.

Theorem D.1 (exponential sample complexity, vanilla LR). Assume the input vectors are i.i.d. with respect to each component. Then Algorithm 3 learns a neural network from a ground-truth residual unit with at least $\mathcal{O}(d \cdot 2^{d+1})$ expected number of samples.

Proof. For each $j \in [d]$, let P_j be the marginal probability of \mathbf{x}_j being positive, i.e. $P_j = P(\mathbf{x}_j > 0) > 0$. Thus, the probability that a sample is positive $P(\mathbf{x} > 0) = \prod_{j \in [d]} P_j$. Then the expected number of sampling trials to obtain d positive samples is $d/\prod_{j \in [d]} P_j$. Similarly, the expected number of sampling trials to obtain d negative samples is $d/\prod_{j \in [d]} [1 - P_j - P(\mathbf{x}_j = 0)]$. Let n be a random natural number s.t. with n samples Table 3: Learning success rate of vanilla LR on residual units with input $\mathbf{x} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0,1)$. The success rates shows an exponentially decreasing trend with the same number of samples. The sample sizes to achieve close to the same rate grows exponentially as the number of dimensions grows linearly.

| d n | le1 | 1e2 | 5e2 | 1e3 | 5e3 | 1e4 | 5e4 | 1e5 |
|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| 4 | 0.002 | 0.137 | 0.999 | 1 | 1 | 1 | 1 | 1 |
| 6 | 0 | 0 | 0.041 | 0.614 | 1 | 1 | 1 | 1 |
| 8 | 0 | 0 | 0 | 0 | 0.579 | 0.998 | 1 | 1 |
| 10 | 0 | 0 | 0 | 0 | 0 | 0.002 | 1 | 1 |
| 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0.001 | 0.322 |

the network has learned. The expected number of samples that guarantees successful learning is the sum of the two expectations

$$\begin{split} \mathbb{E}[\mathbf{n}] &= d \left\{ \frac{1}{\prod_{j \in [d]} P_j} + \frac{1}{\prod_{j \in [d]} [1 - P_j - P(\mathbf{x}_j = 0)]} \right\} \\ &\geq d \left[\frac{1}{\prod_{j \in [d]} P_j} + \frac{1}{\prod_{j \in [d]} (1 - P_j)} \right] \\ &\geq d \cdot 2 \prod_{j \in [d]} \frac{1}{\sqrt{P_j (1 - P_j)}} \\ &\geq d \cdot 2 \prod_{j \in [d]} \frac{1}{(P_j + 1 - P_j)/2} \\ &= d \cdot 2^{d+1}. \end{split}$$

Thus $\mathbb{E}[\mathbf{n}] \geq \mathcal{O}(d \cdot 2^{d+1}).$

With an exponential sample complexity, the time complexity of the LR approach is thereby also exponential because filtering through all the samples costs time with the same complexity as the number of samples, even though the LR itself only costs polynomial time in the number of samples. To summarize, this vanilla LR approach learns exact ground-truth residual units but with exponential complexity in terms of both computational cost and sample size. While this approach to learning a residual unit such as ours is simple and intuitive, we aim to improve on this approach, making full use of all samples available.

Experiments: Sample Efficiency We present the results of the vanilla LR in learning the residual units. As discussed in Section 3, LR requires full-rank linear systems parameterized by samples to learn the exact ground-truth parameters. However, with degenerate linear systems, LR is completely incapable of learning the parameters. Therefore, there exists a hard threshold for the number of samples required by LR that makes the linear equation system full-rank. With such a dichotomous constraint on LR learning, we take the learning success rate among 1000 trials as the metric to evaluate the performance of LR with different sample sizes and number of dimensions.

Table 3 shows the learning success rates with zero-mean Gaussian inputs. For each fixed number of dimensions, it appears there is a hard threshold that switches the learnability of LR. The exponential sample complexity is also reflected there as the number of dimensions grows linearly. Table 4 shows the rates with input mean non-zero, where an overall decline in success rates happens. This observation is explainable because the bottleneck of LR learning is the lower value found between the probability of sampling a positive and a negative vector. A positive mean reduces the probability of the latter, and thereby increases the sample size required.

Table 4: Learning success rate of vanilla LR on residual units with input $\mathbf{x} \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0.1, 1)$. With the non-zero Gaussian mean, the success rates show an overall decline compared with the rates shown in Table 3 for zero-mean Gaussian inputs.

| n d | 1e1 | 1e2 | 5e2 | 1e3 | 5e3 | 1e4 | 5e4 | 1e5 |
|--------|-------|-------|-------|-------|-------|-------|-------|-------|
| 4 | 0.001 | 0.122 | 0.996 | 1 | 1 | 1 | 1 | 1 |
| 6 | 0 | 0 | 0.030 | 0.346 | 1 | 1 | 1 | 1 |
| 8 | 0 | 0 | 0 | 0 | 0.116 | 0.792 | 1 | 1 |
| 10 | 0 | 0 | 0 | 0 | 0 | 0 | 0.619 | 1 |
| 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.001 |

Generalization of Layer 2 Learning Ε

In this appendix, we discuss how layer 2 is learned without satisfying Condition 2.1. We note that the fact that A^* is a scale transformation w.r.t. some components leads to scaling equivalence to our layer 2 objective functional minimizers just as in layer 1. To be more specific, we give a general version of Theorem 4.1 that handles the case without Condition 2.1 and describes the scaling equivalence of G_2 minimizers. See Theorem E.1 for a formal description.

Theorem E.1 (objective minimizer space, layer 2, general). Let G_2 be $\frac{1}{2}\mathbb{E}_{\mathbf{x}}\left[\left\|h\left(\mathbf{x}\right)+\mathbf{x}-C\mathbf{y}\right\|^2\right]$ as a functional, where $C \in \mathbb{R}^{d \times m}$, $h \in \mathbb{C}^0_{\geq 0}$. Then $G_2(C, h)$ reaches its zero minimum iff $CB^* = \operatorname{diag}(k)$ where for each $j \in [d]$,

a) 1/(1+A^{*}_{j,j}) ≤ k_j ≤ 1, if A^{*} is a scale transformation w.r.t. the j-th component.
b) k_j = 1, if A^{*} is not a scale transformation w.r.t. the j-th component.

As in layer 1, the scaling equivalence in layer 2 can also be obtained by assigning A^* as a scale transformation w.r.t. some component and using the properties of the ReLU nonlinearity. See Appendix F.1 for a detailed explanation of Theorem E.1. To obtain the scale factors k and correct a scale-equivalent G_2 minimizer C to the ground-truth, we observe linear models parameterized by the scale factors, where for each $j \in [d]$, $[Cy]_j = k_j x_j$ given that $x_j < 0$. See Theorem E.2 and its proof for justifications of the linear models that are used to compute k.

Theorem E.2 (scale factor, layer 2, general). Assume C is a minimizer of G_1 in the context of Theorem E.1. Then for any $j \in [d]$, the following three propositions are equivalent:

- a) A^* is a scale transformation w.r.t. the *j*-th component.
- b) $[\mathbf{Cy}]_j / \mathbf{x}_j$ is a constant $c_j^{\mathbf{n}}$ given that $\mathbf{x}_j < 0$. c) $[\mathbf{Cy}]_j / \mathbf{x}_j$ is a constant $c_j^{\mathbf{p}}$ given that $\mathbf{x}_j > 0$.

Additionally, if one of the above three propositions is true, then $k_j = c_j^n$.

Proof. For $j \in [d]$, we first prove a) \implies b), c): From $(\mathbf{A}_{j,i}^*)^\top = A_{j,j}^* \cdot \mathbf{e}^{(j)}$ we have

$$\frac{\left[\mathbf{C}\mathbf{y}\right]_{j}}{\mathbf{x}_{j}} = \frac{\left[k_{j}\boldsymbol{e}^{(j)}\right]^{\top}\left[\left(\boldsymbol{A}^{*}\mathbf{x}\right)^{+} + \mathbf{x}\right]}{\mathbf{x}_{j}}$$
$$= \frac{k_{j}\left[\left(\boldsymbol{A}_{j,:}^{*}\mathbf{x}\right)^{+} + \mathbf{x}_{j}\right]}{\mathbf{x}_{j}} = \frac{k_{j}\left[\left(\boldsymbol{A}_{j,:j}^{*}\mathbf{x}_{j}\right)^{+} + \mathbf{x}_{j}\right]}{\mathbf{x}_{j}}$$
$$= \begin{cases} k_{j}, & \mathbf{x}_{j} < 0\\ k_{j}\left(\boldsymbol{A}_{j,:j}^{*} + 1\right), & \mathbf{x}_{j} > 0 \end{cases}.$$

Algorithm 4 Rescale a \hat{G}_2^{NPE} minimizer.

1: **Parameters:** $\varepsilon_{tol} > 0$, LR objective tolerance.

- 2: Input: $\{(\boldsymbol{x}^{(i)}, \boldsymbol{y}^{(i)})\}_{i=1}^{n}$, samples drawn by Eq. (1); $\hat{\boldsymbol{C}}$, a \hat{G}_{2}^{NPE} minimizer.
- 3: Output: \hat{k} , a layer 2 scale factor estimate.
- 4: for each $j \in [d]$ do

5: $\hat{k_j} \leftarrow \operatorname{LR}\left\{x_j^{(i)}, \left[\hat{\boldsymbol{C}}\boldsymbol{y}^{(i)}\right]_j\right\}_{x_i^{(i)} < 0}$

- 6: **if** the LR objective optimal $\hat{R}_j(\hat{k}_j) > \varepsilon_{\text{tol}}$ **then**
- 7: $\hat{k_j} \leftarrow 1$.
- 8: end if
- 9: end for
- 10: return \hat{k} .

Algorithm 5 Learn a ReLU residual unit layer 2.

- 1: Input: $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$, samples drawn by Eq. (1).
- 2: Output: \hat{B} , $\hat{\Xi}$, a layer 2 and $x \mapsto (A^*x)^+$ estimate.
- 3: Go to line 4 if QP, line 5 if LP.
- 4: Solve QP: Eq. (4) and obtain a \hat{G}_2^{NPE} minimizer, denoted by \hat{C} , $\hat{\Xi}$. Go to line 6.
- 5: Solve LP: Eq. (5) and obtain a minimizer \hat{C} , then assign $\hat{\xi}^{(i)} \leftarrow \hat{C}y^{(i)} x^{(i)}$ for each $i \in [n]$.
- 6: Run Algorithm 4 on $\{(\boldsymbol{x}^{(i)}, \boldsymbol{y}^{(i)})\}_{i \in [n]}, \hat{\boldsymbol{C}}$ and obtain $\hat{\boldsymbol{k}}$.

7:
$$\hat{\boldsymbol{B}} \leftarrow \operatorname{LR} \left\{ \operatorname{diag}^{-1}(\hat{\boldsymbol{k}}) \cdot \hat{\boldsymbol{C}} \boldsymbol{y}^{(i)}, \boldsymbol{y}^{(i)} \right\}_{i \in [n]}$$
.
8: $\hat{\boldsymbol{\xi}}^{(i)} \leftarrow \operatorname{diag}^{-1}(\boldsymbol{k}) \left[\hat{\boldsymbol{\xi}}^{(i)} + \boldsymbol{x}^{(i)} \right] - \boldsymbol{x}^{(i)}$ for each $i \in [n]$. {Correct $\hat{\boldsymbol{\Xi}}$ to the function estimation of $(\boldsymbol{A}^* \boldsymbol{x})^+$.}
9: return $\hat{\boldsymbol{B}}, \hat{\boldsymbol{\Xi}}$.

Then we prove $\neg a$) $\implies \neg b$), $\neg c$): $\neg a$) $\implies \exists j' \in [d] \text{ and } j' \neq j \text{ s.t. } A^*_{j,j'} \neq 0.$ Recall

$$\frac{\left[\boldsymbol{C}\mathbf{y}\right]_{j}}{\mathrm{x}_{j}} = \frac{k_{j}\left[\left(\boldsymbol{A}_{j,:}^{*}\mathbf{x}\right)^{+} + \mathrm{x}_{j}\right]}{\mathrm{x}_{j}}$$

The value of $\mathbf{x}_{j'}$ affects the value of $[\mathbf{C}\mathbf{y}]_j/\mathbf{x}_j$ because $A_{j,j'}^* \neq 0$. In fact, from $\operatorname{supp} p(\mathbf{x}) = \mathbb{R}^d$ we have $\operatorname{supp} p(\mathbf{x}_{j'} | \mathbf{x}_j) = \mathbb{R}$, indicating that the value of $[\mathbf{C}\mathbf{y}]_j/\mathbf{x}_j$ can never be kept as a constant when given both $\mathbf{x}_j < 0$ and $\mathbf{x}_j > 0$ because $\mathbf{x}_{j'}$ can be any real number.

With the exact derivations above, we are able to obtain a left inverse of B^* , namely C, that satisfies $CB^* = I_d$. Consider Eq. (1) left multiplied by B^*C

$$B^*C\mathbf{y} = \mathbf{y}.\tag{9}$$

Eq. (9) is also a noiseless unbiased linear model where Cy and y are observable, and thereby B^* is computable due to the easy solvability of its linearity.

Algorithm 5 describes how a residual unit second layer is learned without satisfying Condition 2.1: We first solve the QP/LP and obtain a scaling equivalence to an estimated left inverse of B^* and the function estimate, namely $\hat{\mathbf{C}}$ and $\hat{\mathbf{\Xi}}$. Then we compute the scale factor estimate $\hat{\mathbf{k}}$ by running Algorithm 4, where for each component index $j \in [d]$, we first use a tolerance parameter as a threshold to determine whether ground-truth layer 1 is a scale transformation, and if so, we run LR to estimate the model $[C\mathbf{y}]_j = k_j \mathbf{x}_j$, otherwise the scale factor is directly assigned by 1. Upon correcting $\hat{\mathbf{C}}$ and $\hat{\mathbf{\Xi}}$ by $\hat{\mathbf{k}}$, we obtain a layer 2 estimate $\hat{\mathbf{B}}$ by running LR to estimate the linear model Eq. (9). The strong consistency of our results in this appendix is justified in Appendix H.

