

A Library of Mirrors: Deep Neural Nets in Low Dimensions are Convex Lasso Models with Reflection Features

Emi Zeger*, Yifei Wang*, Aaron Mishkin†, Tolga Ergen*
Emmanuel Candès‡, and Mert Pilanci*
Stanford University

Key words. deep neural networks, convex optimization, LASSO, sparsity

MSC codes. 62M45, 46N10

Abstract. We prove that training certain classes of neural networks on 1-D data is equivalent to solving convex Lasso problems with discrete, explicitly defined dictionary matrices. We consider neural networks with piecewise linear activations and depths ranging from 2 to an arbitrary but finite number of layers. We first show that two-layer networks with piecewise linear activations are equivalent to Lasso models using a discrete dictionary of ramp functions, with breakpoints corresponding to the training data points. In certain general architectures with absolute value or ReLU activations, a third layer surprisingly creates features that reflect the training data about themselves. Additional layers in a *deep narrow network* progressively generate reflections of these reflections. The Lasso representation provides valuable insights into the analysis of globally optimal networks, elucidating their solution landscapes and enabling closed-form solutions in certain minimal regularization cases. Numerical results show that reflections also occur when optimizing standard deep networks using standard non-convex optimizers. Additionally, we demonstrate our theory with autoregressive time series models.

1. Introduction. Training deep neural networks is an important optimization problem. However, the non-convexity of neural nets makes their training challenging. In this paper, we show that for low-dimensional data, e.g., 1-D or 2-D, training a deep neural network can be simplified to solving a convex Lasso problem with an easily constructed dictionary matrix.

Neural networks are used as predictive models for low-dimensional data in applications such as processing time-series data, especially acoustic signals, as well as in “coordinate-based” multi-layer perceptrons (MLP)s (23). They can also be used to predict time series financial data (Section 4). In acoustic signal processing, raw audio data from a single microphone represents a 1-D time series (34; 36) (although it can also be represented as higher-dimensional spectrograms). Coordinate-based MLPs are a class of models used in computer vision and graphics that input low-dimensional positional data, typically 3-D coordinates representing a single location, and output representations of properties such as shape or color at that location (18; 37; 23). Low-dimensional techniques used in coordinate-based MLPs such as Fourier representation features (38) are also useful in physics-informed machine learning problems (23). Neural networks are also used on data assumed to have a low-dimensional structure such as using Bayesian inference to learn low-dimensional representations (9) and uncertainty quantification (42; 43). There has been recent study on the problem of two-layer ReLU neural

*Department of Electrical Engineering

†Department of Computer Science

‡Departments of Statistics and Mathematics

39 networks on 1-D data (21; 22; 26). It has been proved that an optimal two-layer ReLU neural
 40 network precisely interpolates 1-D training data as a piecewise linear function for which the
 41 breakpoints are at the data points (35; 13; 14).

42 However, even for low-dimensional data, the current literature still lacks analysis on the
 43 expressive power and learning capabilities of deeper neural networks with generic activations.
 44 This motivates us to study the optimization of 2 and 3-layer networks with piecewise linear
 45 activations and deeper neural networks with sign and ReLU activations. For 1-D data, we
 46 simplify the training problem by recasting it as a convex Lasso problem, which has been
 47 extensively studied (10; 39; 40).

48 Convex analysis of neural networks was developed in several prior works. As an example,
 49 infinite-width neural networks enable the convexification of the overall model (3; 4; 16). However,
 50 due to the infinite-width assumption, these results do not reflect finite-width neural networks
 51 in practice. Recently, a series of papers (12; 13; 33) developed a convex analytic framework
 52 for the training problem of two-layer neural networks with ReLU activation. As a follow-up
 53 work, a similar approach is used to formulate the training problem for threshold activations
 54 with data in general d -dimensions as a Lasso problem (11). However, the dictionary matrix is
 55 described implicitly and requires high computational complexity to create (11). By focusing on
 56 1-D data, we can provide simple, explicit Lasso dictionaries and consider additional activations.
 57 For example, we analyze networks with sign activation, which is useful in contexts such as
 58 saving memory to meet hardware constraints (7; 24).

59 Throughout the paper, all scalar functions extend to vector and matrix inputs component-
 60 wise. We write vectors as $\mathbf{v} = (v_1, \dots, v_n)$ and denote the set of n -dimensional, real-valued
 61 column and row vectors by \mathbb{R}^n and $\mathbb{R}^{1 \times n}$, respectively. For $L \geq 2$, an L -layer *neural network* for
 62 d -dimensional data is denoted by $f_L(\mathbf{x}; \theta) : \mathbb{R}^{1 \times d} \rightarrow \mathbb{R}$, where $\mathbf{x} \in \mathbb{R}^{1 \times d}$ is an input row vector
 63 and θ is the *parameter set*. The set θ may contain matrices, vectors, and scalars representing
 64 weights and biases. We let $\theta \in \Theta$, where Θ is the *parameter space*. Let $\mathbf{X} \in \mathbb{R}^{N \times d}$ be a *data*
 65 *matrix* consisting of N *training samples* $\mathbf{x}_1, \dots, \mathbf{x}_N \in \mathbb{R}^{1 \times d}$. We consider regression tasks and
 66 call $\mathbf{y} \in \mathbb{R}^N$ the *target vector*. The (non-convex) neural net *training problem* is

$$67 \quad (1.1) \quad \min_{\theta \in \Theta} \frac{1}{2} \|f_L(\mathbf{X}; \theta) - \mathbf{y}\|_2^2 + \frac{\beta}{\tilde{L}} \|\theta_w\|^{\tilde{L}}$$

68 where $\beta > 0$ is a regularization coefficient for a subset of parameters $\theta_w \subset \theta$ that incur a weight
 69 penalty when training. We denote $\|\theta_w\|^{\tilde{L}} = \sum_{q \in \theta_w} \|q\|_2^{\tilde{L}}$, which penalizes the total network
 70 weight. The results can be generalized to other p -norm regularizations as well. \tilde{L} is the *effective*
 71 *regularized depth*, defined to be the usual depth L for ReLU, leaky ReLU, and absolute value
 72 activations, but defined to be just 2 for threshold and sign activations. This is motivated by
 73 the property that neurons with sign or threshold activations are invariant to the magnitude of
 74 their neuron weights, so it does not make sense to penalize the inner weights (Remark D.1). A
 75 central element of this paper is the *Lasso problem*

$$76 \quad (1.2) \quad \min_{\mathbf{z}, \xi} \frac{1}{2} \|\mathbf{A}\mathbf{z} + \xi \mathbf{1} - \mathbf{y}\|_2^2 + \tilde{\beta} \|\mathbf{z}\|_1$$

77 where \mathbf{z} is a vector, $\xi \in \mathbb{R}$, $\mathbf{1}$ is a vector of ones, and $\tilde{\beta} > 0$ (whose relation to β in (1.1) is

78 defined in Section 3). \mathbf{A} is called the *dictionary matrix*, with columns $\mathbf{A}_i \in \mathbb{R}^N$.

79 A neural net is *trained* by searching for θ that optimizes (1.1). A neural net $f_L(\mathbf{x}; \theta)$ is
80 called *optimal* if θ is a global minimizer of (1.1). Unfortunately, training is complicated by the
81 non-convexity of the optimization problem. However, for data of dimension $d=1$, we reformulate
82 the training problem (1.1) into the equivalent but simpler Lasso problem (1.2), where \mathbf{A} is
83 a fixed matrix that is constructed based on the training data \mathbf{X} and neural net architecture.
84 Moreover, the dictionary matrix columns \mathbf{A}_i correspond to piecewise linear functions $f_i : \mathbb{R} \rightarrow \mathbb{R}$
85 we call *features* that satisfy the following:

- 86 1. For all $n = 1, \dots, N$, $f_i(x_n) = \mathbf{A}_{n,i}$.
- 87 2. If \mathbf{z}, ξ are Lasso solutions, then $f(x) = \sum_{i=1}^n z_i f_i(x) + \xi$ is equal to an optimal network.

88

89 The set of features for a given network is a *dictionary*. We call a collection of dictionaries a
90 *library* when referring to the features from multiple depths. We explicitly provide the elements
91 of \mathbf{A} , making it straightforward to build and solve the convex Lasso problem instead of solving
92 the non-convex training problem. This reformulation allows for exploiting fast Lasso solvers
93 based on the proximal gradient method and Least Angle Regression (LARS) (10), (31). The
94 Lasso representation also elucidates the solution path of neural networks. The *solution path*
95 for the Lasso or training problem is the map from $\beta \in (0, \infty)$ to the solution set. The Lasso
96 solution path is well understood (39; 40; 10; 31), providing insight into the solution path of the
97 ReLU training problem (28).

98 Whereas in the training problem (1.1), the quality of the neural net fit to the data is measured
99 by the l_2 loss as $\frac{1}{2} \|f_L(\mathbf{X}; \theta) - \mathbf{y}\|_2^2$, our results generalize to a wide class of convex loss functions
100 $\mathcal{L}_{\mathbf{y}} : \mathbb{R}^N \rightarrow \mathbb{R}$. With a general loss function, (1.1) becomes $\min_{\theta \in \Theta} \mathcal{L}_{\mathbf{y}}(f_L(\mathbf{X}; \theta)) + \frac{\beta}{L} \|\theta_w\|^{\tilde{L}}$. This
101 is shown to be equivalent to the generalization of (1.2), namely $\min_{\mathbf{z}, \xi} \mathcal{L}_{\mathbf{y}}(\mathbf{A}\mathbf{z} + \xi\mathbf{1} - \mathbf{y}) + \tilde{\beta} \|\mathbf{z}\|_1$.
102 The Lasso problem selects solutions \mathbf{z} that generalize well by penalizing their total weight in l_1
103 norm (39). The l_1 norm typically selects a small number of elements in \mathbf{z} to be nonzero. The
104 Lasso equivalence demonstrates that neural networks can learn a sparse representation of the
105 data by selecting certain features to fit \mathbf{y} .

106 This paper is organized as follows. Section 2 defines various neural network architectures.
107 Section 3 describes our main theoretical result: neural networks are solutions to Lasso problems.
108 One important consequence of the Lasso equivalence is that deep networks with ReLU and
109 absolute value activations learn geometric reflections in the data, where the *reflection* of scalars
110 a about b is $R_{(a,b)} = 2b - a$. The next sections describe applications of the Lasso equivalence.
111 Section 4 uses the Lasso formulation to improve training of neural networks that predict financial
112 time-series. Appendix E examines the relationship between the entire set of optimal neural
113 nets given by the training problem versus the Lasso problem, while Appendix C applies the
114 Lasso model to examine neural net behavior under minimum regularization, finding closed-form
115 solutions. Appendix D presents experiments that support our theoretical results, and shows
116 examples where neural networks trained with Adam naturally exhibit Lasso features.

117 **1.1. Contributions.** Our contributions can be summarized as follows:

- 118 • Training neural networks with symmetrized or deep narrow architectures on 1-D data is
119 equivalent to solving Lasso problems with finite, explicit and fixed dictionaries of basis
120 signals that grow richer with depth (Theorems 3.11, 3.6). We identify dictionaries for

- 121 various architectures in closed form (Lemma B.2).
- 122 • Features with reflections of training data appear in libraries for certain 3-layer architec-
- 123 tures with ReLU (Theorem 3.5) or deep narrow networks with absolute value activation
- 124 (Theorem 3.2). Experimentally, training these networks using the Adam optimizer leads
- 125 to the same reflection features that we prove and matches our theoretical results on the
- 126 global optima (Appendix D).
- 127 • After depth 3, the libraries freeze for deep narrow networks with ReLU activation that
- 128 have the same number of biased neurons in the middle and final layers (Lemma 3.7).
- 129 But, the libraries grow when the width expands to twice as many neurons in the
- 130 middle layer as the final layer for networks with ReLU activation, enabling greater
- 131 representation power (Theorem 3.11).

132 **1.2. Notation.** Assume 1-D training data is ordered as $x_1 > x_2 > \dots > x_N$. The indicator

133 function of a logical statement z is $\mathbf{1}\{z\}$. Let $[n] = \{1, 2, \dots, n\}$. The number of nonzero

134 elements in a vector \mathbf{z} is $\|\mathbf{z}\|_0$. Let $\mathbf{1}, \mathbf{0} \in \mathbb{R}^N$ be the all-ones and all-zeros vectors, respectively.

135 A network that uses ReLU activation is a "ReLU network" or "ReLU-activated network," and

136 a feature in a Lasso problem that is equivalent to a ReLU network is a "ReLU feature." Similar

137 terminology holds for other activations and architectures.

138 **2. Neural net architectures.** This section is devoted to defining neural net terminology

139 and notation to be used throughout the rest of the paper. Let $L \geq 2$ be the depth of a neural

140 network (which has $L-1$ hidden layers). We assume the activation $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ is *piecewise*

141 *linear around 0*, i.e., of the form

142

$$143 \quad (2.1) \quad \sigma(x) = \begin{cases} a^-x + b^- & \text{if } x < 0 \\ a^+x + b^+ & \text{otherwise} \end{cases}$$

144 for some $a^-, a^+, b^-, b^+ \in \mathbb{R}$. As shorthand, "piecewise linear" will mean "piecewise linear

145 around 0." We focus on the piecewise linear activations of ReLU, leaky ReLU, absolute value,

146 sign, and threshold functions. The ReLU activation is $\sigma(x) = (x)_+ := \max\{x, 0\}$, and absolute

147 value activation is $\sigma(x) = |x|$. The leaky ReLU generalizes ReLU and absolute value as

148 $\sigma(x) = (a^+\mathbf{1}\{x > 0\} + a^-\mathbf{1}\{x < 0\})x$ where $a^+ \neq a^-$. ReLU, leaky ReLU and absolute value

149 activations will be referred to as "continuous piecewise linear." The threshold activation is

150 $\sigma(x) = \mathbf{1}\{x \geq 0\}$, and the sign activation is $\sigma(x) = \text{sign}(x)$, where $\text{sign}(x)$ is -1 if $x < 0$, and 1

151 if $x \geq 0$. Note $\text{sign}(0) = 1$.

152 For $\mathbf{Z} \in \mathbb{R}^{n \times m}$, $\mathbf{s} \in \mathbb{R}^m$, let $\sigma_{\mathbf{s}}(\mathbf{Z}) = \sigma(\mathbf{Z})\text{Diag}(\mathbf{s})$. When the columns of $\sigma(\mathbf{Z})$ are neuron

153 outputs, the i^{th} column of $\sigma_{\mathbf{s}}(\mathbf{Z}) \in \mathbb{R}^{N \times m}$ represents a neuron scaled by an *amplitude parameter*

154 $s_i \in \mathbb{R}$. Amplitude parameters are (trainable) parameters for only the sign and threshold

155 activations, and are to be ignored by interpreting them as 1 for ReLU, leaky ReLU, and absolute

156 value activations.

157 Next, we define some neural net architectures. The parameter set is partitioned into

158 $\theta = \theta_w \cup \theta_b \cup \{\xi\}$, where θ_b is a set of *internal bias* terms, and ξ is an *external* bias term. We

159 define the elements of each parameter set below. We define neural nets by their output on row

160 vectors $\mathbf{x} \in \mathbb{R}^{1 \times d}$. Their outputs then extend to matrix inputs $\mathbf{X} \in \mathbb{R}^{N \times d}$ row-wise.

161 **2.1. Standard networks.** The following is a commonly studied neural net architecture. Let
 162 $L \geq 2$, the number of layers. Let $m_1 = d$ and $m_l \in \mathbb{N}$ for $l \in [L] - \{1\}$ which are the number of
 163 neurons in each layer. For $l \in [L - 1]$, let $\mathbf{W}^{(l)} \in \mathbb{R}^{m_l \times m_{l+1}}$, $\mathbf{s}^{(l)} \in \mathbb{R}^{m_{l+1}}$, $\mathbf{b}^{(l)} \in \mathbb{R}^{1 \times m_{l+1}}$, $\xi \in \mathbb{R}$,
 164 which denote weights, amplitude parameters, internal biases, and external bias, respectively.
 165 Let $\mathbf{X}^{(1)} = \mathbf{x} \in \mathbb{R}^{1 \times d}$ be the input to the neural net and $\mathbf{X}^{(l+1)} \in \mathbb{R}^{1 \times m_{l+1}}$ be viewed as the
 166 inputs to layer $l + 1$, defined by

$$167 \quad (2.2) \quad \mathbf{X}^{(l+1)} = \sigma_{\mathbf{s}^{(l)}} \left(\mathbf{X}^{(l)} \mathbf{W}^{(l)} + \mathbf{b}^{(l)} \right).$$

168 Let $\boldsymbol{\alpha} \in \mathbb{R}^{m_L}$, which is the vector of final layer coefficients. A *standard neural network* is
 169 $f_L(\mathbf{x}; \theta) = \xi + \mathbf{X}^{(L)} \boldsymbol{\alpha}$. The regularized and bias parameter sets are
 170 $\theta_w = \{ \boldsymbol{\alpha}, \mathbf{W}^{(l)}, \mathbf{s}^{(l)} : l \in [L-1] \}$ and $\theta_b = \{ \mathbf{b}^{(l)} : l \in [L-1] \}$, respectively.

171 Analyzing the training problem for standard networks has traditionally been challenging
 172 due to provably non-zero duality gaps (46). However, parallel architectures can be examined
 173 through the lens of convex analysis, and they can also be converted into standard architectures,
 174 which significantly streamlines the analytical process. Furthermore, employing parallel branches
 175 of neural networks has proven especially effective in enhancing modern mixture-of-experts
 176 structures. By changing the architecture to a parallel structure defined next, we show that the
 177 training problem simplifies to the Lasso problem. These alternative architectures allow neural
 178 nets to be reconstructed more tractably from a Lasso solution than with a standard network.

179 **2.2. Parallel networks.** A parallel network is a linear combination of certain standard
 180 networks in parallel, as we now define. Each standard network is called a *parallel unit*. Let
 181 $L \geq 2$, $m_0 = d$, $m_{L-1} = 1$ and $m_l \in \mathbb{N}$ for $l \in [L] - \{L-1\}$. For $i \in [m_L]$, $l \in [L-1]$,
 182 let $\mathbf{W}^{(i,l)} \in \mathbb{R}^{m_{l-1} \times m_l}$, $\mathbf{s}^{(i,l)} \in \mathbb{R}^{m_l}$, $\mathbf{b}^{(i,l)} \in \mathbb{R}^{1 \times m_l}$, $\xi \in \mathbb{R}$, which are the weights, amplitude
 183 parameters, and biases of the i^{th} parallel unit.¹ Let $\hat{\mathbf{X}}^{(i,1)} = \mathbf{x} \in \mathbb{R}^{1 \times d}$ be the input to the
 184 neural net and $\hat{\mathbf{X}}^{(i,l+1)} \in \mathbb{R}^{1 \times m_l}$ be viewed as the input to layer $l + 1$ in unit i , defined by

$$185 \quad (2.3) \quad \hat{\mathbf{X}}^{(i,l+1)} = \sigma_{\mathbf{s}^{(i,l)}} \left(\hat{\mathbf{X}}^{(i,l)} \mathbf{W}^{(i,l)} + \mathbf{b}^{(i,l)} \right).$$

186 Let $\boldsymbol{\alpha} \in \mathbb{R}^{m_L}$. A *parallel neural network* is $f_L(\mathbf{x}; \theta) = \xi + \sum_{i=1}^{m_L} \hat{\mathbf{X}}^{(i,L)} \alpha_i$. In the parallel
 187 architecture, m_L is the number of neurons in the final layer and for $i \in [m_L]$, we define the
 188 disjoint unions $\theta_w = \bigcup_{i \in [m_L]} \theta_w^{(i)}$ and $\theta_b = \bigcup_{i \in [m_L]} \theta_b^{(i)}$. The regularized and bias parameter sets
 189 are $\theta_w^{(i)} = \{ \alpha_i, \mathbf{s}^{(i,l)}, \mathbf{W}^{(i,l)} : l \in [L-1] \}$, $\theta_b^{(i)} = \{ \mathbf{b}^{(i,l)} : l \in [L-1] \}$, for $i \in [m_L]$. We view
 190 parallel units as functions $\hat{\mathbf{X}}^{(i,L)} : \mathbb{R}^{1 \times d} \rightarrow \mathbb{R}$ and with abuse of notation write $\hat{\mathbf{X}}^{(i,L)}(\mathbf{x}) \in \mathbb{R}$
 191 as the output of $\hat{\mathbf{X}}^{(i,L)}$ evaluated on a sample \mathbf{x} . For a training dataset $\mathbf{X} \in \mathbb{R}^{N \times d}$, we
 192 denote the evaluations of the functions $\hat{\mathbf{X}}^{(i,L)}$ on the training data as $\hat{\mathbf{X}}^{(i,L)}(\mathbf{X})$. The condition
 193 that $m_{L-1} = 1$ is required for parallel networks to output a scalar. Unlike parallel networks,

¹While we define the first dimension of $\mathbf{W}^{(1)}$ to be $m_1 = d$ for standard networks, it is $m_0 = d$ for parallel networks. This offset is due to the requirement that the second dimension of $\mathbf{W}^{(i,L-1)}$ must be 1 rather than m_L as for standard networks.

194 standard networks do not require $m_{L-1} = 1$. Each parallel unit of a parallel network (with
 195 $m_L \geq 1$) is a special case of a standard network (with $m_L = 1$). Therefore a parallel network
 196 with just one parallel unit is a special case of a standard network (with $m_L = 1$.) This paper
 197 primarily focuses on parallel architectures. However, a parallel network can be converted into a
 198 standard network, and a 2-layer network is both a standard and parallel network, as shown in
 199 [Appendix A](#).

200 **3. Main results.** In this section, we show that for 1-D data and certain architectures,
 201 non-convex deep neural net training problems are *equivalent* to Lasso problems, that is, their
 202 optimal values are the same, and given a Lasso solution, we can reconstruct a neural net that is
 203 optimal in the training problem. We say that a neural network *model* is equivalent to a Lasso
 204 model to mean that the optimal models are convertible to each other.

205 Unless otherwise stated, for the rest of this paper assume the data is 1-D. Define $\tilde{\beta} = \frac{\beta}{2}$ in
 206 (1.1),(1.2) for 3-layer symmetrized networks (defined below). For all other networks, let $\tilde{\beta} = \beta$.

207 We focus on 2-layer networks with piecewise linear activations and deep networks with
 208 continuous piecewise linear activations. A *deep narrow network* is a parallel network where
 209 the number of neurons in each parallel path is equal to 1, i.e., $m_1 = \dots = m_{L-1} = 1$. For a 2-layer
 210 network, the parallel and deep narrow networks are the same as the standard network, as
 211 shown in [Appendix A.2](#). A deep narrow network has m_L parallel units. Each l^{th} parallel
 212 unit has $L - 1$ non-linear hidden layers: $l = 1, \dots, L - 1$. Each l^{th} hidden layer has a single
 213 neuron with a weight $\mathbf{W}^{(i,l)} \in \mathbb{R}$ and bias $\mathbf{b}^{(i,l)} \in \mathbb{R}$. Therefore the total number of neurons
 214 is m_L units $\times (L - 1)$ layers per unit $\times (1)$ neuron per layer $= m_L(L - 1)$. For example, a
 215 2-layer network has m_L neurons.

216 A *symmetrized 3-layer network* is a parallel network with continuous piecewise linear
 217 activation where $m_1 = 2$ (which means each i^{th} unit out of m_3 parallel units has 2-D weight
 218 vectors $\mathbf{W}^{(i,1)}$ and $\mathbf{W}^{(i,2)}$ for the first and second layer, respectively) and the parameter space
 219 Θ enforces the constraint $|\mathbf{W}_1^{(i,l)}| = |\mathbf{W}_2^{(i,l)}|$ for $l \in [2], i \in [m_3]$ where $\mathbf{W}_j^{(i,l)}$ is the j^{th} element
 220 of $\mathbf{W}^{(i,l)}$, which is a column/row vector. Therefore a symmetrized 3-layer network has $2m_3$
 221 neurons in the "middle" layer. A symmetrized network extends the expressibility of a deep
 222 narrow network by expanding its width. For other architectures, see the supplementary material
 223 (47).

224 In this section, we focus on the architectures discussed above to derive explicit, simple
 225 features that provide intuition towards understanding the representation power of neural
 226 networks. We reformulate the training problem for 2-layer networks and for deeper depths with
 227 absolute value, ReLU and leaky ReLU activations into a convex Lasso problem. In general,
 228 our convexification approach and proofs provide a framework to analyze neural networks with
 229 more arbitrary widths. We begin for with convex formulations for networks with absolute value
 230 activation, as their convex models are relatively simple. Then we extend these formulations
 231 for the traditional ReLU activation. Finally, we generalize to other various piecewise linear
 232 activations that have more complex convex formulations by formally defining a library of
 233 features. Proofs are deferred to [Appendix D.3](#).

234 **3.1. Absolute value activation.** We first discuss networks with absolute value activation,
 235 as the symmetry of $|x|$ significantly simplifies the features. However, absolute value activation

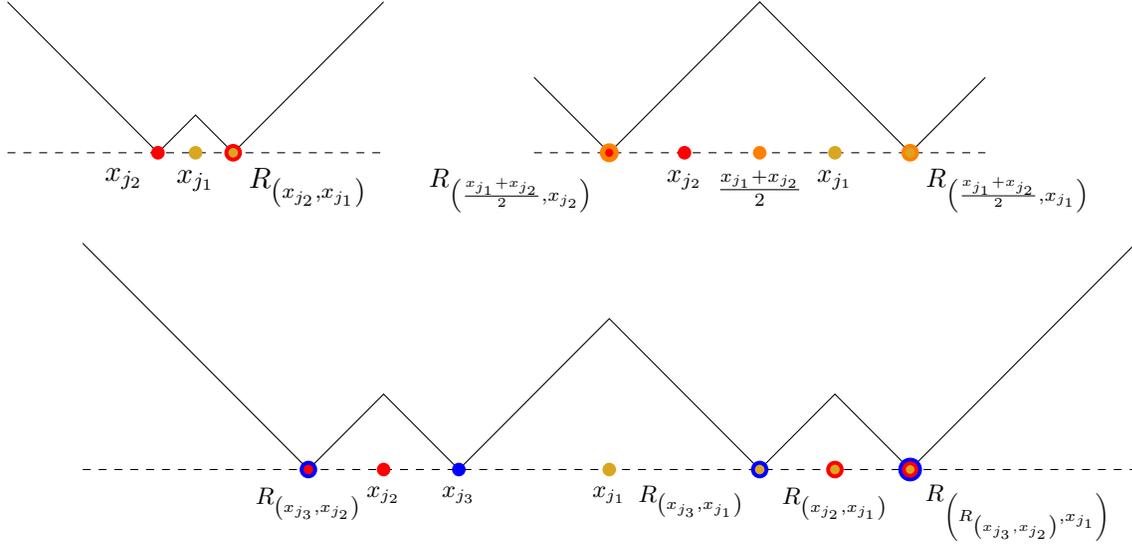


Figure 3.1: Example features are drawn as defined in Definition 3.10, not including reversed directions, for deep narrow networks with absolute value activation. The features drawn are $\hat{\mathbf{X}}^{(3)}(x) = ||x - x_{j_1}| - |x_{j_2} - x_{j_1}||$ where $x_{j_1} = 3, x_{j_2} = 2$ (top left) and $x_{j_1} = 8, x_{j_2} = 4$ (top right) for 3 layers and $\hat{\mathbf{X}}^{(3)}(x) = ||x - x_{j_1}| - |x_{j_2} - x_{j_1}| - ||x_{j_3} - x_{j_1}| - |x_{j_2} - x_{j_1}||$ where $x_{j_1} = 7, x_{j_2} = 4, x_{j_3} = 5$, for 4 layers. Top row: 3-layer features. The top left feature contains a breakpoint at the reflection of x_{j_2} (red) across x_{j_1} (gold), which is denoted as $R(x_{j_2}, x_{j_1})$ (red encircling gold). Other breakpoints are colored similarly. The top right feature has breakpoints at averages of training data and their reflections about other training data. Bottom row: an example of a 4-layer feature, which contains a double reflection of x_{j_3} (blue) reflected across x_{j_2} (red), then reflected across x_{j_1} (gold), which is denoted as $R(R(x_{j_3}, x_{j_2}), x_{j_1})$ (blue encircling red, encircling gold). All lines have slopes ± 1 , and $x_{j_1}, x_{j_2}, x_{j_3}$ are training data.

