

1 **A STRUCTURE-GUIDED GAUSS-NEWTON METHOD**
2 **FOR SHALLOW ReLU NEURAL NETWORK ***

3 ZHIQIANG CAI[†], TONG DING[†], MIN LIU[‡], XINYU LIU[†], AND JIANLIN XIA[†]

4 **Abstract.** In this paper, we propose a structure-guided Gauss-Newton (SgGN) method for solving least squares
5 problems using a shallow ReLU neural network. The method effectively takes advantage of both the least squares
6 structure and the neural network structure of the objective function. By categorizing the weights and biases of the
7 hidden and output layers of the network as nonlinear and linear parameters, respectively, the method iterates back
8 and forth between the nonlinear and linear parameters. The nonlinear parameters are updated by a damped Gauss-
9 Newton method and the linear ones are updated by a linear solver. Moreover, at the Gauss-Newton step, a special
10 form of the Gauss-Newton matrix is derived for the shallow ReLU neural network and is used for efficient iterations.
11 It is shown that the corresponding mass and Gauss-Newton matrices in the respective linear and nonlinear steps
12 are symmetric and positive definite under reasonable assumptions. Thus, the SgGN method naturally produces an
13 effective search direction without the need of additional techniques like shifting in the Levenberg-Marquardt method
14 to achieve invertibility of the Gauss-Newton matrix. The convergence and accuracy of the method are demonstrated
15 numerically for several challenging function approximation problems, especially those with discontinuities or sharp
16 transition layers that pose significant challenges for commonly used training algorithms in machine learning.

17 **Key words.** structure-guided Gauss-Newton method, neural network, least squares, mass matrix, Hessian
18 matrix, positive definiteness

19 **AMS subject classifications.** 65D15, 65K10

20 **1. Introduction.** When a neural network (NN) is used as a model for least-squares data
21 fitting, the procedure for determining the values of the NN parameters is a high-dimensional non-
22 convex optimization problem. This optimization problem tends to be computationally intensive
23 and complicated. Popular and widely used optimization algorithms (iterative solvers) in machine
24 learning are generally first-order gradient-based methods (see, e.g., survey papers [3, 11, 24]) because
25 of their moderate computational cost per iteration and their ease for implementation. However, their
26 efficiency depends heavily on hyper-parameters, the learning rate is cumbersome to tune, and the
27 methods usually converge slowly. Moreover, they exhibit the so-called plateau phenomenon in many
28 training tasks (see, e.g., [1]).

29 Recently, there has been a lot of interest in using second-order methods such as BFGS [4, 10,
30 12, 22] to solve optimization problems that arise from NN-based machine learning applications. For
31 attractive features and recent progress in overcoming the challenges of second-order methods, see
32 survey papers [3, 11, 24]. The Gauss-Newton (GN) method is a popular and widely used optimization
33 technique for solving general least-squares problems, as described in [9, 21]. Originating from the
34 classical Newton’s method, this approach makes use of the least squares structure and approximates
35 the Hessian matrix by its principal part. In recent years, the GN method has found practical appli-
36 cations in the realm of machine learning. The Kronecker-factored approximate curvature (KFAC)
37 method [17] enhances the optimization process of neural networks by leveraging the Kronecker-
38 factored approximation of the curvature matrix to mitigate computational challenges associated
39 with second-order optimization techniques. Building upon this foundation, the GN method was
40 advanced in KFRA [2] for deep learning with various practical aspects addressed, such as the com-
41 putational efficiency, the recursion relationship between layers, and step-size choices. The authors

*Submitted to the editors DATE.

Funding: This work of Zhiqiang Cai and Min Liu was supported in part by the National Science Foundation under grant DMS-2110571. The work of Jianlin Xia was supported in part by the National Science Foundation under grant DMS-2111007.

[†]Department of Mathematics, Purdue University, West Lafayette, IN (caiz@purdue.edu, ding158@purdue.edu,, liu1957@purdue.edu, xiaj@purdue.edu).

[‡]School of Mechanical Engineering, Purdue University, West Lafayette, IN (liu66@purdue.edu)

42 demonstrated the potential of the GN method to provide a more reliable and efficient alternative
 43 to conventional first-order optimization techniques for training deep neural networks.