F Exact Derivation of Objective Functional Minimizers

We give the exact derivation of the objective functional minimizers in the next two subsections. We begin with layer 2 and continue to layer 1.

F.1 Layer 2

Lemma 4.2 is proved as follows.

Proof. " \Leftarrow ": Since $h(\mathbf{x}) = C\mathbf{y} - \mathbf{x} \ge 0$, random vector $h(\mathbf{x}) + \mathbf{x} - C\mathbf{y}$ is always a zero vector, which implies $G_2(\mathbf{C}, h) = 0$. Hence " \Leftarrow " holds.

" \implies ": Since r.v. $\|h(\mathbf{x}) + \mathbf{x} - C\mathbf{y}\|^2 \ge 0$ we have

$$G_2(\boldsymbol{C}, h) = 0 \implies \lambda\left(\left\{\boldsymbol{x} \in \mathbb{R}^d \mid \|\boldsymbol{h}(\boldsymbol{x}) + \boldsymbol{x} - \boldsymbol{C}\boldsymbol{y}\|^2 > 0\right\}\right) = 0, \tag{10}$$

where λ is the Lebesgue measure on \mathbb{R}^d .

Proof by contradiction: Assume that $\exists x' \in \mathbb{R}^d$, $\exists i \in [d]$, s.t. $(Cy' - x')_i < 0$, where y' is the corresponding network output. Let $f(x) = (Cy - x)_i$ which is continuous on \mathbb{R}^d . Therefore for $\epsilon = -f(x') > 0$, $\exists \delta > 0$, $\forall x \in B(x; \delta)$ i.e. $||x - x'|| < \delta$,

$$\begin{aligned} |f(\boldsymbol{x}) - f(\boldsymbol{x}')| &< \epsilon \implies 2f(\boldsymbol{x}') < f(\boldsymbol{x}) < 0 \implies [\boldsymbol{x} - \boldsymbol{C}\boldsymbol{y}]_i > 0 \\ \implies [h(\boldsymbol{x}) + \boldsymbol{x} - \boldsymbol{C}\boldsymbol{y}]_i > 0 \implies \|h(\boldsymbol{x}) + \boldsymbol{x} - \boldsymbol{C}\boldsymbol{y}\|^2 > 0. \end{aligned}$$

Since $\lambda(B(\boldsymbol{x};\delta)) = \frac{\pi^{d/2}}{\Gamma(d/2+1)}\delta^d > 0$ and $B(\boldsymbol{x};\delta)$ is a subset of the measured set in Eq. (10), we have a contradiction. Thus $C\boldsymbol{y} - \boldsymbol{x} \ge 0$ holds in a pointwise manner in \mathbb{R}^d , indicating $\{\boldsymbol{x} \in \mathbb{R}^d \mid h(\boldsymbol{x}) \neq C\boldsymbol{y} - \boldsymbol{x}\}$ must be a null set. With *h*'s continuity, *h* must be $\boldsymbol{x} \mapsto C\boldsymbol{y} - \boldsymbol{x}$. Therefore " \Longrightarrow " holds.

Lemma F.1. $Cy - x \ge 0$ holds for any $x \in \mathbb{R}^d$ and its corresponding output y only if CB^* is a diagonal matrix.

Proof. Let $D = CB^*$. We rewrite $Cy - x \ge 0$ as

$$D\left[\left(\boldsymbol{A}^{*}\boldsymbol{x}\right)^{+}+\boldsymbol{x}\right]-\boldsymbol{x}\geq0.$$
(11)

For further use, we substitute x by -x and the resulting inequality $D\left[\left(-A^*x\right)^+ - x\right] + x \ge 0$ still holds. Added by Eq. (11) we have

$$\boldsymbol{D}\left(\left|\boldsymbol{A}_{k,:}^{*}\cdot\boldsymbol{x}\right|\right)_{d\times 1}\geq 0.$$
(12)

Proof by contradiction: Assume **D** is not diagonal, then $\exists i \neq j$, such that $D_{i,j} \neq 0$. Consider the following two cases:

a) $D_{i,j} > 0$. We take the *i*-th row of Eq. (11) as follows

$$\sum_{k=1}^{d} D_{i,k} \left[\left(\boldsymbol{A}_{k,:}^{*} \cdot \boldsymbol{x} \right)^{+} + x_{k} \right] \geq x_{i}.$$

Let $x_{-j} = 0$ and $x_j < 0$. Then $\sum_{k=1}^{d} D_{i,k} \left[\left(\boldsymbol{A}_{k,:}^* \cdot \boldsymbol{x} \right)^+ + x_k \right] = D_{i,j} \cdot x_j < 0 \implies \bot$. b) $D_{i,j} < 0$. We take the *i*-th row of Eq. (12) as follows

$$\sum_{k=1}^{d} D_{i,k} \left| \boldsymbol{A}_{k,:}^{*} \cdot \boldsymbol{x} \right| \ge 0.$$

Let $\boldsymbol{x} = (\boldsymbol{A}^*)^{-1} \boldsymbol{v}$, where $\boldsymbol{v} = \boldsymbol{e}^{(j)}$. Then $\sum_{k=1}^d D_{i,k} \left| \boldsymbol{A}_{k,:}^* \cdot \boldsymbol{x} \right| = D_{i,j} < 0 \implies \bot$.

With Lemma 4.2 and Lemma F.1, we prove Theorem E.1 as follows.

Proof. With Lemma 4.2, we only need to prove $\forall x \in \mathbb{R}^d, Cy - x \ge 0 \iff CB^* = \operatorname{diag}(k)$. " \Leftarrow ": We have

$$Cy - x = CB^* \left[\left(A^* x \right)^+ + x \right] - x = \operatorname{diag}(k) \left[\left(A^* x \right)^+ + x \right] - x.$$
(13)

For the *i*-th row of Eq. (13), consider the following two cases:

a) $A_{i,i}^* \cdot \boldsymbol{x} = A_{i,i}x_i$, i.e. A^* is a scale transformation w.r.t. the *i*-th row. With $\frac{1}{1+A_{i,i}^*} \leq k_i \leq 1$ we have

$$[Cy - x]_i = k_i \left[\left(A_{i,i}^* x_i \right)^+ + x_i \right] - x_i = k_i \left(A_{i,i}^* x_i \right)^+ + (k_i - 1) x_i.$$

For $x_i \ge 0$, $[Cy - x]_i = (k_i A_{i,i}^* + k_i - 1) x_i \ge 0$. For $x_i < 0$, $[Cy - x]_i = (k_i - 1) x_i \ge 0$. b) $A_{i,:}^* \cdot x \ne A_{i,i} x_i$. With $k_i = 1$ we have

$$[\boldsymbol{C}\boldsymbol{y}-\boldsymbol{x}]_{i}=k_{i}\left[\left(\boldsymbol{A}_{i,:}^{*}\boldsymbol{x}\right)^{+}+x_{i}\right]-x_{i}=\left(\boldsymbol{A}_{i,:}^{*}\boldsymbol{x}\right)^{+}\geq0.$$

Hence " \Leftarrow " holds.

" \implies ": Let $D = CB^*$. With Lemma F.1, D is diagonal. Consider the following two cases:

a) $\mathbf{A}_{i,:}^* \cdot \mathbf{x} = A_{i,i} x_i$. The *i*-th inequality can be written as

$$[Cy - x]_{i} = D_{i,i} \left[\left(A_{i,i}^{*} x_{i} \right)^{+} + x_{i} \right] - x_{i}$$

= $D_{i,i} \left(A_{i,i}^{*} x_{i} \right)^{+} + (D_{i,i} - 1) x_{i} \ge 0.$ (14)

Proof by contradiction: we need to find $\boldsymbol{x} \in \mathbb{R}^d$ which contradicts with Eq. (14) in the following three cases:

- a) If $D_{i,i} \leq 0$, then, $x_i > 0 \implies D_{i,i} \left(A_{i,i}^* x_i\right)^+ + \left(D_{i,i} 1\right) x_i < 0 \implies \bot$. b) If $D_{i,i} > 1$, then, $x_i < 0 \implies \bot$.

c) If
$$0 < D_{i,i} < \frac{1}{1+A_{i,i}^*}$$
, then $\exists a > 0$, s.t. $D_{i,i} = \frac{1}{1+A_{i,i}^*+a}$. Letting $x_i > 0$, we have

$$D_{i,i} \left(A_{i,i}^* x_i \right)^+ + \left(D_{i,i} - 1 \right) x_i = \frac{\left(A_{i,i}^* x_i \right)^+}{1 + A_{i,i}^* + a} - \frac{\left(A_{i,i}^* + a \right) x_i}{1 + A_{i,i}^* + a} < 0 \implies \bot.$$

Hence $\frac{1}{1+A_{i,i}^*} \le D_{i,i} \le 1$.

b) $A_{i,i}^* \cdot x \neq A_{i,i}^* x_i$, i.e. $\exists j \neq i$, s.t. $A_{i,j}^* > 0$. The *i*-th inequality can be written as

$$[Cy - x]_{i} = D_{i,i} \left[(A_{i,:}^{*}x)^{+} + x_{i} \right] - x_{i}$$

= $D_{i,i} (A_{i,:}^{*}x)^{+} + (D_{i,i} - 1) x_{i} \ge 0.$ (15)

Proof by contradiction: we need to find $x \in \mathbb{R}^d$ which contradicts with Eq. (15) in the following three cases:

a) If $D_{i,i} \leq 0$, then, $x_i > 0 \Longrightarrow D_{i,i} \left(\boldsymbol{A}_{i,:}^* \boldsymbol{x} \right)^+ + (D_{i,i} - 1) x_i < 0 \Longrightarrow \bot$. b) If $D_{i,i} > 1$, then, $x_i < 0 \land x_{-i} \leq 0 \Longrightarrow \bot$. c) If $0 < D_{i,i} < 1$, then, $x_i > 0 \land x_j \leq -\frac{A_{i,i}^*}{A_{i,j}^*} x_i \land x_k \leq 0 \implies \bot$, where $k \neq i, j$. Hence $D_{i,i} = 1$.

Hence $D = CB^* = \text{diag}(k)$, and thereby " \implies " holds.

F.2 Layer 1

The proof of Lemma 5.3 is similar to that of Lemma 4.2 because the two lemmas follow the same idea, which is to link objective functional minimization with always-hold inequalities. With Lemma 5.3, we prove Theorem 5.1 as follows.

Proof. With Lemma 5.3, we only need to prove $\forall \boldsymbol{x} \in \mathbb{R}^d$, $(\boldsymbol{A}^*\boldsymbol{x})^+ - \boldsymbol{A}\boldsymbol{x} \ge 0 \iff \forall i \in [d], \boldsymbol{A}_{i,:} = k_i \boldsymbol{A}_{i,:}^*$, where $0 \le k_i \le 1$. This is equivalent to what it is for a single row, i.e. $\forall \boldsymbol{x} \in \mathbb{R}^d$, $(\boldsymbol{a}^*^\top \boldsymbol{x})^+ - \boldsymbol{a}^\top \boldsymbol{x} \ge 0 \iff \boldsymbol{a} = k\boldsymbol{a}^*$, where $0 \le k \le 1$.

" \Leftarrow ": The case where $a^* = 0$ is clear. If a^* is not a zero vector and $a = ka^*$, we have

a) If $\boldsymbol{a}^{*\top}\boldsymbol{x} \ge 0$, $(\boldsymbol{a}^{*\top}\boldsymbol{x})^{+} - \boldsymbol{a}^{\top}\boldsymbol{x} = \boldsymbol{a}^{*\top}\boldsymbol{x} - k\boldsymbol{a}^{*\top}\boldsymbol{x} = (1-k)\boldsymbol{a}^{*\top}\boldsymbol{x} \ge 0$. b) If $\boldsymbol{a}^{*\top}\boldsymbol{x} < 0$, $(\boldsymbol{a}^{*\top}\boldsymbol{x})^{+} - \boldsymbol{a}^{\top}\boldsymbol{x} = -k\boldsymbol{a}^{*\top}\boldsymbol{x} \ge 0$.

Hence " \Leftarrow " holds.

" \implies ": If $a^* = 0$, a must be a zero vector, otherwise let x = a, then $-a^{\top}x < 0 \implies \bot$. If a^* is not a zero vector, consider two cases below:

a) If $\boldsymbol{a} \neq k\boldsymbol{a}^*$ where $k \geq 0$. Let $\boldsymbol{x} = \frac{\|\boldsymbol{a}^*\|}{\|\boldsymbol{a}\|} \cdot \boldsymbol{a} - \boldsymbol{a}^*$, then \boldsymbol{x} is not a zero vector, and

$$\begin{array}{c} \boldsymbol{a}^{*^{\top}}\boldsymbol{x} = \|\boldsymbol{a}^{*}\|^{2}\left(\cos\theta - 1\right) < 0, \\ \boldsymbol{a}^{\top}\boldsymbol{x} = \|\boldsymbol{a}\| \|\boldsymbol{a}^{*}\|\left(1 - \cos\theta\right) > 0 \end{array} \right\} \implies \left(\boldsymbol{a}^{*^{\top}}\boldsymbol{x}\right)^{+} - \boldsymbol{a}^{\top}\boldsymbol{x} < 0 \implies \bot$$

where θ denotes the angle between a and a^* .

b) If $\boldsymbol{a} = k\boldsymbol{a}^*$ where k > 1, then let $\boldsymbol{x} = \boldsymbol{a}^*$ we have

$$\left(\boldsymbol{a^{*}}^{\mathsf{T}}\boldsymbol{x}\right)^{\mathsf{+}} - \boldsymbol{a}^{\mathsf{T}}\boldsymbol{x} = (1-k) \left\|\boldsymbol{a}^{*}\right\|^{2} < 0 \implies \bot$$

Hence $\boldsymbol{a} = k\boldsymbol{a}^*$ where $0 \le k \le 1$ if \boldsymbol{a}^* is not zero, and thereby " \Longrightarrow " holds.