236 models for two-layer networks are equivalent to ReLU activation models, as long as a skip
 237 connection is present. A 2-layer network with a *skip connection* is $f_L^{\text{skip}}(x; \theta) = f_L(x; \theta) + \omega x$
 238 if $x \in \mathbb{R}$ (or more generally, $f_L^{\text{skip}}(\mathbf{x}; \theta) = f_L(\mathbf{x}; \theta) + \mathbf{x}\omega$) where $\omega \in \mathbb{R}^d$ is a trainable parameter
 239 in θ .

240 **Lemma 3.1.** *The training problem for a 2-layer network with skip connection and ReLU*
 241 *activation remains equivalent if the activation is changed to absolute value, and there is a map*
 242 *between the solutions for either activation.*

243 By Lemma 3.1, the absolute value activation is of interest to analyze as it can map to ReLU
 244 when skip connections are incorporated. Next we define some terms used to state our results.

245 For a piecewise linear function $f : \mathbb{R} \rightarrow \mathbb{R}$, x is a *breakpoint* of f if f changes slope or is
 246 discontinuous at x . A breakpoint is a "kink" in the graph of f . For $a, b, c \in \mathbb{R}$, *double reflection*
 247 is a reflection of a reflection, i.e., is of the form $R(R(a, c), b)$ or $R(b, R(a, c))$. The *generalized*

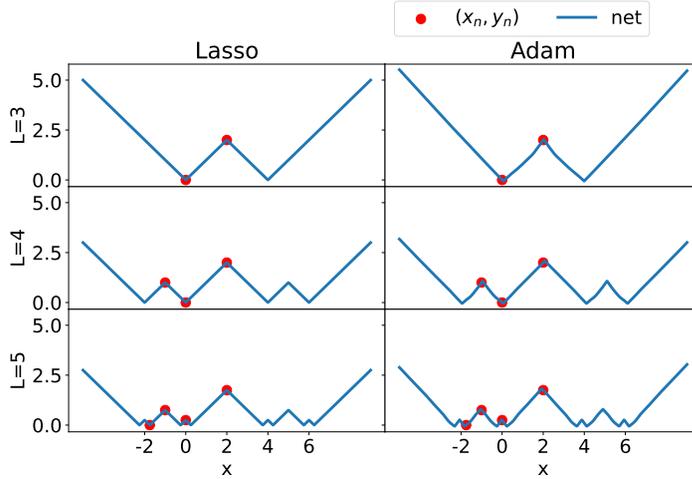


Figure 3.2: Lasso and Adam-trained deep narrow networks with absolute value activation. For $L=3$, the breakpoint at 4 is not a training point; it is the reflection of $x_2=0$ across $x_1=2$. For $L=4$, the breakpoint at 6 is not a training point; it is $x_2=0$ reflected about $x_3=-1$ to -2 (which not a training point) and then reflected across $x_1=2$. Similarly the 5-layer network contains more complex reflections.

248 reflection of a and c about b is $a + c - b = R_{(b, \frac{a+c}{2})}$, the reflection of b about the average of a
 249 and c . When $a = c$, the generalized reflection of a and c about b is the reflection of a about b . A
 250 breakpoint that is of the form $R_{(x_{j_1}, x_{j_2})}$ for training data x_{j_1}, x_{j_2} is called a *reflection breakpoint*
 251 and a feature with a reflection breakpoint is called a *reflection feature*; and similarly for double
 252 reflections. Networks with absolute value activation can be modeled as Lasso problems with
 253 reflection features, as stated next.

254 **Theorem 3.2 (Lasso equivalent of deep absolute value networks).** *A deep narrow network of*
 255 *arbitrary depth with $\sigma(x) = |x|$ is equivalent² to a Lasso model with a finite set of features. Its*
 256 *dictionary matrix for 2 layers is $\mathbf{A}_{i,j} = |x_i - x_j|$. For 3 and 4 layers, its library includes features*
 257 *whose i^{th} element is $||x_i - x_{j_1}| - |x_{j_2} - x_{j_1}||$ and $|||x_i - x_{j_1}| - |x_{j_2} - x_{j_1}|| - ||x_{j_3} - x_{j_1}| - |x_{j_2} - x_{j_1}|||$,*
 258 *respectively, over all training samples $x_i, x_{j_1}, x_{j_2}, x_{j_3}$. A similar pattern holds for deeper*
 259 *networks. The 2-layer features have breakpoints exactly at training data. The libraries for 3*
 260 *and 4 layers additionally include reflection and double reflection features, respectively.*

261 **Theorem 3.2** implies that an absolute value network learns to model data with a discrete and
 262 fixed dictionary of features. **Figure 3.1** plots these features for $L=3$ and $L=4$. It illustrates that
 263 the breakpoints of the $L=3$ features occur at training data reflections. The bottom row plots
 264 a possible feature for 4 layers, which shows breakpoints at reflections and double reflections.
 265 Even with just one neuron per layer (per path for parallel networks), adding a third layer in
 266 neural networks with absolute value activation adds features to the dictionary, and adds new

²See Section 3 for the definition of model equivalence.

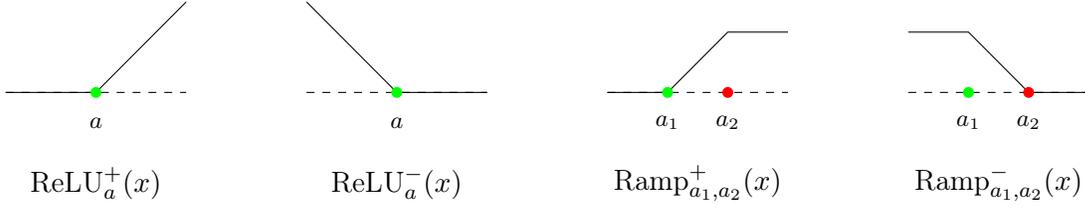


Figure 3.3: Examples of capped ramp functions in Definition 3.3. If $a, a_1, a_2 \in \{x_1, \dots, x_N\}$, these functions are Lasso features for 3-layer deep narrow networks with ReLU activation.

267 locations of breakpoint to those features, namely at reflections $R_{(x_i, x_j)}$ of data points about
 268 themselves. Adding a fourth layer creates double reflections. In Figure 3.2, standard 3, 4, and
 269 5-layer networks are trained with Adam. We make $\beta=10^{-7}$ close to zero, and also solve the
 270 Lasso problem as $\beta \rightarrow 0$ (Appendix C). The Lasso solutions are used to initialize a small subset
 271 of the neurons in the non-convex training. The Adam-trained networks closely match the Lasso
 272 solutions. Moreover, the 5-layer networks suggest that features continue to gain more complex
 273 reflections with depth. Simulation details are given in Appendix D.

274 Our approach lays a foundation which future work can use to enumerate complete feature
 275 libraries for $L \geq 4$ layers. Reflection features allow neural networks to fit functions with
 276 breakpoints at locations in between data points. The reflection breakpoints for absolute
 277 value networks suggest that they can learn geometric structures or symmetries from the data.
 278 Moreover, as the depth increases, this dictionary expands, which deepens its representation
 279 power.

280 **3.2. ReLU activation.** With intuition from absolute value networks, we next discuss ReLU
 281 networks. Since reflection features appear in 3-layer networks with absolute value activation
 282 and $|x| = \frac{(x)_+ + (-x)_+}{2}$, we might expect reflections to also appear in 3-layer ReLU networks with
 283 twice as many neurons in the middle layer as the absolute value network. This is indeed
 284 the case, as shown below. First, we define parameterized families of scalar input, real-valued
 285 functions that will be used to describe features for ReLU networks.

286 **Definition 3.3.** Let $a_1 \in [-\infty, \infty), a_2 \in (-\infty, \infty]$. The capped ramp functions are

$$287 \quad \text{Ramp}_{a_1, a_2}^+(x) = \begin{cases} 0 & \text{if } x \leq a_1 \\ x - a_1 & \text{if } a_1 \leq x \leq a_2 \\ a_2 - a_1 & \text{if } x \geq a_2 \end{cases}, \quad \text{Ramp}_{a_1, a_2}^-(x) = \begin{cases} a_2 - a_1 & \text{if } x \leq a_1 \\ a_2 - x & \text{if } a_1 \leq x \leq a_2 \\ 0 & \text{if } x \geq a_2 \end{cases}$$

288 provided that $a_1 \leq a_2$, and otherwise $\text{Ramp}_{a_1, a_2}^+ = \text{Ramp}_{a_1, a_2}^- = 0$. In particular, the ramp
 289 functions are $\text{ReLU}_a^+(x) = \text{Ramp}_{a, \infty}^+ = (x - a)_+$ and $\text{ReLU}_a^-(x) = \text{Ramp}_{-\infty, a}^- = (a - x)_+$.

290

291 In Definition 3.3, the parameters a, a_1, a_2 are the breakpoints of ramp and capped ramp
 292 functions. This is illustrated in Figure 3.3. Using these features, we state our results on Lasso
 293 equivalence for ReLU networks.

294 **Theorem 3.4 (deep narrow ReLU network representation capability stagnates).** *A deep narrow*
 295 *network of arbitrary depth with ReLU activation is equivalent to a Lasso model with a finite*
 296 *set of features. Its library contains only ramps and capped ramps, and beyond 3 layers, ReLU*
 297 *features do not change as the network deepens.*

298 In contrast to absolute value activation, for deep narrow networks, the ReLU library never
 299 gains reflection features. However, a symmetrized ReLU architecture creates reflections.

300 **Theorem 3.5 (wider ReLU networks do not stagnate and generate reflections).** *A three-layer*
 301 *symmetrized network with ReLU activation is equivalent to a Lasso model with a finite set of*
 302 *features, including those with breakpoints at reflections.*

303 **Figure 3.4** plots ReLU features. A 2-layer ReLU network learns ramp features. Extending
 304 the depth to one more layer without changing the width creates a deep narrow 3-layer network,
 305 which additionally learns capped ramp features. Extending both the depth and width creates a
 306 3-layer symmetrized network, which learns an additional host of features, including generalized
 307 reflection features shown in the bottom row of the figure. The generalized reflections $x_j+x_k-x_i$
 308 are reflections when $j=k$. Reflection features appear in Adam-trained ReLU networks as
 309 predicted by the Lasso formulation, as shown in **Figure 3.5**. Details of this simulation are given
 310 in **Appendix D**.

311 In addition to ReLU and absolute value networks, leaky ReLU networks are also equivalent
 312 to Lasso models. Moreover, we can upper bound the size of the libraries.

313 **Theorem 3.6.** *A deep narrow network of any depth $L \geq 2$ and piecewise linear activation is*
 314 *equivalent to a Lasso problem with a finite set of features. The number of features is $O(N^2)$ for*
 315 *ReLU activation and $O(N^{L-1}2^L L!)$ for leaky ReLU and absolute value activations.*

316 In particular, the number of features is finite and at most polynomial in the number of training
 317 points and exponential in the depth. However, the number of features is an overestimate and
 318 in fact saturates for certain activations and architectures, as seen in the next result.

319 **Lemma 3.7.** *Training a deep narrow ReLU network with an arbitrary number of layers*
 320 *($L \geq 2$) is a Lasso problem where features are ReLU or capped ramp functions with breakpoints*
 321 *at data points. The number of features is $O(N^2)$.*

322 **Lemma 3.7** is a consequence of **Theorem 3.6** and **Theorem 3.4** but is stated to emphasize
 323 the property that the number of features is capped for deep narrow ReLU networks. Having
 324 only one neuron per layer limits the expressibility of ReLU networks. In contrast to absolute
 325 value, ReLU networks rely more heavily on wider layers for expressibility.

326 **3.3. General piecewise linear activations.** In the next definition, recall a^+, a^- (2.1) are
 327 the positive and negative slopes for an activation σ .

328 **Definition 3.8.** *For $\alpha, \beta, a^+, a^- \in \mathbb{R}$ such that $a^+ \neq a^-$, let the normalized midpoint be*

$$329 \quad (3.1) \quad m_{\alpha, \beta} = -\frac{a^+ \max\{\alpha, \beta\} - a^- \min\{\alpha, \beta\}}{a^+ - a^-}.$$

330 Note $m_{\alpha, \beta}$ is the midpoint $\frac{\alpha+\beta}{2}$ between α and β if $\sigma(x)=|x|$. If $\sigma(x)=(x)_+$, or if $\sigma(x)=|x|$ and
 331 $\alpha=\beta$, then $m_{\alpha, \beta}=\alpha$. We use the normalized midpoint to define bias parameters for features.

| depth L | activation σ | features | reflections? | $\xi=0?$ | number of features | results |
|-----------|----------------------|-------------------------------------|---------------------|----------|---|--------------------------------|
| 2 | ReLU, leaky ReLU | $\sigma(x-x_n),$ $\sigma(x_n-x)$ | no | yes | $2N$ | Theorem 3.11 Corollary 3.15 |
| 2 | absolute value | $\sigma(x-x_n)$ | no | yes | N | Theorem 3.11 Corollary 3.15 |
| 2 | threshold | $\sigma(x-x_n),$ $\sigma(x_n-x)$ | no | no | $2N$ | Theorem 3.11 Corollary 3.15 |
| 2 | sign | $\sigma(x-x_n)$ | no | no | N | Theorem 3.11 Corollary 3.15 |
| 3+ | ReLU | deep library | no | yes | $O(N^2)$ | Lemma 3.7 |
| 3 | ReLU, symmetrized | includes deep library | yes, generalized | yes | $O(2^{2(L-1)}N^3)$ for deep library | Theorem 3.5 Theorem 3.13 |
| 3 | absolute value | deep library | yes | yes | $O(N^2)$ | Theorem 3.2, Theorem 3.11 |
| 4+ | absolute value | includes deep library | yes, multilevel | yes | $O(N^{L-1}L!)$ for deep library | Theorem 3.2, Theorem 3.6 |

Table 3.1: Properties of equivalent Lasso problems for different architectures. Multilevel reflections include reflections of reflections. By Definition 3.10, the size of the deep library is $O(2^L N^{L-1})$. For absolute value, there is no 2^L since all weights can be 1.

332 Recall that $\hat{\mathbf{X}}^{(1,L)}$ (2.3) denotes a parallel unit, which is a special case of a standard network
333 (Subsection 2.2). Let $\hat{\mathbf{X}}^{(l)} = \hat{\mathbf{X}}^{(1,l)}$, $\mathbf{W}^{(l)} = \mathbf{W}^{(1,l)}$, and $\mathbf{b}^{(l)} = \mathbf{b}^{(1,l)}$. The output $\hat{\mathbf{X}}^{(l)}$ of each layer
334 can be interpreted as a feature extracted by that layer. This motivates the following definition.

335 **Definition 3.9.** A bias parameter $\mathbf{b}^{(l)}$ is a data feature bias if $\mathbf{b}^{(l)} = -\hat{\mathbf{X}}^{(l)}(\mathbf{x}_0)\mathbf{W}^{(l)}$ for some
336 column vector $\mathbf{x}_0 \in \{x_1, \dots, x_N\}^{m_l}$. The bias parameter $\mathbf{b}^{(1)}$ is a first-layer midpoint feature
337 bias if $\mathbf{b}^{(1)} = -m_{x_{j_1}, x_{j'_1}} \mathbf{W}^{(1)}$ for some $j_1, j'_1 \in [N]$ and it is in a 3-layer deep narrow network
338 with a non-monotone, continuous piecewise linear activation (such as $\sigma(x) = |x|$.)

339 These bias parameters are used to define a deep library.

340 **Definition 3.10 (Deep Library).** Consider a 3-layer symmetrized or L -layer deep narrow
341 network. A feature is a function of the form $\hat{\mathbf{X}}^{(L)}(x)$ with data or midpoint feature biases and
342 all elements of $\mathbf{W}^{(1)}$ being ± 1 . The deep library is the set of output vectors $\hat{\mathbf{X}}^{(L)}(\mathbf{X}) \in \mathbb{R}^N$ of
343 all possible features on the training data.

344 The deep library contains a finite number of standard networks evaluated on the training
345 data. The next result shows that a neural network learns features in the deep library.

346 **Theorem 3.11 (complete Lasso libraries for general activations).** Consider a deep narrow
347 L -layer network where $L \in \{2, 3\}$ and with activation σ which is ReLU, leaky ReLU, absolute
348 value, sign, or threshold if $L=2$, and ReLU, leaky ReLU or absolute value if $L=3$. Consider

349 a Lasso problem whose dictionary is the deep library and where $\xi=0$ if σ is sign or threshold.
 350 Suppose (\mathbf{z}^*, ξ^*) is a solution, and let $m^* = \|\mathbf{z}^*\|_0$. This Lasso problem is equivalent to the
 351 training problem for the network, provided $m_L \geq m^*$.

352 The critical number of neurons satisfies $m^* \leq 2N$ for 2-layer networks and in general is at
 353 most the number of features. The notion of equivalence between optimization problems is
 354 that they have the same optimal value and their solutions can be found from each other.
 355 **Definition B.3** defines a map to reconstruct a neural network from a Lasso solution, and
 356 **Lemma B.4** shows that this is optimal. The general reconstructed network is equivalent to the
 357 function $f(x) = \sum_i z_i^* f_i(x) + \xi^*$ where $f_i(x) = \hat{\mathbf{X}}^{(i,L)}(x)$ is the feature corresponding to the
 358 i^{th} column \mathbf{A}_i of the dictionary. In the simplest case when $L=2$, a feature is of the form
 359 $f_i(x) = \sigma(x\mathbf{W}^{(i,1)} + \mathbf{b}^{(i,1)}) = \sigma((x - x_n)\mathbf{W}^{(i,1)})$ with a data feature bias $\mathbf{b}^{(i,1)} = -\mathbf{W}^{(i,1)}x_n$
 360 where $\mathbf{W}^{(1)} = \pm 1$, and $f_i(\mathbf{X})$ is a vector in the deep library. **Lemma B.2** explicitly describes
 361 features for various networks. While **Theorem 3.2** describes a subset of the library for $L \geq 3$
 362 layers, **Theorem 3.11** defines the full library for 3 layers.

363 **Remark 3.12.** *The dictionary for an architecture discussed in **Theorem 3.11** is a superset of*
 364 *any dictionary with the same architecture but shallower depth.*

365 The next theorem generalizes **Theorem 3.5** to leaky ReLU activations.

366 **Theorem 3.13.** *The training problem for a 3-layer symmetrized network with monotone*
 367 *activations such as ReLU is equivalent to a Lasso problem with solution (\mathbf{z}^*, ξ^*) and whose*
 368 *dictionary contains a deep library with $O(2^{2(L-1)}N^3)$ features, provided $m_L \geq m^*$, where*
 369 *$m^* = \|\mathbf{z}^*\|_0$.*

370 **Theorem 3.13** states that the deep library is a sub-dictionary for symmetrized networks.
 371 Finding the full dictionary for a symmetrized network is an area of future work. **Theorem 3.11**
 372 shows that instead of training a neural network with a non-convex problem and reaching a
 373 possibly local optimum, we can simply solve a straightforward Lasso problem whose convexity
 374 guarantees that gradient descent approaches global optimality. In previous work (11), a similar
 375 Lasso formulation is developed for networks with threshold activation but requires up to 2^N
 376 features of length N in the dictionary for a 2-layer network. In contrast, with 1-D data,
 377 **Theorem 3.11** shows that at most $2N^2$ features are needed for a 2-layer network.

378 **Remark 3.14.** *Note that when the network is 2 layers, the equivalent Lasso dictionary only*
 379 *contains features with breakpoints at training data, leading to a prediction with breakpoints*
 380 *only at data locations. In contrast, when the network has 3 layers there can be breakpoints at*
 381 ***reflections** of data points with respect to other data points due to the reflection features. As a*
 382 *result, for activations such as absolute value and ReLU with symmetrized networks, the sequence*
 383 *of dictionaries as the network gets deeper converges to a richer library that includes reflections.*

384 Our approach lays the foundation to further analyze the evolution of feature libraries over
 385 expanding depth and widths as an area of future work. For 2-layer networks, the dictionary
 386 (**Theorem 3.11**) is simple, as described next.

387 **Corollary 3.15 (2-layer libraries).** *Let $\mathbf{A}_+, \mathbf{A}_- \in \mathbb{R}^{N \times N}$ with $(\mathbf{A}_+)_{i,n} = \sigma(x_i - x_n)$, $(\mathbf{A}_-)_{i,n} =$*
 388 *$\sigma(x_n - x_i)$. We can write the dictionary matrix for 2-layer networks as $\mathbf{A} = \mathbf{A}_+$ for absolute*
 389 *value and sign activations, and $\mathbf{A} = [\mathbf{A}_+, \mathbf{A}_-] \in \mathbb{R}^{N \times 2N}$ for ReLU, leaky ReLU, and threshold*

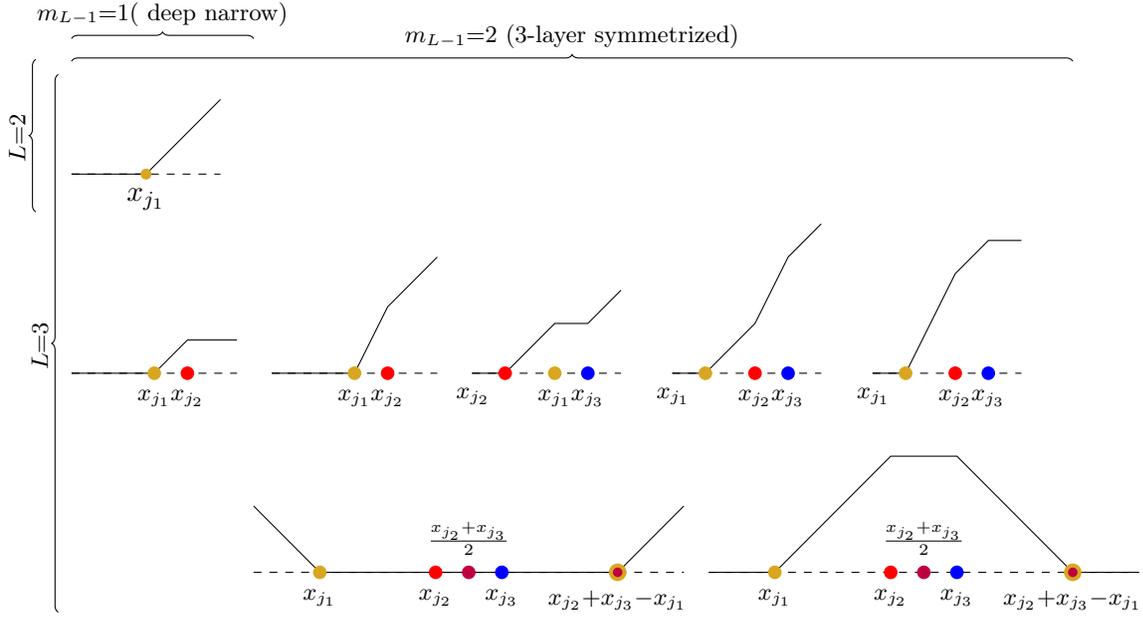


Figure 3.4: Example ReLU features, excluding reverse directions. Bottom row: 3-layer symmetrized ReLU features with breakpoints at generalized reflections of x_{j_1} (gold) across x_{j_2} (red) and x_{j_3} (blue) $x_{j_2+x_{j_3}-x_{j_1}}$ (gold encircling purple). Lines have slopes $\pm 2, \pm 1$ or 0.

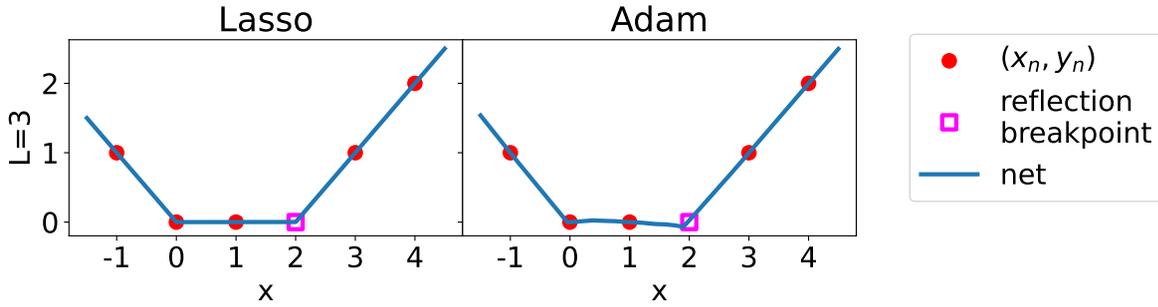


Figure 3.5: Lasso and Adam-trained 3-layer symmetrized ReLU networks. Most crucially, the breakpoint at 2 is not a training point; it is the reflection of $x_4 = 0$ across $x_3 = 1$.

390 activations.

391 In Corollary 3.15, \mathbf{A}_+ and \mathbf{A}_- contain features $\hat{\mathbf{X}}^{(2)}(\mathbf{X})$ where $\mathbf{W}^{(1)} = 1$ and $\mathbf{W}^{(1)} = -1$,
 392 respectively (Definition 3.10). Figure 3.6 illustrates the features for \mathbf{A}_+ for the ReLU and
 393 threshold activations. Using the notation of Definition 3.10, the 3-layer deep narrow absolute
 394 value features are of the form $\hat{\mathbf{X}}^{(3)}(x) = ||x - x_{j_1}| - |x_{j_2} - x_{j_1}||$, and for 4 layers, include features
 395 of the form $\hat{\mathbf{X}}^{(4)}(x) = |||x - x_{j_1}| - |x_{j_2} - x_{j_1}| - ||x_{j_3} - x_{j_1}| - |x_{j_2} - x_{j_1}|||$. These features are plotted

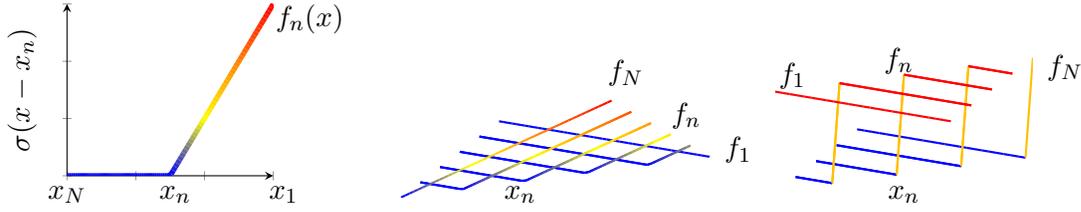


Figure 3.6: The features $f_n(x) = \sigma(x - x_n)$ define the 2-layer dictionary $\mathbf{A}_+ \in \mathbb{R}^{N \times N}$ as $\mathbf{A}_{+,i,n} = f_n(x_i)$. An individual feature $f_n(x)$ for ReLU activation is plotted in the left figure, and it is plotted alongside all other $f_1(x), \dots, f_N(x)$ in the middle figure. The right figure plots $f_1(x), \dots, f_N(x)$ for threshold activation.