44 The purpose of this paper is to design and investigate a novel structure-guided Gauss-Newton
 45 (SgGN) iterative method for solving least squares optimization problems using shallow ReLU NN.
 46 The method utilizes both the least squares structure and the ReLU NN structure. Following the
 47 study of the set of approximating functions generated by shallow ReLU NN and the establishment
 48 of the linearly independence of neurons, we set up a least-squares optimization problem and the
 49 corresponding algebraic systems for the stationary points. To quickly and accurately solve the
 50 optimization problem, we categorize the NN parameters into linear parameters (the weights \mathbf{c} and
 51 bias $\boldsymbol{\alpha}$ of the output layer) and nonlinear parameters (the weights \mathbf{w} and bias b of the hidden layer),
 52 denoted by $(\mathbf{c}, \boldsymbol{\alpha})$ and $\mathbf{r} = (b, \mathbf{w}^T)^T$, respectively. Such a classification leads to a natural block
 53 iterative process for solving the optimization problem by iterating back and forth between $(\mathbf{c}, \boldsymbol{\alpha})$
 54 and \mathbf{r} . $(\mathbf{c}, \boldsymbol{\alpha})$ and \mathbf{r} are updated by a linear solver and a damped Gauss-Newton nonlinear iterative
 55 solver, respectively.

56 At each SgGN iteration, the linear solver involves a mass matrix $\mathcal{A}(\mathbf{r})$ defined in (3.4) below.
 57 The nonlinear Gauss-Newton iterative solver is based on the following newly derived form of a
 58 *Gauss-Newton matrix* for shallow ReLU NN:

$$59 \quad (D(\mathbf{c}) \otimes I_{d+1})\mathcal{H}(\mathbf{r})(D(\mathbf{c}) \otimes I_{d+1}),$$

60 where d is the input data dimension, I_{d+1} is the order- $(d+1)$ identity matrix, $D(\mathbf{c})$ is a diagonal
 61 matrix consisting of the linear parameter \mathbf{c} , and $\mathcal{H}(\mathbf{r})$ depends on \mathbf{r} and is referred as the *layer*
 62 *Gauss-Newton matrix*. A specific form of $\mathcal{H}(\mathbf{r})$ is given in (4.2). Both $\mathcal{A}(\mathbf{r})$ and $\mathcal{H}(\mathbf{r})$ are functions
 63 of the nonlinear parameter \mathbf{r} and are independent of the linear parameters $(\mathbf{c}, \boldsymbol{\alpha})$.

64 Theoretically, we show that both $\mathcal{A}(\mathbf{r})$ and $\mathcal{H}(\mathbf{r})$ are *symmetric and positive definite* under the
 65 condition that all neurons are linearly independent (see Lemmas 3.1 and 4.2). Hence, a significant
 66 distinction between the SgGN method and the usual Gauss-Newton method is that there is no need
 67 to use additional techniques like shifting in the Levenberg-Marquardt method [14, 16] to achieve
 68 *invertibility* of the Gauss-Newton matrix.

69 The SgGN method provides an innovative way to effectively take advantage of both the qua-
 70 dratic structure and the NN structure in least-squares optimization problems arising from shallow
 71 ReLU NN approximations. The natural positive definiteness of $\mathcal{A}(\mathbf{r})$ and $\mathcal{H}(\mathbf{r})$ means many effi-
 72 cient direct/iterative linear solvers may be used for updating linear parameters (see (4.4)) and for
 73 computing search directions of nonlinear parameters (see (4.5)).

74 The SgGN method works for both continuous and discrete least-squares approximations. Its
 75 convergence and accuracy are demonstrated numerically for several one and two dimensional prob-
 76 lems, especially those with discontinuities or sharp transition layers that pose significant challenges
 77 for commonly used training algorithms in machine learning such as BFGS [4, 10, 12, 22] and Adam
 78 [13]. The loss curves for all test problems clearly show that the SgGN method significantly outper-
 79 forms those methods in terms of the convergence and accuracy. This conclusion is further enhanced
 80 by examining the ability of the methods in effectively moving the breaking hyper-planes (points for
 81 one dimension and lines for two dimensions). The breaking hyper-planes are determined by the non-
 82 linear parameters (weights and biases of the hidden layer). A data fitting application is also tested
 83 to show that the SgGN method can be naturally extended to discrete least-squares optimization,
 84 which makes it useful for many data science applications.