G QP Convexity

The LPs in the main paper are trivially convex. So in this appendix, we only justify the convexity of our QPs: We first prove the convexity of single-sample objectives, then the convexity of the empirical objectives with nonparametric estimation, i.e. \hat{G}_1^{NPE} and \hat{G}_2^{NPE} is obtained by the convexity of convex function summations.

Lemma G.1. Suppose $f(u) = \frac{1}{2} \|Tu - b\|^2$ where u is a real matrix. Then f is convex w.r.t. u.

Lemma G.1 is easily obtained since the Hessian $f''(\mathbf{T}) = \mathbf{T}^{\top}\mathbf{T}$ is positive semidefinite. In the following, we demonstrate and justify the convexity of the QPs of both layers by rewriting their single-sample objectives into the formulation of f and summing them without loss of convexity.

Theorem G.2. *QP: Eq.* (4), and *QP: Eq.* (7) are convex optimization problems.

Proof. First of all, constraints of both QPs are trivially linear and convex. Thus, we only need to justify the convexity of the two empirical objectives, \hat{G}_1^{NPE} and \hat{G}_2^{NPE} (see Eqs. (4) and (7)). Consider the single-sample version of \hat{G}_1^{NPE} , namely $g_1^{\text{NPE}}(\boldsymbol{A}, \boldsymbol{\phi}; \boldsymbol{x}, \boldsymbol{h}) = \frac{1}{2} \|\boldsymbol{\phi} + \boldsymbol{A}\boldsymbol{x} - \boldsymbol{h}\|^2$, which, in the formulation of f, can be rewritten with

$$egin{aligned} m{T} = egin{bmatrix} m{x}^{ op} & & 1 & & \ & m{x}^{ op} & & 1 & & \ & m{x}^{ op} & & 1 & & \ & & \ddots & & \ddots & \ & & m{x}^{ op} & & & \ddots & & \ & & m{x}^{ op} & & & 1 \end{bmatrix}, \ m{u} = egin{bmatrix} m{A}_{1,:}^{ op} \ m{A}_{2,:}^{ op} \ m{B}_{2,:}^{ op} \ m{A}_{d,:}^{ op} \ m{A}_{d,:}^$$

which guarantees the convexity of g_1^{NPE} w.r.t. \boldsymbol{A} and $\boldsymbol{\phi}$ by Lemma G.1. For $g_2^{\text{NPE}}(\boldsymbol{C},\boldsymbol{\xi};\boldsymbol{x},\boldsymbol{y}) = \frac{1}{2} \|\boldsymbol{\xi} + \boldsymbol{x} - \boldsymbol{C}\boldsymbol{y}\|^2$, we have

$$egin{aligned} egin{aligned} egin{aligne} egin{aligned} egin{aligned} egin{aligned} egin$$

which guarantees the convexity of g_2^{NPE} w.r.t. C and $\boldsymbol{\xi}$ by Lemma G.1. Now we consider the summation. Taking layer 1 as an example, by definition we have

$$\hat{G}_{1}^{\text{NPE}}(\boldsymbol{A}, \boldsymbol{\Phi}) = \sum_{i \in [n]} g_{1}^{\text{NPE}}(\boldsymbol{A}, \boldsymbol{\phi}^{(i)}; \boldsymbol{x}^{(i)}, \boldsymbol{h}^{(i)}).$$

For each $i \in [n]$, equivalently, we take Φ as a variable instead of $\phi^{(i)}$ in g_1^{NPE} , but with only $\phi^{(i)} \in \Phi$ determining the value of g_1^{NPE} . In this sense, g_1^{NPE} is convex w.r.t. A and Φ for each $i \in [n]$. Thus, the sum \hat{G}_1^{NPE} is convex w.r.t. A and Φ . Similarly, \hat{G}_2^{NPE} is convex w.r.t. C and Ξ .

H Strong Consistency

In this appendix, we justify the strong consistency of our estimators for the residual unit layer 1/2 learning and the whole network.

According to Theorem E.1 and Theorem 5.1, the solutions to our objective functionals are continuous sets. In addition, we also present intermediate results in terms of continuous sets, e.g. possible left-inverse matrices for B^* in the results for layer 2. Thus, to analyze the consistency of our learning algorithm, we define distances between sets, so that the convergence of sets can be well defined. Further point convergence results, i.e. layer 1/2 estimator strong consistency, are based on the set convergence we define.

Definition H.1 (deviation). Let $\mathbb{A} \subseteq \mathbb{M}$ and $\mathbb{B} \subseteq \mathbb{M}$ be two non-empty sets from a metric space (\mathbb{M}, d) . The deviation of set \mathbb{A} from the set \mathbb{B} , denoted by $D(\mathbb{A}, \mathbb{B})$, is

$$D(\mathbb{A},\mathbb{B}) = \sup_{\boldsymbol{a}\in\mathbb{A}} d(\boldsymbol{a},\mathbb{B}) = \sup_{\boldsymbol{a}\in\mathbb{A}} \inf_{\boldsymbol{b}\in\mathbb{B}} d(\boldsymbol{a},\boldsymbol{b}),$$

where sup and inf represent supremum and infimum, respectively.

Definition H.2 (Hausdorff distance). Let $\mathbb{A} \subseteq \mathbb{M}$ and $\mathbb{B} \subseteq \mathbb{M}$ be two non-empty sets from a metric space (\mathbb{M}, d) . The Hausdorff distance between \mathbb{A} and \mathbb{B} , denoted by $D_{\mathrm{H}}(\mathbb{A}, \mathbb{B})$, is

$$D_{\mathrm{H}}(\mathbb{A}, \mathbb{B}) = \max\{D(\mathbb{A}, \mathbb{B}), D(\mathbb{B}, \mathbb{A})\}.$$

Remark. In the following, we use the Frobenius norm to define the distance between matrices, i.e. $d(\mathbf{X}, \mathbf{Y}) = \|\mathbf{X} - \mathbf{Y}\|_{\mathrm{F}}$.

H.1 Layer 2

In this subsection, we prove the strong consistency of the layer 2 estimator.

H.1.1 Objective Minimizer Space Estimator

For simplicity of notation, we use $\hat{\mathbb{S}}_n$ to denote our layer 2 QP/LP solution space by n random samples¹³ as a random set

$$\hat{\mathbb{S}}_n \coloneqq \{ \boldsymbol{C} \in \mathbb{R}^{d \times m} : \boldsymbol{C} \mathbf{y}^{(i)} - \mathbf{x}^{(i)} \ge 0, \, \forall i \in [n] \},$$
(16)

 $^{^{13}\}mathrm{Here},$ we take samples as random variables for empirical analysis.

and \mathbb{S}^* to denote the value space of C in Lemma 4.2 which minimizes layer 2 objective functional

$$\mathbb{S}^* := \{ \boldsymbol{C} \in \mathbb{R}^{d \times m} : \boldsymbol{C}\boldsymbol{y} - \boldsymbol{x} \ge 0, \, \forall \, \boldsymbol{x} \in \mathbb{R}^d \text{ and its corresponding output } \boldsymbol{y} \}.$$

For further use, we name $\hat{\mathbb{P}}_n$ as the set of n sampled inputs which define $\hat{\mathbb{S}}_n$, i.e. $\hat{\mathbb{P}}_n := {\mathbf{x}^{(i)}}_{i \in [n]}$ where each $\mathbf{x}^{(i)}$ is the same random variable in Eq. (16). By the definitions above, we describe the strong consistency of our QP/LP as Lemma H.1.

Lemma H.1 (QP/LP strong consistency, layer 2). $D_{\mathrm{H}}(\hat{\mathbb{S}}_n, \mathbb{S}^*) \xrightarrow{a.s.} 0 \text{ as } n \to \infty.$

Proof. First we prove that $D_{\mathrm{H}}(\hat{\mathbb{S}}_n, \mathbb{S}^*) \xrightarrow{\mathrm{p}} 0$ as $n \to \infty$. Recall Theorem E.1, $C\mathbf{y}-\mathbf{x} \ge 0$ only if $CB^* = \mathrm{diag}(\mathbf{k})$. We inherit the notation as defining $\mathbf{D} = CB^*$. Theorem E.1 is based on Lemma F.1, and we prove them by raising points that show contradiction, i.e. violate the inequality that holds in a pointwise fashion:

- a) In the proof of Lemma F.1, we use d points: $-e^{(i)}$, for $i \in [d]$, to make $\sum_{k=1}^{d} D_{i,k} \left[\left(\mathbf{A}_{k,:}^{*} \cdot \mathbf{x} \right)^{+} + x_k \right] < x_i$, so that $D_{i,j}$ $(i \neq j)$ cannot be positive; and another d points: $(\mathbf{A}^{*})^{-1} e^{(i)}$, to show $\sum_{k=1}^{d} D_{i,k} \left| \mathbf{A}_{k,:}^{*} \cdot \mathbf{x} \right| < 0$, so that $D_{i,j}$ $(i \neq j)$ cannot be negative. For each $-e^{(i)}$ we use here, since the violations follow strict inequalities, we know there exists a neighborhood of $-e^{(i)}$, $\mathbb{N}_i = \mathbb{N}(-e^{(i)})$, such that $\forall \mathbf{z} \in \mathbb{N}_i$, $\sum_{k=1}^{d} D_{i,k} \left[\left(\mathbf{A}_{k,:}^{*} \cdot \mathbf{z} \right)^{+} + z_k \right] < z_i$. We can similarly find such neighborhood of each $(\mathbf{A}^{*})^{-1} e^{(i)}$ that the strict inequality holds within the neighborhood respectively. We index them as \mathbb{N}_{d+1} to \mathbb{N}_{2d} .
- b) In the proof of Theorem E.1, we further construct d points: for each $i \in [d]$, we take a point x such that $x_i > 0 \land x_j \leq -\frac{A_{i,i}^*}{A_{i,j}^*} x_i \land x_k \leq 0$, where $k \neq i, j$. This counterexample shows $[Cy x]_i < 0$, and eliminates the possibility of $0 < D_{i,i} < 1$ when $A_{i,:}^*$ is not a scale transformation. We can similarly find neighborhood of each point and index them as \mathbb{N}_{2d+1} to \mathbb{N}_{3d} . Note that we omit some cases in the proof of Theorem E.1, because the first 2d points are sufficient to use in those cases to show contradiction.

In the sampling procedure, if we sample at least one point in each neighborhood \mathbb{N}_i , Theorem E.1 assures the solution we get $\hat{\mathbf{C}}_n$ would lie in the true optimal set \mathbb{S}^* . The probability that the sampling procedure "omits" any of the neighborhoods is

$$P\left(\hat{\mathbb{P}}_{n}\bigcap\mathbb{N}_{1}=\emptyset \text{ or } \hat{\mathbb{P}}_{n}\bigcap\mathbb{N}_{2}=\emptyset \text{ or } \dots \text{ or } \hat{\mathbb{P}}_{n}\bigcap\mathbb{N}_{3d}=\emptyset\right) \leq \sum_{i=1}^{3d} P\left(\hat{\mathbb{P}}_{n}\bigcap\mathbb{N}_{i}=\emptyset\right)$$

$$\leq 3d[1-\min_{i\in[3d]}P(\mathbb{N}_{i})]^{n}.$$
(17)

Since the measure on each neighborhood $P(\mathbb{N}_i) = \int_{\boldsymbol{x} \in \mathbb{N}_i} p(\boldsymbol{x}) > 0$,

$$P\left(\hat{\mathbb{P}}_n \bigcap \mathbb{N}_1 = \emptyset \text{ or } \hat{\mathbb{P}}_n \bigcap \mathbb{N}_2 = \emptyset \text{ or } \dots \text{ or } \hat{\mathbb{P}}_n \bigcap \mathbb{N}_{3d} = \emptyset\right) \to 0, \text{ as } n \to \infty.$$

Here we obtain $D_{\mathrm{H}}(\hat{\mathbb{S}}_n, \mathbb{S}^*) \xrightarrow{\mathrm{P}} 0$ as $n \to \infty$. Now we take the infinite sum over the both sides of Eq. (17)

$$\sum_{n \in [\infty]} P\left(\hat{\mathbb{P}}_n \bigcap \mathbb{N}_1 = \emptyset \text{ or } \hat{\mathbb{P}}_n \bigcap \mathbb{N}_2 = \emptyset \text{ or } \dots \text{ or } \hat{\mathbb{P}}_n \bigcap \mathbb{N}_{3d} = \emptyset\right)$$

$$\leq 3d \sum_{n \in [\infty]} \left[1 - \min_{i \in [3d]} P(\mathbb{N}_i)\right]^n = 3d \left[\frac{1}{\min_{i \in [3d]} P(\mathbb{N}_i)} - 1\right] < +\infty.$$

By the Borel-Cantelli lemma (Borel, 1909), $D_{\mathrm{H}}(\hat{\mathbb{S}}_n, \mathbb{S}^*) \xrightarrow{\mathrm{a.s.}} 0$ as $n \to \infty$.