396 in Figure 3.1, which highlights their reflection and double reflection breakpoints. Figure 3.1 also
 397 illustrates breakpoints at averages and reflections of averages of training data, which correspond
 398 to midpoint feature biases in Definition 3.9. Figure 3.4 shows a subset of the ReLU library, with
 399 generalized reflection features for 3-layer symmetrized networks. Figure 3.4 illustrates general
 400 feature shapes not including mirrored directions. In Figure B.1, we choose distinct training
 401 samples $x_i, x_j, x_k \in \{-1, 0, 2\}$ and numerically compute all possible deep library features for
 402 3-layer, symmetrized ReLU networks using Definition 3.10. Since the features $\hat{\mathbf{X}}^{(L)}(x)$ are
 403 continuous with respect to x_i, x_j, x_k , the features for non-distinct x_i, x_j, x_k can be extrapolated
 404 by merging adjacent training points. Figure B.1 shows that when exhaustively enumerating over
 405 all possible values of data feature biases and ± 1 weights, there are redundant features, and so
 406 the upper bound on features in Theorem 3.13 can be an overestimate. The numerically plotted
 407 features are consistent with Figure 3.4. In addition to graphical representations, Lemma B.2 in
 408 Appendix B gives examples of simple, explicit expressions for features defined in Definition 3.10.

409 **3.4. Discussion: extensions, limitations and open problems.** Appendix B.1 gives explicit
 410 examples of reconstructed networks and discusses the equivalence of the Lasso model and
 411 reconstructed network on test points. We note that the optimal solutions may not be unique.
 412 Appendix E discusses the relation between the solution sets of the Lasso and non-convex
 413 training problems with respect to a specified straightforward reconstruction. Lemma E.2 shows
 414 that the optimal Lasso fit $\hat{\mathbf{y}} = \mathbf{A}\mathbf{z}^* + \xi^*\mathbf{1}$, the linear fit $\mathbf{A}\mathbf{z}^*$ and the optimal bias ξ^* are
 415 unique, and the set of all optimal Lasso solutions can be characterized by an *equicorrelation*
 416 *set*, as described in Proposition E.3. Figure E.1 analyzes the generalization properties of neural
 417 networks reconstructed from different Lasso solutions, and shows that their performance can
 418 vary widely. Further analysis of all possible reconstruction maps between the models and the
 419 span of all stationary solutions of the non-convex model is an area of future work.

420 There is an analogous result to the convex equivalence in Theorem 3.11 for deep and wide
 421 sign and threshold networks which is provided in the supplementary material (47). Further
 422 extending the Lasso formulations for wider, more general neural network architectures is an
 423 area of future work. A limitation of this work is the computationally large number of features
 424 as the network grows larger and deeper. In practice, given a formulation of Lasso features, one
 425 can subsample features to reduce the computational complexity. Subsampling Lasso problem

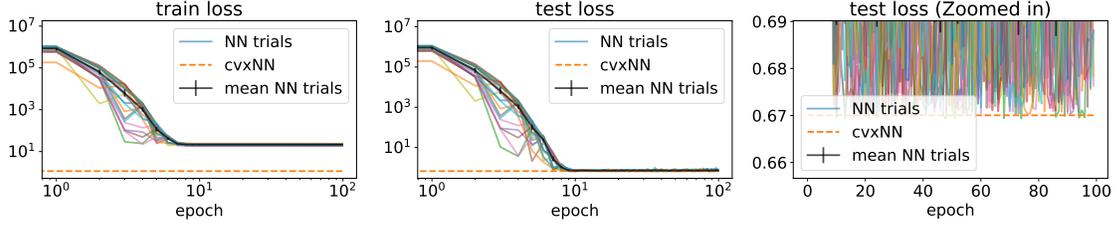


Figure 4.1: Comparison of autoregressive models $x_t = f(x_{t-1}; \theta) + \epsilon_t$ using convex (cvxNN) and non-convex (NN) training with 2-layer neural networks and the linear model AR(1) to forecast minutely 2017 Bitcoin prices (1). The non-linear models outperform the linear AR(1) model. Moreover, SGD generally underperforms in training and test loss compared to the convex model which is guaranteed to find a global optimum of the NN objective. The non-black, solid colored curves plot each NN trial, and the black curve plots their mean, with the vertical lines indicating one standard deviation of the data above and below the mean.

426 representations and their approximation guarantees has been studied in (15; 25; 44; 45) and is
 427 shown to maintain improved performance over training the non-convex network. Future work
 428 can adapt these approaches to the networks discussed in this paper. Regardless of the number
 429 of features, Lasso problems are fully transparent and intuitive as convex models that perform
 430 feature selection, in contrast to non-convex, black-box neural networks. There is extensive
 431 literature on the training and generalization behavior of Lasso models.

432 **4. Application: Time-series modeling.** In this section, we apply the Lasso problem for
 433 neural networks to an autoregression problem. Suppose at times $1, \dots, T + 1$ we observe data
 434 points $x_1, \dots, x_{T+1} \in \mathbb{R}$ that follow the time-series model

$$435 \quad (4.1) \quad x_t = f(x_{t-1}; \theta) + \epsilon_t,$$

436 where $f : \mathbb{R} \rightarrow \mathbb{R}$ is parameterized by some parameter θ and $\epsilon_t \sim \mathcal{N}(0, \sigma^2)$ represents observation
 437 noise. The parameter θ is unknown, and the goal is find θ that best fits the model (4.1) to the
 438 data x_1, \dots, x_{T+1} . For example, the *auto-regressive model with lag 1* (AR(1)) is a linear model
 439 $f(x; \theta) = ax$ where $\theta = \{a\}$ is chosen as a solution to $\min_{\theta \in \Theta} \sum_{t=1}^T (f(x_t; \theta) - x_{t+1})^2$. For a more
 440 expressive model, instead of $f(x; \theta) = ax$ suppose we use a 2-layer neural network

$$441 \quad (4.2) \quad f_2^{\text{NN}}(x; \theta) = \sum_{i=1}^m |xw_i + b_i| \alpha_i.$$

442 The parameter set is $\theta = \{w_i, b_i, \alpha_i\}_{i=1}^m$. Consider a training problem of the form

$$443 \quad (4.3) \quad \min_{\theta \in \Theta} \frac{1}{T} \sum_{t=1}^T \mathcal{L}(f_2^{\text{NN}}(x_t; \theta) - x_{t+1}) + \frac{\beta}{2} \|\theta_w\|_2^2.$$

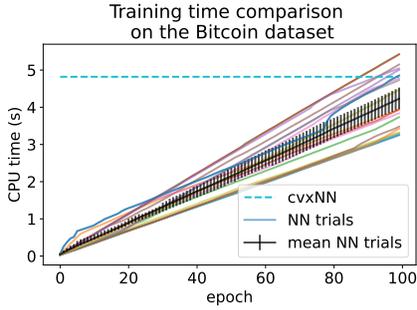


Figure 4.2: Training time for 2-layer neural networks trained on 2017 minutely Bitcoin data in Figure 4.1. The non-black, solid colored curves plot the time for each NN trial, and the black curve plots their mean, with the vertical lines indicating one standard deviation of the data above and below the mean.

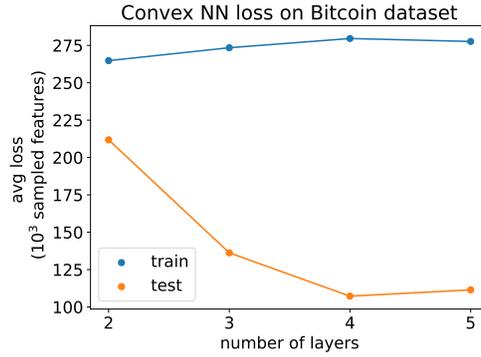


Figure 4.3: Mean squared loss over $T = 10^3$ time samples of Bitcoin data (BTC-2017min) for autoregressive deep narrow networks (one for each depth) trained on subsampled Lasso problems with 10^3 features.

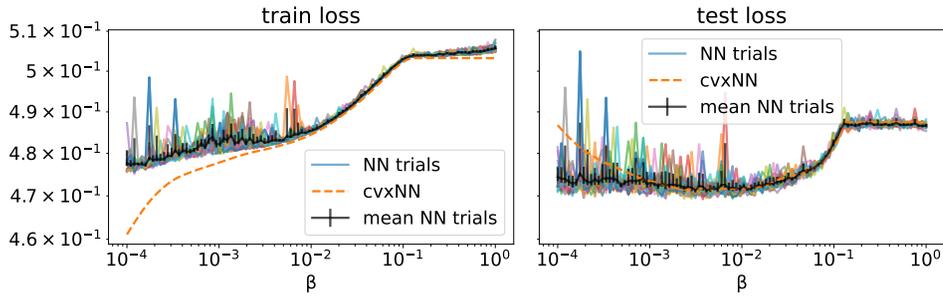


Figure 4.4: The regularization path for planted data with $\sigma^2 = 1$, $p = 5$ planted neurons. The non-black, solid colored curves plot each NN trial, and the black curve plots their mean, with the vertical lines indicating one standard deviation of the data above and below the mean.

444 If $\mathcal{L}(x)=x^2$ then (4.3) is an *autoregression* problem, where the network represents predictors
 445 of x_{t+1} from x_t . Problem (4.3) can be solved by converting it to an equivalent Lasso problem
 446 (1.2) where $A_{i,j}=|x_i-x_j|$ and $y_i = x_{i+1}$. Figure 4.1 shows that using the Lasso model to
 447 predict Bitcoin price reaches a lower loss than training with the non-convex model. In our
 448 experiments, we define the training loss as the objective of the training problem (4.3), which is
 449 the sum of the performance loss and the regularization term $\frac{\beta}{2}\|\theta_w\|_2^2$. We define the test loss as
 450 just the performance loss $\frac{1}{T}\sum_{t=1}^T \mathcal{L}(f_2^{\text{NN}}(x'_t; \theta) - x'_{t+1})$ evaluated on test data x'_t , without the
 451 regularization term. Therefore we note that the training loss can be higher than the test loss,

452 as seen in Figure 4.1. We plot the training loss defined in this way to test our theory that the
 453 training loss is globally optimized using the convex training. We plot the test loss without the
 454 regularization term to show the practical performance of the network on unseen data and how
 455 well the regularization improved generalization. In Figure 4.2, for each trial of training the
 456 non-convex network, we measure the time that each epoch takes and plot the cumulative time.
 457 As shown in Figure 4.2, the non-convex model takes less time to train than the convex model
 458 until it converges to a solution, but the non-convex model converges to a suboptimal network.

459 We also investigate training on planted data. We build a network $f_2^{\text{NN}}(x; \theta)$ (4.2) with p
 460 known, or *planted* neurons. We use this network to generate training samples x_1, \dots, x_{T+1} based
 461 on (4.1) with $f(x; \theta) = f_2^{\text{NN}}(x; \theta)$ where $x_1 \sim \mathcal{N}(0, \sigma^2)$. Using the same model $f_2^{\text{NN}}(x; \theta)$, we
 462 also generate test samples $x_1^{\text{test}}, \dots, x_{T+1}^{\text{test}}$ in an analogous way, independently of the training data.
 463 This allows us to generate independent testing and training sets with the same distributions,
 464 whereas in the real Bitcoin data, the testing samples subsequently follow the training samples
 465 and may still have time-based correlations with the training data. We use $T = 1000$ time
 466 samples. Then, we try to recover the planted neurons based on only the training samples
 467 by solving the NN AR/QR training problems. The *regularization path* is the optimal neural
 468 net’s performance loss as a function of the regularization coefficient β . Figure 4.4 plots the
 469 regularization path for $\sigma^2=1$ in (4.1) and $p=10$ neurons in the planted model. The regularization
 470 path is the optimal loss as a function of the regularization parameter β . Training the convex
 471 model finds an optimal training loss, while training the non-convex model may find a suboptimal
 472 loss. The loss versus β curves taken by the non-convex models, which are the jagged curves in
 473 Figure 4.4, are therefore noisy estimates of the true regularization path taken by the convex
 474 model, which is the smoother curve in Figure 4.4. More details are found in Appendix E.2.

475 In Figure 4.3 and Figure 4.5, deep narrow neural networks of various depths with $\sigma(x) = |x|$
 476 are trained using the convex Lasso model on Bitcoin price data (minutely price in 2017). To
 477 allow a comparison where the training complexity is the same for all networks, each network
 478 uses 10^3 Lasso features that are randomly sampled from the deep library as follows. Since data
 479 feature biases enable reflection features, for simplicity and efficiency, the experiment uses deep
 480 library features with only data feature biases, and omits midpoint features. The experiment
 481 randomly and uniformly samples a data feature bias $\mathbf{b}^{(1)} = -x_{j_1}$ out of $j_1 \in [N]$, which is
 482 used to create a 2-layer Lasso feature. Then another training sample x_{j_2} is randomly selected,
 483 and a 3-layer feature is created with $\mathbf{b}^{(1)} = -x_{j_1}, \mathbf{b}^{(2)} = -\hat{\mathbf{X}}^{(2)}(x_{j_2})$ using the same x_{j_1} as for
 484 the 2-layer feature. Similarly another sample x_{j_3} is randomly selected and a 4-layer feature is
 485 constructed from x_{j_3} and the prior samples with $\mathbf{b}^{(1)} = x_{j_1}, \mathbf{b}^{(2)} = \hat{\mathbf{X}}^{(2)}(x_{j_2}), \mathbf{b}^{(3)} = \hat{\mathbf{X}}^{(3)}(x_{j_3})$.
 486 A similar procedure creates a 5th-layer feature. In this process, the feature selected for layer L
 487 in this process will either be the same feature as for layer $L - 1$ if $j_{L-1} = j_{L-2}$, or be the feature
 488 for $L - 1$ with an additional breakpoint. The experiment repeats this process 10^3 times to
 489 create a single subsampled Lasso training problem with 10^3 sampled features for each network
 490 depth L . Figure 4.3 and Figure 4.5 show the results for a single example network trained from
 491 such a subsampled Lasso. The experiment repeats this sampling procedure 200 times to create
 492 200 subsampled Lasso problems for each depth. Figure 4.6 and Figure 4.7 show results from
 493 the 200 networks trained for each depth.

494 Note that an alternative sampling procedure could be first enumerating all features and

495 then randomly sampling from the entire feature set, which would create independent samples
496 from different depths. However, since enumerating all features can be computationally intensive,
497 the sampling is done in the sequential procedure described above. The dictionaries of the
498 networks are not subsets of each other, and so the Lasso training loss does not necessarily
499 decrease with depth, as seen in Figure 4.3.

500 Figure 4.3 and Figure 4.5 are consistent with Figure 4.6 and Figure 4.7, and show that
501 while deeper networks can have higher training error as shown in Figure 4.6, they generally
502 have better test error, at least up to a certain depth. The test performance of the networks
503 in Figure 4.7 improves until $L = 4$ but slightly worsens for $L = 5$ (not shown in Figure 4.5).
504 Figure 4.6 also shows that while 5-layer networks can perform well on test data depending on
505 the subsampled features, they generally have high variance in test loss and large test errors.
506 The improved test performance with limited depth suggests that certain deeper features, which
507 have richer levels of reflections, tend to generalize data better.

508 This performance loss when $L = 5$ may be due to the network overfitting to overly complex
509 features after the network is already deep enough that its features sufficiently capture the
510 relationships in the data. It may also occur if the 10^3 subsampled features are no longer
511 representative of the entire library, whose size increases with depth: the small fraction of
512 subsampled features may be missing out on key features for the data. The number of subsampled
513 features was set to 10^3 to reduce computational complexity, which is a practical limitation in
514 this experiment. Increasing the number of sampled features may increase the performance of
515 the 5-layer networks.

516 In our experiments, the Lasso-trained deep narrow networks of depth greater than 2 do not
517 outperform deep networks that are trained with the non-convex training problem, unlike the
518 experiment in Figure 4.1 where 2-layer convex networks outperformed non-convex networks.
519 This may be due to insufficient Lasso feature sampling and hyperparameter choices in the
520 non-convex program. A practical limitation affecting these results is insufficient sampling;
521 addressing this issue will be the subject of future work.

522 **5. Conclusion.** Our results show that certain deep neural networks with a variety of
523 activation functions trained on 1-D data with weight regularization can be recast as convex
524 Lasso models with simple dictionary matrices. This provides critical insight into their solution
525 path as the weight regularization changes. The Lasso problem also provides a fast way to train
526 neural networks for 1-D data. Moreover, the understanding of the neural networks through
527 Lasso models could also be used to explore designing better neural network architectures.

528 We proved that reflection features can emerge in the Lasso dictionary for certain absolute
529 value and ReLU networks when the depth is 3 or deeper. This leads to predictions that have
530 breakpoints at reflections of data points about other data points. In contrast, for networks of
531 depth 2, the breakpoints are only located at a subset of training data. We believe that this
532 mechanism enables deep neural networks to generalize to the unseen by encoding a geometric
533 regularity prior.

534 Our analysis of various architectures provides a foundation for studying more complex
535 network topologies. The 1-D results can extend to sufficiently structured or low rank data in
536 higher dimensions. Generalizing to higher dimensions is also an area of future work. Building
537 on a similar theme, (32) showed that the structure of hidden neurons can be expressed through

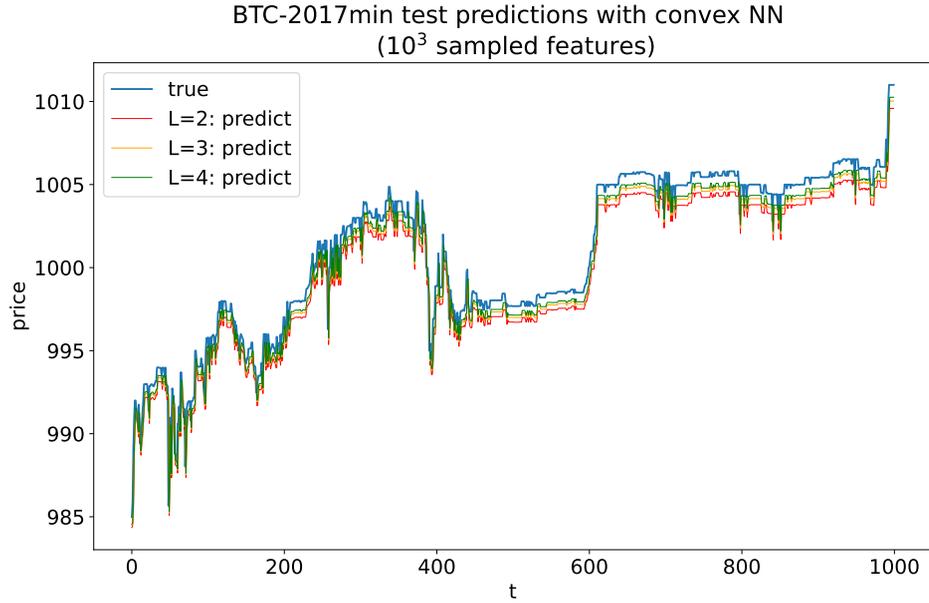


Figure 4.5: Predictions of example autoregressive neural networks (one for each depth) trained using the convex Lasso problem with 10^3 sampled features on predicting $T = 10^3$ BTC-2017min samples. The network predictions improve as the depth increases from $L = 2$ to $L = 4$. For each network of depth L , the regularization is $\beta = \frac{1}{T}$ in the training problem (4.3) for 2 layers and (1.1) for more general L layers.

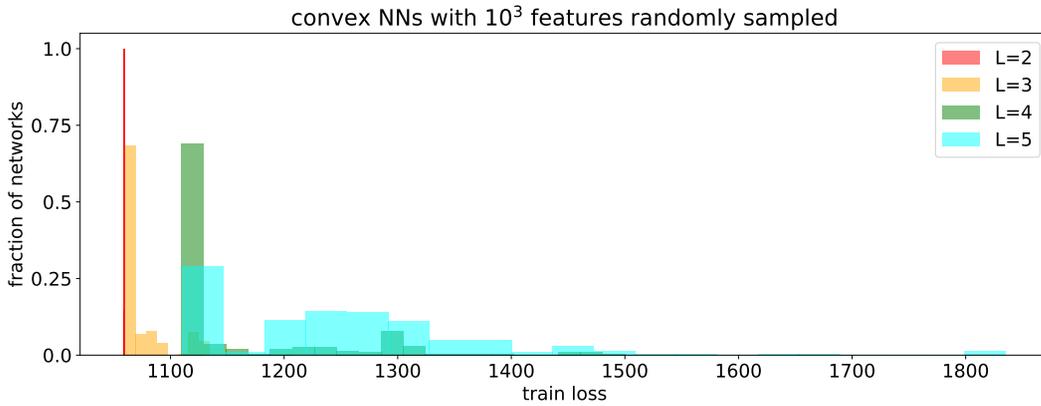


Figure 4.6: Histogram of training loss of deep narrow networks (200 networks for each depth L) trained on the Lasso problem with randomly sampled 10^3 features for autoregression on BTC-2017min with $T = 10^3$ samples. trained on the Lasso problem with randomly sampled 10^3 features for autoregression on BTC-2017min with $T = 10^3$ samples.

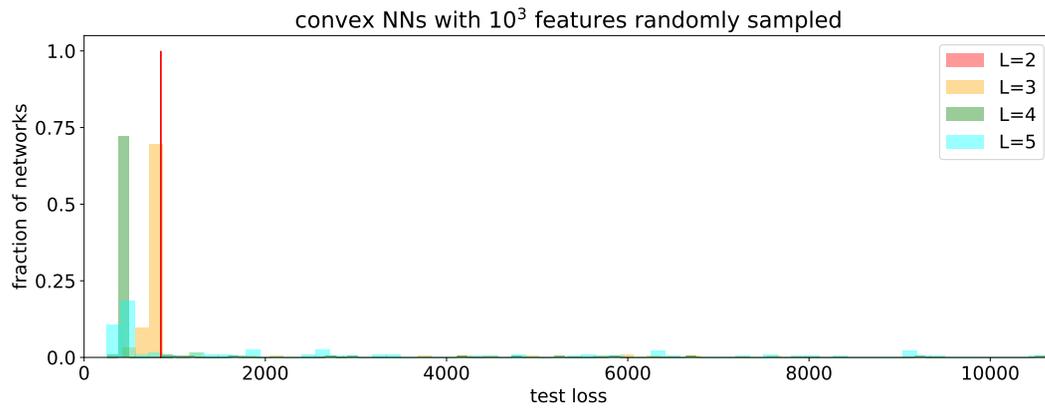


Figure 4.7: Histogram of training loss of deep narrow networks (200 networks for each depth L) trained on the Lasso problem with randomly sampled 10^3 features for autoregression on BTC-2017min with $T = 10^3$ samples. The test loss appears to improve with depth. For visual clarity, the loss is only shown up to a value that contains $\geq 80\%$ of the 4-layer samples, $\geq 90\%$ of the 3-layer samples, and all of the 2-layer samples (using the full dictionary.)

538 convex optimization and Clifford's Geometric Algebra. The techniques developed in this paper
 539 can be combined with the Clifford Algebra to develop higher-dimensional analogues of the
 540 results.

541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584

References.

- [1] Kaggle, <https://www.kaggle.com/datasets/prasoonkottarathil/btcinusd>.
- [2] S. ARORA, N. COHEN, W. HU, AND Y. LUO, *Implicit regularization in deep matrix factorization*, Advances in neural information processing systems, 32 (2019).
- [3] F. BACH, *Breaking the curse of dimensionality with convex neural networks*, JMLR, 18 (2017), pp. 629–681.
- [4] Y. BENGIO, N. ROUX, P. VINCENT, O. DELALLEAU, AND P. MARCOTTE, *Convex neural networks*, in Advances in Neural Information Processing Systems, vol. 18, MIT Press, 2005.
- [5] J. M. BORWEIN AND A. S. LEWIS, *Convex Analysis and Nonlinear Optimization: Theory and Examples*, Springer, 2000.
- [6] S. BOYD AND L. VANDENBERGHE, *Convex optimization*, Cambridge university press, 2004.
- [7] A. BULAT AND G. TZIMIROPOULOS, *XNOR-Net++: Improved binary neural networks*, ArXiv, abs/1909.13863 (2019).
- [8] M.-H. CHEN AND B. W. SCHMEISER, *General hit-and-run monte carlo sampling for evaluating multidimensional integrals*, Operations Research Letters, 19 (1996), pp. 161–169.
- [9] P. CHEN, K. WU, J. CHEN, T. O’LEARY-ROSEBERRY, AND O. GHATTAS, *Projected Stein variational newton: A fast and scalable Bayesian inference method in high dimensions*, Advances in Neural Information Processing Systems, 32 (2019).
- [10] B. EFRON, T. HASTIE, I. JOHNSTONE, AND R. TIBSHIRANI, *Least angle regression*, The Annals of statistics, 32 (2004), pp. 407–499.
- [11] T. ERGEN, H. I. GULLUK, J. LACOTTE, AND M. PILANCI, *Globally optimal training of neural networks with threshold activation functions*, arXiv:2303.03382, (2023).
- [12] T. ERGEN AND M. PILANCI, *Convex geometry of two-layer ReLU networks: Implicit autoencoding and interpretable models*, PMLR, 26–28 Aug. 2020, pp. 4024–4033.
- [13] T. ERGEN AND M. PILANCI, *Convex geometry and duality of over-parameterized neural networks*, The Journal of Machine Learning Research, 22 (2021), pp. 9646–9708.
- [14] T. ERGEN AND M. PILANCI, *Revealing the structure of deep neural networks via convex duality*, in ICML, PMLR, 2021, pp. 3004–3014.
- [15] T. ERGEN AND M. PILANCI, *The convex landscape of neural networks: Characterizing global optima and stationary points via lasso models*, ArXiv, abs/2312.12657 (2023), <https://api.semanticscholar.org/CorpusID:266374487>.
- [16] C. FANG, Y. GU, W. ZHANG, AND T. ZHANG, *Convex formulation of overparameterized deep neural networks*, arXiv:1911.07626, (2019).
- [17] S. FEIZI, H. JAVADI, J. M. ZHANG, AND D. TSE, *Porcupine neural networks: (almost) all local optima are global*, ArXiv, abs/1710.02196 (2017).
- [18] K. GENOVA, F. COLE, D. VLASIC, A. SARNA, W. T. FREEMAN, AND T. A. FUNKHOUSER, *Learning shape templates with structured implicit functions*, 2019 IEEE/CVF International Conference on Computer Vision (ICCV), (2019), pp. 7153–7163, <https://api.semanticscholar.org/CorpusID:118642996>.
- [19] G. GIDEL, F. BACH, AND S. LACOSTE-JULIEN, *Implicit regularization of discrete gradient dynamics in linear neural networks*, Advances in Neural Information Processing Systems, 32 (2019).