85 The paper is organized as follows. Section 2 introduces the set of approximating functions
 86 generated by shallow ReLU neural networks and establishes linear independence of neurons. The
 87 least-squares optimization problem and the corresponding nonlinear algebraic system for stationary
 88 points are described in Section 3. In Section 4, the structure of the Gauss-Newton matrix for the
 89 nonlinear parameters is derived and the resulting SgGN method is proposed. Section 5 presents the
 90 SgGN method for discrete least-squares optimization. Numerical results are given in Section 6 for

91 various function approximations and data science applications, following by some conclusions and
 92 discussions in Section 7.

93 **2. Shallow ReLU Neural Network.** This section describes shallow ReLU NN as a set of
 94 continuous piece-wise linear functions from \mathbb{R}^d to \mathbb{R} and discusses some analytical and geomet-
 95 rical properties of the network. Here the output dimension is restricted to one for simplicity of
 96 presentation since the extension of materials covered by the paper to higher output dimensions is
 97 straightforward.

98 ReLU refers to the rectified linear activation function defined by

$$99 \quad (2.1) \quad \sigma(t) = \max\{0, t\} = \begin{cases} t, & t > 0, \\ 0, & t \leq 0. \end{cases}$$

100 Its first- and second-order derivatives are the Heaviside (unit) step and the Dirac delta functions
 101 given by

$$102 \quad (2.2) \quad H(t) = \sigma'(t) = \begin{cases} 1, & t > 0, \\ 0, & t < 0 \end{cases} \quad \text{and} \quad \delta(t) = \sigma''(t) = H'(t) = \begin{cases} \infty, & t = 0, \\ 0, & t \neq 0, \end{cases}$$

103 respectively.

104 Let Ω be a connected, bounded open domain in \mathbb{R}^d . For any $\mathbf{x} = (x_1, \dots, x_d)^T \in \Omega \subset \mathbb{R}^d$,
 105 by appending 1 to the inhomogeneous (x_1, \dots, x_d) -coordinates, we have the following homogeneous
 106 coordinates:

$$107 \quad \mathbf{y}^T = (1, \mathbf{x}^T) = (1, x_1, \dots, x_d).$$

108 A standard shallow ReLU neural network with n neurons may be viewed as the set of continuous
 109 piece-wise linear functions from $\Omega \subset \mathbb{R}^d$ to \mathbb{R} as follows:

$$110 \quad (2.3) \quad \hat{\mathcal{M}}_n(\Omega) = \left\{ \sum_{i=1}^n c_i \sigma(\mathbf{r}_i \cdot \mathbf{y}) + \alpha_0 : \mathbf{x} \in \Omega, c_i \in \mathbb{R}, \mathbf{r}_i^T = (b_i, \mathbf{w}_i^T), b_i \in \mathbb{R}, \mathbf{w}_i \in \mathcal{S}^{d-1}, \alpha_0 \in \mathbb{R} \right\},$$

111 where $\mathbf{c} = (c_1, \dots, c_n)^T$ and α_0 are the respective output weight and bias, $\mathbf{r}_i = (b_i, \mathbf{w}_i^T) \in \mathbb{R}^{d+1}$ is
 112 for the parameters of the i^{th} neuron with b_i and \mathbf{w}_i the respective bias and weight of the neuron,
 113 and \mathcal{S}^{d-1} is the unit sphere in \mathbb{R}^d . The constraint that the weight of each neuron belongs to the
 114 unit sphere is a consequence of normalization for the ReLU activation function (see [15]) and may
 115 narrow down the solution set of a given approximating problem. For convenience, denote

$$116 \quad (2.4) \quad \mathbf{r}^T = (\mathbf{r}_1^T, \dots, \mathbf{r}_n^T) = (b_1, \mathbf{w}_1^T, \dots, b_n, \mathbf{w}_n^T).$$

117 Notice that the ReLU activation function $\sigma(t)$ is a continuous piece-wise linear function with
 118 one *breaking point* at $t = 0$. Hence each ridge function $\sigma(\mathbf{r}_i \cdot \mathbf{y}) = \sigma(\mathbf{w}_i \cdot \mathbf{x} + b_i)$ is a continuous
 119 piece-wise linear function with a *breaking hyper-plane* (see [5, 15]):

$$120 \quad (2.5) \quad \mathcal{P}_i(\mathbf{r}_i) = \{\mathbf{x} \in \Omega \subset \mathbb{R}^d : \mathbf{w}_i \cdot \mathbf{x} + b_i = 0\}.$$