H.1.2 Scale Factor Estimator

To avoid ambiguity, we use $n_{\rm sf}$ to denote the number of samples used in Algorithm 4. The samples pairs are $\{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^{n_{\rm sf}}$. Without loss of generality, the following discussion focuses on some fixed index $j \in [d]$. In Algorithm 4, we plug in our estimator $\hat{\mathbf{C}}_n$ and use LR to estimate k_j given that $\mathbf{x}_j^{(i)} < 0$

$$\mathbf{k}_{n_{\rm sf}}(\hat{\mathbf{C}}_n) = \operatorname*{arg\,min}_k \quad \frac{1}{2n_{\rm sf}} \sum_{i \in [n_{\rm sf}]} \left\| [\hat{\mathbf{C}}_n \mathbf{y}^{(i)}]_j - k \mathbf{x}_j^{(i)} \right\|^2 = \frac{\sum_{i \in [n_{\rm sf}]} \mathbf{x}_j^{(i)} [\hat{\mathbf{C}}_n \mathbf{y}^{(i)}]_j}{\sum_{i \in [n_{\rm sf}]} \left(\mathbf{x}_j^{(i)} \right)^2}.$$
 (18)

We first give the strong consistency of layer 2 scale factor estimator for as $n_{\rm sf} \to \infty$, as described in Lemma H.2. Lemma H.2 (scale factor estimator strong consistency, layer 2). Suppose A^* is a scale transformation w.r.t. the *j*-th component, and $k_{n_{sf}}(\hat{\mathbf{C}}_n)$ is the n_{sf} -sample estimator of k_j via LR: Eq. (18) given that $\mathbf{x}_j^{(i)} < 0$. Define sets

$$\mathbb{U}_{n_{sf},n}\coloneqq \{ \mathrm{k}_{n_{sf}}(oldsymbol{C}):oldsymbol{C}\in \hat{\mathbb{S}}_n \}, \ and \ \mathbb{U}^*\coloneqq \left[rac{1}{1+A_{j,j}^*}, 1
ight].$$

Then $\lim_{n_{sf}\to\infty}\lim_{n\to\infty}D_{\mathrm{H}}(\mathbb{U}_{n_{sf},n},\mathbb{U}^*)\stackrel{a.s.}{=} 0.$

Proof. Following our notation, Theorem E.1 and Theorem E.2 ensure that if \mathbf{A}^* is a scale transformation w.r.t. the *j*-th component, for any C belonging to the true optimal set \mathbb{S}^* , $\mathbf{k}_{n_{\mathrm{sf}}}(C) \in \left[\frac{1}{1+A_{j,j}^*}, 1\right]$. And the "iff" statement strengthens that $\mathbb{U}^* = \{\mathbf{k}_{n_{\mathrm{sf}}}(C) : C \in \mathbb{S}^*\}$ for any $n_{\mathrm{sf}} \in \mathbb{Z}^+$. Note that since $\mathbb{S}^* \subset \hat{\mathbb{S}}_n$, we have $\mathbb{U}^* \subset \mathbb{U}_{n_{\mathrm{sf}},n}$. We only need to prove $D(\mathbb{U}_{n_{\mathrm{sf}},n}, \mathbb{U}^*) \xrightarrow{\mathrm{a.s.}} 0$ as $n \to \infty$.

 $\forall \hat{\mathbf{C}}_n \in \hat{\mathbb{S}}_n \text{ and } \forall \mathbf{C} \in \mathbb{S}^*,$

$$\begin{split} \left| \mathbf{k}_{n_{sf}}(\hat{\mathbf{C}}_{n}) - \mathbf{k}_{n_{sf}}(C) \right| \\ &= \frac{1}{\sum_{i \in [n_{sf}]} \left(\mathbf{x}_{j}^{(i)} \right)^{2}} \left| \sum_{i \in [n_{sf}]} \mathbf{x}_{j}^{(i)} \left[(e^{(j)})^{\top} \left(\hat{\mathbf{C}}_{n} - C \right) B^{*} \left[\left(A^{*} \mathbf{x}^{(i)} \right)^{+} + \mathbf{x}^{(i)} \right] \right] \right| \\ &\leq \frac{1}{\sum_{i \in [n_{sf}]} \left(\mathbf{x}_{j}^{(i)} \right)^{2}} \left[\sum_{i \in [n_{sf}]} \left| \mathbf{x}_{j}^{(i)} \right| \left\| (e^{(j)})^{\top} \left(\hat{\mathbf{C}}_{n} - C \right) B^{*} \right\|_{2} \left(\left\| A^{*} \mathbf{x}^{(i)} \right\|_{2} + \left\| \mathbf{x}^{(i)} \right\|_{2} \right) \right] \\ &\leq \frac{1}{\sum_{i \in [n_{sf}]} \left(\mathbf{x}_{j}^{(i)} \right)^{2}} \left[\sum_{i \in [n_{sf}]} \left| \mathbf{x}_{j}^{(i)} \right| \left\| \hat{\mathbf{C}}_{n} - C \right\|_{F} \left\| B^{*} \right\|_{F} \left(\left\| A^{*} \right\|_{F} + 1 \right) \left\| \mathbf{x}^{(i)} \right\|_{2} \right] \\ &= \frac{\sum_{i \in [n_{sf}]} \left[\left| \mathbf{x}_{j}^{(i)} \right| \left\| B^{*} \right\|_{F} \left(\left\| A^{*} \right\|_{F} + 1 \right) \left\| \mathbf{x}^{(i)} \right\|_{2} \right]}{\sum_{i \in [n_{sf}]} \left(\mathbf{x}_{j}^{(i)} \right)^{2}} \left\| \hat{\mathbf{C}}_{n} - C \right\|_{F}. \end{split}$$

Then,

$$\sup_{\hat{\mathbf{C}}_{n}\in\hat{\mathbb{S}}_{n}} \inf_{\boldsymbol{C}\in\mathbb{S}^{*}} \left| \mathbf{k}_{n_{\mathrm{sf}}}(\hat{\mathbf{C}}_{n}) - \mathbf{k}_{n_{\mathrm{sf}}}(\boldsymbol{C}) \right|$$

$$\leq \frac{\sum_{i\in[n_{\mathrm{sf}}]} \left[\left| \mathbf{x}_{j}^{(i)} \right| \left\| \boldsymbol{B}^{*} \right\|_{\mathrm{F}} \left(\left\| \boldsymbol{A}^{*} \right\|_{\mathrm{F}} + 1 \right) \left\| \mathbf{x}^{(i)} \right\|_{2} \right]}{\sum_{i\in[n_{\mathrm{sf}}]} \left(\mathbf{x}_{j}^{(i)} \right)^{2}} \sup_{\hat{\mathbf{C}}_{n}\in\hat{\mathbb{S}}_{n}} \inf_{\boldsymbol{C}\in\mathbb{S}^{*}} \left\| \hat{\mathbf{C}}_{n} - \boldsymbol{C} \right\|_{\mathrm{F}}$$

which implies

$$D(\mathbb{U}_{n_{\rm sf},n},\mathbb{U}) \le \frac{\sum_{i \in [n_{\rm sf}]} \left[\left| \mathbf{x}_{j}^{(i)} \right| \| \boldsymbol{B}^{*} \|_{\rm F} \left(\| \boldsymbol{A}^{*} \|_{\rm F} + 1 \right) \left\| \mathbf{x}^{(i)} \right\|_{2} \right]}{\sum_{i \in [n_{\rm sf}]} \left(\mathbf{x}_{j}^{(i)} \right)^{2}} D(\hat{\mathbb{S}}_{n}, \mathbb{S}^{*}).$$
(19)

Take the $n_{\rm sf} \to \infty$ limit over both sides of Eq. (19). With the strong law of large numbers¹⁴ we have

$$\lim_{n_{\mathrm{sf}}\to\infty} D(\mathbb{U}_{n_{\mathrm{sf}},n},\mathbb{U}) \leq \frac{\mathbb{E}\left[\left|\mathbf{x}_{j}^{(i)}\right| \left\|\boldsymbol{B}^{*}\right\|_{\mathrm{F}}\left(\left\|\boldsymbol{A}^{*}\right\|_{\mathrm{F}}+1\right)\left\|\mathbf{x}^{(i)}\right\|_{2}\right]}{\mathbb{E}\left[\mathbf{x}_{j}^{(i)}\right]^{2}} D(\hat{\mathbb{S}}_{n},\mathbb{S}^{*}), \, \mathrm{w.p.} \, 1.$$

Since $D(\hat{\mathbb{S}}_n, \mathbb{S}^*) \xrightarrow{\text{a.s.}} 0$ as $n \to \infty$, we have $D(\mathbb{U}_{n_{\text{sf}}, n}, \mathbb{U}^*) \xrightarrow{\text{a.s.}} 0$ as $n \to \infty$ then $n_{\text{sf}} \to \infty$.

H.1.3 Layer 2 Weights Estimator

In Algorithm 5, we solve \boldsymbol{B} via LR. Let $\hat{\mathbf{z}}^{(i)} = \operatorname{diag}^{-1}(\hat{\mathbf{k}}) \cdot \hat{\mathbf{C}}_n \mathbf{y}^{(i)} \in \mathbb{R}^d$, where $\hat{\mathbf{k}}$ is obtained through Algorithm 4 with input $\hat{\mathbf{C}}_n$. Assume we are using sample size of n_w to do the LR. The optimization problem is:

$$\min_{\boldsymbol{B}} \sum_{i \in [n_{w}]} \left\| \mathbf{y}^{(i)} - \boldsymbol{B} \hat{\mathbf{z}}^{(i)} \right\|^{2}.$$
(20)

Now we present the strong consistency of layer 2 estimator, as described in Theorem H.3.

Theorem H.3 (strong consistency, layer 2). Suppose $\hat{\mathbf{B}}_{n_{sf}}$ is the solution to Eq. (20). Then $\hat{\mathbf{B}}_{n_w} \xrightarrow{a.s.} \mathbf{B}^*$ as $n, n_{sf}, n_w \to \infty$.

Proof. Let $\boldsymbol{\beta}$ denote vec \boldsymbol{B} (flattening \boldsymbol{B} into a vector), then $\boldsymbol{B}\hat{\mathbf{z}}^{(i)} = (\hat{\mathbf{z}}^{(i)})^{\top} \otimes \boldsymbol{I}_m$. Here the operation \otimes denotes the Kronecker product. Then we can define an equivalent optimization problem

$$\min_{\boldsymbol{\beta}} \frac{1}{2n_{w}} \sum_{i \in [n_{w}]} \left\| \mathbf{y}^{(i)} - \left[\left(\hat{\mathbf{z}}^{(i)} \right)^{\top} \otimes \boldsymbol{I}_{m} \right] \boldsymbol{\beta} \right\|^{2}.$$

Taking the derivatives of β , we obtain

$$-2\sum_{i\in[n_{w}]}\left[\left[\left(\hat{\mathbf{z}}^{(i)}\right)^{\top}\otimes\boldsymbol{I}_{m}\right]^{\top}\left(\mathbf{y}^{(i)}-\left[\left(\hat{\mathbf{z}}^{(i)}\right)^{\top}\otimes\boldsymbol{I}_{m}\right]\boldsymbol{\beta}\right)\right]=0.$$

Then the optimal solution $\hat{\beta}_{n_{w}}$ of this optimization can be written in closed form:

$$\begin{split} \hat{\boldsymbol{\beta}}_{n_{w}} &= \left[\sum_{i \in [n_{w}]} \left[\left(\hat{\mathbf{z}}^{(i)} \right)^{\top} \otimes \boldsymbol{I}_{m} \right]^{\top} \left[\left(\hat{\mathbf{z}}^{(i)} \right)^{\top} \otimes \boldsymbol{I}_{m} \right] \right]^{-1} \left(\sum_{i \in [n_{w}]} \left[\left(\hat{\mathbf{z}}^{(i)} \right)^{\top} \otimes \boldsymbol{I}_{m} \right]^{\top} \mathbf{y}^{(i)} \right) \\ &= \left[\sum_{i \in [n_{w}]} \left[\hat{\mathbf{z}}^{(i)} \otimes \boldsymbol{I}_{m} \right] \left[\left(\hat{\mathbf{z}}^{(i)} \right)^{\top} \otimes \boldsymbol{I}_{m} \right] \right]^{-1} \left(\sum_{i \in [n_{w}]} \left[\hat{\mathbf{z}}^{(i)} \otimes \boldsymbol{I}_{m} \right] \mathbf{y}^{(i)} \right) \\ &= \left[\sum_{i \in [n_{w}]} \left[\hat{\mathbf{z}}^{(i)} \left(\hat{\mathbf{z}}^{(i)} \right)^{\top} \right] \otimes \boldsymbol{I}_{m} \right]^{-1} \left(\left[\sum_{i \in [n_{w}]} \hat{\mathbf{z}}^{(i)} \otimes \boldsymbol{I}_{m} \right] \mathbf{y}^{(i)} \right) \\ &= \left[\left[\sum_{i \in [n_{w}]} \hat{\mathbf{z}}^{(i)} \left(\hat{\mathbf{z}}^{(i)} \right)^{\top} \right] \otimes \boldsymbol{I}_{m} \right]^{-1} \left(\left[\sum_{i \in [n_{w}]} \hat{\mathbf{z}}^{(i)} \otimes \boldsymbol{I}_{m} \right] \mathbf{y}^{(i)} \right). \end{split}$$

 14 Here we assume that the Kolmogorov's strong law assumption on moments (Sen & Singer, 1994) is met as is commonly done in similar analysis.