- 585 [20] S. GUNASEKAR, B. E. WOODWORTH, S. BHOJANAPALLI, B. NEYSHABUR, AND N. SRE-
586 BRO, *Implicit regularization in matrix factorization*, Advances in neural information
587 processing systems, 30 (2017).
- 588 [21] N. JOSHI, G. VARDI, AND N. SREBRO, *Noisy interpolation learning with shallow univariate*
589 *relu networks*, arXiv.2307.15396, (2023).
- 590 [22] K. KARHADKAR, M. MURRAY, H. TSERAN, AND G. MONTÚFAR, *Mildly overpa-*
591 *rameterized relu networks have a favorable loss landscape*, arXiv:2305.19510, (2023),
592 <https://arxiv.org/abs/2305.19510>.
- 593 [23] G. E. KARNIADAKIS, I. G. KEVREKIDIS, L. LU, P. PERDIKARIS, S. WANG,
594 AND L. YANG, *Physics-informed machine learning*, Nature Reviews Physics, 3
595 (2021), pp. 422–440, <https://doi.org/10.1038/s42254-021-00314-5>, <https://doi.org/10.1038/s42254-021-00314-5>.
- 596 [24] M. KIM AND P. SMARAGDIS, *Bitwise neural networks for efficient single-channel source*
597 *separation*, in 2018 IEEE International Conference on Acoustics, Speech and Signal
598 Processing (ICASSP), 2018, pp. 701–705.
- 600 [25] S. KIM AND M. PILANCI, *Convex relaxations of relu neural networks approximate global*
601 *optima in polynomial time*, ICML, (2024).
- 602 [26] G. KORNOWSKI, G. YEHUDAI, AND O. SHAMIR, *From tempered to benign overfitting in*
603 *ReLU neural networks*, arXiv:2305.15141, (2023), <https://arxiv.org/abs/2305.15141>.
- 604 [27] A. MISHKIN AND M. PILANCI, *The solution path of the group lasso*, 2022, <https://api.semanticscholar.org/CorpusID:259504228>.
- 605 [28] A. MISHKIN AND M. PILANCI, *Optimal sets and solution paths of ReLU networks*, in
606 International Conference on Machine Learning, ICML 2023, PMLR, 2023.
- 607 [29] A. MISHKIN AND M. PILANCI, *Optimal sets and solution paths of relu networks*, in Pro-
608 ceedings of the 40th International Conference on Machine Learning, ICML’23, JMLR.org,
609 2023.
- 610 [30] B. NEYSHABUR, *Implicit regularization in deep learning*, arXiv preprint arXiv:1709.01953,
611 (2017).
- 612 [31] M. R. OSBORNE, B. PRESNELL, AND B. A. TURLACH, *A new approach to variable*
613 *selection in least squares problems*, IMA journal of numerical analysis, 20 (2000), pp. 389–
614 403.
- 615 [32] M. PILANCI, *From complexity to clarity: Analytical expressions of deep neural network*
616 *weights via Clifford’s geometric algebra and convexity*, arXiv:2309.16512, (2023).
- 617 [33] M. PILANCI AND T. ERGEN, *Neural networks are convex regularizers: Exact polynomial-*
618 *time convex optimization formulations for two-layer networks*, in Proceedings of the 37th
619 International Conference on Machine Learning, vol. 119, 13–18 July 2020, pp. 7695–7705.
- 620 [34] H. PURWINS, B. LI, T. VIRTANEN, J. SCHLÜTER, S.-Y. CHANG, AND T. SAINATH, *Deep*
621 *learning for audio signal processing*, IEEE Journal of Selected Topics in Signal Processing,
622 13 (2019), pp. 206–219.
- 623 [35] P. SAVARESE, I. EVRON, D. SOUDRY, AND N. SREBRO, *How do infinite width bounded*
624 *norm networks look in function space?*, Annual Conference on Learning Theory, (2019),
625 pp. 2667–2690.
- 626 [36] J. SERRÀ, S. PASCUAL, AND C. S. PERALES, *Blow: a single-scale hyperconditioned flow*
627 *for non-parallel raw-audio voice conversion*, Advances in Neural Information Processing
628

- 629 Systems, 32 (2019).
- 630 [37] K. O. STANLEY, *Compositional pattern producing networks: A novel abstraction of*
631 *development*, Genetic Programming and Evolvable Machines, 8 (2007), p. 131–162, <https://doi.org/10.1007/s10710-007-9028-8>, <https://doi.org/10.1007/s10710-007-9028-8>.
- 632
- 633 [38] M. TANCIK, P. P. SRINIVASAN, B. MILDENHALL, S. FRIDOVICH-KEIL, N. RAGHAVAN,
634 U. SINGHAL, R. RAMAMOORTHI, J. T. BARRON, AND R. NG, *Fourier features let*
635 *networks learn high frequency functions in low dimensional domains*, in Proceedings of the
636 34th International Conference on Neural Information Processing Systems, NIPS '20, Red
637 Hook, NY, USA, 2020, Curran Associates Inc.
- 638 [39] R. TIBSHIRANI, *Regression shrinkage and selection via the Lasso*, Journal of the Royal
639 Statistical Society: Series B (Methodological), 58 (1996), pp. 267–288, [https://arxiv.org/](https://arxiv.org/abs/https://rss.onlinelibrary.wiley.com/doi/pdf/10.1111/j.2517-6161.1996.tb02080.x)
640 [abs/https://rss.onlinelibrary.wiley.com/doi/pdf/10.1111/j.2517-6161.1996.tb02080.x](https://arxiv.org/abs/https://rss.onlinelibrary.wiley.com/doi/pdf/10.1111/j.2517-6161.1996.tb02080.x).
- 641 [40] R. J. TIBSHIRANI, *The lasso problem and uniqueness*, Electronic Journal of Statistics, 7
642 (2013), pp. 1456–1490.
- 643 [41] S. VAITER, C. DELEDALLE, G. PEYRÉ, J. FADILI, AND C. DOSSAL, *The degrees of*
644 *freedom of the group lasso for a general design*, CoRR, abs/1212.6478 (2012).
- 645 [42] Y. WANG, P. CHEN, AND W. LI, *Projected Wasserstein gradient descent for high-*
646 *dimensional Bayesian inference*, SIAM/ASA Journal on Uncertainty Quantification, 10
647 (2022), pp. 1513–1532.
- 648 [43] Y. WANG, P. CHEN, M. PILANCI, AND W. LI, *Optimal neural network approximation of*
649 *Wasserstein gradient direction via convex optimization*, arXiv:2205.13098, (2022).
- 650 [44] Y. WANG, S. KIM, P. CHU, I. SUBRAMANIAM, AND M. PILANCI, *Randomized geometric*
651 *algebra methods for convex neural networks*, ArXiv, abs/2406.02806 (2024), [https://api.](https://api.semanticscholar.org/CorpusID:270258553)
652 [semanticscholar.org/CorpusID:270258553](https://api.semanticscholar.org/CorpusID:270258553).
- 653 [45] Y. WANG, J. LACOTTE, AND M. PILANCI, *The hidden convex optimization landscape*
654 *of regularized two-layer relu networks: an exact characterization of optimal solutions*, in
655 International Conference on Learning Representations, 2021.
- 656 [46] T. E. YIFEI WANG AND M. PILANCI, *Parallel deep neural networks have zero duality*
657 *gap*, ICLR, (2023s).
- 658 [47] E. ZEGER, Y. WANG, T. E. AARON MISHKIN, E. CAND'ES, AND M. PILANCI, *A library*
659 *of mirrors (technical report)*, arXiv preprint:2403.01046, (2024).

660 **Appendix A. Detailed results for Section 2.**

661 **A.1. Parallel to standard architecture conversion.** Given a parallel network, a standard
 662 network can be constructed as follows. Let $\mathbf{W}^{(1)} = [\mathbf{W}^{(1,1)} \dots \mathbf{W}^{(m_1,1)}]$. For $l \geq 1$, let
 663 $\mathbf{b}^{(l)} = (\mathbf{b}^{(1,l)} \dots \mathbf{b}^{(m_l,l)})$. For $l > 1$, let $\mathbf{W}^{(l)} = \text{blockdiag}(\mathbf{W}^{(1,l)} \dots \mathbf{W}^{(m_l,l)})$. And let α, ξ be
 664 the same in the standard network as the parallel one.

665 **A.2. 2-layer network equivalences.** A 2-layer network is a standard network $f_2(\mathbf{x}; \theta) =$
 666 $\xi + \mathbf{X}^{(2)}\alpha$ by setting $L = 2, m_1 = d$ in (2.2) as $\mathbf{X}^{(2)} = \sigma(\mathbf{X}^{(1)}\mathbf{W}^{(1)} + \mathbf{b}^{(1)})$ where $\mathbf{W}^{(1)} \in$
 667 $\mathbb{R}^{m_1 \times m_2}, \mathbf{b}^{(1)} \in \mathbb{R}^{1 \times m_2}$. In other words, $\mathbf{W}^{(1)}$ contains $\mathbf{W}^{(i,1)}$ as columns, and $\mathbf{b}^{(1)}$ contains
 668 $\mathbf{b}^{(i,1)}$ as elements.

669 A 2-layer network is a special case of a parallel network $f_2(\mathbf{x}; \theta) = \xi + \sum_{i=1}^{m_2} \hat{\mathbf{X}}^{(i,2)}\alpha_i$ by set-
 670 ting $L = 2, m_0 = d, m_1 = m_{L-1} = 1$ in (2.3) where $\hat{\mathbf{X}}^{(i,2)} = \sigma(\hat{\mathbf{X}}^{(i,1)}\mathbf{W}^{(i,1)} + \mathbf{b}^{(i,1)})$ for
 671 $\mathbf{W}^{(i,1)} \in \mathbb{R}^{m_0 \times m_1}, \mathbf{b}^{(i,1)} \in \mathbb{R}^{1 \times m_1}$. A 2-layer network is also a special case of a deep narrow
 672 network where $m_1 = m_{L-1} = 1$.

673

674 **Appendix B. Detailed results for Section 3 .** Proofs are defered to [Appendix D.3](#).

675 **Definition B.1.** Define the function $\mathcal{W}_{\alpha,\beta} : \mathbb{R} \rightarrow \mathbb{R}$ parameterized by $\alpha, \beta \in \mathbb{R}$ as

$$676 \quad \mathcal{W}_{\alpha,\beta}(x) = \begin{cases} \alpha - x & \text{if } x \leq \alpha \\ x - \alpha & \text{if } \alpha \leq x \leq \beta \\ R_{(\alpha,\beta)} - x & \text{if } \beta \leq x \leq R_{(\alpha,\beta)} \\ x - R_{(\alpha,\beta)} & \text{if } x \geq R_{(\alpha,\beta)}. \end{cases}$$

677 [Figure 3.4](#) plots all ReLU features $\hat{\mathbf{X}}^{(L)}(x)$ for an example dataset. As shown in the figure,
 678 symmetrized ReLU network features contain generalized reflections of the form $x_i + x_j - x_k$. We
 679 now describe features for other activations.

680 **Lemma B.2 (Examples of the Deep Library).** Consider a deep narrow network of depth
 681 $L \in \{2, 3\}$.

682 **If $L = 2$:** for $\mathbf{b}^{(1)} = -x_j \mathbf{W}^{(1)}$,

$$683 \quad \hat{\mathbf{X}}^{(L)}(x) = \begin{cases} \sigma(x_{j_1} - x) & \text{if } \mathbf{W}^{(1)} = -1 \\ \sigma(x - x_{j_1}) & \text{if } \mathbf{W}^{(1)} = 1 \end{cases}$$

684 **If $L = 3$:**

685 if $\sigma(x) = \text{ReLU}(x)$ and $m_1 = 1$: for $\mathbf{b}^{(1)} = -x_{j_1} \mathbf{W}^{(1)}$,
 686 if $\mathbf{W}^{(2)} = 1$:

$$687 \quad \hat{\mathbf{X}}^{(2)}(x) = \begin{cases} \text{ReLU}^-_{\min\{x_{j_1}, x_{j_2}\}}(x) & \text{if } \mathbf{W}^{(1)} = -1 \\ \text{ReLU}^+_{\max\{x_{j_1}, x_{j_2}\}}(x) & \text{if } \mathbf{W}^{(1)} = 1 \end{cases}$$

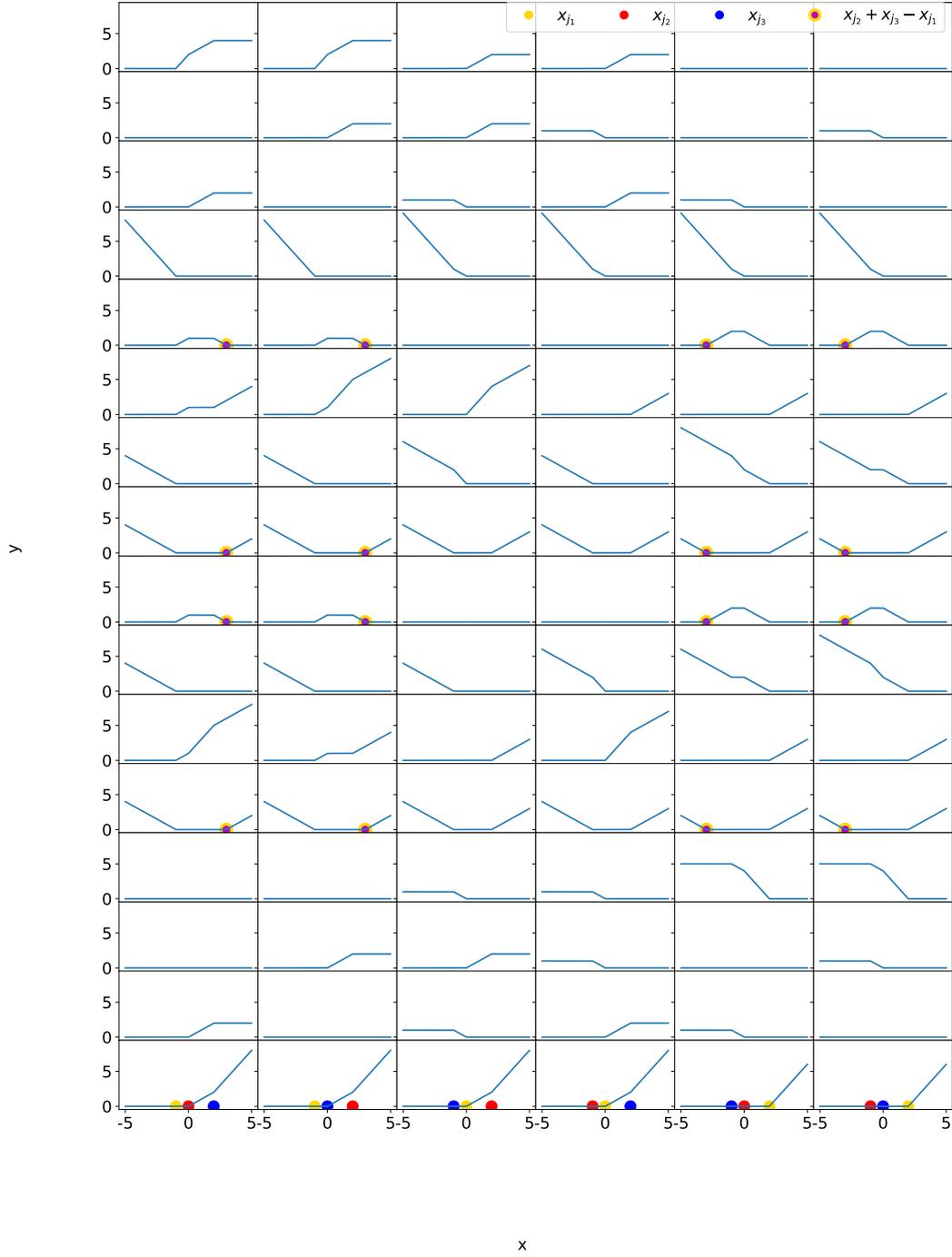


Figure B.1: Figure for [Appendix B](#). Deep library features for a 3-layer symmetrized ReLU network. Each row corresponds to a different set of weights $\mathbf{W}^{(1)} \in \{-1, 1\}^{1 \times 2}$, $\mathbf{W}^{(2)} \in \{-1, 1\}^{2 \times 1}$. Each column corresponds to a different ordering of x_i, x_j, x_k . The generalized reflections of x_{j_1} (gold) across x_{j_2} (red) and x_{j_3} (blue) are depicted by gold encircling purple.

688 if $\mathbf{W}^{(2)} = -1$:

$$689 \quad \hat{\mathbf{X}}^{(L)}(x) = \begin{cases} \text{Ramp}_{x_{j_2}, x_{j_1}}^+(x) & \text{if } \mathbf{W}^{(1)} = -1 \\ \text{Ramp}_{x_{j_1}, x_{j_2}}^-(x) & \text{if } \mathbf{W}^{(1)} = 1. \end{cases}$$

690 if $\sigma(x) = |x|$ and $m_1 = 1$: for $a \in \left\{x_{j_1}, \frac{x_{j_1} + x_{j_2}}{2}\right\}$,

$$691 \quad \hat{\mathbf{X}}^{(L)}(x) = \begin{cases} \mathcal{W}_{\min\{x_{j_2}, R(x_{j_2}, a)\}, a}(x) & \text{if } \mathbf{b}^{(1)} = -a\mathbf{W}^{(1)} \\ \mathcal{W}_{\min\{R(a, x_{j_1}), R(a, x_{j_2})\}, a}(x) & \text{if } \mathbf{b}^{(1)} = -a\mathbf{W}^{(1)}. \end{cases}$$

692 **B.1. Reconstruction results.** In this section, let (\mathbf{z}^*, ξ^*) be a solution to the Lasso problem.
 693 We give a map to efficiently and explicitly reconstruct an optimal neural net from (\mathbf{z}^*, ξ^*) by
 694 leveraging the structure of the deep library (Definition 3.10). Proofs are deferred to Appen-
 695 dix D.3.

696 **Lasso model on future test points:** The Lasso problem can be interpreted as finding a
 697 linear model $f : \mathbb{R}^m \rightarrow \mathbb{R}$, $f(\mathbf{a}) = \mathbf{z}^T \mathbf{a} + \xi$ such that for given samples $\mathbf{a}_1, \dots, \mathbf{a}_N \in \mathbb{R}^m$ and
 698 labels $y_1, \dots, y_N \in \mathbb{R}$, we have $f(\mathbf{a}_i) \sim y_i$. Now suppose the future test point is $x \in \mathbb{R}$. For
 699 simplicity, consider a 2-layer network with $\sigma(x) = x$. Since $\sigma(x) = x$, the Lasso dictionary is a
 700 simple $n \times n$ dictionary with $A_{i,j} = \sigma(x_i - x_j)$. For $\mathbf{a} = (\sigma(x - x_1), \dots, \sigma(x - x_N))$, the Lasso
 701 model and reconstructed network have the same value as $f(\mathbf{a}) = \mathbf{z}^T \mathbf{a} + \xi = \sum_{i=1}^N z_i \sigma(x - x_i) + \xi$.
 702
 703

704 **Definition B.3 (Reconstructed parameters).** The reconstructed parameters for a parallel
 705 network are constructed as follows. For each i^{th} column \mathbf{A}_i of the dictionary matrix such that
 706 $z_i^* \neq 0$, let $\hat{\mathbf{X}}^{(i,L)}$ be the parallel unit $\hat{\mathbf{X}}^{(L)}$ corresponding to that column in the deep library. Let
 707 $\boldsymbol{\alpha} = \mathbf{z}^*$ and $\xi = \xi^*$. For sign and threshold activation, let all amplitude parameters be 1. Finally,
 708 unscale parameters (Definition D.17).

709 Definition B.3 reconstructs parameters using the columns of the deep library, which are defined
 710 in Definition 3.9 and used to construct the Lasso problem. A reconstructed network is a
 711 network with reconstructed parameters. Some examples of optimal reconstructed parameters
 712 and networks are given below.
 713

714 **Reconstruction Examples.** Let \mathbf{z}^*, ξ^* be a solution to the Lasso problem. Let $\mathcal{J} =$
 715 $\{j : z_j^* \neq 0\}$. The following are some examples of optimal reconstructed networks using
 716 Definition B.3.

717 For simplicity, let $\sigma(x) = |x|$. For a 2-layer network, a set of optimal parameters is

$$718 \quad (\text{B.1}) \quad \alpha_j = \text{sign}(z_j^*) \sqrt{|z_j^*|}, \mathbf{W}^{(j,1)} = \sqrt{|z_j^*|}, \mathbf{b}^{(j,1)} = -x_j \sqrt{|z_j^*|}, \xi = \xi^*.$$

719 The corresponding optimal reconstructed network is

$$720 \quad (\text{B.2}) \quad f(x) = \sum_{j \in \mathcal{J}} \sigma \left(x \mathbf{W}^{(j,1)} + \mathbf{b}^{(j,1)} \right) \alpha_j + \xi$$

$$721 \quad (\text{B.3}) \quad = \sum_{j \in \mathcal{J}} \sigma \left(x \sqrt{|z_j^*|} - x_j \sqrt{|z_j^*|} \right) \text{sign}(z_j^*) \sqrt{|z_j^*|} + \xi^*$$

$$722 \quad (\text{B.4}) \quad = \sum_{j \in \mathcal{J}} z_j^* \sigma(x - x_j) + \xi^*.$$

723 By optimal parameter scaling ([Definition D.17](#)), the optimal weights and biases have a factor
724 of $\sqrt{|z_j^*|}$ in [\(B.1\)](#), which are plugged in to the neural network form [\(B.2\)](#) to yield [\(B.3\)](#), which
725 simplifies to [\(B.4\)](#). The reconstruction for general activations is similar.

726 A 3-layer deep narrow network has dictionary elements

727 $A_{i,j} = \sigma(\sigma(x_i - x_{j_1}) - \sigma(x_{j_2} - x_{j_1}))$ where the columns are indexed by $j = (j_1, j_2)$. A
728 set of optimal parameters is $\alpha_j = \text{sign}(z_j^*) |z_j^*|^{\frac{1}{3}}$, $\mathbf{W}^{(j,2)} = |z_j^*|^{\frac{1}{3}}$, $\mathbf{W}^{(j,1)} = |z_j^*|^{\frac{1}{3}}$, $\mathbf{b}^{(j,2)} =$
729 $-\sigma(x_{j_2} - x_{j_1}) |z_j^*|^{\frac{2}{3}}$, $\mathbf{b}^{(j,1)} = -x_{j_1} |z_j^*|^{\frac{1}{3}}$, $\xi = \xi^*$. The corresponding optimal reconstructed
730 network simplifies to

$$731 \quad f(x) = \sum_{j \in \mathcal{J}} z_j^* \sigma(\sigma(x - x_{j_1}) - \sigma(x_{j_2} - x_{j_1})) + \xi^*.$$

732 For a 3-layer symmetrized network, a set of optimal parameters is $\alpha_j = \text{sign}(z_j^*) |2z_j^*|^{\frac{1}{3}}$,

$$733 \quad \mathbf{W}_1^{(j,2)} = s_{j_4} \frac{1}{\sqrt{2}} |2z_j^*|^{\frac{1}{3}}, \mathbf{W}_2^{(j,2)} = s_{j_5} \frac{1}{\sqrt{2}} |2z_j^*|^{\frac{1}{3}}, \mathbf{W}_1^{(j,1)} = s_{j_6} \frac{1}{\sqrt{2}} |2z_j^*|^{\frac{1}{3}}, \mathbf{W}_2^{(j,1)} = s_{j_7} \frac{1}{\sqrt{2}} |2z_j^*|^{\frac{1}{3}},$$

$$734 \quad \mathbf{b}_1^{(j,1)} = -x_{j_1} \mathbf{W}_1^{(j,1)} \frac{|2z_j^*|^{\frac{1}{3}}}{\sqrt{2}}, \mathbf{b}_2^{(j,1)} = -x_{j_2} \mathbf{W}_2^{(j,1)} \frac{|2z_j^*|^{\frac{1}{3}}}{\sqrt{2}},$$

$$735 \quad \mathbf{b}^{(j,2)} = - \left(- \left(\sigma \left((x_{j_3} - x_{j_1}) W_1^{(j,1)} \right) W_1^{(j,2)} - \sigma \left((x_{j_3} - x_{j_2}) W_2^{(j,1)} \right) W_2^{(j,2)} \right) \right) \frac{|2z_j^*|^{\frac{2}{3}}}{2}, \xi = \xi^*,$$

736 where $s_{j_4}, \dots, s_{j_7} \in \{-1, 1\}$ and the multi-index is $j = (j_1, \dots, j_7)$. The corresponding optimal
737 reconstructed network simplifies to

$$738 \quad f(x) = \sum_{j \in \mathcal{J}} z_j^* \sigma \left(\left(\left(\sigma \left((x - x_{j_1}) W_1^{(j,1)} \right) W_1^{(j,2)} + \sigma \left((x - x_{j_2}) W_2^{(j,1)} \right) W_2^{(j,2)} \right) \right) \right. \\ \left. - \left(\sigma \left((x_{j_3} - x_{j_1}) W_1^{(j,1)} \right) W_1^{(j,2)} - \sigma \left((x_{j_3} - x_{j_2}) W_2^{(j,1)} \right) W_2^{(j,2)} \right) \right) + \xi^*.$$

739 General 3-layer, wide networks are out of the scope of the paper.

740 **Lemma B.4.** *A reconstructed parallel network is optimal in the training problem.*

741 **Appendix C. Solution sets of Lasso under minimal regularization.** One of the insights
742 that the Lasso formulation provides is that under minimal regularization, certain neural nets
743 perfectly interpolate the data. Proofs in this section are deferred to [Appendix D.4](#).

744 **Corollary C.1.** *For the ReLU, absolute value, sign, and threshold networks with $L = 2$ layers,*
745 *and sign-activated deeper networks, if $m_L \geq m^*$, then $f_L(\mathbf{X}; \theta) \rightarrow \mathbf{y}$ as $\beta \rightarrow 0$.*

746 In [Corollary C.1](#), m^* depends on L and the activation and is defined in [Theorem 3.11](#). The
747 Lasso equivalence and reconstruction also shed light on optimal neural network structure as
748 regularization decreases. The *minimum (l_1) norm subject to interpolation* version of the Lasso
749 problem is

$$750 \quad (\text{C.1}) \quad \min_{\mathbf{z}, \xi} \|\mathbf{z}\|_1 \text{ s.t. } \mathbf{A}\mathbf{z} + \xi\mathbf{1} = \mathbf{y}.$$

751 Assume \mathbf{A} has at least as many columns as rows. Loosely speaking, as $\beta \rightarrow 0$, if \mathbf{A} has full
752 rank, the Lasso problem (1.2) "approaches" the minimum norm problem (C.1), where $\xi = 0$ for
753 sign and threshold activations (27; 29). The rest of this section describes the solution sets of
754 (C.1) for certain networks. Recall that the training samples are ordered. For 2-layer networks
755 with absolute value activation, the dictionary \mathbf{A} is $N \times N$ where $A_{i,j} = |x_i - x_j|$, and the next
756 result analyzes the first and last elements z_1^* and z_N^* of a solution to the Lasso problem.