121 For the set $\hat{\mathcal{M}}_n(\Omega)$ of neural network functions defined in (2.3) with fixed weights and biases \mathbf{r} ,
 122 there are n breaking hyper-planes $\{\mathcal{P}_i(\mathbf{r}_i)\}_{i=1}^n$. Together with the boundary of the domain Ω , these
 123 hyper-planes form a *physical partition*, denoted by $\mathcal{K}(\mathbf{r})$, of Ω [8, 15]. This partition $\mathcal{K}(\mathbf{r})$ consists
 124 of irregular, polygonal sub-domains of Ω (see Figures 6.6(e) and 6.6(i) below for some examples).
 125 Each function in $\hat{\mathcal{M}}_n(\Omega)$ is then a continuous piece-wise linear function with respect to $\mathcal{K}(\mathbf{r})$.

126 Below, we discuss linear independence of some ridge functions for fixed parameter \mathbf{r} in (2.4).
 127 To this end, let $\sigma_0(\mathbf{x}) = 1$ and for $i = 1, \dots, n$, let

$$128 \quad \sigma_i(\mathbf{x}) = \sigma(\mathbf{r}_i \cdot \mathbf{y}) \quad \text{and} \quad H_i(\mathbf{x}) = H(\mathbf{r}_i \cdot \mathbf{y}),$$

129 where σ and H are the ReLU activation and Heaviside step functions given in (2.1) and (2.2),
 130 respectively. It is already known from Lemma 2.1 in [15] that the set of functions $\{\sigma_i(\mathbf{x})\}_{i=0}^n$ is
 131 linearly independent if the hyper-planes $\{\mathcal{P}_i(\mathbf{r}_i)\}_{i=1}^n$ are distinct. We further have the following
 132 result.

133 **LEMMA 2.1.** *For fixed \mathbf{r} in (2.4), assume that the hyper-planes $\{\mathcal{P}_i(\mathbf{r}_i)\}_{i=1}^n$ are distinct. Then*
 134 *the set of functions $\{H_i(\mathbf{x}), x_1 H_i(\mathbf{x}), \dots, x_d H_i(\mathbf{x})\}_{i=1}^n$ is linearly independent.*

135 *Proof.* For each $i = 1, \dots, n$, the linear independence of $\{1, x_1, \dots, x_d\}$ implies that the set of
 136 functions

$$137 \quad \{H_i(\mathbf{x}), x_1 H_i(\mathbf{x}), \dots, x_d H_i(\mathbf{x})\} = H_i(\mathbf{x})\{1, x_1, \dots, x_d\}$$

138 is linearly independent. Now, the linear independence of $\{H_i(\mathbf{x}), x_1 H_i(\mathbf{x}), \dots, x_d H_i(\mathbf{x})\}_{i=1}^n$ follows
 139 from the assumption on the distinct hyper-planes. \square

140 We conclude this section in discussing a possible restriction on the biases of all neurons in
 141 $\hat{\mathcal{M}}_n(\Omega)$. For each \mathbf{r}_i , there are two ridge functions:

$$142 \quad \sigma(\mathbf{w}_i \cdot \mathbf{x} + b_i) \quad \text{and} \quad \sigma(-\mathbf{w}_i \cdot \mathbf{x} - b_i).$$

143 They share the same breaking hyper-plane $\mathcal{P}_i(\mathbf{r}_i)$ and are linearly independent. Nevertheless, the
 144 following identity is well known:

$$145 \quad \sigma(\mathbf{w}_i \cdot \mathbf{x} + b_i) - \sigma(-\mathbf{w}_i \cdot \mathbf{x} - b_i) = \mathbf{w}_i \cdot \mathbf{x} + b_i$$

146 for all $\mathbf{x} \in \mathbb{R}^d$. This identity implies that clearly only one of the two ridge functions is needed in
 147 $\hat{\mathcal{M}}_n(\Omega)$, if $\hat{\mathcal{M}}_n(\Omega)$ contains all linear functions

$$148 \quad \mathbb{P}_1 = \text{span} \{\phi_i(\mathbf{x})\}_{i=0}^d \quad \text{with} \quad \phi_0(\mathbf{x}) = 1, \quad \phi_i(\mathbf{x}) = x_i.$$