We inherit the notation from the last two subsections. By Lemma H.1, $d(\hat{\mathbf{C}}_n, \mathbb{S}^*) \xrightarrow{\text{a.s.}} 0$ as $n \to \infty$. Thus $\forall \varepsilon > 0, \exists N$ such that $\forall n \ge N, d(\hat{\mathbf{C}}_n, \mathbb{S}^*) \le \varepsilon$ w.p. 1, i.e. $\exists \mathbf{C}_n \in \mathbb{S}^*$ s.t. $d(\hat{\mathbf{C}}_n, \mathbf{C}_n) \le \varepsilon$ w.p. 1. Then by Lemma H.2, $\exists K > 0$, for ε that is small enough, $\exists N_{\text{sf}}$ such that $\forall n_{\text{sf}} \ge N_{\text{sf}}$, we have $\left| \mathbf{k}_{n_{\text{sf}}}(\hat{\mathbf{C}}_n) - \mathbf{k}_{n_{\text{sf}}}(\mathbf{C}_n) \right| \le K\varepsilon$ w.p.1. For simplicity, we omit the under-script n_{sf} of $\mathbf{k}_{n_{\text{sf}}}$ in the following discussion. In fact,

$$\begin{aligned} \left| \hat{\mathbf{z}}_{j}^{(i)} - \mathbf{z}_{j}^{(i)} \right| &= \left| \frac{1}{\mathbf{k}(\hat{\mathbf{C}}_{n})} \left(\mathbf{e}^{(i)} \right)^{\top} \hat{\mathbf{C}}_{n} \mathbf{y}^{(i)} - \frac{1}{\mathbf{k}(\mathbf{C}_{n})} \left(\mathbf{e}^{(i)} \right)^{\top} \mathbf{C}_{n} \mathbf{y}^{(i)} \right| \\ &= \left| \frac{1}{\mathbf{k}(\hat{\mathbf{C}}_{n})\mathbf{k}(\mathbf{C}_{n})} \left[\mathbf{k}(\mathbf{C}_{n}) \left(\mathbf{e}^{(i)} \right)^{\top} \hat{\mathbf{C}}_{n} - \mathbf{k}(\hat{\mathbf{C}}_{n}) \left(\mathbf{e}^{(i)} \right)^{\top} \mathbf{C}_{n} \right] \mathbf{y}^{(i)} \right| \\ &\leq \left| \frac{1}{\mathbf{k}(\hat{\mathbf{C}}_{n})\mathbf{k}(\mathbf{C}_{n})} \right| \left| \mathbf{k}(\mathbf{C}_{n}) \left(\mathbf{e}^{(i)} \right)^{\top} \hat{\mathbf{C}}_{n} \mathbf{y}^{(i)} - \mathbf{k}(\mathbf{C}_{n}) \left(\mathbf{e}^{(i)} \right)^{\top} \mathbf{C}_{n} \mathbf{y}^{(i)} \right| \\ &+ \left| \frac{1}{\mathbf{k}(\hat{\mathbf{C}}_{n})\mathbf{k}(\mathbf{C}_{n})} \right| \left| \mathbf{k}(\mathbf{C}_{n}) \left(\mathbf{e}^{(i)} \right)^{\top} \mathbf{C}_{n} \mathbf{y}^{(i)} - \mathbf{k}(\hat{\mathbf{C}}_{n}) \left(\mathbf{e}^{(i)} \right)^{\top} \mathbf{C}_{n} \mathbf{y}^{(i)} \right| \\ &\leq \frac{(1 + \mathbf{A}_{j,j}^{*})^{2}}{1 - K\varepsilon \left(1 + \mathbf{A}_{j,j}^{*} \right)} \left(\left\| \mathbf{y}^{(i)} \right\| \left\| \hat{\mathbf{C}}_{n} - \mathbf{C}_{n} \right\|_{\mathrm{F}} + \left| \mathbf{k}(\hat{\mathbf{C}}_{n}) - \mathbf{k}(\mathbf{C}_{n}) \right| \left\| \mathbf{C}_{n} \mathbf{B}^{*} \right\|_{\mathrm{F}} \left(\left\| \mathbf{A}^{*} \right\|_{\mathrm{F}} + 1 \right) \left\| \mathbf{x}^{(i)} \right\| \right) \\ &\leq \frac{(1 + \mathbf{A}_{j,j}^{*})^{2}}{1 - K\varepsilon \left(1 + \mathbf{A}_{j,j}^{*} \right)^{2}} \left(\left\| \mathbf{y}^{(i)} \right\| + dK \left(\left\| \mathbf{A}^{*} \right\|_{\mathrm{F}} + 1 \right) \left\| \mathbf{x}^{(i)} \right\| \right) \varepsilon. \end{aligned}$$

Then,

$$\begin{split} \left\| \hat{\mathbf{z}}^{(i)} \left(\hat{\mathbf{z}}^{(i)} \right)^{\top} - \mathbf{z}^{(i)} \left(\mathbf{z}^{(i)} \right)^{\top} \right\|_{\mathrm{F}} \\ &= \left\| \hat{\mathbf{z}}^{(i)} \left(\hat{\mathbf{z}}^{(i)} \right)^{\top} - \mathbf{z}^{(i)} \left(\hat{\mathbf{z}}^{(i)} \right)^{\top} + \mathbf{z}^{(i)} \left(\hat{\mathbf{z}}^{(i)} \right)^{\top} - \mathbf{z}^{(i)} \left(\mathbf{z}^{(i)} \right)^{\top} \right\|_{\mathrm{F}} \\ &\leq \left\| \left[\hat{\mathbf{z}}^{(i)} - \mathbf{z}^{(i)} \right] \left(\hat{\mathbf{z}}^{(i)} \right)^{\top} \right\|_{\mathrm{F}} + \left\| \mathbf{z}^{(i)} \left[\left(\hat{\mathbf{z}}^{(i)} \right)^{\top} - \left(\mathbf{z}^{(i)} \right)^{\top} \right] \right\|_{\mathrm{F}} \\ &\leq \left\| \left[\hat{\mathbf{z}}^{(i)} - \mathbf{z}^{(i)} \right] \left[\hat{\mathbf{z}}^{(i)} - \mathbf{z}^{(i)} \right]^{\top} \right\|_{\mathrm{F}} + 2 \left\| \mathbf{z}^{(i)} \left[\left(\hat{\mathbf{z}}^{(i)} \right)^{\top} - \left(\mathbf{z}^{(i)} \right)^{\top} \right] \right\|_{\mathrm{F}} \\ &\leq \sum_{j \in [d]} \left(\hat{\mathbf{z}}^{(i)}_{j} - \mathbf{z}^{(i)}_{j} \right)^{2} + 2 \left\| \mathbf{z}^{(i)} \right\| \left\| \hat{\mathbf{z}}^{(i)} - \mathbf{z}^{(i)} \right\|. \end{split}$$

It follows that $\left\| \hat{\mathbf{z}}^{(i)} \left(\hat{\mathbf{z}}^{(i)} \right)^{\top} - \mathbf{z}^{(i)} \left(\mathbf{z}^{(i)} \right)^{\top} \right\|_{\mathrm{F}}$ is also bounded by $\mathcal{O}(\varepsilon)$. With similar techniques we can prove

$$\left\|\sum_{i\in[n_{w}]}\hat{\mathbf{z}}^{(i)}\left(\hat{\mathbf{z}}^{(i)}\right)^{\top}-\sum_{i\in[n_{w}]}\mathbf{z}^{(i)}\left(\mathbf{z}^{(i)}\right)^{\top}\right\|_{F}\leq\mathcal{O}(\varepsilon),$$

and

$$\left\| \hat{\mathbf{z}}^{(i)} \otimes \boldsymbol{I}_m - \mathbf{z}^{(i)} \otimes \boldsymbol{I}_m \right\|_{\mathrm{F}} \leq \mathcal{O}(\varepsilon).$$

Denote $\left[\sum_{i \in [n_w]} \mathbf{z}^{(i)} \left(\mathbf{z}^{(i)}\right)^{\top}\right]$ as \boldsymbol{P} and $\sum_{i \in [n_w]} \mathbf{z}^{(i)}$ as \boldsymbol{Q} . Substitute $\mathbf{z}^{(i)}$ with $\hat{\mathbf{z}}^{(i)}$ in the above expression we have \hat{P} and \hat{Q} . Hence,

$$\begin{split} \left\| \hat{\boldsymbol{\beta}}_{n_{w}} - \boldsymbol{\beta}^{*} \right\|_{F} &= \left\| \begin{bmatrix} \hat{\boldsymbol{P}} \otimes \boldsymbol{I}_{m} \end{bmatrix}^{-1} \left(\begin{bmatrix} \hat{\boldsymbol{Q}} \otimes \boldsymbol{I}_{m} \end{bmatrix} \mathbf{y}^{(i)} \right) - \begin{bmatrix} \boldsymbol{P} \otimes \boldsymbol{I}_{m} \end{bmatrix}^{-1} \left(\begin{bmatrix} \boldsymbol{Q} \otimes \boldsymbol{I}_{m} \end{bmatrix} \mathbf{y}^{(i)} \right) \right\|_{F} \\ &\leq \left\| \begin{bmatrix} \hat{\boldsymbol{P}} \otimes \boldsymbol{I}_{m} \end{bmatrix}^{-1} \left(\begin{bmatrix} \hat{\boldsymbol{Q}} \otimes \boldsymbol{I}_{m} \end{bmatrix} \mathbf{y}^{(i)} \right) - \begin{bmatrix} \hat{\boldsymbol{P}} \otimes \boldsymbol{I}_{m} \end{bmatrix}^{-1} \left(\begin{bmatrix} \boldsymbol{Q} \otimes \boldsymbol{I}_{m} \end{bmatrix} \mathbf{y}^{(i)} \right) \right\|_{F} \\ &+ \left\| \begin{bmatrix} \hat{\boldsymbol{P}} \otimes \boldsymbol{I}_{m} \end{bmatrix}^{-1} \left(\begin{bmatrix} \boldsymbol{Q} \otimes \boldsymbol{I}_{m} \end{bmatrix} \mathbf{y}^{(i)} \right) - \begin{bmatrix} \boldsymbol{P} \otimes \boldsymbol{I}_{m} \end{bmatrix}^{-1} \left(\begin{bmatrix} \boldsymbol{Q} \otimes \boldsymbol{I}_{m} \end{bmatrix} \mathbf{y}^{(i)} \right) \right\|_{F} \\ &= \left\| \begin{bmatrix} \hat{\boldsymbol{P}} \otimes \boldsymbol{I}_{m} \end{bmatrix}^{-1} \right\|_{F} \left\| \begin{bmatrix} \left(\hat{\boldsymbol{Q}} - \boldsymbol{Q} \right) \otimes \boldsymbol{I}_{m} \end{bmatrix} \mathbf{y}^{(i)} \right\|_{F} + \left\| \begin{bmatrix} \hat{\boldsymbol{P}} \otimes \boldsymbol{I}_{m} \end{bmatrix}^{-1} - \begin{bmatrix} \boldsymbol{P} \otimes \boldsymbol{I}_{m} \end{bmatrix}^{-1} \right\|_{F} \\ & \left\| \begin{bmatrix} \boldsymbol{Q} \otimes \boldsymbol{I}_{m} \end{bmatrix} \mathbf{y}^{(i)} \right\|_{F} \end{split}$$

In the first part, by the triangle inequality,

$$\left\|\left[\hat{\boldsymbol{P}}\otimes\boldsymbol{I}_{m}\right]^{-1}\right\|_{\mathrm{F}}\leq\left\|\left[\hat{\boldsymbol{P}}\otimes\boldsymbol{I}_{m}\right]^{-1}-\left[\boldsymbol{P}\otimes\boldsymbol{I}_{m}\right]^{-1}\right\|_{\mathrm{F}}+\left\|\left[\boldsymbol{P}\otimes\boldsymbol{I}_{m}\right]^{-1}\right\|_{\mathrm{F}}.$$

So we only need to prove $\left\| \left[\hat{\boldsymbol{P}} \otimes \boldsymbol{I}_m \right]^{-1} - \left[\boldsymbol{P} \otimes \boldsymbol{I}_m \right]^{-1} \right\|_{\mathrm{F}} \leq \mathcal{O}(\varepsilon)$ to claim $\left\| \hat{\boldsymbol{\beta}}_{n_{\mathrm{w}}} - \boldsymbol{\beta}^* \right\|_{\mathrm{F}} \leq \mathcal{O}(\varepsilon)$. Denote $\hat{P} - P = \Delta P$. From Eq. (21), we know every entry of $\Delta P \otimes I_m$ can be bounded by $\mathcal{O}(\varepsilon)$. By simple calculation we have (

$$(\boldsymbol{P} + \Delta \boldsymbol{P})^{-1} = \boldsymbol{P}^{-1} - \boldsymbol{P}^{-1} \Delta \boldsymbol{P} \boldsymbol{P}^{-1} + \mathcal{O}(\varepsilon^2).$$

Then we have

$$\boldsymbol{P}^{-1} - \hat{\boldsymbol{P}}^{-1} = \boldsymbol{P}^{-1} \Delta \boldsymbol{P} \boldsymbol{P}^{-1} + \mathcal{O}(\varepsilon^2) = \mathcal{O}(\varepsilon).$$

H.2 Layer 1

In this subsection, we justify the strong consistency of layer 1 objective functional minimizer estimator in detail, i.e. the layer 1 QP/LP solution space. We will omit the detailed proof of Algorithm 2 line 4 to 7 strong consistency since it is similar to the proof of Lemma H.2. Additionally, we also omit the strong consistency of the $x \mapsto (A^*x)^+$ function estimator because it can be directly obtained by Lemmas H.1 and H.2 and the continuous mapping theorem.

We use a new optimization problem equivalent to the optimization of G_1 . Before that, we first define the equivalence between two optimization problems as follows.

Definition H.3. Let opt1 and opt2 be two optimization problems, with f_1 , f_2 as the respective objective functions. Then opt1 and opt2 are said to be equivalent if given a feasible solution to opt1, namely x_1 , a feasible solution to opt2 is uniquely corresponded, namely x_2 , such that $f_1(x_1) = f_2(x_2)$, and vice versa.

The new optimization problem and its equivalence to the optimization of G_1 is described in Lemma H.4. **Lemma H.4.** The optimization of G_1 (Eq. (6)) is equivalent to

$$\min_{\boldsymbol{A}} f(\boldsymbol{A}) = \frac{1}{2} \mathbb{E}_{\mathbf{x}} \left[\left\| (\boldsymbol{A}\mathbf{x} - \mathbf{h})^{+} \right\|^{2} \right].$$
(22)

Proof. To see this, suppose A_1 is one optimal solution to Eq. (22), then we can construct $r_1(\mathbf{x}) = (A_1\mathbf{x} - \mathbf{h})^+$ so that $G_1(\mathbf{A}_1, r_1) = f(\mathbf{A}_1)$ and the optimality implies

$$\min_{\boldsymbol{A},r} G_1(\boldsymbol{A},r) \le \min_{\boldsymbol{A}} f(\boldsymbol{A}).$$

On the other hand, suppose (\mathbf{A}_2, r_2) is an optimum of G_1 . Let $r_3(\mathbf{x}) = (\mathbf{h} - \mathbf{A}_2 \mathbf{x})^+$, then $\forall \mathbf{x} \in \mathbb{R}^d$ and \mathbf{h} be the corresponding hidden output, if $[\mathbf{h} - \mathbf{A}_2 \mathbf{x}]_j \ge 0$, then $[r_3(\mathbf{x}) + \mathbf{A}_2 \mathbf{x} - \mathbf{h}]_j = 0$, otherwise $[r_3(\mathbf{x}) + \mathbf{A}_2 \mathbf{x} - \mathbf{h}]_j^2 = [\mathbf{A}_2 \mathbf{x} - \mathbf{h}]^2 \le [h_2(\mathbf{x}) + \mathbf{A}_2 \mathbf{x} - \mathbf{h}]_j^2$ since h_2 is nonnegative. So we know that

$$\min_{\mathbf{A}_{2}} G_{1}(\mathbf{A}, r) = G_{1}(\mathbf{A}_{2}, r_{2}) = G_{1}(\mathbf{A}_{2}, r_{3}) = f(\mathbf{A}_{2}).$$

From the optimality, we further have

$$\min_{\boldsymbol{A},r} G_1(\boldsymbol{A},r) \geq \min_{\boldsymbol{A}} f(\boldsymbol{A}).$$

From the simple calculation above, we can see that one optimal solution to Eq. (22) has a one-to-one correspondence to an optimal solution to G_1 .