757 [Proposition C.2](#). *Let $L = 2$. Suppose σ is the absolute value activation. Let \mathbf{z}^* be a solution*
758 *to (C.1). Then, we have $z_1^* z_N^* \leq 0$. Moreover, the entire solution set of (C.1) for \mathbf{z}^* is*

$$759 \quad (\text{C.2}) \quad \left\{ \mathbf{z}^* + t \text{sign}(z_1^*)(1, 0, \dots, 0, 1)^T \mid -|z_1^*| \leq t \leq |z_N^*| \right\}.$$

760 [Proposition C.3](#). *For a 2-layer network with sign activation and $\beta \geq 0$, the Lasso problem*
761 *(1.2) has a unique solution. Furthermore, the minimum norm solution in (C.1) is $\mathbf{z}^* = \mathbf{A}^{-1}\mathbf{y}$.*

762 Given an optimal bias term ξ^* , if \mathbf{A} is invertible, then $\mathbf{z}^* = \mathbf{A}^{-1}(\mathbf{y} - \xi^*\mathbf{1})$ is optimal in
763 (C.1). [Appendix D.5](#) explicitly finds \mathbf{A}^{-1} for some activation functions. The structure of \mathbf{A}^{-1}
764 suggests the behavior of neural networks under minimal regularization: sign-activated neural
765 networks act as difference detectors, while neural networks with absolute value activation, whose
766 subgradient is the sign activation, act as a second-order difference detectors (see [Remark D.23](#)).
767 The next result shows that threshold-activated neural networks are also difference detectors,
768 but for the special case of positive, nonincreasing y_n . An example of such data is cumulative
769 revenue, e.g. $y_n = \sum_{i=1}^n r_i$ where r_i is the revenue in dollars earned on day i . Since the data is
770 assumed to be ordered, $i < j$ implies $x_i \geq x_j$, and hence for nonincreasing y_n , we have $y_i \geq y_j$.

771 [Proposition C.4](#). *Let $L = 2$. Suppose σ is threshold activation and $y_1 \geq \dots \geq y_N \geq 0$. Then*

$$772 \quad z_n^* = \begin{cases} y_n - y_{n-1} & \text{if } n \leq N - 1 \\ y_N & \text{if } n = N \\ 0 & \text{else} \end{cases}$$

773 *is the unique solution to the minimum norm problem (C.1).*

774 The next result gives a lower bound on the optimal value of the minimum weight problem
775 for ReLU networks. If we can find \mathbf{z} with a l_1 -norm that meets the lower bound and a ξ
776 such that $\mathbf{A}\mathbf{z} + \xi\mathbf{1} = \mathbf{y}$, then we know \mathbf{z}, ξ is optimal. In this section, for $n \in [N - 1]$, let
777 $\mu_n = \frac{y_n - y_{n+1}}{x_n - x_{n+1}}$ be the slope between the n^{th} and $n + 1^{\text{th}}$ data points. Let $\mu_N = 0$.

778 [Lemma C.5](#). *The optimal value $\|\mathbf{z}^*\|_1$ of the minimum norm problem (C.1) for deep narrow*
779 *networks with ReLU or absolute value activation is at least $\max_{n \in [N-1]} |\mu_n|$.*

780 When $L=2$, the next result gives a solution to the min-norm problem. For $i \in [N]$, let $(z_+)_i$
781 and $(z_-)_i$ be the Lasso variable corresponding to the features $\text{ReLU}_{x_i}^+$ and $\text{ReLU}_{x_i}^-$, respectively.
782 In other words, \mathbf{z}_+ corresponds to \mathbf{A}_+ , and \mathbf{z}_- corresponds to \mathbf{A}_- as defined in [Corollary 3.15](#).

783 **Lemma C.6.** *The min-norm problem (C.1) for $L=2, \sigma=\text{ReLU}$ has optimal value $\|\mathbf{z}^*\|_1 =$
784 $\sum_{n=1}^{N-1} |\mu_n - \mu_{n+1}|$. An optimal solution is $(z_+)_{n+1} = \mu_n - \mu_{n+1}$ for $n \in [N-1]$, $\mathbf{z}_- = \mathbf{0}$, and $\xi = y_N$.*

785 **Lemma C.7.** *For a 3-layer symmetrized ReLU network and training data as shown in
786 [Figure 3.2](#), the optimal value $\|\mathbf{z}^*\|_1$ of the minimum norm problem (C.1) is at least 1.*

787 **Appendix D. Numerical results.** The following simulations support our theoretical results.

788 In [Figure 3.2](#), a deep narrow, absolute value network is trained. In [Figure 3.5](#), a 3-layer
789 symmetrized ReLU network is trained. Both neural nets have standard architecture, $m_L=100$,
790 and are trained with Adam with the non-convex problem. The learning rate is $5(10^{-3})$, weight
791 decay is 10^{-4} , and $\beta=10^{-7} \approx 0$. When $\beta \rightarrow 0$, the Lasso problem approaches the minimum norm
792 problem (C.1). For each $L \in \{3, 4, 5\}$, the neural net reconstructed from the single feature in
793 the left plot with corresponding Lasso parameter $z_i^* = 1$ and $\xi^* = 0$ is optimal in the minimum
794 norm problem (C.1) by [Lemma C.5](#) and [Lemma C.7](#). This network is used to initialize a subset
795 of the neurons in the non-convex training. All other weights are initialized randomly according
796 to Pytorch defaults. The figures show that the networks trained with the non-convex problem
797 closely match the Lasso solutions. The networks exhibit breakpoints at data points and their
798 reflections not in the training data. The standard architecture in the non-convex model shows
799 the applicability of the Lasso formulation. SGD gives similar results as Adam.

800 **Supplementary Material.**

801 **Definitions and preliminaries.**

802 **D.1. Activation function.** A function f is *bounded* if there is $M \geq 0$ with $|f(x)| \leq M$ for
 803 all x . If $\sigma(x)$ is piecewise linear around zero, $\sigma(x)$ is bounded if and only if $a^- = a^+ = 0$, e.g.
 804 $\sigma(x)$ is a threshold or sign activation. We call f *symmetric* if it is an even or odd function, for
 805 example absolute value. The activation $\sigma(x)$ is defined to be *homogeneous* if for any $a \geq 0$,
 806 $\sigma(ax) = a\sigma(x)$. Homogeneous activations include ReLU, leaky ReLU, and absolute value, and
 807 don't have amplitude parameters. We say $\sigma(x)$ is *sign-determined* if its value depends only on
 808 the sign of its input and not its magnitude. Threshold and sign activations are sign-determined.

809 **D.2. Effective depth.**

810 **Remark D.1.** *Plugging (2.3) into itself shows that in parallel network, for $l \in [L-2]$,*

811
$$\hat{\mathbf{X}}^{(i,l+2)} = \sigma_{\mathbf{s}^{(i,l+1)}} \left(\sigma \left(\hat{\mathbf{X}}^{(i,l)} \mathbf{W}^{(i,l+1)} + \mathbf{b}^{(i,l)} \right) \mathbf{s}^{(i,l)} \mathbf{W}^{(i,l+1)} + \mathbf{b}^{(i,l+1)} \right).$$

812 *The inner parameters of a parallel network are $\mathbf{s}^{(i,l)}$ for $l \leq L-2$ and $\mathbf{W}^{(i,l)}$ for $l \leq L-1$. Suppose*
 813 *σ is sign-determined. Then $f_L(\mathbf{X}; \theta)$ is invariant to the value of the inner parameters, so*
 814 *they would be driven to 0 by weight regularization. We define the minimum value in (1.1) as*
 815 *an infimum which is approached as the norms of the inner parameters approach 0. Therefore*
 816 *the inner parameters are not regularized (the effective depth is 2), and we optimize for their*
 817 *directions rather than their magnitudes.*

818 **Parallel networks with data dimension $d \geq 1$.**

819
 820 For convenience, we will omit ξ when writing $f_L(\mathbf{X}; \theta)$. This does not change the training
 821 problem because we can write (1.1) as $\min_{\theta - \{\xi\}} \frac{\beta}{L} \|\theta_w\|_2^2 + \min_{\xi} \{\mathcal{L}_y(f_L(\mathbf{X}) - \xi \mathbf{1} + \xi \mathbf{1})\}$ and apply the
 822 change of variables/functions $\theta' = \theta - \{\xi\}$, $f_L(\mathbf{X}; \theta)' = f_L(\mathbf{X}; \theta) - \xi \mathbf{1}$ and $\mathcal{L}_y(\mathbf{z}') = \min_{\xi} L_y(\mathbf{z} + \xi \mathbf{1})$.
 823 The loss function absorbs ξ while preserving its convexity. We begin by assuming the data is
 824 d -dimensional and consider the general *training problem*

825
$$\min_{\theta \in \Theta} \mathcal{L}_y(f_L(\mathbf{X}; \theta)) + \frac{\beta}{L} r(\theta)$$

826 with a general regularization of the form $r(\theta) = \sum_{i=1}^{m_L} \sum_{\theta^{(i,l)} \in \theta_w^{(i)}} (r^{(i,l)}(\theta^{(i,l)}))^{\tilde{L}}$, where $\theta^{(i,l)}$ is a
 827 parameter such as $\mathbf{W}^{(i,l)}$ and $r^{(i,l)}$ is a regularization function such as $r^{(i,l)}(\mathbf{W}^{(i,l)}) = \|\mathbf{W}^{(i,l)}\|_p$.
 828 Assume $r^{(i,l)}$ is nonnegative and positively homogeneous, i.e., for any $\alpha \geq 0$, $r^{(i,l)}(\alpha \mathbf{W}) =$
 829 $\alpha r^{(i,l)}(\mathbf{W}) \geq 0$. The training problem (1.1) is a special case of the general training problem.

830 **Definition D.2.** *A simplified neural network with sign-determined activation is*

$$\begin{aligned}
f_L(\mathbf{x}; \theta) &= \xi + \sum_{i=1}^{m_L} \hat{\mathbf{X}}^{(i,L)} \alpha_i s^{(i,L-1)} \\
831 \quad (\text{D.1}) \quad \hat{\mathbf{X}}^{(i,l+1)} &= \sigma \left(\hat{\mathbf{X}}^{(i,l)} \mathbf{W}^{(i,l)} + \mathbf{b}^{(i,l)} \right), l \in [L-1] \\
\theta_w^{(i)} &= \left\{ \alpha_i, s^{(i,L-1)} \right\}, \theta_b^{(i)} = \left\{ \mathbf{b}^{(i,l)}, \mathbf{W}^{(i,l)} : l \in [L-1] \right\} \text{ for } i \in [m_L]
\end{aligned}$$

832 for a parallel network.

833 In other words, a simplified network has amplitude parameters only in the last layer.

834 **Lemma D.3.** *The simplified neural network is equivalent to the original neural network.*

835 *Proof.* By [Remark D.1](#), it suffices to only regularize the outermost $\tilde{L}=2$ layers. For parallel
836 networks, apply a change of variables $\mathbf{W}^{(i,l)'} = \mathbf{s}^{(i,l-1)} \mathbf{W}^{(i,l)}$ for $2 \leq l \leq L-1$. This removes $\mathbf{s}^{(i,l)}$
837 from θ for $l \leq L-2$. ■

838 Let $\mathcal{D} = \{L-\tilde{L}+1, \dots, L\}$. By [Lemma D.3](#), the training problem's regularization is

$$839 \quad (\text{D.2}) \quad r(\theta) = \sum_{i=1}^{m_L} \sum_{l \in \mathcal{D}} \left(r^{(i,l)} \left(\theta_w^{(i,l)} \right) \right)^{\tilde{L}}$$

840 where $\theta^{(i,L)} = \alpha_i$ and $\theta^{(i,l)} = \mathbf{W}^{(i,l)}$ for $l < L$ in parallel networks with homogeneous σ , $\theta^{(i,L)} = \alpha_i$
841 and $\theta^{(i,L-1)} = s^{(i,L-1)}$ in parallel networks with sign-determined σ . Similar to the parallel
842 network in [Subsection 2.2](#), extend the standard (2.2) network definitions row-wise to the cases
843 where the input is $\mathbf{X} \in \mathbb{R}^{N \times d}$.

844 **Lemma D.4.** *Let $\tilde{\mathbf{X}}^{(i)} \in \mathbb{R}^N$ be $\hat{\mathbf{X}}^{(i,L)}$ for a parallel network, where the input is $\mathbf{X} \in \mathbb{R}^{N \times d}$.
845 The training problem is equivalent to*

$$846 \quad (\text{D.3}) \quad \min_{\theta \in \Theta: r_{(i,l)}(\theta^{(i,l)})=1, l \in \mathcal{D}-\{L\}} \mathcal{L}_y \left(\sum_{i=1}^{m_L} \alpha_i \tilde{\mathbf{X}}^{(i)} \right) + \beta \sum_{i=1}^{m_L} r_{(i,L)}(\alpha_i)$$

847 *Proof.* By the AM-GM inequality on (D.2), a lower bound on the training problem is

$$848 \quad (\text{D.4}) \quad \min_{\theta \in \Theta} \mathcal{L}_y(f_L(\mathbf{X}; \theta)) + \beta \sum_{i=1}^{m_L} \prod_{l \in \mathcal{D}} r_{(i,l)}(\theta^{(i,l)}).$$

849 Consider the minimization problem

$$850 \quad (\text{D.5}) \quad \min_{\theta \in \Theta: r_{(i,l)}(\theta^{(i,l)})=1, l \in \mathcal{D}-\{L\}} \mathcal{L}_y(f_L(\mathbf{X}; \theta)) + \beta \sum_{i=1}^{m_L} \prod_{l \in \mathcal{D}} r_{(i,l)}(\theta^{(i,l)})$$

851 Problem (D.5) is an upper bound on (D.4). Given optimal $\{\theta^{(i,l)}\}$ in (D.4), the rescaled pa-
852 rameters $\theta^{(i,l)'} = \theta^{(i,l)} / r_{(i,l)}(\theta^{(i,l)})$ for $l \in \mathcal{D}-\{L\}$ and $\theta^{(i,L)'} = \theta^{(i,L)} \prod_{l \in \mathcal{D}} r_{(i,l)}(\theta^{(i,l)})$ (and rescaled

853 bias parameters) achieve the same objective in (D.5). Hence (D.5) and (D.4) are equivalent.
 854 Given optimal $\{\theta^{(i,l)}\}$ in (D.5), the rescaled parameters $\theta^{(i,l)'} = |\theta^{(i,l)}|^{\frac{1}{L}} \theta^{(i,l)}$ (and rescaled bias
 855 parameters) achieve the same objective in the training problem, which is therefore equivalent
 856 to (D.5). Simplifying (D.5) gives (D.3). ■

857 **Lemma D.4** applies to networks without any weight constraints, i.e., it excludes 3-layer
 858 symmetrized networks. A L -layer *symmetrized network* is a parallel network with homoge-
 859 neous activation such that $m_{L-3}=1$ and the elements of $\mathbf{W}^{(i,l)}$ have the same magnitude
 860 for $l \in \{L-2, L-1\}$. In a symmetrized network, the last two layer's weights are vectors:
 861 $\mathbf{W}^{(i,L-2)} \in \mathbb{R}^{1 \times m_{L-2}}$ and $\mathbf{W}^{(i,L-1)} \in \mathbb{R}^{m_{L-2}}$. The constraint on the weight magnitudes is en-
 862 coded in Θ . A 3-layer symmetrized network and a deep narrow network are special cases of a
 863 symmetrized network.

864 **Lemma D.5.** Let $r^{(i,l)}(\mathbf{W}^{(i,l)}) = \|\mathbf{W}^{(i,l)}\|_2$ for $l \in \{L-2, L-1\}$. The rescaled problem for
 865 a symmetrized network is equivalent to

(D.6)

$$866 \quad \min_{\theta \in \Theta: r^{(i,l)}(\mathbf{W}^{(i,l)})=1 \text{ for } l \in [L-3]; |\mathbf{W}_j^{(i,l)}|=1 \text{ for } l \in \{L-2, L-1\}, j \in [m_{L-3}]} \mathcal{L}_y(f_L(\mathbf{X}; \theta)) + \tilde{\beta} \sum_{i=1}^{m_L} r_{(i,L)}(\alpha_i)$$

867 where $\tilde{\beta} = \frac{\beta}{m_{L-2}}$.

868 *Proof.* Since $m_{L-3} = 1$, we have $\mathbf{W}^{(i,L-2)} \in \mathbb{R}^{1 \times m_{L-2}}$ and $\mathbf{W}^{(i,L-1)} \in \mathbb{R}^{m_{L-2}}$. The constraint
 869 states that for $l \in \{L-2, L-1\}$, $|\mathbf{W}_1^{(i,l)}| = \dots = |\mathbf{W}_{m_{L-2}}^{(i,l)}|$. For $l \in \{L-2, L-1\}$, given $\mathbf{W}^{(i,l)}$, $\mathbf{b}^{(i,l)}$
 870 and α_i apply a change of variables $\mathbf{W}^{(i,l)'} = \sqrt{m_{L-2}} \mathbf{W}^{(i,l)}$ and $\mathbf{b}^{(i,l)'} = \sqrt{m_{L-2}}^{l+3-L} \mathbf{b}^{(i,l)}$ and
 871 $\alpha_i' = \frac{1}{m_{L-2}} \alpha_i$ to the parameters in (D.3) to arrive at (D.6). ■

872 Henceforth, assume $r_{(i,L)}(\alpha_i) = |\alpha_i|$.

873 **Definition D.6.** Define the rescaled training problem as

$$874 \quad (\text{D.7}) \quad \min_{\theta \in \Theta} \mathcal{L}_y \left(\sum_{i=1}^{m_L} \alpha_i \tilde{\mathbf{X}}^{(i)} \right) + \tilde{\beta} \sum_{i=1}^{m_L} |\alpha_i|.$$

875 If the network is symmetrized, $\tilde{\beta} = \frac{\beta}{m_{L-2}}$ and Θ includes the constraints that $r^{(i,l)}(\mathbf{W}^{(i,l)})$
 876 $= 1$ for $l \in [L-3]$ and $|\mathbf{W}_j^{(i,l)}| = 1$ for $l \in \{L-2, L-1\}$. Otherwise, $\tilde{\beta} = \beta$ and $r_{(i,l)}(\theta^{(i,l)}) = 1$ for
 877 $l \in \mathcal{D} - \{L\}$. $\tilde{\mathbf{X}}^{(i)}$ is defined as in Lemma D.4.

878 For 3-layer networks with l_2 regularization ($r^{(i,l)}(\mathbf{W}^{(i,l)}) = \|\mathbf{W}^{(i,l)}\|_2$), Θ constrains the
 879 absolute value of all elements of all inner layer weights to be 1 in the rescaled training problem.
 880 The rescaled training problem is equivalent to both symmetrized and non-symmetrized networks.

881 **Lemma D.7.** The rescaled training problem (D.7) is equivalent to the training problem for
 882 all the architectures and activations discussed above.

883 *Proof.* Follows from Lemma D.4 and Lemma D.5. ■

884 **Lemma D.8.** *A lower bound on the rescaled training problem is the dual problem*

$$885 \quad (\text{D.8}) \quad \max_{\lambda \in \mathbb{R}^N} -\mathcal{L}_{\mathbf{y}}^*(\lambda) \quad \text{s.t.} \quad \max_{\theta \in \Theta} \left| \lambda^T \tilde{\mathbf{X}} \right| \leq \tilde{\beta},$$

886 where $\tilde{\mathbf{X}} = \tilde{\mathbf{X}}^{(1)}$ and $f^*(\mathbf{x}) := \max_{\mathbf{z}} \{\mathbf{z}^T \mathbf{x} - f(\mathbf{x})\}$ is the convex conjugate of f .

887 *Proof.* Find the dual of (D.3), by rewriting (D.3) as

$$888 \quad (\text{D.9}) \quad \min_{\theta \in \Theta} \mathcal{L}_{\mathbf{y}}(\mathbf{z}) + \tilde{\beta} \|\boldsymbol{\alpha}\|_1, \quad \text{s.t.} \quad \mathbf{z} = \sum_{i=1}^{m_L} \alpha_i \tilde{\mathbf{X}}^{(i)}.$$

889 The Lagrangian of problem (D.9) is $L(\lambda, \theta) = \mathcal{L}_{\mathbf{y}}(\mathbf{z}) + \tilde{\beta} \|\boldsymbol{\alpha}\|_1 - \lambda^T \mathbf{z} + \sum_{i=1}^{m_L} \lambda^T \tilde{\mathbf{X}}^{(i)} \alpha_i$. Minimize
890 the Lagrangian over \mathbf{z} and $\boldsymbol{\alpha}$ and use Fenchel duality (6). The dual of (D.9) is

$$891 \quad (\text{D.10}) \quad \max_{\lambda \in \mathbb{R}^N} -\mathcal{L}_{\mathbf{y}}^*(\lambda) \quad \text{s.t.} \quad \max_{\theta \in \Theta} \left| \lambda^T \tilde{\mathbf{X}}^{(i)} \right| \leq \tilde{\beta}, i \in [m_L].$$

892 In the tree and parallel nets, $\tilde{\mathbf{X}}^{(i)}$ is of the same form for all $i \in [m_L]$. So the m_L constraints
893 in (D.10) collapse to a single constraint. Then we can write (D.10) as (D.8). ■

894 For a parallel network, the dual problem (D.8) is

$$895 \quad (\text{D.11}) \quad \max_{\lambda \in \mathbb{R}^N} -\mathcal{L}_{\mathbf{y}}^*(\lambda) \quad \text{s.t.} \quad \max_{\theta \in \Theta} \left| \lambda^T \hat{\mathbf{X}}^{(L)} \right| \leq \tilde{\beta},$$

896 where $\hat{\mathbf{X}}^{(1)} = \mathbf{X}$. Henceforth, all regularizations are l_2 -norm: e.g., for a parallel network,
897 $r^{(i,l)}(\mathbf{W}^{(i,l)}) = \|\mathbf{W}^{(i,l)}\|_2$, denoting the square root of sum of squares of $\mathbf{W}^{(i,l)}$'s elements.

898 **Deep narrow and symmetrized networks with $d=1$.** In this section, we assume the data is
899 1-D and find the maximizers of $\left| \lambda^T \hat{\mathbf{X}}^{(L)} \right|$ in the dual constraint (D.11). To this end, assume the
900 elements of $\mathbf{W}^{(l)}$ are ± 1 . Note that $\mathbf{b}^{(L-1)}$ is a scalar and $\mathbf{X}^{(l+1)} = \sigma(\mathbf{X}^{(l)} \mathbf{W}^{(l)} + \mathbf{b}^{(l)} \mathbf{1}) \in \mathbb{R}^N$.
901 The next remark refers to the leaky ReLU slopes a^+ and a^- defined in Section 2. Let $\mathcal{K}(f)$
902 and $\mathcal{Z}(f)$ be the sets of breakpoints and zeros of a function f , respectively.

903 **Remark D.9.** Let $b = \mathbf{b}^{(L-1)}$, $a_n = \hat{\mathbf{X}}_n^{(L-1)} \mathbf{W}^{(L-1)}$, $g_n(b) = \sigma(a_n + b)$, and $g(b) = \sum_{n=1}^N \lambda_n g_n(b)$
904 $= \lambda^T \mathbf{X}^{(L)}$. Let $\mathcal{I} = \bigcup_{n=1}^N \mathcal{K}(g_n) \supset \mathcal{K}(g)$. Then $g(b) = \sum_{n=1}^N \lambda_n a(a_n + b) = ab \sum_{n=1}^N \lambda_n +$
905 $\sum_{n=1}^N \lambda_n a_n a$ for b large enough (with $a = a^+$) and for b small enough (with $a = a^-$). So $a^- = a^+ = 0$
906 or $\sum_{n=1}^N \lambda_n = 0$ if and only if g is bounded, if and only if g has a (finite) maximizer and minimizer.
907 In this case, assuming g is not a constant function, \mathcal{I} contains a maximizer and minimizer of g .

908 Henceforth, assume $\lambda^T \mathbf{1} = 0$ if $a^- \neq 0$ or $a^+ \neq 0$. Suppose $\mathbf{b}^{(l)}, \dots, \mathbf{b}^{(L-1)}$ are scalar-valued. We
909 call $\mathbf{b}^{(l)} = \mathbf{b}^{(l)*}$ optimal if $\mathbf{b}^{(l)*} \in \arg \max_{\mathbf{b}^{(l)}} \max_{\mathbf{b}^{(l+1)}} \dots \max_{\mathbf{b}^{(L-1)}} \left| \sum_{n=1}^N \lambda_n \hat{\mathbf{X}}_n^{(L)} \right|$. The next
910 result refers to the normalized midpoint defined in (3.1).