149 Hence, we may assume that $b_i \in \mathbb{R}_0^+ = [0, +\infty)$ for all i to further narrow down the solution set.

150 In general, $\hat{\mathcal{M}}_n(\Omega)$ does not contain \mathbb{P}_1 . Nevertheless, this may be resolved by either expanding
 151 $\hat{\mathcal{M}}_n(\Omega)$ to contain \mathbb{P}_1 or fixing the weights and biases of d neurons in $\hat{\mathcal{M}}_n(\Omega)$ such that the cor-
 152 responding d hyper-planes are linearly independent and do not intersect with the domain Ω . For
 153 convenience, in this paper we choose the former. That is, we use the following set of approximating
 154 functions:

$$155 \quad \mathcal{M}_n(\Omega) = \left\{ \sum_{i=1}^n c_i \sigma_i(\mathbf{x}) + \sum_{i=0}^d \alpha_i \phi_i(\mathbf{x}) : \mathbf{x} \in \Omega, c_i \in \mathbb{R}, \alpha_i \in \mathbb{R}, b_i \in \mathbb{R}_0^+, \mathbf{w}_i \in \mathcal{S}^{d-1} \right\}.$$

156 The SgGN method developed in this paper can be applied to the standard shallow ReLU neural
 157 network $\hat{\mathcal{M}}_n(\Omega)$ directly without enforcing $b_i \in \mathbb{R}_0^+$.

158 **3. Continuous Least-Squares Optimization Problems.** Given a function $u(\mathbf{x})$ defined on
 159 Ω , the best least-squares approximation to u in $\mathcal{M}_n(\Omega)$ is to find $u_n(\mathbf{x}) \in \mathcal{M}_n(\Omega)$ such that

$$160 \quad (3.1) \quad \mathcal{J}_\mu(u_n) = \min_{v \in \mathcal{M}_n(\Omega)} \mathcal{J}_\mu(v) \quad \text{with} \quad \mathcal{J}_\mu(v) = \frac{1}{2} \|v - u\|_\mu^2$$

161 where $\|\cdot\|_\mu$ is the weighted $L^2(\Omega)$ norm given by

$$162 \quad \|f\|_\mu = \left(\int_\Omega \mu(\mathbf{x}) f^2(\mathbf{x}) d\mathbf{x} \right)^{1/2}.$$

163 The corresponding weighted $L^2(\Omega)$ inner product is denoted by $\langle \cdot, \cdot \rangle_\mu$.

164 Below, we use the optimality condition to derive the corresponding systems of algebraic equa-
165 tions. To this end, let

$$166 \quad u_n(\mathbf{x}) = \sum_{i=1}^n c_i \sigma_i(\mathbf{x}) + \sum_{i=0}^d \alpha_i \phi_i(\mathbf{x})$$

167 be a solution of (3.1). Then the *linear parameter*

$$168 \quad \hat{\mathbf{c}} = (\mathbf{c}, \boldsymbol{\alpha})^T = (c_1, \dots, c_n, \alpha_0, \dots, \alpha_d)^T$$

169 and the *nonlinear parameter* \mathbf{r} defined in (2.4) are critical points of the loss function $\mathcal{J}_\mu(u_n)$. That
170 is, $\hat{\mathbf{c}}$ and \mathbf{r} satisfy the following systems of algebraic equations

$$171 \quad (3.2) \quad \nabla_{\hat{\mathbf{c}}} \mathcal{J}_\mu(u_n) = \mathbf{0} \quad \text{and} \quad \nabla_{\mathbf{r}} \mathcal{J}_\mu(u_n) = \mathbf{0},$$

172 where $\nabla_{\hat{\mathbf{c}}}$ and $\nabla_{\mathbf{r}}$ denote the gradients with respect to the respective parameters $\hat{\mathbf{c}}$ and \mathbf{r} .

173 In the following, we derive specific forms of the algebraic equations in (3.2). Notice that

$$174 \quad \nabla_{\hat{\mathbf{c}}} u_n(\mathbf{x}) = (\sigma_1, \dots, \sigma_n, \phi_0, \dots, \phi_d)^T.$$

175 Let

$$176 \quad \begin{cases} a_{ij}(\mathbf{r}) = \langle \sigma_j, \sigma_i \rangle, \quad f_i(\mathbf{r}) = \langle u, \sigma_i \rangle, & \text{for } i, j = 1, \dots, n, \\ b_{ij} = \langle \phi_j, \phi_i \rangle, \quad g_i = \langle u, \phi_i \rangle, & \text{for } i, j = 0, 1, \dots, d, \\ c_{ij}(\mathbf{r}) = \langle \phi_j, \sigma_i \rangle, & \text{for } i = 1, \dots, n, \quad j = 0, 1, \dots, d, \end{cases}$$