Similarly, the empirical version of the two problems are equivalent, which indicates their consistency in the empirical estimation being equivalent. In the following, we justify the strong consistency of empirical Eq. (22) instead of G_1

$$\min_{\boldsymbol{A}} \hat{f}_n(\boldsymbol{A}) = \frac{1}{2n} \sum_{i \in [n]} \left\| \left(\boldsymbol{A} \boldsymbol{x}^{(i)} - \boldsymbol{h}^{(i)} \right)^+ \right\|^2.$$

Denote $\mathbb{T}^* := \{ \operatorname{diag}(\mathbf{k}) \cdot \mathbf{A}^* \mid 0 \leq k_j \leq 1, j \in [d] \}$ as the true optimal solution set, and $\hat{\mathbb{T}}_n$ as the optimal solution set corresponding to the *n*-sample problem. In the following, we justify four conditions in a row that hold for f to derive the strong consistency of its optimal solution estimator.

Lemma H.5. Let $\hat{\mathbb{T}}_{n'}$ be the layer 1 QP/LP solution space by n' samples. Then there exists a compact set \mathbb{C} determined by \mathbf{A}^* , namely $\mathbb{C}(\mathbf{A}^*)$, s.t. $\hat{\mathbb{T}}_{n'} \subset \mathbb{C}(\mathbf{A}^*)$ w.p. 1 as $n \to \infty$.

Proof. $\forall l \in [d]$, let a^* be the *l*-th row of A^* , and \hat{a}_l be the *l*-th row of $\hat{A}_{n'}$. We'd like first to prove that the set

$$\mathbb{T}_{n'}^{l} = \{ \boldsymbol{a}_{l} : \boldsymbol{a}_{l} \text{ is the } l\text{-th row of } \boldsymbol{A}, \text{ where } \boldsymbol{A} \in \mathbb{T}_{n'} \}$$

is compact w.p.1.

Suppose n' > d, and among the n' samples, we classify them into two folds. To avoid ambiguity, let $\boldsymbol{u}^{(i)}$ be the points such that $(\boldsymbol{a}^*)^{\top} \boldsymbol{u}^{(i)} > 0$, $i \in [q]$; and $\boldsymbol{v}^{(j)}$ be the points such that $(\boldsymbol{a}^*)^{\top} \boldsymbol{v}^{(j)} < 0$, $j \in [n-q]$. From the analysis of Theorem 5.1, we have $(\boldsymbol{a}^*)^{\top} \boldsymbol{u}^{(i)} \ge \hat{\boldsymbol{a}}^{\top} \boldsymbol{u}^{(i)}$, $\forall i \in [q]$ and $\hat{\boldsymbol{a}}^{\top} \boldsymbol{v}^{(j)} \le 0$, $\forall j \in [n-q]$. It follows that we can rewrite $\hat{\mathbb{T}}_{n'}^l$ as a polyhedron

$$\hat{\mathbb{T}}_{n'}^l = \{ \boldsymbol{a} \in \mathbb{R}^d : \boldsymbol{a}^\top u^{(i)} \leq (\boldsymbol{a}^*)^\top \boldsymbol{u}^{(i)}, \boldsymbol{a}^\top \boldsymbol{v}^{(j)} \leq 0 \}.$$

We are going to show the polyhedron $\hat{\mathbb{T}}_{n'}^l$ is bounded by contradiction. If it is not bounded, then $\exists d \in \mathbb{R}^d$, $d \neq \mathbf{0}$ and $\tilde{a} \in \hat{\mathbb{T}}_{n'}^l$, such that $\forall \lambda > 0$, $\tilde{a} + \lambda d \in \hat{\mathbb{T}}_{n'}^l$. Then

$$\left(ilde{oldsymbol{a}}+\lambdaoldsymbol{d}
ight)^{ op}oldsymbol{u}^{(i)}= ilde{oldsymbol{a}}^{ op}oldsymbol{u}^{(i)}+\lambdaoldsymbol{d}^{ op}oldsymbol{u}^{(i)}\leq\left(oldsymbol{a}^*
ight)^{ op}oldsymbol{a}^{(i)}$$

Similarly,

$$(\tilde{\boldsymbol{a}} + \lambda \boldsymbol{d})^{\top} \boldsymbol{v}^{(i)} = \tilde{\boldsymbol{a}}^{\top} \boldsymbol{v}^{(j)} + \lambda \boldsymbol{d}^{\top} \boldsymbol{v}^{(j)} \le 0.$$

From the definition, λ can be arbitrarily big, then $d^{\top} u^{(i)} \leq 0$, $\forall i \in [q]$, and $d^{\top} v^{(j)} \leq 0$, $\forall j \in [n-q]$.

Since we know span $\{u^{(i)}\} = \mathbb{R}^d$ w.p. 1, then there \exists some i^* such that $d^\top u^{(i^*)} < 0$, w.p. 1. (Otherwise if $d^\top u^{(i)} = 0$ for all $i \in [q]$, then either span $\{u^{(i)}\} \neq \mathbb{R}^d$ or d = 0.) Under our assumption that, $\hat{\mathbb{T}}_{n'}^l$ is not bounded, we know the following system (w.r.t x) has a feasible solution w.p. 1

$$\begin{bmatrix} -\boldsymbol{U} \\ -\boldsymbol{V} \end{bmatrix} \boldsymbol{x} \ge 0, \, \boldsymbol{x}^{\top} \boldsymbol{u}^{(i^*)} < 0, \tag{I}$$

where every row of \boldsymbol{U} and \boldsymbol{V} is $(\boldsymbol{u}^{(i)})^{\top}$ and $(\boldsymbol{v}^{(j)})^{\top}$ respectively. By Farkas' Lemma (Farkas, 1902), the system

$$[-\boldsymbol{U}^{\top}, -\boldsymbol{V}^{\top}] \cdot \boldsymbol{x} = \boldsymbol{u}^{(i^*)}, \, \boldsymbol{x} \ge 0 \tag{II}$$

is not feasible (w.p. 1). We claim that $u^{(i^*)}$ lies in the conic hull of $-v^{(j)}$'s w.p. 1. So that the second system actually has a feasible solution and thus it raises the contradiction.

Denote the conic hull as

$$\mathbb{H} = \left\{ \boldsymbol{t} \in \mathbb{R}^d : \boldsymbol{t} = \sum_{j \in [n-q]} \lambda_j \left(-\boldsymbol{v}^{(j)} \right), \, \lambda_j \ge 0 \text{ for } \forall j \in [n-q] \right\}.$$

Now suppose $\boldsymbol{u}^{(i^*)} \notin \mathbb{H}$, by the supporting hyperplane theorem (Luenberger, 1997), $\exists \boldsymbol{b} \in \mathbb{R}^d$, $\boldsymbol{b} \neq \boldsymbol{0}$, such that $\boldsymbol{b}^\top \boldsymbol{u}^{(i^*)} \leq \boldsymbol{b}^\top \boldsymbol{t}$ for $\forall \boldsymbol{t} \in \mathbb{H}$. Then by definition,

$$- \boldsymbol{v}^{(j)} \in \{ \boldsymbol{t} : \boldsymbol{b}^{\top} \boldsymbol{t} \leq \boldsymbol{b}^{\top} \boldsymbol{u}^{(i^*)} \}, ext{ for } orall j \in [n-q].$$

Denote the hyperplane $\mathbb{J} = \{ \boldsymbol{t} : \boldsymbol{b}^{\top} \boldsymbol{t} \leq \boldsymbol{b}^{\top} \boldsymbol{u}^{(i^*)} \}$, then

$$P\left(\boldsymbol{u}^{(i^*)} \notin \mathbb{H}\right) \leq P\left(-\boldsymbol{v}^{(j)} \in \mathbb{J}, \text{ for } \forall j \in [n-q]\right) = P\left(-\boldsymbol{v}^{(j)} \in \mathbb{J}\right)^{n-q}$$

Since we have in this case a geometric sequence, we know its infinite sum is bounded. By Borel-Cantelli lemma (Borel, 1909), we conclude that $u^{(i^*)} \in \mathbb{H}$ w.p. 1. Then system II is feasible w.p. 1. So that $\hat{\mathbb{T}}_{n'}^l$ is compact w.p. 1.

Now we prove that there exists a compact set $\mathbb{C}(\mathbf{A}^*)$, s.t. $\hat{\mathbb{T}}_{n'} \subset \mathbb{C}(\mathbf{A}^*)$ w.p. 1 as $n \to \infty$. Similarly, we focus on the analysis of one row. As discussed above, $\hat{\mathbb{T}}_{n'}^l$ is compact w.p. 1. Let $\hat{\mathbb{W}}_{n'}^1$ be the set of all $\mathbf{u}^{(i)}$ sampled in estimating $\hat{\mathbb{T}}_{n'}$, and $\hat{\mathbb{W}}_{n'}^2$ be the set of all $\mathbf{v}^{(j)}$ sampled. Now for another set $\hat{\mathbb{T}}_{n''}$, similarly define sample point sets $\hat{\mathbb{W}}_{n''}^1$ and $\hat{\mathbb{W}}_{n''}^2$. We claim that $\forall \mathbf{u}^{(i)} \in \hat{\mathbb{W}}_{n'}^1$, $\mathbf{u}^{(i)}$ lies in the conic hull of $\hat{\mathbb{W}}_{n''}^1$. Actually, this part of the proof is very similar to the way we prove $\mathbf{u}^{(i^*)} \in \mathbb{H}$ w.p. 1, so we will omit it here.

 $\forall \boldsymbol{a} \in \hat{\mathbb{T}}_{n''}, \text{ let } \boldsymbol{x}^{(1)} \text{ and } \boldsymbol{x}^{(2)} \text{ be two different points in } \hat{\mathbb{W}}_{n''}^{1}. \text{ Then } \boldsymbol{a}^{\top}(\lambda_{1}\boldsymbol{x}^{(1)} + \lambda_{2}\boldsymbol{x}^{(2)}) \leq (\boldsymbol{a}^{*})^{\top}(\lambda_{1}\boldsymbol{x}^{(1)} + \lambda_{2}\boldsymbol{x}^{(2)}) \text{ for } \forall \lambda_{1} \geq 0 \text{ and } \lambda_{2} \geq 0. \text{ This simple calculation reveals } \boldsymbol{a}^{\top}\boldsymbol{u}^{(i)} \leq (\boldsymbol{a}^{*})^{\top}\boldsymbol{u}^{(i)} \text{ for } \forall \boldsymbol{u}^{(i)} \in \hat{\mathbb{W}}_{n'}^{1} \text{ (from the claim we made). This implies that } \hat{\mathbb{T}}_{n''} \subset \mathbb{C}(\boldsymbol{A}^{*}) \text{ w.p. 1 as } n \to \infty, \text{ too.}$

Lemma H.6. The minimizer space of $f(\mathbf{A})$, i.e. $\mathbb{T}^* = \{ \operatorname{diag}(\mathbf{k}) \cdot \mathbf{A}^* \mid 0 \leq k_j \leq 1, j \in [d] \}$, is contained in $\mathbb{C}(\mathbf{A}^*)$.

Lemma H.7. $f(\mathbf{A})$ is finite valued and continuous on $\mathbb{C}(\mathbf{A}^*)$.

Lemmas H.6 and H.7 are easily obtained by the formulation of f (see Eq. (22)) and Lemma H.5.

Lemma H.8 (uniform a.s. convergence). $\hat{f}_n(A) \xrightarrow{a.s.} f(A)$ as $n \to \infty$, uniformly in $A \in \mathbb{C}(A^*)$.

Proof. Name single-sample objective $g(\boldsymbol{x}, \boldsymbol{A}) = \frac{1}{2} \left\| (\boldsymbol{A}\boldsymbol{x} - \boldsymbol{h})^+ \right\|^2$. The uniform a.s. convergence is guaranteed by the uniform law of large numbers (Jennrich, 1969):

- a) By Lemma H.5, $\mathbb{C}(\mathbf{A}^*)$ is a compact set.
- b) g is continuous w.r.t. A by its formulation and measurable over x at each $A \in \mathbb{C}(A^*)$.
- c) In fact,

$$g(\boldsymbol{x}, \boldsymbol{A}) = \frac{1}{2} \left\| (\boldsymbol{A}\boldsymbol{x} - \boldsymbol{h})^{+} \right\|^{2} \\ \leq \|\boldsymbol{A}\boldsymbol{x}\|^{2} + \|\boldsymbol{A}^{*}\boldsymbol{x}\|^{2} \\ \leq (\|\boldsymbol{A}\|_{\mathrm{F}} + \|\boldsymbol{A}^{*}\|_{\mathrm{F}}) \|\boldsymbol{x}\|^{2}.$$

Since $A \in \mathbb{C}(A^*)$ is in a compact set,

$$g(oldsymbol{x},oldsymbol{A}) \leq \left[\sup_{oldsymbol{A} \in \mathbb{C}(oldsymbol{A}^*)} \left\|oldsymbol{A}
ight\|_{\mathrm{F}} + \left\|oldsymbol{A}^*
ight\|_{\mathrm{F}}
ight] \left\|oldsymbol{x}
ight\|^2.$$

Thus the dominating function exists.¹⁵

By Lemmas H.5 to H.8, all of the conditions are satisfied in (Shapiro et al., 2014, Thm. 5.3). Thus, we have the strong consistency of layer 1 objective optima estimator as described in Lemma H.9.