911 **Lemma D.10.** Let σ be leaky ReLU. Let $\alpha, \beta \in \mathbb{R}$ with $\alpha \neq \beta$. Let $f(x) = \sigma(x + \alpha) - \sigma(x + \beta)$.
912 For $s \in \{-, +\}$, let $g_s(x) = \sigma(sx)$. Then, if σ is monotone,

$$913 \quad (\text{D.12}) \quad \mathcal{K}(g_s(f)) \subset \{-\alpha, -\beta\}.$$

914 Otherwise if σ is not monotone,

$$915 \quad (\text{D.13}) \quad \mathcal{K}(g_s(f)) \subset \{-\alpha, -\beta, m_{\alpha, \beta}\}.$$

916 *Proof.* Since $f_{\alpha, \beta}$ is piecewise linear,

$$917 \quad (\text{D.14}) \quad \mathcal{K}(g_s(f)) \subset \mathcal{K}(f) \cup \mathcal{Z}(f).$$

918 Moreover if $f(x)$ has the same sign for all x , then observe that

$$919 \quad (\text{D.15}) \quad \mathcal{K}(g_s(f)) = \mathcal{K}(f).$$

920 Let $\mathcal{G}(f) = \{(x, f(x)) : x \in \mathcal{K}(f)\}$. We find that

$$921 \quad (\text{D.16}) \quad \mathcal{G}(f) = \left\{ \left(-\max\{\alpha, \beta\}, a^-(\alpha - \beta) \right), \left(-\min\{\alpha, \beta\}, a^+(\alpha - \beta) \right) \right\}.$$

922 In particular, $\mathcal{K}(f) = \{-\alpha, -\beta\}$. Observe that f has constant value (an in particular, sign)
 923 beyond its breakpoints. If σ is monotone, then $\text{sign}(a^+) = \text{sign}(a^-)$, so (D.16) implies
 924 $\text{sign}(f(x)) = \text{sign}(a^+(\alpha - \beta))$ for all x , so (D.15) implies (D.12). If σ is not monotone, linearly
 925 interpolating between the points in $\mathcal{G}(f)$ (D.16) shows that f changes sign exactly when
 926 $x = m_{\alpha, \beta}$, so (D.14) implies (D.13). ■

927 **Lemma D.11.** *Let σ be piecewise linear if $L=2$ and leaky ReLU otherwise.*

928 *Let $l \in \{L-1, L-2\}$. Consider the bias $\mathbf{b}^{(l)}$. Let $m_{L-2}=1$. If σ is monotone there is either a*
 929 *data feature bias that is optimal. If σ is not monotone there is either a data feature bias or a*
 930 *midpoint feature bias $b^{(L-2)} = m_{\mathbf{X}_n^{(L-2)}, \mathbf{X}_n^{(L-1)}} \mathbf{W}^{(L-2)}$ (if $l=L-2$) that is optimal.*

931 *Proof.* Remark D.9 implies that $\bigcup_{n=1}^N \mathcal{K}(\hat{\mathbf{X}}_n^{(L)})$ contains an optimal $\mathbf{b}^{(L-1)}$. The break-
 932 point of $\hat{\mathbf{X}}_n^{(L)} = \sigma(\mathbf{X}_n^{(L-1)} \mathbf{W}^{(L-1)} + \mathbf{b}^{(L-1)})$ as a function of $\mathbf{b}^{(L-1)}$ is $\mathbf{b}^{(L-1)} = -\mathbf{X}_n^{(L-1)} \mathbf{W}^{(L-1)}$.
 933 Therefore for some $n^{(L-1)} \in [N]$, the data feature bias $\mathbf{b}^{(L-1)} = -\mathbf{X}_{n^{(L-1)}}^{(L-1)} \mathbf{W}^{(L-1)}$ is optimal.
 934 Plugging in this optimal $\mathbf{b}^{(L-1)}$ back into $\hat{\mathbf{X}}_n^{(L)}$ gives

$$935 \quad \lambda^T \mathbf{X}^{(L)} = \sum_{n=1}^N \lambda_n \sigma \left(\left(\hat{\mathbf{X}}_n^{(L-1)} - \hat{\mathbf{X}}_{n^{(L-1)}}^{(L-1)} \right) \mathbf{W}^{(L-1)} \right).$$

936 For $L > 2$, expanding $\hat{\mathbf{X}}^{(L-1)}$ gives
 (D.17)

$$937 \quad \lambda^T \mathbf{X}^{(L)} = \sum_{n=1}^N \lambda_n \sigma \left(\left(\sigma \left(\hat{\mathbf{X}}_n^{(L-2)} \mathbf{W}^{(L-2)} + \mathbf{b}^{(L-2)} \right) - \sigma \left(\hat{\mathbf{X}}_{n^{(L-1)}}^{(L-2)} \mathbf{W}^{(L-2)} + \mathbf{b}^{(L-2)} \right) \right) \mathbf{W}^{(L-1)} \right).$$

938 Since $m_{L-2}=1$, we have $\mathbf{b}^{(L-2)} \in \mathbb{R}$. Applying Lemma D.10 with $x = \mathbf{b}^{(L-2)}$, $\alpha = \hat{\mathbf{X}}_n^{(L-2)} \mathbf{W}^{(L-2)}$,
 939 $\beta = \hat{\mathbf{X}}_{n^{(L-1)}}^{(L-2)} \mathbf{W}^{(L-2)}$, $s = \mathbf{W}^{(L-1)}$ and a similar argument as Remark D.9 gives the result. ■

940 **Lemma D.12.** *Let $L \geq 2$. Consider a deep narrow ReLU network with L layers. For every*
 941 *$l \in \{2, \dots, L\}$, a data feature bias is optimal for layer $L-l+1$, and with this optimal bias, if*
 942 *$l < L$, there exists $N_1, N_2 \in [N]$ such that for all $n \in [N]$, either $\hat{\mathbf{X}}_n^{(L)} = 0$ or*

943 (D.18)
$$\hat{\mathbf{X}}_n^{(L)} = \pm \left(\sigma \left(\hat{\mathbf{X}}_{N_1}^{(L-l)} \mathbf{W}^{(L-l)} + \mathbf{b}^{(L-l)} \right) - \sigma \left(\hat{\mathbf{X}}_{N_2}^{(L-l)} + \mathbf{b}^{(L-l)} \right) \right).$$

944 *Proof.* We prove by induction on l .

945 Base case: suppose $l = 2$. Then the claim holds by (D.17) in the proof of Lemma D.11.

946 Now, suppose the claim holds for $l=k \in \{2, \dots, L-1\}$. Applying argument similar to the proof

947 of Lemma D.11 to (D.18) shows that a data feature bias $\mathbf{b}^{(L-k)} = -\hat{\mathbf{X}}_{n^{(L-k)}}^{(L-k)} \mathbf{W}^{(L-k)}$ is optimal

948 for $\mathbf{b}^{(L-l+1)}$ when $l=k+1$. Plugging in this optimal $\mathbf{b}^{(L-k)}$ into $\hat{\mathbf{X}}^{(L)}$ in (D.18) for $l=k$ and

949 expanding $\hat{\mathbf{X}}^{(L-k)} = \sigma \left(\hat{\mathbf{X}}^{(L-k-1)} \mathbf{W}^{(L-k-1)} + \mathbf{b}^{(L-k-1)} \right)$ shows that for some $s^{(1)}, s^{(2)} \in \{-1, 1\}$,

950 either $\hat{\mathbf{X}}_n^{(L)} = 0$ or

$$\begin{aligned} \hat{\mathbf{X}}_n^{(L)} &= \sigma \left(s^{(1)} \left(\sigma \left(s^{(2)} \left(\sigma \left(\hat{\mathbf{X}}_{n'}^{(L-k-1)} \mathbf{W}^{(L-k-1)} + \mathbf{b}^{(L-k-1)} \right) - \sigma \left(\hat{\mathbf{X}}_{n^{(L-2)}}^{(L-k-1)} \mathbf{W}^{(L-3)} + \mathbf{b}^{(L-k-1)} \right) \right) \right) \right. \\ &\quad \left. - \sigma \left(s^{(2)} \left(\sigma \left(\hat{\mathbf{X}}_{n^{(L-1)}}^{(L-k-1)} \mathbf{W}^{(L-k-1)} + \mathbf{b}^{(L-k-1)} \right) - \sigma \left(\hat{\mathbf{X}}_{n^{(L-2)}}^{(L-k-1)} \mathbf{W}^{(L-k-1)} + \mathbf{b}^{(L-k-1)} \right) \right) \right) \right) \\ 951 &\in \left\{ s^{(2)} \left(\sigma \left(\hat{\mathbf{X}}_{n'}^{(L-k-1)} \mathbf{W}^{(L-k-1)} + \mathbf{b}^{(L-k-1)} \right) - \sigma \left(\hat{\mathbf{X}}_{n^{(L-2)}}^{(L-k-1)} + \mathbf{b}^{(L-k-1)} \right) \right), \right. \\ &\quad \left. - s^{(2)} \left(\sigma \left(\hat{\mathbf{X}}_{n^{(L-1)}}^{(L-k-1)} \mathbf{W}^{(L-k-1)} + \mathbf{b}^{(L-k-1)} \right) - \sigma \left(\hat{\mathbf{X}}_{n^{(L-2)}}^{(L-k-1)} + \mathbf{b}^{(L-k-1)} \right) \right), \right. \\ &\quad \left. s^{(2)} \left(\sigma \left(\hat{\mathbf{X}}_{n'}^{(L-k-1)} \mathbf{W}^{(L-k-1)} + \mathbf{b}^{(L-k-1)} \right) - \sigma \left(\hat{\mathbf{X}}_{n^{(L-1)}}^{(L-k-1)} + \mathbf{b}^{(L-k-1)} \right) \right), 0 \right\}. \end{aligned}$$

952 So (D.18) holds for $l=k+1$. Finally if (D.18) holds for $l=L-1$ then by a similar argument

953 as the proof of Lemma D.11, a data feature bias for $\mathbf{b}^{(L-l+1)}$ is optimal when $l=k+1=L$. By

954 induction, the result holds. \blacksquare

955 **Lemma D.13.** Let $L \geq 2$. Consider a deep narrow ReLU network with L layers. For

956 $l \in [L-1]$, suppose $\mathbf{b}^{(l)}$ is a data feature. For every $l \in \{2, \dots, L-1\}$, there exist $N_1, N_2 \in [N]$

957 such that for all x ,

958 (D.19)
$$\hat{\mathbf{X}}^{(l)}(x) \in \left\{ \text{ReLU}_{x_{N_1}}^\pm(x), \text{Ramp}_{x_{N_1}, x_{N_2}}^\pm(x) \right\}.$$

959 *Proof.* We prove by induction on l . Base case: for $l = 2$, $\hat{\mathbf{X}}^{(2)}(x) = \sigma(x \mathbf{W}^{(1)} + \mathbf{b}^{(1)}) =$

960 $\sigma((x - x_{n^{(1)}}) \mathbf{W}^{(1)}) = \text{ReLU}_{x_{n^{(1)}}}^{\mathbf{W}^{(1)}}(x)$. Now, suppose the claim holds for $l = 1, \dots, k < L-1$. Then

961 $\hat{\mathbf{X}}^{(k+1)}(x) = \sigma \left(\hat{\mathbf{X}}^{(k)} \mathbf{W}^{(k)} + \mathbf{b}^{(k)} \right) = \sigma \left(\left(\hat{\mathbf{X}}_n^{(k)} - \hat{\mathbf{X}}_{n^{(k)}}^{(k)} \right) \mathbf{W}^{(k)} \right)$ so either

962 there exist $N'_1, N'_2 \in \{N_1, n^{(k)}\}$, $s \in \{+, -\}$ such that for all x ,

963 $\hat{\mathbf{X}}^{(k+1)}(x) = \sigma \left(\left(\text{ReLU}_{x_{N_1}}^\pm(x) - \text{ReLU}_{x_{N_1}}^\pm(x_{n^{(k)}}) \right) \mathbf{W}^{(k)} \right) \in \left\{ \text{ReLU}_{x_{N'_1}}^s(x), \text{Ramp}_{x_{N'_1}, x_{N'_1}}^s(x) \right\}$, or

964 there exist $N'_1, N'_2 \in \{N_1, N_2, n^{(k)}\}$, $s \in \{+, -\}$ such that for all x ,

965 $\hat{\mathbf{X}}^{(k+1)}(x) = \sigma \left(\left(\text{Ramp}_{x_{N_1}, x_{N_2}}^\pm(x) - \text{Ramp}_{x_{N_1}, x_{N_2}}^\pm(x_{n^{(k)}}) \right) \mathbf{W}^{(k)} \right) = \text{Ramp}_{x_{N'_1}, x_{N'_2}}^s(x)$. By induc-

966 tion, the result holds. \blacksquare

967 **Lemma D.14.** Consider a deep narrow network. Let $L \geq 2$. For all $l \in [L-1]$, there exist

968 $n^{(L-l)} \in [N]$ and a l -tuple (a_0, a_1, \dots, a_l) such that for all $l \in [L-1]$,

969 (D.20)
$$\mathbf{b}^{(L-l)*} = - \sum_{i=1}^l a_i \hat{\mathbf{X}}_n^{(L-l)} \mathbf{W}^{(L-l)}$$

970 *is optimal. Note that optimal $\mathbf{b}^{(1)}, \dots, \mathbf{b}^{(L-1)}$ can be found sequentially, by computing $\mathbf{b}^{(l)*}$*
 971 *as a function of $\mathbf{b}^{(l-1)*}$ and so on. Moreover, there are $O(2^L L!)$ options for $\mathbf{b}^{(1)*}$. Additionally,*
 972 *for all $l \in [L-1]$, under such optimal $\mathbf{b}^{(1)}, \dots, \mathbf{b}^{(l-1)}$, for all $n \in [N]$,*

973 (D.21)
$$\hat{\mathbf{X}}_n^{(L-l+1)} = \sigma \left(\hat{\mathbf{X}}_n^{(L-l)'} \mathbf{W}^{(L-l)} \right)$$

974 *where*

975 (D.22)
$$\hat{\mathbf{X}}_n^{(L-l)'} = a_0 \hat{\mathbf{X}}_n^{(L-l)} + \sum_{i=1}^l a_i \hat{\mathbf{X}}_n^{(L-l)}.$$

976 *Proof.* We prove (D.20), (D.21) and (D.22) by strong induction on l .

977 Base case: suppose $l = 1$. Then (D.20), (D.21), (D.22) hold by Lemma D.11 and its proof.

978 Now, suppose (D.20), (D.21) and (D.22) hold for $l=1, \dots, k < L-1$. By Equation (D.21) for
 979 $l=k$, $\hat{\mathbf{X}}_n^{(L-k+1)}(x)$ is locally a linear combination of the $k+1$ terms $\hat{\mathbf{X}}_n^{(L-k+1)}, \hat{\mathbf{X}}_n^{(L-k+1)}, \dots$
 980 $, \hat{\mathbf{X}}_n^{(L-k+1)}$, and hence so is $\hat{\mathbf{X}}^{(L)}$. So the breakpoints of $\hat{\mathbf{X}}^{(L)}$ as a function of $\mathbf{b}^{(L-(k+1))}$
 981 are linear combinations of the $k+1$ terms, which proves (D.20) for $l=k+1$. Plugging in this
 982 breakpoint $\mathbf{b}^{(L-k-1)}$ into $\hat{\mathbf{X}}_n^{(L-k)}$ proves (D.21) and (D.22) for $l=k+1$. Therefore (D.20), (D.21)
 983 and (D.22) hold for all $l \in [L-1]$.

984 Now we prove that the number of options for $\mathbf{b}^{(1)}$ is $O(2^L L!)$. By Equation (D.22) for
 985 $l=L-1$, as a function of $\mathbf{b}^{(1)}$, $\hat{\mathbf{X}}_j^{(L-1)'}$ has breakpoints at $-\hat{\mathbf{X}}_i^{(L-l-1)} \mathbf{W}^{(L-l-1)}$ for $i \in \{n, n^{(L-1)},$
 986 $\dots, n^{(L-l)}\}$, which totals to $l+1=L$ breakpoints. A piecewise linear function f has up to plus
 987 or minus one as many zeros as breakpoints, and the zeros become breakpoints in $\sigma(f)$. So by
 988 Equation (D.21), $\hat{\mathbf{X}}_n^{(L-l+1)}$ as a function of $\mathbf{b}^{(1)}$ has breakpoints at $O(2L)$ places. Then by
 989 Equation (D.22) for $l=L-2$, $\hat{\mathbf{X}}_n^{(L-l+1)'}$ as a function of $\mathbf{b}^{(1)}$ has breakpoints at $O(2(L-1)(L))$
 990 places. And by Equation (D.21) for $l=L-2$, $\hat{\mathbf{X}}_n^{(L-k+2)}$ as a function of $\mathbf{b}^{(1)}$ has breakpoints
 991 at $O(2^2 L(L-1))$ places. By repeating this argument for $l=L-1, \dots, 1$, $\hat{\mathbf{X}}_n^{(L)}$ as a function of
 992 $\mathbf{b}^{(1)}$ has breakpoints at $O(2^L L!)$ places. ■

993 Recall the deep library defined in Definition 3.10. Note that any \mathbf{a} in the deep library is of
 994 the form $\mathbf{a} = \hat{\mathbf{X}}^{(L)}(\mathbf{X})$, with $\mathbf{a}_n = \hat{\mathbf{X}}^{(L)}(x_n)$.

995 **Lemma D.15.** *For $L \in \{2, 3\}$ the maximization constraint in (D.11) is equivalent to: for all*
 996 *vectors \mathbf{a} in the deep library,*

997 (D.23)
$$|\lambda^T \mathbf{a}| \leq \tilde{\beta}$$

$$\mathbf{1}^T \lambda = 0 \text{ if } a^+ \neq 0 \text{ or } a^- \neq 0.$$

998 *Proof.* Lemma D.11 gives (D.23). Now, if the activation is symmetric, then $\hat{\mathbf{X}}^{(L)}$ is invariant
999 under the sign of the components of $\mathbf{W}^{(1)}$. Next, recall $x_1 > \dots > x_N$ and $\text{sign}(0) = 1$. If
1000 the activation is sign, then for all $n \in [N - 1]$, with $\mathbf{W}^{(1)} = 1$ and $\mathbf{b}^{(1)} = -x_n$ we have
1001 $\hat{\mathbf{X}}^{(L=2)}(\mathbf{X}) = \sigma(\mathbf{X} - x_n) = (\mathbf{1}_{1:n}^T, -\mathbf{1}_{n+1:N}^T) = -\sigma(x_{n+1} - \mathbf{X})$; while for $\mathbf{W}^{(1)} = -1$ and
1002 $\mathbf{b}^{(1)} = x_{n+1}$ we have $\hat{\mathbf{X}}^{(L=2)}(\mathbf{X}) = \sigma(x_{n+1} - \mathbf{X})$. And for $\mathbf{W}^{(1)} = 1$ and $\mathbf{b}^{(1)} = -x_N$, we have
1003 $\hat{\mathbf{X}}^{(2)}(\mathbf{X}) = \mathbf{1}$ while for $\mathbf{W}^{(1)} = -1$ and $\mathbf{b}^{(1)} = x_1$ we have $\hat{\mathbf{X}}^{(2)}(\mathbf{X}) = \mathbf{1}$. So for $L = 2$ with
1004 symmetric or sign activation, (D.23) is unchanged if $\mathbf{W}^{(1)} \in \{-1, 1\}$ is restricted to be 1. ■

1005 **Lemma D.16.** *Let \mathbf{A} be a matrix whose set of columns is the deep library. Replace the*
1006 *maximization constraint in (D.11) with (D.23). The dual of (D.11) then is*

$$1007 \text{ (D.24)} \quad \min_{\mathbf{z}, \xi \in \mathbb{R}} \mathcal{L}_{\mathbf{y}}(\mathbf{A}\mathbf{z} + \xi\mathbf{1}) + \tilde{\beta}\|\mathbf{z}\|_1, \quad \text{where } \xi = 0 \text{ if } c_1 = c_2 = 0.$$

1008 *Proof.* Problem (D.11) can be written as

$$1009 \text{ (D.25)} \quad - \min_{\lambda \in \mathbb{R}^N} \mathcal{L}_{\mathbf{y}}^*(\lambda) \quad \text{s.t.} \quad \lambda^T \mathbf{1} = 0 \text{ if } a^- \neq 0 \text{ or } a^+ \neq 0, \text{ and } |\lambda^T \mathbf{A}| \leq \tilde{\beta} \mathbf{1}^T.$$

1010 The Lagrangian of the negation of (D.25) with bidual variables \mathbf{z}, ξ is

$$1011 \text{ (D.26)} \quad L(\lambda, \mathbf{z}, \xi) = \mathcal{L}_{\mathbf{y}}^*(\lambda) - \lambda^T(\mathbf{A}\mathbf{z} + \xi\mathbf{1}) - \tilde{\beta}\|\mathbf{z}\|_1, \text{ where } \xi = 0 \text{ if } a^- = a^+ = 0.$$

1012 Equation (D.26) holds because the constraint $|\lambda^T \mathbf{A}| \leq \tilde{\beta} \mathbf{1}^T$ i.e., $\lambda^T \mathbf{A} - \tilde{\beta} \mathbf{1}^T \leq \mathbf{0}^T, -\lambda^T \mathbf{A} -$
1013 $\tilde{\beta} \mathbf{1}^T \leq \mathbf{0}^T$, appears in the Lagrangian as $\lambda^T \mathbf{A} (\mathbf{z}^{(1)} - \mathbf{z}^{(2)}) - \tilde{\beta} \mathbf{1}^T (\mathbf{z}^{(1)} + \mathbf{z}^{(2)})$ with bidual
1014 variables $\mathbf{z}^{(1)}, \mathbf{z}^{(2)}$, which are combined into one bidual variable $\mathbf{z} = \mathbf{z}^{(1)} - \mathbf{z}^{(2)}$. This makes
1015 $\mathbf{z}^{(1)} + \mathbf{z}^{(2)} = \|\mathbf{z}\|_1$. Changing variables $\mathbf{z}' = -\mathbf{z}, \xi' = -\xi$ gives (D.26). Since $\mathcal{L}^{**} = \mathcal{L}$ (5),
1016 $\inf_{\lambda} L(\lambda, \mathbf{z}, \xi) = -\mathcal{L}_{\mathbf{y}} - \tilde{\beta}\|\mathbf{z}\|_1$ and negating its maximization gives (D.24). ■

1017 **Definition D.17 (Parameter unscaling).** *For symmetrized networks, change $\alpha'_i = 2\alpha_i$. Let*
1018 $\gamma_i = |\alpha_i|^{\frac{1}{L}}$. *For ReLU, absolute value, or leaky ReLU, parameter unscaling is the following*
1019 *transformation. For symmetrized networks, first change variables as $\mathbf{W}^{(i,l)'} = \frac{1}{\sqrt{2}} \mathbf{W}^{(i,l)}$ and*
1020 $\mathbf{b}^{(i,l)'} = \frac{1}{\sqrt{2}} \mathbf{b}^{(i,l)}$ *for $l \in [2]$. Then, for all architectures, change variables as $q' = \text{sign}(q)\gamma_i$ for*
1021 $q \in \theta_w^{(i)}$. *For parallel networks, change variables as $\mathbf{b}^{(i,l)'} = \mathbf{b}^{(i,l)} (\gamma_i)^l$. For sign and threshold*
1022 *activations, parameter unscaling is the transformation $\alpha'_i = \text{sign}(\alpha_i)\sqrt{|\alpha_i|}$ and $s^{(i,L-1)'} = \sqrt{|\alpha_i|}$*
1023 *for parallel nets.*

1024 D.3. Proofs of results in Section 3 .

1025 *Proof of Lemma 3.7.* The Lasso equivalence follows from a similar argument as the proof
1026 of Theorem 3.11. The Lasso features follow from Lemma D.12 and Lemma D.13 with $l = L$. ■

1027 *Proof of Theorem 3.6.* By the proof of Lemma D.14 and a similar argument as Theo-
1028 rem 3.11, the training problem is equivalent to a Lasso problem. Moreover Lemma D.14 and
1029 choice of $n^{(1)}, \dots, n^{(L-1)} \in [N]$ give the number of possible $\mathbf{b}^{(1)}$ for given $\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(L-1)}$ as
1030 $O(N^{L-1} 2^L L!)$. Note that by (D.21) and (D.22), $\mathbf{b}^{(1)}$ determines all other $\mathbf{b}^{(l)}$. Finally, there
1031 are 2^L possibilities of $\mathbf{W}^{(l)} \in \{-1, 1\}$. The ReLU result follows from Lemma 3.7. ■

1032 *Proof of Theorem 3.11.* By Lemma D.16, problem (D.24) is a lower bound on the training
1033 problem (1.1), where the Lasso features are all vectors in the deep library. Let (\mathbf{z}^*, ξ^*) be a Lasso
1034 solution. From Definition 3.10, it can be seen that $\mathbf{A}\mathbf{z}^* + \xi^* = \sum_i z_i^* \mathbf{A}_i + \xi^* = \sum_i \alpha_i \hat{\mathbf{X}}^{(i,L)}(\mathbf{X}) + \xi^*$
1035 $= f_L(\mathbf{X}; \theta)$ and $\|\mathbf{z}^*\|_1 = \|\boldsymbol{\alpha}\|_1$. Therefore a reconstructed neural net achieves the same objective
1036 in the rescaled training problem, as (\mathbf{z}^*, ξ^*) does in the Lasso objective. Parameter unscal-
1037 ing (Definition D.17) makes them achieve the same objective in the training problem (see
1038 Remark D.20). Therefore the Lasso problem in Theorem 3.11 is equivalent to the training
1039 problem. ■

1040 *Proof of Lemma B.4.* By Theorem 3.11 and its proof, the reconstructed neural net achieves
1041 the same objective as the Lasso problem, which is equivalent to the training problem. So the
1042 reconstructed neural net is optimal. ■

1043 *Proof of Theorem 3.13.* By considering each component of $\mathbf{b}^{(l)}$ separately in the proof of
1044 Lemma D.11, we see that optimal bias parameters for a symmetrized network can include data
1045 features. The Lasso equivalence follows from a similar argument as the proof of Theorem 3.11.
1046 In the deep library, each $\mathbf{W}^{(1)}, \mathbf{W}^{(2)}$ consists of 2 elements that are ± 1 , and given $\mathbf{W}^{(1)}, \mathbf{W}^{(2)}$,
1047 each of the 2 elements of $\mathbf{b}^{(1)}$ as well as the scalar $\mathbf{b}^{(2)}$ have N choices. So the number of
1048 possible features in the deep library is at most $2^{2(L-1)}N^3 = 2^4N^3$. ■

1049 *Proof of Theorem 3.5.* Follows from Theorem 3.13 and evaluating the features. See proof
1050 of Figure 3.4 below. ■

1051 *Proof of Theorem 3.2.* For $L \in \{2, 3\}$, apply Theorem 3.6, Theorem 3.11 where $\sigma(x) = |x|$
1052 and Lemma B.2. For $L=4$, continue the same line of argument as Lemma D.11 by plugging
1053 in $\mathbf{b}^{(L-2)} = -\hat{\mathbf{X}}_{n^{(L-2)}}^{(L-2)} \mathbf{W}^{(L-2)}$ into (D.17), expanding $\hat{\mathbf{X}}^{(L-2)} = \sigma\left(\hat{\mathbf{X}}^{(L-3)} \mathbf{W}^{(L-3)} + \mathbf{b}^{(L-3)}\right)$, and
1054 observing that the breakpoints of $\lambda^T \hat{\mathbf{X}}^{(L)}$ as a function of $\mathbf{b}^{(L-3)}$ include data features, and
1055 hence these can be optimal bias parameters. Plugging in these data feature biases in to
1056 $\hat{\mathbf{X}}^{(4)}$ gives the features and reflections. Continue repeating a similar argument for $L > 4$.
1057 The reconstruction and Lasso equivalence follows from a similar argument as the proof of
1058 Theorem 3.11. ■

1059 *Remark D.18.* In the proof of Theorem 3.2 for $L \geq 4$, we only used one of the possible
1060 forms of an optimal $\mathbf{b}^{(L-2)}$ specified by Lemma D.11 for absolute value activation, and we only
1061 specified one of the possible forms of a breakpoint in $\lambda^T \hat{\mathbf{X}}^{(L)}$ as a function of $\mathbf{b}^{(L-3)}$. So a
1062 4-layer network may have additional features not specified.

1063 *Proof of Theorem 3.4.* Follows from Lemma 3.7 and Theorem 3.11 where $\sigma(x) = (x)_+$ is
1064 monotone. ■

1065 *Proof of Corollary 3.15.* Follows from Theorem 3.11 for $L = 2$. ■

1066 **Remark D.19.** For $b_1, b_2 \geq 0$,

$$1067 \quad ||x - b_1| - b_2| = \begin{cases} b_1 - b_2 - x & \text{if } x \leq b_1 - b_2 \\ x - (b_1 - b_2) & \text{if } b_1 - b_2 \leq x \leq b_1 \\ b_1 + b_2 - x & \text{if } b_1 \leq x \leq b_1 + b_2 \\ x - (b_1 + b_2) & \text{if } x \geq b_1 + b_2 \end{cases}$$

$$= \mathcal{W}_{b_1 - b_2, b_1}.$$

1068 *Proof of Lemma B.2, Figure 3.1, Figure 3.4.* Follows from direct computation of Equation (2.2) and application of Remark D.19. In particular, for ReLU symmetrized networks, 1069 it can be verified that for $\mathbf{W}^{(1)} = (-1, 1)$, $\mathbf{W}^{(2)} = (-1, -1)$; $\mathbf{W}^{(1)} = (-1, 1)$, $\mathbf{W}^{(2)} = (1, 1)$; 1070 $\mathbf{W}^{(1)} = (1, -1)$, $\mathbf{W}^{(2)} = (-1, -1)$; and $\mathbf{W}^{(1)} = (1, -1)$, $\mathbf{W}^{(2)} = (1, 1)$, the deep library features 1071 can contain reflections. Note ReLU symmetrized features are consistent with Figure B.1. ■ 1072

1073 **Remark D.20.** For sign-determined activations, by Remark D.1, the inner weights can be 1074 unregularized. So, reconstructed parameters (as defined in Definition B.3) that are unscaled 1075 according to Definition D.17 achieve the same objective in the training problem as the optimal 1076 value of the rescaled problem.