177 and let

$$178 \quad \mathcal{A}_{11}(\mathbf{r}) = [a_{ij}(\mathbf{r})]_{n \times n}, \quad \mathcal{A}_{12}(\mathbf{r}) = [c_{ij}(\mathbf{r})]_{n \times d}, \quad \mathcal{A}_{22} = [b_{ij}]_{d \times d}, \quad \mathbf{f}_1(\mathbf{r}) = [f_i(\mathbf{r})]_{n \times 1}, \quad \mathbf{f}_2 = [g_i]_{d \times 1}.$$

179 It is easy to see that the first equation in (3.2) implies

$$180 \quad (3.3) \quad \mathcal{A}(\mathbf{r}) \hat{\mathbf{c}} = \mathbf{f}(\mathbf{r}),$$

181 where $\mathcal{A}(\mathbf{r})$ and $\mathbf{f}(\mathbf{r})$ are the *mass matrix* and the right-hand side vector given by

$$182 \quad (3.4) \quad \mathcal{A}(\mathbf{r}) = \begin{bmatrix} \mathcal{A}_{11}(\mathbf{r}) & \mathcal{A}_{12}(\mathbf{r}) \\ \mathcal{A}_{12}^T(\mathbf{r}) & \mathcal{A}_{22} \end{bmatrix} \quad \text{and} \quad \mathbf{f}(\mathbf{r}) = \begin{bmatrix} \mathbf{f}_1(\mathbf{r}) \\ \mathbf{f}_2 \end{bmatrix},$$

183 respectively.

184 LEMMA 3.1. Under the assumption of Lemma 2.1, additionally assume that for all $i = 1, \dots, n$

$$185 \quad (3.5) \quad b_i \geq 0 \quad \text{and} \quad \mathcal{P}_i(\mathbf{r}_i) \cap \Omega \neq \emptyset,$$

186 where $\mathcal{P}_i(\mathbf{r}_i)$ is the breaking hyper-plane defined in (2.5). Then the mass matrix $\mathcal{A}(\mathbf{r})$ is symmetric
187 and positive definite.

188 *Proof.* Obviously, $\mathcal{A}(\mathbf{r})$ is symmetric. For any $\hat{\mathbf{c}} = (\mathbf{c}^T, \boldsymbol{\alpha}^T)^T \in \mathbb{R}^{n+d+1}$, it is easy to see that

189
$$\hat{\mathbf{c}}^T \mathcal{A}(\mathbf{r}) \hat{\mathbf{c}} = \left\| \sum_{i=1}^n c_i \sigma_i(\mathbf{x}) + \sum_{i=0}^d \alpha_i \phi_i(\mathbf{x}) \right\|_{\mu}^2.$$

190 Hence, in order to show the positive definiteness of $\mathcal{A}(\mathbf{r})$, it suffices to prove that

191
$$\{\sigma_i(\mathbf{x})\}_{i=1}^n \cup \{\phi_i(\mathbf{x})\}_{i=0}^d$$

192 is linearly independent. This can be done through proof by induction. To do so, first notice
 193 that the second assumption in (3.5) implies that $\sigma_i(\mathbf{x})$ vanishes in a non-empty subdomain Ω_i
 194 of Ω . The linear independence of $\{\sigma_1(\mathbf{x})\} \cup \{\phi_i(\mathbf{x})\}_{i=0}^d$ is a direct consequence of the fact that
 195 $\sigma_1(\mathbf{x})|_{\Omega_1} \equiv 0$ and the linear independence of $\{\phi_i(\mathbf{x})\}_{i=0}^d$ in $\Omega \setminus \Omega_1$. Similarly, the linear independence
 196 of $\{\sigma_i(\mathbf{x})\}_{i=1}^k \cup \{\phi_i(\mathbf{x})\}_{i=0}^d$ follows from the fact that $\sigma_k(\mathbf{x})|_{\Omega_k} \equiv 0$ and the assumption on the linear
 197 independence of $\{\sigma_i(\mathbf{x})\}_{i=1}^{k-1} \cup \{\phi_i(\mathbf{x})\}_{i=0}^d$. This completes the proof of the lemma. \square

198 Next, we calculate $\nabla_{\mathbf{r}} \mathcal{J}_{\mu}(u_n)$. To simplify expression of formulas, we use the Kronecker product,
 199 denoted by \otimes , of two matrices. Let