Lemma H.9 (QP/LP strong consistency, layer 1). $D_{\mathrm{H}}(\hat{\mathbb{T}}_{n'}, \mathbb{T}^*) \xrightarrow{a.s.} 0 \text{ as } n \to \infty.$

Similar to Lemma H.2, we have the strong consistency of the layer 1 scale factor estimator as described in Lemma H.10.

Lemma H.10 (scale factor estimator strong consistency, layer 1). Let $k_{n'_{sf}}(\hat{\mathbf{C}}_n)$ be the n'_{sf} -sample estimator of k_j via LR: Algorithm 2 line 5. given that $h_j^{(i)} > 0$. Define sets

$$\mathbb{V}_{n'_{et},n} \coloneqq \{ \mathrm{k}_{n'_{et}}(\boldsymbol{A}) : \boldsymbol{A} \in \hat{\mathbb{T}}_{n'} \}, \text{ and } \mathbb{V}^* \coloneqq [0,1].$$

Then $\lim_{n'_{sf}\to\infty} \lim_{n'\to\infty} D_{\mathrm{H}}(\mathbb{V}_{n'_{sf},n'},\mathbb{V}^*) \stackrel{a.s.}{=} 0.$

Remark. In case $k_j = 0$, suppose the algorithm finds a solution over a continuous distribution with [0, 1] as support and the probability that it finds a solution with scale factor 0 is 0.

With Theorem 5.2 and the continuous mapping theorem, the strong consistency of layer 1 estimation is guaranteed.

Theorem H.11 (strong consistency, layer 1). Suppose $\hat{\mathbf{A}}_{n'}$ is scaled by $\hat{\mathbf{k}}_{n'_{sf}}$. Then $\hat{\mathbf{A}}_{n'} \xrightarrow{a.s.} \mathbf{A}^*$ as $n', n'_{sf} \to \infty$.

By Theorems H.3 and H.11, Theorem 6.3 is guaranteed by the continuous mapping theorem.

I Sample Complexity Results

In this appendix, we begin with an intuition that justifies how our core QP/LP approaches eliminate the exponential sample efficiency compared to the vanilla LR approach. In layer 2 learning, the optimization of C is done to find a feasible point in the space determined by n inequalities each of which corresponds to a sample. Taking the *j*-th row of the inequality, i.e. $C_{j,:} \mathbf{y} \ge x_j$. Every time we get a new sample $\mathbf{x}^{(i)}, \mathbf{y}^{(i)}$ in,

- The inequality $C_{j,:} y^{(i)} \ge x_j^{(i)}$ eliminates the solution space of $C_{j,:}$ by one of the spaces divided by plane $C_{j,:} y^{(i)} = x_j^{(i)}$. This property guarantees fast convergence speed at early phase since $C_{j,:}$'s feasibility starts from \mathbb{R}^d .
- In fact, when ReLU does not activate A*x⁽ⁱ⁾ at the *j*-th row, i.e. A_j^{*}, x⁽ⁱ⁾ ≤ 0, the theoretical solution of C_j, that C_j, B* = (0,...,0,1,0,...,0) (a 1 at index *j*) directly lies on the separating plane C_j, y⁽ⁱ⁾ = x_j⁽ⁱ⁾:

$$C_{j,:} y^{(i)} = C_{j,:} B^* \left[\left(A^* x^{(i)} \right)^+ + x^{(i)} \right] = x_j^{(i)}$$

This property remarkably speeds up further constraints on the solution space of $C_{j,:}$ to the correct estimate and is with high probability, since it directly depends on the sign of $A^* \mathbf{x}^{(i)}$.

 $^{^{15}\}text{Here,}$ we assume $\mathbb{E}[\|\mathbf{x}\|^2] < +\infty$ as is commonly done in empirical analysis.

With the vanilla LR approach mentioned in Section 3 all the dimensions need to have the correct sign at the same time. This is the reason for the exponential complexity of this approach. Our main algorithm only requires one dimension to get the correct sign for a single sample, which avoids the exponential sample size expectation. The convergence speed is determined by the actual probability of each dimension not getting activated by ReLU, and such probability is determined by specified input distribution.

Roadmap Our sample complexity results follow several intermediate results. We first demonstrate that the constraints that define our optimization objective give a bounded polytope (Theorem I.2). This happens with a high probability, given a set of samples. Once we have shown that, we then show that the optimization problem constraints we use define a polytope with a small diameter with a high probability given a sufficient number of samples, where the true network parameters exist. The boundnesses of the constraint polytope is required for this result.

Details We follow the above intuitive explanation with a more rigorous analysis that describes the sample complexity of recovering $C_{j,:}$ for $j \in [d]$.

Let us assume we have n samples of $\boldsymbol{y}^{(i)}$, $i \in [n]$. We denote by $\boldsymbol{F} \in \mathbb{R}^{n \times d}$ such that $F_{ij} = -y_j^{(i)}$. We assume n > d. We denote by $L_j = \mathbb{E}[\mathbf{x}_j^2]$.

Condition I.1 (Haar condition on samples). F satisfies the Haar condition, i.e. that every submatrix of size $d \times d$ is invertible, almost surely.

We will assume that Condition I.1 holds for the rest of the discussion. For a discussion of the Haar condition, see Schmidt & Mattheiss (1977).

Fix $j \in [d]$. The first step in our process is to show that with high likelihood, the inequalities $C_{j,:}y^{(i)} \ge x_j^{(i)}$ for $i \in [n]$, or equivalently $C_{j,:}(-y^{(i)}) \le (-x_j^{(i)})$, or equivalently $FC_{j,:} \le -x_j^{(i)}$ define a bounded polytope with high probability.

Schmidt & Mattheiss (1977) describe a necessary and sufficient condition for this set of inequalities, which describe an intersection of half-spaces, to be bounded. More specifically, under Condition I.1, such an intersection is unbounded if and only if the subspace spanned by the columns of F intersected with the negative quadrant of \mathbb{R}^n is non-empty.

This means that this halfspace intersection is unbounded if and only if there exists $\alpha \in \mathbb{R}^d$, $\alpha \neq 0$, such that $F_{i,:}\alpha \leq 0$ for $i \in [n]$.

For a given $\alpha \in \mathbb{R}^d$, $\alpha \neq 0$, let $h_{\alpha}(\boldsymbol{y}) = \operatorname{sgn}(\sum_{i=1}^d \alpha_i \boldsymbol{y}_i)$ where $\operatorname{sgn}(z)$ for $z \in \mathbb{R}$ is 0 if $z \leq 0$ and 1 otherwise. Define:

$$\overline{P}(\alpha) = \mathbb{E}[h_{\alpha}(\boldsymbol{y})],$$
$$\overline{P}_{n}(\alpha) = \frac{1}{n} \sum_{i=1}^{n} h_{\alpha}(\boldsymbol{y}^{(i)}).$$

A well-known result of VC-theory (Vapnik, 1999), specialized to linear separators, states that: **Theorem I.1** (Vapnik 1999). For any $t > \sqrt{2/n}$, it holds that:

$$P\left(\sup_{\alpha \in \mathbb{R}^d} |\overline{P}_n(\alpha) - \overline{P}(\alpha)| \ge t\right) \le 4\left(\frac{2en}{d+1}\right)^{(d+1)} \exp(-nt^2/8).$$

We continue with the assuming the following separability condition:

Condition I.2 (separability of **y**). There exists a > 0 such that for any $\alpha \in \mathbb{R}^d$, $\alpha \neq 0$, $\overline{P}(\alpha) > a$.

Theorem I.2 (unboundedness of halfspace intersection). Under Condition I.2 and Condition I.1, the halfspace intersection defined by \mathbf{F} for identifying $C_{j,:}$ is bounded with probability

$$P\left(\mathbf{F} \text{ defines unbounded intersection}\right) \leq 4\left(\frac{2en}{d+1}\right)^{(d+1)}\exp(-na^2/8)$$

 $if n > \max\{2/a^2, d\}.$

Proof. For any $\alpha \in \mathbb{R}^d$, $\alpha \neq 0$, it holds that if $\overline{P}_m(\alpha) = 0$ then for all $i \in [n]$, $\operatorname{sgn}(\sum_{j=1}^d \alpha_j \boldsymbol{y}_j^{(i)}) = 0$. Each such case means that for this specific α , we found a span of the columns of \boldsymbol{F} such that its *i*th coordinate is 0, and as such, a case in which \boldsymbol{F} is unbounded.

This means that if for all $\alpha \in \mathbb{R}^d$, it holds that $\overline{P}_n(\alpha) > 0$, then the *n* samples define a bounded halfspace intersection.

In addition, if $\overline{P}_n(\alpha) = 0$, then

$$\overline{P}_n(\alpha) - \overline{P}(\alpha) = -\overline{P}(\alpha), \tag{23}$$

therefore,

 $|\overline{P}_n(\alpha) - \overline{P}(\alpha)| = \overline{P}(\alpha).$

Since $a < \overline{P}(\alpha)$, it holds that if Eq. (23) holds then:

$$|\overline{P}_n(\alpha) - \overline{P}(\alpha)| > a.$$

Therefore,

$$\begin{split} P\left(\boldsymbol{F} \text{ defines unbounded intersection}\right) &\leq P\left(\exists \alpha \neq 0, \, \overline{P}_n(\alpha) = 0\right) \\ &\leq P\left(\sup_{\alpha} |P_n(\alpha) - P(\alpha)| > a\right) \\ &\leq 4\left(\frac{2en}{d+1}\right)^{(d+1)} \exp(-na^2/8), \end{split}$$

where the last inequality is the result of Theorem I.1.

Theorem I.2 gives a well-behaving sample complexity for the *n* samples to define halfspace intersection that is bounded, when solving for $C_{j,:}$. Note that the theorem does not depend on *j*, because the use of \mathbf{x}_j is not necessary for the boundedness result of Schmidt & Mattheiss (1977).

We now further show that when the halfspace intersection is bounded, then the diameter r of the bounded space is smaller (where diameter is defined as the maximal distance between any two points).

The diameter is bounded by the maximal distance between all pairs of vertices of the intersection (the intersection is a bounded polytope). Let $\mathbb{K} \subseteq \{v \mid \mathbf{F}v \leq -\mathbf{x}_{j}^{(:)}\}$ be the set of these vertices. The size $|\mathbb{K}|$ is bounded from above by $O(n^{d/2})$ (by McMullen's Upper Bound theorem; Toth et al. 2017).

Condition I.3 (equality saturation for vertices). For every $v \in \mathbb{K}$, there are *d* equalities that are satisfied, in the form of $\mathbf{F}_{I,:}v = -x_i^{(I)}$, such that $I \subset [n]$ and |I| = d.

Condition I.4 (bounded input distribution). For every $j \in [d]$, $|x_j| \leq b$ for some b > 0.

Lemma I.3 (Hoeffding's inequality). Let Z_1, \ldots, Z_n be n random variables such that $Z_i \in [a_i, b_i]$ almost surely. Then, for all t > 0:

$$P\left(\sum_{i=1}^{n} Z_i - \sum_{i=1}^{n} \mathbb{E}[Z_i] \ge t\right) \le \exp\left(-\frac{2t^2}{\sum_{i=1}^{n} (b_i - a_i)^2}\right).$$

Lemma I.4. Assume Condition I.1, Condition I.3 and Condition I.4 are satisfied.

In addition, let $\sigma^* = \mathbb{E}[\sigma]$ where σ is the smallest singular value of $\mathbf{F}_{I,:}$ for $I = \{1, \ldots, d\}$. Let $\varepsilon > \frac{\sqrt{db^2}}{\sigma^* \nu}$, where $0 \leq \nu = \frac{\log(\delta/4)}{\sigma^*} + 1 \leq 1$ for any $1 > \delta > 0$. Define:

$$L(\varepsilon, d, \sigma^*, b, L_j, \delta) = \sqrt[d]{\frac{\delta}{4D^2}} \exp\left(\frac{2((\sigma^*)^2 \nu^2 \varepsilon^2 - dL_j)^2}{d^2 b^2}\right)$$

For any set of samples $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$, $i \in [n]$ such that the halfspace intersection is bounded, it holds that with probability at least $1-\delta > 0$ the diameter of the halfspace intersection is smaller than ε if the following inequality is satisfied:

$$n \ge L(\varepsilon, d, \sigma^*, b, \delta) \ge d.$$

Proof. The diameter of the intersection is not larger than the maximal distance between two pairs of vertices $||v - u||, u, v \in \mathbb{K}$. By Condition I.3, let I, J be the two subsets of integers such that |I| = |J| = d and

$$egin{aligned} &m{F}_{I,:} u = -m{x}_j^{(I)}, \ &m{F}_{J,:} v = -m{x}_j^{(J)}. \end{aligned}$$

By the Haar assumption on F, the inverse of F exists, and it holds that:

$$||u - v|| = ||\boldsymbol{F}_{I,:}^{-1} \boldsymbol{x}_{j}^{(I)} - \boldsymbol{F}_{J,:}^{-1} \boldsymbol{x}_{j}^{(J)}|| \le ||\boldsymbol{F}_{I,:}^{-1}||^{*} \cdot ||\boldsymbol{x}_{j}^{(I)}|| + ||\boldsymbol{F}_{J,:}^{-1}||^{*} \cdot ||\boldsymbol{x}_{j}^{(J)}||.$$

Consider $||F_{I,:}^{-1}||^*$, the operator norm of $F_{I,:}^{-1}$ which equals $1/\sigma_I$ where σ_I is the smallest singular value of $F_{I,:}$. Similarly, we define σ_J . Therefore:

$$||u - v|| \le ||\boldsymbol{x}_{j}^{(I)}|| / \sigma_{I} + ||\boldsymbol{x}_{j}^{(J)}|| / \sigma_{J}.$$

The event that the diameter r is larger than ε can be bounded by (A(I, J)) is the event $\sigma_I \ge c, \sigma_J \ge c$ for some c):

$$P(r \ge \varepsilon) \le p(r \ge \varepsilon, A(I, J)) + P(\sigma_I \le c \land \sigma_J \le c).$$

In addition,

$$\begin{split} P\left(r \geq \varepsilon, A(I,J)\right) &\leq P\left(\max_{u,v \in \mathbb{K}} ||u-v|| \geq \varepsilon, A(I,J)\right) \\ &\leq P\left(\max_{u,v} ||\boldsymbol{x}_{j}^{(I)}||/c + ||\boldsymbol{x}_{j}^{(J)}||/c \geq \varepsilon\right) \\ &\leq \sum_{u,v \in \mathbb{K}} P\left(||\boldsymbol{x}_{j}^{(I)}||/c + ||\boldsymbol{x}_{j}^{(J)}||/c \geq \varepsilon\right), \end{split}$$

We know that $|\mathbb{K}| \leq D(n^{d/2})$ for some constant D (Toth et al., 2017), therefore $P(r \geq \varepsilon, A(I, J))$ can be further bounded by

$$D^2 n^d \cdot P\left(||\boldsymbol{x}_j^{(I)}||/c + ||\boldsymbol{x}_j^{(J)}||/c \ge \varepsilon\right).$$

Since the distribution of $x_j^{(I)}$ and $x_j^{(J)}$ (and the singular values) is identical, $P(r \ge \varepsilon, A(I, J))$ can be further bounded by

$$2D^2 n^d \cdot P\left(||\boldsymbol{x}_j^{(I)}||/c \ge \varepsilon\right),\tag{24}$$

for any I, for example $I = \{1, \ldots, d\}$.