1077 *Proof of Lemma 3.1.* Since $|x| = 2(x)_+ + x$, we have $\sum_{i=1}^{m_2} |\mathbf{X}\mathbf{W}^{(i,1)} + \mathbf{b}^{(i,1)}| \alpha_i + \mathbf{X}\omega + \xi$ 1078 $= \sum_{i=1}^{m_2} 2(\mathbf{X}\mathbf{W}^{(i,1)} + \mathbf{b}^{(i,1)})_+ \alpha_i + \mathbf{X}(\omega - \sum_{i=1}^{m_2} \mathbf{W}^{(i,1)} \alpha_i) - \sum_{i=1}^{m_2} \mathbf{b}^{(i,1)} \alpha_i + \xi$. Given a solution 1079 $\mathbf{W}^{(i,l)*}, \mathbf{b}^{(i,l)*}, \alpha_i^*, \xi^*, \omega^*$ to the training problem for absolute value activation with skip con- 1080 nection, the parameters $\mathbf{W}^{(i,l)} = \mathbf{W}^{(i,l)*}, \mathbf{b}^{(i,l)*} = \mathbf{b}^{(i,l)*}, \alpha_i = 2\alpha_i^*, \xi = \xi^* - \sum_{i=1}^{m_2} \mathbf{b}^{(i,l)*} \alpha_i^*, \omega = \omega^* - \sum_{i=1}^{m_2} \mathbf{W}^{(i,1)*} \alpha_i^*$ 1081 achieve the same objective in the training problem for ReLU activation with 1082 skip connection. Conversely, since $(x)_+ = \frac{|x| + x}{2}$, we have $\sum_{i=1}^{m_2} (\mathbf{X}\mathbf{W}^{(i,1)} + \mathbf{b}^{(i,1)})_+ \alpha_i + \mathbf{X}\omega + \xi$ 1083 $= \frac{1}{2} \sum_{i=1}^{m_2} |\mathbf{X}\mathbf{W}^{(i,1)} + \mathbf{b}^{(i,1)}| \alpha_i + \mathbf{X}(\omega + \frac{1}{2} \sum_{i=1}^{m_2} \mathbf{W}^{(i,1)} \alpha_i) + \frac{1}{2} \sum_{i=1}^{m_2} \mathbf{b}^{(i,1)} \alpha_i + \xi$. Given a solution 1084 $\mathbf{W}^{(i,l)*}, \mathbf{b}^{(i,l)*}, \alpha_i^*, \xi^*, \omega^*$ to the training problem for ReLU activation with skip connection, the 1085 parameters $\mathbf{W}^{(i,l)} = \mathbf{W}^{(i,l)*}, \mathbf{b}^{(i,l)*} = \mathbf{b}^{(i,l)*}, \alpha_i = \frac{1}{2} \alpha_i^*, \xi = \xi^* + \frac{1}{2} \sum_{i=1}^{m_2} \mathbf{b}^{(i,l)*} \alpha_i^*,$ 1086 $\omega = \omega^* + \frac{1}{2} \sum_{i=1}^{m_2} \mathbf{W}^{(i,1)*} \alpha_i^*$ achieve the same objective in the training problem for absolute value 1087 activation with skip connection. Therefore the problems achieve the same optimal value and 1088 we have given a map between the solutions. ■

1089 **Solution sets of Lasso under minimal regularization.**

1090 **Remark D.21.** For each optimal \mathbf{z}^* of the Lasso problem, minimizing the objective over ξ 1091 gives the optimal bias term as $\xi^* = (\mathbf{y} - \mathbf{1}\mathbf{1}^T \mathbf{y}) - (\mathbf{A} - \mathbf{1}\mathbf{1}^T \mathbf{A}) \mathbf{z}^*$.

1092 **Remark D.22.** For a neural net with $L = 2$ layers and sign activation, by Theorem 3.11 the 1093 Lasso problem has an objective function $f(\mathbf{z}) = \frac{1}{2} \|\mathbf{A}\mathbf{z} - \mathbf{y}\|_2^2 + \beta \|\mathbf{z}\|_1$ where $\mathbf{A} \in \mathbb{R}^{N \times N}$. By 1094 Lemma D.25, \mathbf{A} is full rank, which makes f strongly convex. Therefore the Lasso problem has a 1095 unique solution \mathbf{z}^* (6). Moreover, for any Lasso problem, \mathbf{z}^* satisfies the subgradient condition 1096 $\mathbf{0} \in \delta f(\mathbf{z}) = \mathbf{A}^T (\mathbf{A}\mathbf{z}^* - \mathbf{y}) + \beta \partial \|\mathbf{z}^*\|_1$. Equivalently,

$$1097 \quad \frac{1}{\beta} \mathbf{A}_n^T (\mathbf{A}\mathbf{z}^* - \mathbf{y}) \in \begin{cases} \{-\text{sign}(\mathbf{z}_n^*)\} & \text{if } \mathbf{z}_n^* \neq 0 \\ [-1, 1] & \text{if } \mathbf{z}_n^* = 0 \end{cases}, \quad n \in [N].$$

1098 **D.4. Proofs of results in Appendix C.** Let $\mathbf{e}^{(n)} \in \mathbb{R}^N$ be the n^{th} canonical basis vector,
 1099 that is $\mathbf{e}_i^{(n)} = \mathbf{1}\{i = n\}$.

1100 *Proof of Proposition C.2.* We analyze the solution set of $\mathbf{A}\mathbf{z} + \xi\mathbf{1} = \mathbf{y}$. We note that
 1101 $(I - \mathbf{1}\mathbf{1}^T/N)\mathbf{A}\mathbf{z} = (I - \mathbf{1}\mathbf{1}^T/N)\mathbf{y}$. As \mathbf{z}^* is optimal in (C.1), this implies that $(I -$
 1102 $\mathbf{1}\mathbf{1}^T/N)\mathbf{A}\mathbf{z}^* = (I - \mathbf{1}\mathbf{1}^T/N)\mathbf{y}$. This implies that $(I - \mathbf{1}\mathbf{1}^T/N)\mathbf{A}(\mathbf{z} - \mathbf{z}^*) = \mathbf{0}$. As $x_1 > x_2 >$
 1103 $\dots > x_N$, we have $\mathbf{A}(\mathbf{e}^{(1)} + \mathbf{e}^{(N)}) \propto \mathbf{1}$. As \mathbf{A} is invertible by Lemma D.26 in Appendix D.5,
 1104 this implies that there exists $t \in \mathbb{R}$ such that $\mathbf{z} - \mathbf{z}^* = t(\mathbf{e}^{(1)} + \mathbf{e}^{(N)})$. It is impossible to have
 1105 $z_1^* z_N^* > 0$ from the optimality of \mathbf{z}^* . Otherwise, by taking $t = -\text{sign}(z_1^*) \min\{|z_1^*|, |z_N^*|\}$, we have
 1106 $\|\mathbf{z}\|_1 = \|\mathbf{z}^*\|_1 - 2 \min\{|z_1^*|, |z_N^*|\} < \|\mathbf{z}^*\|_1$. Therefore, we have $z_1^* z_N^* \leq 0$. We can reparameterize
 1107 $\mathbf{z} = \mathbf{z}^* + t \text{sign}(z_1^*)(\mathbf{e}^{(1)} + \mathbf{e}^{(N)})$. It is easy to verify that for t such that $-|z_1^*| \leq t \leq |z_N^*|$, we
 1108 have $\|\mathbf{z}\|_1 = \|\mathbf{z}^*\|_1$, while for other choice of t , we have $\|\mathbf{z}\|_1 > \|\mathbf{z}^*\|_1$. Therefore, the solution
 1109 set of (C.1) is given by (C.2). ■

1110 *Proof of Proposition C.3.* Follows from Remark D.22 describing the Lasso objective. ■

1111 *Proof of Proposition C.4.* By Lemma D.27, for $n \in [N - 1]$, $z_{+n}^* = y_n - y_{n-1} \geq 0$ and
 1112 $z_{+N}^* = y_N \geq 0$. So \mathbf{z}^* achieves an objective value of $\|\mathbf{z}^*\|_1 = y_1$ in (C.1). Now let \mathbf{z}
 1113 be any solution to (C.1). Then $\mathbf{A}\mathbf{z} = \mathbf{y}$. Since the first row of \mathbf{A} is $[\mathbf{1}^T, \mathbf{0}^T]$, we have
 1114 $y_1 = (\mathbf{A}\mathbf{z})_1 = \mathbf{1}^T \mathbf{z}_+ \leq \|\mathbf{z}_+\|_1 \leq \|\mathbf{z}\|_1 \leq \|\mathbf{z}^*\|_1 = y_1$. So $\|\mathbf{z}_+\|_1 = \|\mathbf{z}\|_1 = y_1$, leaving $\mathbf{z}_- = \mathbf{0} = \mathbf{z}_-^*$.
 1115 Therefore $\mathbf{z}_+ = \mathbf{A}^{-1} \mathbf{y} = \mathbf{z}_+^*$. Applying Lemma D.27 gives the result. ■

1116 *Proof of Corollary C.1.* By Lemma D.25, Lemma D.26, Lemma D.27 and Lemma D.28,
 1117 the dictionary matrix for the 2-layer net is full rank for sign, absolute value, threshold and
 1118 ReLU activations. The dictionary matrices for deeper nets with sign activation are also full
 1119 rank by Remark 3.12. Let $\mathbf{u} = \mathbf{A}^T(\mathbf{A}\mathbf{z}^* - \mathbf{y})$. By Remark D.22, as $\beta \rightarrow 0$, we have $\mathbf{u} \rightarrow \mathbf{0}$, so
 1120 $\mathbf{A}\mathbf{z} - \mathbf{y} = (\mathbf{A}\mathbf{A}^T)^{-1} \mathbf{A}\mathbf{u} \rightarrow \mathbf{0}$. So as $\beta \rightarrow 0$, the optimal Lasso objective approaches 0, and by
 1121 and Theorem 3.11, so does the training problem. So $f_L(\mathbf{X}; \theta) \xrightarrow{\beta \rightarrow 0} \mathbf{y}$. ■

1122 *Proof of Lemma C.7.* It can be verified that as shown in Figure 3.4, the Lasso features for
 1123 a symmetrized network all have slope magnitude 0, 1, or 2. However, only monotone features
 1124 contain a segment with slope magnitude 2, and the training data (x_n, y_n) in Figure 3.5 is not
 1125 monotone. There is a "left branch" consisting of $\{(x_5, y_5) = (-1, 1), (x_4, y_4) = (0, 0)\}$ and a
 1126 "right branch" consisting of $\{(x_2, y_2) = (3, 1), (x_1, y_1) = (4, 2)\}$. Let \mathbf{z}^*, ξ^* be a Lasso solution
 1127 that fits the data exactly: $\mathbf{A}\mathbf{z}^* + \xi^* = \mathbf{y}$. Let \mathcal{I}^- (\mathcal{I}^+) be the set of indices i where the i^{th} feature
 1128 is monotone and has negative (positive) slope over the left (right) branch. Let \mathcal{I}^0 be the set of
 1129 indices corresponding to features that are not monotone. Let m^- and m^+ be the magnitude
 1130 of the slopes of the left and right branch, respectively. Then $2 \sum_{i \in \mathcal{I}^-} |z_i^*| + \sum_{i \in \mathcal{I}^0} |z_i^*| \geq m^-$
 1131 and $2 \sum_{i \in \mathcal{I}^+} |z_i^*| + \sum_{i \in \mathcal{I}^0} |z_i^*| \geq m^+$. Note $\mathcal{I}^-, \mathcal{I}^+, \mathcal{I}^0$ are all pairwise disjoint. Therefore
 1132 $\|\mathbf{z}^*\|_1 \geq \frac{m^+ + m^-}{2} = 1$. ■

1133 As seen from the proof, Lemma C.7 can be generalized to other training data examples.

1134 *Proof of Lemma C.5.* Since ReLU and absolute value activations have slopes ± 1 or 0, and
 1135 the weights of deep narrow networks are ± 1 , the features $\hat{\mathbf{X}}^{(L)}(x)$ have slopes ± 1 or 0. Observe
 1136 $\mathbf{y} = f_L(\mathbf{X}; \theta) = \xi^* \mathbf{1} + \mathbf{A}\mathbf{z}^* = \xi^* \mathbf{1} + \sum_i z_i^* \mathbf{A}_i$. For any $n \in [N - 1]$, $|\mu_n| = \left| \frac{f_L(\mathbf{X}; \theta)_{n+1} - f_L(\mathbf{X}; \theta)_n}{x_{n+1} - x_n} \right| =$
 1137 $\left| \frac{\sum_i z_i^* (\mathbf{A}_{n+1, i} - \mathbf{A}_{n, i})}{x_{n+1} - x_n} \right| = \left| \sum_i z_i^* \frac{\hat{\mathbf{X}}^{(L)}(x_{n+1}) - \hat{\mathbf{X}}^{(L)}(x_n)}{x_{n+1} - x_n} \right|$

$$\leq \sum_i |z_i^*| \left| \frac{\tilde{\mathbf{X}}^{(L)}(x_{n+1}) - \tilde{\mathbf{X}}^{(L)}(x_n)}{x_{n+1} - x_n} \right| \leq \sum_i |z_i^*| = \|\mathbf{z}^*\|. \quad \blacksquare$$

Proof of Lemma C.6. Let $n \in [N - 1]$. Let $\mathcal{S}_n^+ = \{i \in [N] : i > n, (z_+)_i \neq 0\}$, $\mathcal{S}_n^- = \{i \in [N] : i \leq n, (z_-)_i \neq 0\}$. Observe $\mathcal{S}_{n+1}^+ = \mathcal{S}_n^+ - \{n+1\}$ if $(z_+)_{n+1} \neq 0$ and $\mathcal{S}_{n+1}^+ = \mathcal{S}_n^+$ otherwise. Similarly, $\mathcal{S}_{n+1}^- = \mathcal{S}_n^- \cup \{n+1\}$ if $(z_-)_{n+1} \neq 0$ and $\mathcal{S}_{n+1}^- = \mathcal{S}_n^-$ otherwise. Now, $\text{ReLU}_{x_i}^+$ has slope 1 after x_i and $\text{ReLU}_{x_i}^-$ has slope 1 before x_i , so $\mu_n = \sum_{i \in \mathcal{S}_n^+} (z_+)_i + \sum_{i \in \mathcal{S}_n^-} (z_-)_i$. Therefore,

$$|\mu_n - \mu_{n+1}| = |-(z_+)_{n+1} + (z_-)_{n+1}| \leq |(z_+)_{n+1}| + |(z_-)_{n+1}|. \text{ So } \sum_{n=1}^{N-1} |\mu_n - \mu_{n+1}| \leq \|\mathbf{z}^*\|_1 - |(z_+)_{1}| - |(z_-)_{1}| \leq \|\mathbf{z}^*\|_1.$$

Now, for any $n \in [N - 1]$, $(\mathbf{A}\mathbf{z} + \xi \mathbf{1})_n - (\mathbf{A}\mathbf{z} + \xi \mathbf{1})_{n+1} = \sum_{i=1}^N (z_+)_i (\mathbf{A}_{+n,i} - \mathbf{A}_{+n+1,i}) = \sum_{i=1}^N (z_+)_i ((x_n - x_i)_+ - (x_{n+1} - x_i)_+) = \sum_{i=n+1}^N (\mu_{i-1} - \mu_i) (x_n - x_{n+1}) = (x_n - x_{n+1}) \mu_n = y_n - y_{n+1}$. And $(\mathbf{A}\mathbf{z} + \xi \mathbf{1})_N = \xi + \sum_{i=1}^N (z_+)_i (x_N - x_i)_+ = y_N + \sum_{i=1}^N (z_+)_i (0) = y_N$. So $\mathbf{A}\mathbf{z} + \xi \mathbf{1} = \mathbf{y}$. And $\|\mathbf{z}\|_1 = \sum_{n=1}^{N-1} |\mu_n - \mu_{n+1}|$, which exactly hits the lower bound on $\|\mathbf{z}^*\|_1$. Therefore \mathbf{z}, ξ is optimal. \blacksquare

Remark D.23. By Lemma D.26, the absolute value network “de-biases” the target vector, normalizes it by the interval lengths $x_i - x_{i+1}$, and applies \mathbf{E} (which contains the difference matrix Δ) twice, acting as a second-order difference detector. By Lemma D.25, the sign network’s dictionary inverse contains Δ just once, acting as a first-order difference detector.

Inverses of 2-layer dictionary matrices.

In this section, we consider the 2-layer dictionary matrix \mathbf{A} as defined in Corollary 3.15. Define the finite difference matrix

$$\Delta = \begin{pmatrix} 1 & -1 & 0 & \cdots & 0 & 0 \\ 0 & 1 & -1 & \cdots & 0 & 0 \\ 0 & 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 1 & -1 \\ 0 & 0 & 0 & \cdots & 0 & 1 \end{pmatrix} \in \mathbb{R}^{N-1 \times N-1}.$$

Multiplying a matrix on its right by Δ subtracts its consecutive rows. Define the diagonal matrix $\mathbf{D} \in \mathbb{R}^{N \times N}$ by $\mathbf{D}_{i,i} = \frac{1}{x_i - x_{i+1}}$ for $i \in [N - 1]$ and $\mathbf{D}_{N,N} = \frac{1}{x_1 - x_N}$. For $n \in \mathbb{N}$, let

$$\mathbf{A}_n^{(s)} = \begin{pmatrix} 1 & 1 & 1 & \cdots & 1 \\ -1 & 1 & 1 & \cdots & 1 \\ -1 & -1 & 1 & \cdots & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ -1 & -1 & -1 & \cdots & 1 \end{pmatrix}, \mathbf{A}_n^{(t)} = \begin{pmatrix} 1 & 1 & 1 & \cdots & 1 \\ 0 & 1 & 1 & \cdots & 1 \\ 0 & 0 & 1 & \cdots & 1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \end{pmatrix} \in \mathbb{R}^{n \times n}.$$

Remark D.24. The dictionary matrices for sign and threshold activation satisfy $\mathbf{A} = \mathbf{A}_N^{(s)}$ and $\mathbf{A}_+ = \mathbf{A}_N^{(t)}$, respectively.

D.5. Results about 2-layer dictionary matrices.

Lemma D.25. The dictionary matrix for $\sigma(x) = \text{sign}(x)$ has inverse $\mathbf{A}^{-1} = \frac{1}{2} \left(\begin{array}{c|c} & 0 \\ \Delta & \vdots \\ & 0 \\ \hline & -1 \\ 1 & 0 & \cdots & 0 & 1 \end{array} \right).$

1165 *Proof.* Multiplying the two matrices (see Remark D.24) gives the identity. ■

1166 **Lemma D.26.** *The dictionary matrix \mathbf{A} for absolute value activation has inverse $\mathbf{A}^{-1} =$*
 1167 $\frac{1}{2}\mathbf{PEDE}$, where $\mathbf{P} = \left(\begin{array}{cccc|c} 0 & 0 & \cdots & 0 & 1 \\ & & & & 0 \\ & & & & \vdots \\ & & & & 0 \\ & & & & -1 \\ & & & & 0 \end{array} \right)$, $\mathbf{E} = \left(\begin{array}{cccc|c} & & & & 0 \\ & & & & \vdots \\ & & & \Delta & 0 \\ & & & & -1 \\ -1 & 0 & \cdots & 0 & -1 \end{array} \right)$.

1168 *Proof.* For $i \in [N-1], j \in [N]$, $\mathbf{A}_{i,j} - \mathbf{A}_{i+1,j} = \begin{cases} (x_j - x_i) - (x_j - x_{i+1}) = x_{i+1} - x_i & \text{if } i > j \\ (x_i - x_j) - (x_{i+1} - x_j) = x_i - x_{i+1} & \text{if } i \leq j. \end{cases}$
 1169 And for all $j \in [N]$, $\mathbf{A}_{1,j} + \mathbf{A}_{N,j} = (x_1 - x_j) + (x_j - x_N) = x_1 + x_N$. Therefore,

1170
$$\mathbf{DEA} = \left(\begin{array}{c|cccc} -1 & & & & \\ \vdots & & & & \\ -1 & & \mathbf{A}_{N-1}^{(s)} & & \\ -1 & -1 & \cdots & -1 & \end{array} \right), \quad \frac{1}{2}\mathbf{EDEA} = \left(\begin{array}{c|cccc} 0 & & & & \\ \vdots & & & & \\ 0 & & & I_{N-1} & \\ 1 & 0 & \cdots & 0 & \end{array} \right).$$

1171 Applying the permutation \mathbf{P} makes $\frac{1}{2}\mathbf{PEDEA} = \mathbf{I}$, so $\mathbf{A}^{-1} = \frac{1}{2}\mathbf{PEDE}$. ■

1172 **Lemma D.27.** *The inverse of \mathbf{A}_+ for threshold activation is $\mathbf{A}_+^{-1} = \left(\begin{array}{cccc|c} & & & & 0 \\ & & & & \vdots \\ & & & \Delta & 0 \\ & & & & -1 \\ 0 & 0 & \cdots & 0 & 1 \end{array} \right)$.*

1173 *Proof.* Multiplying the two matrices (see Remark D.24) gives the identity. ■

1174 **Lemma D.28.** *The submatrix $[(\mathbf{A}_+)_{1:N,2:N}, (\mathbf{A}_-)_{1:N,1}] \in \mathbb{R}^{N \times N}$ of the dictionary matrix for*
 1175 *ReLU activation has inverse $\mathbf{E}_+\mathbf{DE}_-$, where*

1176
$$\mathbf{E}_+ = \left(\begin{array}{c|ccc} & & & 0 \\ & & & \vdots \\ & & \Delta & 0 \\ & & & 1 \\ 0 & 0 & \cdots & 0 & 1 \end{array} \right), \quad \mathbf{E}_- = \left(\begin{array}{c|ccc} & & & 0 \\ & & & \vdots \\ & & \Delta & 0 \\ & & & -1 \\ 0 & 0 & \cdots & 0 & -1 \end{array} \right).$$

1177 *Proof.* For $i \in [N-1], j \in [N]$,

1178
$$(\mathbf{A}_+)_{i,j} - (\mathbf{A}_+)_{i+1,j} = \begin{cases} 0 & \text{if } i \geq j \\ (x_i - x_j) - (x_{i+1} - x_j) = x_i - x_{i+1} & \text{if } i < j \end{cases}$$

1179 and $(\mathbf{A}_-)_{i,1} - (\mathbf{A}_-)_{i+1,1} = (x_1 - x_i) - (x_1 - x_{i+1}) = x_{i+1} - x_i$. Observe that $\mathbf{DE}_-\mathbf{A} =$

1180
$$\left(\begin{array}{c|ccc} & & & -1 \\ & & & \vdots \\ & & \mathbf{A}_{N-1}^{(t)} & -1 \\ & & & -1 \\ 0 & 0 & \cdots & 0 & 1 \end{array} \right)$$
, and applying \mathbf{E}_+ gives \mathbf{I} . ■

1181 **Appendix E. The solution sets of Lasso and the training problem.**

1182 Recall that each 2-layer feature corresponds to $\mathbf{W}^{(1)} \in \{-1, 1\}$, $\mathbf{b}^{(1)} = -\mathbf{W}^{(1)}x_n$ for some
 1183 $n \in [N]$. We next define functions $R^{z \rightarrow \alpha}, R^{\alpha, w, b \rightarrow \theta}, R^{w, b}, R$ that are composed with each other
 1184 to reconstruct 2-layer networks. Each feature in a 2-layer network is of the form $f(x) =$
 1185 $\sigma(x\mathbf{W}^{(1)} + \mathbf{b}^{(1)})$ where $\mathbf{W}^{(1)} \in \{-1, 1\}$ and $\mathbf{b}^{(1)} = -\mathbf{W}^{(1)}x_n$ for some $n \in [N]$.

1186 **Definition E.1.** Consider a 2-layer ReLU, absolute value, or leaky ReLU network. Let
 1187 $R^{z \rightarrow \alpha}(z) = \text{sign}(z) \sqrt{|z|}$. For $w \in \{-1, 1\}$, $b \in \{-x_n : n \in [N]\}$, let $R^{\alpha, w, b \rightarrow \theta}(\alpha) = (\alpha, \alpha w, -\alpha w b)$.
 1188 Let $R^{w, b}(z) = R^{\alpha, w, b \rightarrow \theta}(R^{z \rightarrow \alpha}(z))$. Define the 2-layer reconstruction function $R(\mathbf{z})$ which outputs
 1189 a vector whose i^{th} element is $R^{w, b}(z_i)$, where $(w, b) = (\mathbf{W}^{(1)}, \mathbf{b}^{(1)})$ corresponds to the i^{th} feature.

1190 In **Definition E.1**, we have abused notation where α, z, w in the superscripts of functions are
 1191 only used to identify the function, and are distinct from the arguments to the functions. The
 1192 function $R^{\alpha, w, b \rightarrow \theta}$ is parameterized by w, b . For a 2-layer network, let $w_i = \mathbf{W}^{(i,1)}, b_i = \mathbf{b}^{(i,1)} \in \mathbb{R}$
 1193 and $\mathbf{w}, \mathbf{b}, \boldsymbol{\alpha}$ be vectors stacking together all w_i, b_i, α_i respectively. The reconstructed parameters
 1194 are $(\boldsymbol{\alpha}, \mathbf{w}, \mathbf{b}) = R(\mathbf{z}^*)$ and $\xi = \xi^*$.

1195 For a 2-layer network, let $w_i = \mathbf{W}^{(i,1)}, b_i = \mathbf{b}^{(i,1)} \in \mathbb{R}$ and $\mathbf{w}, \mathbf{b}, \boldsymbol{\alpha}$ be vectors stacking together
 1196 all w_i, b_i, α_i respectively. The reconstructed parameters are $(\boldsymbol{\alpha}, \mathbf{w}, \mathbf{b}) = R(\mathbf{z}^*)$ and $\xi = \xi^*$. We
 1197 have shown that training neural networks on 1-D data is equivalent to fitting a Lasso model.
 1198 Now we develop analytical expressions for all minima of the Lasso problem and its relationship
 1199 to the set of all minima of the training problem. These results, which build on the existing
 1200 literature for convex reformulations (28) as well as characterizations of the Lasso (10), illustrate
 1201 that the Lasso model provides insight into non-convex networks. We focus on 2-layer models
 1202 with ReLU, leaky ReLU and absolute value activations, although our results can be extended
 1203 to other architectures by considering the corresponding neural net reconstruction. Proofs are
 1204 deferred to **Appendix E.1**.

1205 We start by characterizing the set of global optima to the Lasso problem (1.2). Suppose
 1206 (\mathbf{z}^*, ξ^*) is a solution to the convex training problem. We first extend classical uniqueness results
 1207 (41; 40) for the model fit to the Lasso with unregularized bias.