200
$$\mathbf{H}(\mathbf{x}) = (H_1(\mathbf{x}), \dots, H_n(\mathbf{x}))^T.$$

201 For $i, j = 1, \dots, n$, the fact that

202
$$\nabla_{\mathbf{r}_i} \sigma_j(\mathbf{x}) = \begin{cases} \mathbf{0}, & i \neq j, \\ H_j(\mathbf{x}) \mathbf{y}, & i = j \end{cases}$$

203 implies

204 (3.6)
$$\nabla_{\mathbf{r}} u_n(\mathbf{x}) = (D(\mathbf{c}) \otimes I_{d+1}) (\mathbf{H}(\mathbf{x}) \otimes \mathbf{y}),$$

205 where $D(\mathbf{c}) = \text{diag}(c_1, \dots, c_n)$ is the diagonal matrix with the i^{th} -diagonal element c_i .

206 Let

207 (3.7)
$$\mathbf{G}(\mathbf{c}, \boldsymbol{\alpha}, \mathbf{r}) = \int_{\Omega} \mu(\mathbf{x}) (u_n(\mathbf{x}) - u(\mathbf{x})) \mathbf{H}(\mathbf{x}) \otimes \mathbf{y} \, d\mathbf{x}.$$

208 It is easy to check that the second equation in (3.2) becomes

209 (3.8)
$$\nabla_{\mathbf{r}} \mathcal{J}_{\mu}(u_n) = (D(\mathbf{c}) \otimes I_{d+1}) \mathbf{G}(\mathbf{c}, \boldsymbol{\alpha}, \mathbf{r}) = 0.$$

210 (3.3) and (3.8) define two algebraic systems that may be used to solve for the linear parameter
 211 \mathbf{c} and the nonlinear parameter \mathbf{r} , respectively. For convenience, (3.3) and (3.8) are referred to as
 212 the *linear and nonlinear problems*, respectively.

213 **4. A structure-guided Gauss-Newton (SgGN) method.** In this section, we introduce
 214 our SgGN method for solving the minimization problem in (3.1), guided by both the least squares
 215 structure and the ReLU NN structure. The method utilizes a block structure of the systems of
 216 algebraic equations given in (3.3) and (3.8) and iterates back and forth between the linear parameter
 217 \mathbf{c} by solving (3.3) and the nonlinear parameter \mathbf{r} by using the Gauss-Newton method.

218 To efficiently apply the Gauss-Newton method, below we derive a special ‘‘Gauss-Newton’’
 219 matrix by making use of the neural network structure. To this end, let $\delta_i(\mathbf{x}) = \delta(\mathbf{r}_i \cdot \mathbf{y})$ for

220 $i = 1, \dots, n$, where $\delta(t)$ is the Dirac delta function defined in (2.2). Denote the $n \times n$ diagonal
 221 matrix with the i^{th} -diagonal element $\delta_i(\mathbf{x})$ by

$$222 \quad \Lambda(\mathbf{x}) = \text{diag}(\delta_1(\mathbf{x}), \dots, \delta_n(\mathbf{x})).$$

223 For $i, j = 1, \dots, n$, it is easy to check that

$$224 \quad \nabla_{\mathbf{r}_i} H_j(\mathbf{x}) = \begin{cases} \mathbf{0}, & i \neq j, \\ \delta_j(\mathbf{x}) \mathbf{y}, & i = j. \end{cases}$$

225 This implies

$$226 \quad (4.1) \quad \nabla_{\mathbf{r}} \mathbf{H}^T(\mathbf{x}) = \Lambda(\mathbf{x}) \otimes \mathbf{y}.$$

227 Now we introduce the following $n(d+1) \times n(d+1)$ matrix:

$$228 \quad (4.2) \quad \mathcal{H}(\mathbf{r}) = \int_{\Omega} \mu [\mathbf{H}\mathbf{H}^T] \otimes [\mathbf{y}\mathbf{y}^T] \, d\mathbf{x},$$

229 where we write $\mu(\mathbf{x})$ as μ and $\mathbf{H}(\mathbf{x})$ as \mathbf{H} for the ease of notation. Lemma 4.1 below shows that
 230 $\mathcal{H}(\mathbf{r})$ is the principal part of the Gauss-Newton matrix and is referred as the *layer Gauss-Newton*
 231 *matrix* in this paper.