We can use any value for $c = \mathbb{E}[\sigma_I] \cdot \nu$ to be smaller than σ_I , and the inequality still holds. Therefore, Eq. (24) can be further bounded by (assuming $\nu < 1$):

$$2D^2n^d \cdot \left(P\left(||\boldsymbol{x}_j^{(I)}|| / \mathbb{E}[\sigma_I]\nu \ge \varepsilon \right) \right).$$

This can be further bounded by:

$$2D^2 n^d \cdot \left(P\left(||\boldsymbol{x}_j^{(I)}||^2 \ge \mathbb{E}[\sigma_I]^2 \nu^2 \varepsilon^2 \right) \right),$$
(25)

where we also square the norm of $\boldsymbol{x}_{j}^{(I)}$ (and its bounding term).

The term $P(\sigma_I \leq \mathbb{E}[\sigma_I] \cdot \nu \wedge \sigma_J \leq \mathbb{E}[\sigma_J] \cdot \nu)$ for small ν is going to be close to 0, more specifically, it measures the probability of the complement of the event A(I, J). This probability is at most twice the probability that $\exp(-\sigma_I) \geq \exp(-\nu \cdot \mathbb{E}[\sigma_I])$, which by Markov's inequality is smaller than $\mathbb{E}[\exp(-\sigma_I)] / \exp(-\nu \cdot \mathbb{E}[\sigma_I]) \leq \exp(\mathbb{E}[(\nu-1) \cdot \sigma_I]) = \exp(((\nu-1)\sigma^*))$ (by Jensen's inequality). Therefore, the probability of the complement of A(I, J) is bounded by:

$$P\left(\sigma_I \le \nu \sigma^* \land \sigma_J \le \sigma^*\right) \le 2\exp((\nu - 1) \cdot \sigma^*).$$
(26)

Consider that $||\boldsymbol{x}_{j}^{(I)}||^{2}$ is the sum of *d* i.i.d. samples \boldsymbol{x}_{j}^{2} . Under Condition I.4, we assume $|\boldsymbol{x}_{j}^{2}| \leq b$. Therefore, by Hoeffding's inequality, we can get an upper bound on the probability term in Eq. (25) to hold:

$$P\left(||\boldsymbol{x}_{j}^{(I)}||^{2} - d \cdot \mathbb{E}[\boldsymbol{x}_{j}^{2}] \geq (\sigma^{*})^{2} \nu^{2} \varepsilon^{2} - d \cdot \mathbb{E}[\boldsymbol{x}_{j}^{2}]\right) \leq \exp\left(\frac{-2((\sigma^{*})^{2} \nu^{2} \varepsilon^{2} - d \cdot \mathbb{E}[\boldsymbol{x}_{j}^{2}])^{2}}{db^{2}}\right),$$
(27)

remembering it is stated in the theorem that $\sigma^* \nu \varepsilon \geq \sqrt{db^2}$ and the fact that $\sqrt{db^2} \geq \sqrt{d \cdot \mathbb{E}[x_j^2]}$ by Condition I.4. Hence, $\mathbb{E}[\sigma_I]^2 \nu^2 \varepsilon^2 - d \cdot \mathbb{E}[x_j^2] > 0$, and the use of Hoeffding's inequality is allowed.

Let $\delta > 0$, then taking Eq. (27) with the union bound constant from K and setting it to $\delta/2$ we get:

$$2D^2 n^d \exp\left(\frac{-2((\sigma^*)^2 \nu^2 \varepsilon^2 - dL_j)^2}{db^2}\right) \le \delta/2.$$

We can now get an upper bound on n:

$$n \leq \sqrt[d]{\frac{\delta}{4D^2}} \exp\left(\frac{2((\sigma^*)^2\nu^2\varepsilon^2 - dL_j)^2}{d^2b^2}\right) = L(\varepsilon, d, \sigma^*, b, L_j, \delta).$$

We want, in addition, the complement of A(I, J) to have probability at most $\delta/2$, so therefore, we choose $\nu = \frac{\log(\delta/4)}{\sigma^*} + 1$ according to Eq. (26).

Noting that the diameter of the sphere gets smaller as n increases (because there are more halfspaces that possibly intersect into a smaller space).

Note that Lemma I.4 requires $\delta \ge \exp(-\sigma^*)$ in its statement (this can be inferred from the relationship between ν and δ). To have an arbitrary confidence $1 - \delta$ in the diameter being smaller than ε , we can run the algorithm k times with an n as needed. The probability of all of them having diameter larger than ε is smaller than δ^k .

Lemma I.4 describes under which condition we get a small error for $C_{j,:}$ if we fix j. Using a union bound, we bound the probability that for any $j \in [d]$, the corresponding feasible halfspace intersection space is small. This will require adding a factor of d to δ , and a factor of d to n. We overcome this by making sure that the probability for each j having a diameter larger than ε is smaller than δ/d , and we use dn samples for a choice of n satisfied by Lemma I.4. Using a union bound, we can show that the probability that for any j the diameter is larger than ε is smaller than δ .

In addition, note that Lemma I.4 works only when we assume the halfspace intersection is bounded. This is the case where we can use the diameter bound through the set of vertices. To take this into account, we also use Theorem I.2. This leads to the following result regarding the sample complexity of our algorithm.

Theorem I.5 (sample complexity of learning ReLU two-layered networks (layer 2)). Assume Condition I.1, Condition I.2, Condition I.3 and Condition I.4 are satisfied.

In addition, let $\sigma^* = \mathbb{E}[\sigma]$ where σ is the smallest singular value of $\mathbf{F}_{I,:}$ for $I = \{1, \ldots, d\}$. Let $\varepsilon > \frac{\sqrt{db^2}}{\sigma^* \nu}$, where $0 \le \nu = \frac{\log(\delta/8d)}{\sigma^*} + 1 \le 1$ for any $1 > \frac{\delta}{2d} > 0$.

Define:

$$L^*(\varepsilon, d, \sigma^*, b, L_j, \delta) = \min_j \sqrt[d]{\frac{\delta}{8dD^2}} \exp\left(\frac{2((\sigma^*)^2\nu^2\varepsilon^2 - dL_j)^2}{d^2b^2}\right)$$

and define:

$$L_u(a,\delta,d) = \frac{(d+1)\log 6 + \log 4 - \log \frac{\delta}{2}}{a^2/8 - 1/e}$$

For any set of samples $(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})$, $i \in [n]$, it holds that with probability at least $1 - \delta > 0$ the diameter of the halfspace intersection is smaller than ε if the following inequalities hold:

$$n \ge d \cdot L^*(\varepsilon, d, \sigma^*, b, \delta),$$
$$L^*(\varepsilon, d, \sigma^*, b, \delta) \ge \max\left\{d, L_u(a, \delta, d), 2/a^2\right\}.$$

Proof. Let U be the event that the diameter is larger than ε for any $j \in [d]$. Let V be the event that F represents bounded space. The event U is the one we need to show has small probability.

By Theorem I.2, we know that if $n \ge L_u(a, \delta, d)$ then **F** is unbounded with probability smaller than $\delta/2$. Therefore, $P(V) \ge 1 - \delta/2$.

From Lemma I.4 and the analysis before this theorem statement, we know that for n as required by this theorem statement, any time \mathbf{F} represents an bounded half-space intersection, the probability of the diameter this intersection (for any j) being larger than ϵ is smaller than $\delta/2$. Therefore, $P(U \mid V) \leq \delta/2$.

Hence,

$$P(U) = P(U \mid V) p(V) + P(U \mid \neg V) (1 - P(V)) \le P(U \mid V) + (1 - P(V)) \le \delta/2 + \delta/2 = \delta.$$

Table 5: Prediction errors (RMSE) as p (the fraction of nonnegative entries in A) increases (n = 512, d = 16). Each set of rows describes an experiment with a different number of scale vector rows in A (s).

| s | Alg. | p = 0 | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | 1.0 |
|----|---------------|--------|--------|--------|---------|----------|----------|----------|----------|----------|----------|----------|
| 0 | SGD | 0.4353 | 0.4135 | 0.3844 | 0.3511 | 0.3364 | 0.3255 | 0.3310 | 0.3480 | 0.3792 | 0.4158 | 0.4417 |
| | LP | 0.0923 | 0.0991 | 5.2210 | 18.3364 | ∞ | ∞ | ∞ | ∞ | ∞ | ∞ | 55.4519 |
| | \mathbf{QP} | 0.0492 | 0.0583 | 0.5445 | 10.0328 | ∞ | ∞ | ∞ | ∞ | ∞ | ∞ | 0.4129 |
| 2 | SGD | 0.4252 | 0.4058 | 0.3777 | 0.3510 | 0.3315 | 0.3216 | 0.3293 | 0.3447 | 0.3704 | 0.4060 | 0.4234 |
| | \mathbf{QP} | 0.0729 | 0.0819 | 0.5072 | 7.8984 | 9.8862 | ∞ | ∞ | ∞ | ∞ | 5.8306 | 0.4297 |
| 4 | SGD | 0.4152 | 0.3962 | 0.3711 | 0.3423 | 0.3286 | 0.3213 | 0.3228 | 0.3382 | 0.3683 | 0.3944 | 0.4118 |
| | \mathbf{QP} | 0.0926 | 0.1070 | 0.4480 | 5.6844 | ∞ | ∞ | ∞ | ∞ | ∞ | 12.0391 | ∞ |
| 8 | SGD | 0.3941 | 0.3823 | 0.3637 | 0.3410 | 0.3289 | 0.3248 | 0.3256 | 0.3428 | 0.3578 | 0.3743 | 0.3874 |
| | \mathbf{QP} | 0.1221 | 0.1527 | 0.2835 | 1.3792 | 8.4902 | 7.2903 | 7.4253 | ∞ | 5.9887 | 8.2727 | 0.5928 |
| 16 | SGD | 0.3820 | 0.3829 | 0.3827 | 0.3812 | 0.3825 | 0.3825 | 0.3827 | 0.3815 | 0.3831 | 0.3818 | 0.3827 |
| | QP | 0.1810 | 0.1841 | 0.1853 | 0.1823 | 0.1833 | 0.1845 | 0.1823 | 0.1823 | 0.1825 | 0.1840 | 0.1803 |

J Nonnegativity of A and the Scale Transformation Condition

Throughout the paper, we assume that the first layer, the matrix \mathbf{A} , has nonnegative entries. As we mention before in Section 8 and Section 7.2, such restrictions on neural network weights have successfully been studied before and shown to be effective in practice.

While the theory of our algorithm does require \mathbf{A} to be nonnegative, we use this condition in a specific part of the proof of Theorem E.1 in Appendix F.1. This theorem is only necessary to prove that the objective functional for layer 2 achieves the correct estimate of the network. It is not necessary for proofs for layer 1 or for the proofs of the transition from the layer objectives into the QP or LP formulation. In that sense, we can aim to alleviate this nonnegativity condition quite in a modular way when addressing this issue. For example, we could parametrize the number of negative entries in \mathbf{A} and make more nuanced arguments (albeit more complex) in the proofs in Appendix F.1.

This condition does not prevent us from running the algorithm as-is with negative entries in **A**. To test our algorithm sensitivity to this assumption, we repeated the experiment in Section 7.1, with d = 16 and n = 512, and checked what error the algorithm gives as a function of $p \in \{k/10 \mid k \in \{0, ..., 10\}\}$ – a fraction of random entries that are set to be negative in **A**. To test the effect of Condition 2.1 (which is alleviated in Appendix E), we also vary the number of rows that are a scaling vector (through parameter $s \in \{0, 2, 4, 8, 16\}$, which denotes the number of such rows in **A**). We repeated each experiment five times with different parameters (and each experiment includes multiple executions with different samples). We experimented both with the LP and the QP objectives for s = 0 and the QP objective for s > 0, and compared them against the SGD algorithm.

The results are given in Table 5. When the assumption about nonnegative entries holds, both the QP algorithm and the LP algorithm (for s = 0) give an error much lower than SGD. This holds even for a small p > 0. However, as p increases, the results of LP and QP degrade quickly and then recover (for the QP) when p = 1.0. This latter result could be due to the symmetry of our input distribution. we also note that in general, the QP algorithm, as expected, is more robust to p > 0 than the LP algorithm. As we increase the number of rows in **A** that are scale vectors, we see that our algorithm becomes much more robust to the number of negative entries in the matrix **A** (with respect to ∞ entries). However, having scale vectors in **A** does seem to make the learning problem more difficult for the QP algorithm.