1208 **Lemma E.2.** Fix $\beta > 0$. Then the optimal model fit $\hat{\mathbf{y}} = \mathbf{A}\mathbf{z}^* + \xi^*\mathbf{1}$ and linear fit $\mathbf{A}\mathbf{z}^*$ are
 1209 the same for all optimal solutions (\mathbf{z}^*, ξ^*) to Lasso problem (1.2). As a result, ξ^* is also unique.

1210 As a result, we can uniquely identify the equicorrelation set \mathcal{E}_β using $\hat{\mathbf{y}}$,

$$1211 \quad \hat{\mathbf{y}} = \mathbf{A}\mathbf{z}^* + \xi^*\mathbf{1}, \quad \mathcal{E}_\beta = \left\{ i : |\mathbf{A}_i^\top(\hat{\mathbf{y}} - \mathbf{y})| = \beta \right\}.$$

1212 The equicorrelation set contains the features maximally correlated with the residual $\hat{\mathbf{y}} - \mathbf{y}$ and
 1213 plays a critical role in the solution set.

1214 **Proposition E.3.** Suppose $\beta > 0$. Then the set of global optima of the Lasso problem (1.2) is
 (E.1)

$$1215 \quad \Phi^*(\beta) = \left\{ (\mathbf{z}, \xi^*) : z_i \neq 0 \Rightarrow \text{sign}(z_i) = \text{sign}(\mathbf{A}_i^\top(\hat{\mathbf{y}} - \mathbf{y})), z_i = 0 \forall i \notin \mathcal{E}_\beta, \mathbf{A}\mathbf{z} = \hat{\mathbf{y}} - \xi^*\mathbf{1} \right\}$$

1216 The solution set $\Phi^*(\beta)$ is polyhedral and its vertices correspond exactly to minimal models,
 1217 i.e. models with the fewest non-zero elements of \mathbf{z} (28). Let R be the reconstruction function
 1218 described in **Definition E.1**. All networks generated from applying R to a Lasso solution are
 1219 globally optimal in the training problem. The next result gives a full description of such
 1220 networks.

1221 **Proposition E.4.** Suppose $\beta > 0$ and the activation is ReLU, leaky ReLU or absolute value.

1222 The set of all 2-layer Lasso-reconstructed networks is

$$1223 \quad (E.2) \quad R(\Phi(\beta)) = \left\{ (\mathbf{w}, \mathbf{b}, \alpha, \xi^*) : \alpha_i \neq 0 \Rightarrow \text{sign}(\alpha_i) = \text{sign} \left(\mathbf{A}_i^\top (\mathbf{y} - \hat{\mathbf{y}}) \right), b_i = -x_i \frac{\tilde{\alpha}_i}{\sqrt{|\alpha_i|}}, \right. \\ \left. w_i = \frac{\tilde{\alpha}_i}{\sqrt{|\alpha_i|}}; \alpha_i = 0 \forall i \notin \mathcal{E}_\beta, f_2(\mathbf{X}; \theta) = \hat{\mathbf{y}} \right\}.$$

1224 **Proposition E.4** shows that all neural nets trained using our Lasso and reconstruction share
1225 the same model fit whose set of active neurons is at most the equicorrelation set. By finding
1226 just *one* optimal neural net that solves Lasso, we can form $R(\Phi^\beta)$ and compute all others.

1227 The min-norm solution path is continuous for the Lasso problem (40). Since the solution
1228 mapping in **Definition B.3**, **Appendix D.3** is continuous, the corresponding reconstructed
1229 neural net path is also continuous as long as the network is sufficiently wide. Moreover, we
1230 can compute this path efficiently using the LARS algorithm (10). This is in contrast to the
1231 under-parameterized setting, where the regularization path is discontinuous (28).

1232 What subset of optimal, or more generally, stationary, points of the non-convex training
1233 problem (1.1) consist of Lasso-generated networks $R(\Phi(\beta))$? First, $R(\Phi(\beta))$ can generate
1234 additional optimal networks through neuron splitting, described as follows. Consider a single
1235 neuron $\alpha \sigma_s(wx + b)$ (where $\alpha, w, b \in \mathbb{R}$), and let $\{\gamma_i\}_{i=1}^n \subset [0, 1]^n$ be such that $\sum_{i=1}^n \gamma_i = 1$.
1236 The neuron can be *split* into n neurons $\{\sqrt{\gamma_i} \alpha \sigma(\sqrt{\gamma_i} wx + \sqrt{\gamma_i} b)\}_{i=1}^n$ (45). For any collection
1237 Θ of parameter sets θ , let $P(\Theta)$ be the collection of parameter sets generated by all possible
1238 neuron splits and permutations of each $\theta \in \Theta$. Next, let $\mathcal{C}(\beta)$ and $\tilde{\mathcal{C}}(\beta)$ be the sets of Clarke
1239 stationary points and solutions to the non-convex training problem (1.1), respectively.

1240 **Proposition E.5.** Suppose $L = 2, \beta > 0$, the activation is ReLU, leaky ReLU or absolute
1241 value and $m^* \leq m \leq 2N$. Let $\Theta^P = \{\theta : \forall i \in [m], \exists j \in [N] \text{ s.t. } b_i = -x_j w_i\}$. Then

$$1242 \quad (E.3) \quad P(R(\Phi(\beta))) = \tilde{\mathcal{C}}(\beta) \cap \Theta^P = \mathcal{C}(\beta) \cap \Theta^P.$$

1243 **Proposition E.5** states that up to neuron splitting and permutation, our Lasso method
1244 gives all stationary points in the training problem satisfying $b_i = -x_i w_i$. Moreover, all such
1245 points are optimal in the training problem, similar to (17).

1246 Since optimal solutions are stationary, a neural net reconstructed from the Lasso model is
1247 in $\tilde{\mathcal{C}}(\beta) \subset \mathcal{C}(\beta)$. However, $\mathcal{C}(\beta) \not\subset \Theta^P$. This is because there may be other neural nets with
1248 the same output on \mathbf{X} as the reconstructed net so that they are all in $\mathcal{C}(\beta)$, but that differ
1249 in the the unregularized parameters \mathbf{b} and ξ , so that they are not in Θ^P . For example, if β
1250 is large enough, the Lasso solution is $\mathbf{z} = \mathbf{0}$ (10), so the reconstructed net will have $\boldsymbol{\alpha} = \mathbf{0}$,
1251 which makes the neural net invariant to \mathbf{b} . In this section, we analyzed the general structure of
1252 the Lasso solution set when $\beta > 0$. Next, we show how this can be leveraged to explore the
1253 generalization performance of different optimal neural networks.

1254 **Generalization properties of networks constructed from different Lasso solutions.** While
1255 all two-layer Lasso-reconstructed networks must have the same model fit for a fixed training
1256 set (**Proposition E.4**), they may have very different predictions on an independent test set. In
1257 particular, some models may generalize better at test time due to the “implicit regularization”
1258 (20) of the rule used to choose a model from the optimal set. Studying implicit regularization for

1259 different problems and selection rules has been a topic of significant interest in the optimization
 1260 community in recent years (30; 2; 19). However, our characterization of the optimal set
 1261 introduces a completely new perspective that we can (at least in principle) compute the
 1262 generalization of all optimal Lasso solutions simultaneously.

1263 We approximate the distribution of test mean-square error (MSE) over the set of optimal
 1264 three-layer, deep narrow ReLU networks reconstructed from Lasso solutions by sampling optimal
 1265 solutions uniformly from the Lasso solution set. Since the Lasso optimal set is polyhedral, it is
 1266 possible to sample optimal Lasso models uniformly at random using a hit-and-run algorithm (8).
 1267 These samples are then reconstructed into non-convex neural network models and evaluated
 1268 on a test set. In this way, we can approximate not just the implicit regularization of specific
 1269 solutions, but also approximate the entire generalization range of the optimal Lasso set.

1270 Figure E.1 shows the distribution of test MSEs for a planted affine model and a planted
 1271 three-layer narrow ReLU network. We estimate each distribution using 5000 hit-and-run
 1272 samples as described above. For each synthetic problem, we generate a set of 5 training
 1273 examples and 100 test examples by sampling $x_i \sim \mathcal{N}(0, 1)$ and then computing $y_i = ax_i + b + \epsilon_i$
 1274 or $y_i = f_3(x_i; \theta) + \epsilon_i$ for the planted affine and planted ReLU models, respectively. Here,
 1275 $\epsilon_i \sim \mathcal{N}(0, 0.1^2)$ and θ is chosen so that $f_3(x_i; \theta)$ depends on only five features from the deep
 1276 library.

1277 We find that test performance of different optimal neural networks varies widely despite the
 1278 models having exactly the same training MSE and model norm. For the planted ReLU problem,
 1279 the variation over the optimal set is even greater than the variation due to the regularization
 1280 strength, implying that implicit regularization may play an even more important role than
 1281 hyper-parameter selection for some problems. We also find that the min-norm solution (marked
 1282 with a red star) tends to the median generalization performance. This surprising observation
 1283 contradicts the common wisdom that min-norm solutions have superior generalization, at least
 1284 for networks reconstructed from Lasso solutions.

1285 The solution sets of Lasso and the training problem.

1286 E.1. Proofs for results in Appendix E.

1287 Lemma E.2. Uniqueness of $\hat{\mathbf{y}}$ is a consequence of strict convexity of the squared-error loss.
 1288 We refer to (40, Lemma 1) for a proof of this result. To see that $\hat{\mathbf{q}} = \mathbf{Az}^*$ is also unique,
 1289 consider minimizing over ξ in (1.2). Since ξ is unregularized, this is a quadratic minimization
 1290 problem with solution

$$1291 \quad \xi^* = \bar{\mathbf{y}} - \frac{1}{n} \langle \mathbf{1}, \mathbf{Az} \rangle$$

$$1292 \quad = \bar{y} - \bar{q},$$

1293 where we have defined $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ and $\bar{q} = \frac{1}{n} \sum_{i=1}^n \hat{q}_i$. Plugging this back into the lasso
 1294 gives an optimization problem purely in \mathbf{z} ,

$$1295 \quad \min_{\mathbf{z}} \frac{1}{2} \left\| \mathbf{Az} - \frac{1}{n} \langle \mathbf{1}, \mathbf{Az} \rangle \mathbf{1} - (\mathbf{y} - \bar{y} \mathbf{1}) \right\|_2^2 + \tilde{\beta} \|\mathbf{z}\|_1.$$

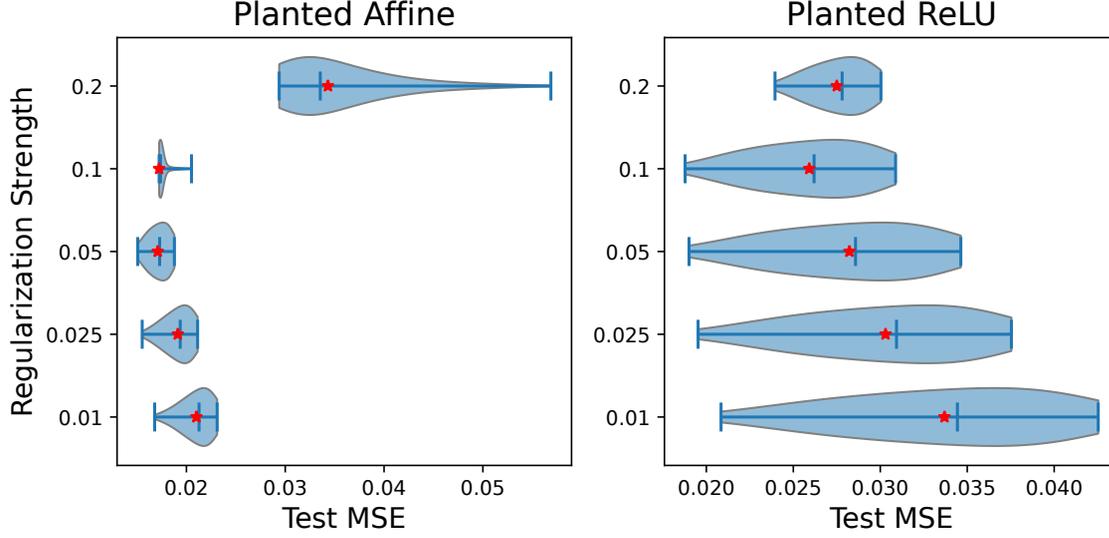


Figure E.1: Distribution of 5000 test MSEs obtained by sampling uniformly from the optimal set for three-layer narrow ReLU networks. We experiment with two synthetic 1-D datasets, one with a planted noisy affine model (left) and one with a planted noisy three-layer narrow ReLU network (right). The red star marks the test MSE of the minimum Euclidean-norm solution to the Lasso reformulation, while the bars mark the minimum, median, and maximum sample test MSEs, respectively.

1296 The squared-error term of the objective can be parameterized in terms of $\hat{\mathbf{q}}$ as

1297
$$\frac{1}{2} \left\| \hat{\mathbf{q}} - \frac{1}{n} \langle \mathbf{1}, \hat{\mathbf{q}} \rangle \mathbf{1} - (\mathbf{y} - \bar{y} \mathbf{1}) \right\|_2^2.$$

1298 Clearly this objective is convex in $\hat{\mathbf{q}}$, but it's not yet obvious if it is strictly convex. The Hessian
1299 of the squared-error is given by,

1300
$$\mathbf{H} = \mathbf{I} - \frac{1}{n} \mathbf{1} \mathbf{1}^\top.$$

1301 For any $\mathbf{w} \in \mathbb{R}^n$, we have,

1302
$$\begin{aligned} \mathbf{w}^\top \mathbf{H} \mathbf{w} &= \|\mathbf{w}\|_2^2 - \frac{1}{n} (\langle \mathbf{1}, \mathbf{w} \rangle)^2 \\ 1303 &\geq \|\mathbf{w}\|_2^2 - \frac{1}{n} \|\mathbf{1}\|_2^2 \|\mathbf{w}\|_2^2 \\ 1304 &= 0, \end{aligned}$$

1305 by Cauchy-Schwarz. Moreover, equality holds if and only if $\mathbf{w} \propto \mathbf{1}$. This implies that the
1306 squared-loss is strictly convex along the cord between $\hat{\mathbf{q}}_1$ and $\hat{\mathbf{q}}_2$ if and only if $\hat{\mathbf{q}}_1, \hat{\mathbf{q}}_2 \not\propto \mathbf{1}$.

1307 Suppose an optimal fit satisfies $\hat{\mathbf{q}} \propto \mathbf{1}$ with $\mathbf{z}^* \neq 0$. Then by setting $\xi' = \xi^* + \hat{\mathbf{q}}_0$ and
1308 $\mathbf{z}' = \mathbf{0}$ we can strictly decrease the lasso objective (since \mathbf{z} is regularized and ξ is not). But this

1309 contradicts optimality of \mathbf{z}^* . We conclude that no optimal linear fit $\hat{\mathbf{q}}$ is proportional to the
1310 one vector. Combining this with the preceding argument shows that the least-squares objective
1311 is strictly convex in $\hat{\mathbf{q}}$. It is now straightforward to show that the optimal $\hat{\mathbf{q}}$ is unique following
1312 a similar argument as (40, Lemma 1). ■

1313 *Proposition E.3.* The result is almost a sub-case of that given by Mishkin and Pilanci
1314 (28) with the exception that the bias parameter ξ , is not regularized. Therefore optimality
1315 conditions do not impose a sign constraint and it is sufficient that $\mathbf{1}^\top(\mathbf{A}\mathbf{z} + \xi\mathbf{1} - \mathbf{y}) = 0$ for ξ
1316 to be optimal. This stationarity condition is guaranteed by $\mathbf{A}\mathbf{z} + \xi\mathbf{1} = \hat{\mathbf{y}}$. In Lemma E.2, we
1317 showed that these optimality conditions admit a unique solution in $\mathbf{A}\mathbf{z}$ and ξ . Thus, we can
1318 conclude ξ^* is the unique optimal output bias and $\hat{\mathbf{y}} - \xi^*\mathbf{1}$ is unique. Now let us look at the
1319 parameters z_i . If $i \notin \mathcal{E}(\beta)$, then $z_i = 0$ is necessary and sufficient from standard results on the
1320 Lasso (40). If $i \in \mathcal{E}(\beta)$ and $z_i \neq 0$, then $\mathbf{A}_i^\top(\hat{\mathbf{y}} - \mathbf{y}) = \beta \text{sign}(z_i)$, which shows that \mathbf{z}_i satisfies
1321 first-order conditions. If $z_i = 0$, then first-order optimality is immediate since $|\mathbf{A}_i^\top(\hat{\mathbf{y}} - \mathbf{y})| \leq \beta$,
1322 holds. Putting these cases together completes the proof. ■

1323 *Proposition E.4.* This result follows from applying the reconstruction in Definition B.3 to
1324 each optimal point in Φ . The reconstruction sets $\alpha_i = \text{sign}(z_i)\sqrt{|z_i|}$. From this we deduce
1325 $\text{sign}(\alpha_i) = \text{sign}(\mathbf{A}_i^\top(\mathbf{y} - \hat{\mathbf{y}}))$. The solution mapping determines the values of w_i and b_i and in
1326 terms of α_i . Finally, the constraint $f_2(\mathbf{X}; \theta) = \hat{\mathbf{y}}$ follows immediately by equality of the convex
1327 and non-convex prediction functions on the training set. ■

1328 *Lemma E.6.* Suppose $L = 2$, and the activation is ReLU, leaky ReLU or absolute value.
1329 Suppose $m^* \leq m \leq 2N$. Since $m \leq 2N$, we can let $\Theta^{\text{Lasso,stat}} = \{\theta : \exists j \in [N] \text{ s.t. } b_i =$
1330 $-x_j w_i\} \subset \Theta$. Let $\theta^* \in \Theta^{\text{Lasso,stat}} \cap C(\beta)$. Then θ^* is a minima of the Lasso problem.

1331 *Proof.* We can ignore ξ since its derivative and reconstruction are straightforward, and it
1332 does not interact with any other parameters. We denote vector-vector operations as being
1333 performed elementwise.

1334 Since $m^* \leq m$, a neural net reconstructed from Lasso is optimal in the training problem.
1335 Let $F(\theta)$ and $F^{\text{Lasso}}(\mathbf{z})$ be the objectives of the non-convex training problem (1.1) and Lasso
1336 (1.2), respectively. The parameters θ are stationary if $\theta \in C(\beta)$, i.e., $0 \in \partial F(\theta)$.

1337 Let $\Theta^{\text{Lasso}} = \{\theta : \exists j \in [N] \text{ s.t. } |w_i| = |\alpha_i|, b_i = -x_j w_i\} \subset \Theta^{\text{Lasso,stat}}$. By a similar argument as
1338 the proof of Theorem 3 in (45), since $0 \in \partial F(\theta^*)$, we have $|w_i| = |\alpha_i|$ for all neurons i . Therefore
1339 $\theta^* \in \Theta^{\text{Lasso}}$. Thus for $\alpha^* \in \theta^*$, observe that $\theta^* = R^{\alpha, w, b \rightarrow \theta}(\alpha^*)$. Let $\tilde{F}(\alpha) = F(R^{\alpha, w, b \rightarrow \theta}(\alpha))$.

1340 Since $0 \in \partial F(\theta^*)$ at $(\alpha^*) = (R^{\alpha, w, b \rightarrow \theta})^{-1}(\theta^*)$, we have $\tilde{F}(\alpha^*) = F(\theta^*)$. Perform the follow-
1341 ing operations elementwise. The chain rule gives $\partial \tilde{F}(\alpha^*) = \partial F(\theta^*) \partial R^{\alpha, w, b \rightarrow \theta}(\alpha^*) \ni \mathbf{0}$. Let
1342 $R^{\alpha \rightarrow z}(\alpha) = \text{sign}(\alpha)\alpha^2$. At $\mathbf{z}^* = R^{\alpha \rightarrow z}(\alpha^*)$, we have $F^{\text{Lasso}}(\mathbf{z}^*) = \tilde{F}(\alpha^*)$. The chain rule gives
1343 $\partial F^{\text{Lasso}}(\mathbf{z}^*) = \partial \tilde{F}(\alpha^*) \partial R(\mathbf{z}^*) \ni \mathbf{0}$. Since the Lasso problem (1.2) is convex, the result holds. ■

1344 *Proof of Proposition E.5.* Observe that $R(\Phi(\beta)) \subseteq \tilde{\mathcal{C}}(\beta) \cap \Theta^{\text{Lasso,stat}} \subseteq \mathcal{C}(\beta) \cap \Theta^{\text{Lasso,stat}} \subseteq$
1345 $R(\Phi(\beta))$, where the first and last subset inequality follow from Theorem 3.11 and Lemma E.6,
1346 respectively. Therefore all subsets in the above expression are equal. Observe that
1347 $P(\Theta^{\text{Lasso,stat}}) = \Theta^P$ and so $P(C(\beta) \cap \Theta^{\text{Lasso,stat}}) = C(\beta) \cap P(\Theta^{\text{Lasso,stat}}) = C(\beta) \cap \Theta^P$ and sim-
1348 ilarly $P(\tilde{\mathcal{C}}(\beta) \cap \Theta^{\text{Lasso,stat}}) = \tilde{\mathcal{C}}(\beta) \cap \Theta^P$. Now apply P to all subsets above. ■

1349 **E.2. Detailed results of Section 4.** Unless otherwise stated, all experiments train for 100
1350 epochs. In Figure 4.3, Figure 4.5, Figure 4.6 and Figure 4.7, the number of time steps is
1351 $T = 1001$, the deep narrow architecture is specified by $m_L = 10$ parallel units, and $\beta = \frac{1}{T}$ in
1352 the training problem (4.3) and (1.1) for more general L layers.

1353 In Figure 4.1, Figure 4.2 and Figure 4.4, we solve the autoregression problems with 20 trials
1354 of stochastic gradient descent (SGD) initializations versus using the Lasso problem (1.2). We
1355 compare against the baseline linear method $f(x; \theta) = ax + b$ (AR1+bias), where we include an
1356 additional bias term b . For each dataset, training and test sets with T data points each are
1357 chosen. In Figure 4.1 and Figure 4.2, the learning rate is 0.05 and $T = 2000$. In Figure 4.4,
1358 the learning rate is 0.1. Figure 4.1 shows that using Lasso(cvxNN) yields a consistent lower
1359 bound on the training loss and demonstrates strong generalization properties, compared to
1360 large fluctuation in the loss curves of the non-convex method (NN). Additional results are in
1361 the supplementary material (47).

1362 Appendix E.2 discusses planted data, for example shown in the regularization path plotted
1363 in Figure 4.4. Figure E.3 shows training curves for networks for planted data with $\sigma^2 = 0.01$
1364 in (4.1) and learning rate 0.1. In Figure 4.4 and Figure E.3, $T = 1000$ time steps is used for
1365 generating each of the planted training and testing data. Figure E.3 shows that using Lasso
1366 gives lower training loss when using planted data. This appears to occur because different trials
1367 of NN (neural net trained directly without Lasso) get stuck into local minima. The global
1368 optimum that cvxNN reaches also enjoys effective generalization properties, as seen by the test
1369 loss. In Figure 4.1, Figure 4.2, Figure 4.4, and Figure E.3, $m = 200$ neurons are trained over
1370 100 epochs and $\beta = 1$ in the training problem (4.3).

1371 Figure E.2 shows that even for the same number of features, the training time increases
1372 with depth. Autoregressive deep narrow neural networks are trained on $T = 1000$ samples
1373 of BTC-2017min data with the Lasso program, with 10^3 features subsampled. The random
1374 subsampling is repeated 40 times according to the procedure described in Section 4.

1375 A *quantile regression* problem is the training problem (4.3) with the quantile loss $\mathcal{L} =$
1376 $L_\tau(z) = 0.5|z| + (\tau - 0.5)z$, where the network $f_2^{\text{NN}}(x_t; \theta)$ models the τ -quantile of the prediction
1377 for x_{t+1} based on the observation x_t . Additional results and comparisons including applying
1378 the same experiments to quantile regression and hourly BTC price (BTC-hourly) are in the
1379 supplementary material (47).

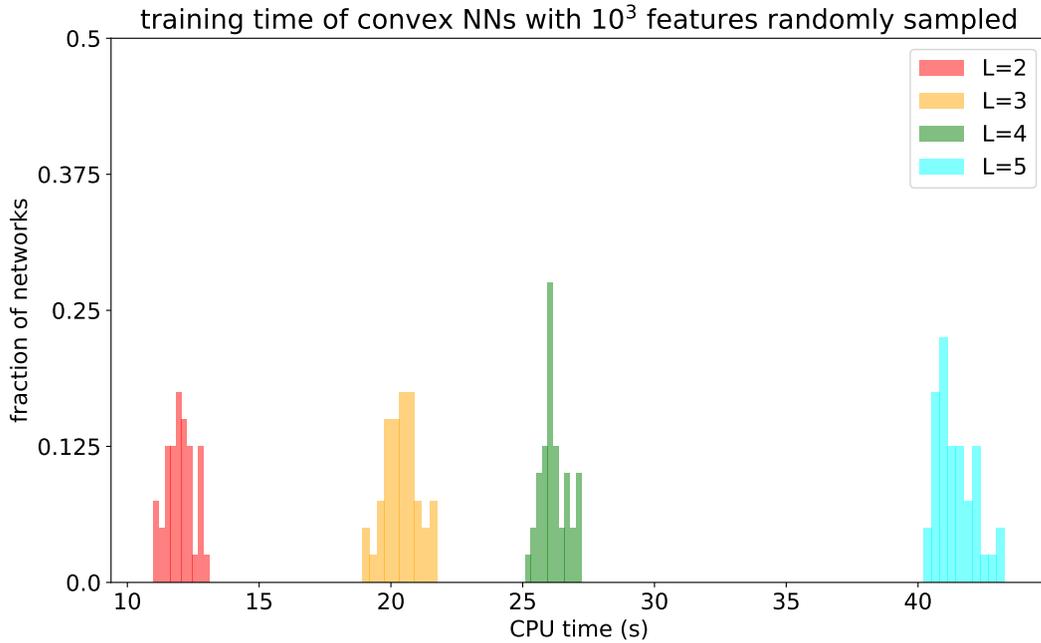


Figure E.2: Histogram of training time of deep narrow networks (40 networks for each depth L) trained on the Lasso problem with randomly sampled 10^3 features for autoregression on BTC-2017min with $T = 10^3$ samples. The training time increases with depth.

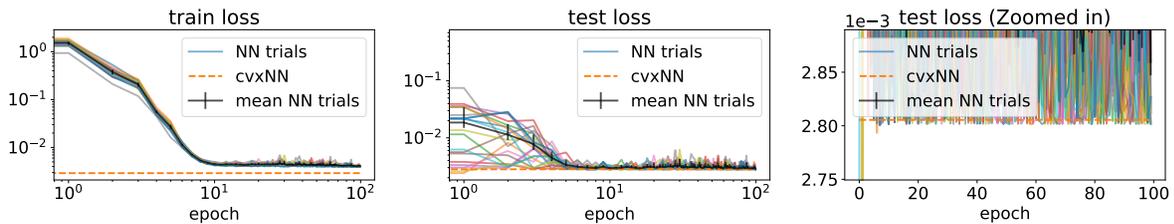


Figure E.3: Loss curves for planted data where $\sigma^2 = 0.01$ in (4.1) and there are $p = 10$ planted neurons. The non-black, solid colored curves plot each NN trial, and the black curve plots their mean, with the vertical lines indicating one standard deviation of the data above and below the mean.