232 LEMMA 4.1. *The Hessian matrix has the form of*

$$233 \quad (4.3) \quad \nabla_{\mathbf{r}} (\nabla_{\mathbf{r}} \mathcal{J}_p(u_n))^T = (D(\mathbf{c}) \otimes I_{d+1}) \mathcal{H}(\mathbf{r}) (D(\mathbf{c}) \otimes I_{d+1}) + \hat{\mathcal{H}}(\mathbf{c}, \mathbf{r}) (D(\mathbf{c}) \otimes I_{d+1}),$$

234 where $\hat{\mathcal{H}}(\mathbf{c}, \mathbf{r})$ is given by

$$235 \quad \hat{\mathcal{H}}(\mathbf{c}, \mathbf{r}) = \int_{\Omega} \mu (u_n(\mathbf{x}) - u(\mathbf{x})) \Lambda(\mathbf{x}) \otimes (\mathbf{y}\mathbf{y}^T) \, d\mathbf{x}.$$

236 *Proof.* It follows from (3.8), the product rule, (3.6), and (4.1) that

$$\begin{aligned} 237 \quad \nabla_{\mathbf{r}} \mathbf{G}(\mathbf{c}, \boldsymbol{\alpha}, \mathbf{r}) &= \int_{\Omega} \mu (\nabla_{\mathbf{r}} u_n) (\mathbf{H} \otimes \mathbf{y})^T \, d\mathbf{x} + \int_{\Omega} \mu (u_n - u) (\nabla_{\mathbf{r}} \mathbf{H}^T) \otimes \mathbf{y}^T \, d\mathbf{x} \\ 238 &= (D(\mathbf{c}) \otimes I_{d+1}) \int_{\Omega} \mu (\mathbf{H} \otimes \mathbf{y}) (\mathbf{H} \otimes \mathbf{y})^T \, d\mathbf{x} + \int_{\Omega} \mu (u_n - u) \Lambda(\mathbf{x}) \otimes (\mathbf{y}\mathbf{y}^T) \, d\mathbf{x}, \end{aligned}$$

239 which, together with (3.8) and the transpose rule of the Kronecker product, implies (4.3). This
 240 completes the proof of the lemma. \square

241 LEMMA 4.2. *Under the assumption of Lemma 2.1, the Gauss-Newton matrix $\mathcal{H}(\mathbf{r})$ is symmetric*
 242 *and positive definite.*

243 *Proof.* Clearly, $\mathcal{H}(\mathbf{r})$ is symmetric. For any $\mathbf{v}^T = (\boldsymbol{\beta}_1^T, \dots, \boldsymbol{\beta}_n^T) \in \mathbb{R}^{n(d+1)}$ with $\boldsymbol{\beta}_i \in \mathbb{R}^{d+1}$, let

$$244 \quad v(\mathbf{x}) = \sum_{i=1}^n (\boldsymbol{\beta}_i^T \mathbf{y}) H_i(\mathbf{x}).$$

245 It is easy to check that

$$246 \quad \mathbf{v}^T \mathcal{H}(\mathbf{r}) \mathbf{v} = \mathbf{v}^T \left(\int_{\Omega} \mu(\mathbf{x}) (\mathbf{H} \otimes \mathbf{y}) (\mathbf{H} \otimes \mathbf{y})^T \, d\mathbf{x} \right) \mathbf{v} = \langle v, v \rangle_{\mu} \geq 0,$$

247 which, together with Lemma 2.1, implies that $\mathcal{H}(\mathbf{r})$ is positive definite. This completes the proof of
 248 the lemma. \square

249 With the Gauss-Newton matrix $\mathcal{H}(\mathbf{r})$ defined in (4.2), Lemma 4.1, and Lemma 4.2, we are ready
 250 to describe one step of the SgGN method. Given the previous iterate $(\hat{\mathbf{c}}^{(k)}, \mathbf{r}^{(k)}) = (\mathbf{c}^{(k)}, \boldsymbol{\alpha}^{(k)}, \mathbf{r}^{(k)})$,
 251 compute the current iterate $(\hat{\mathbf{c}}^{(k+1)}, \mathbf{r}^{(k+1)}) = (\mathbf{c}^{(k+1)}, \boldsymbol{\alpha}^{(k+1)}, \mathbf{r}^{(k+1)})$ as follows.

252 (i) Compute the linear parameter $\hat{\mathbf{c}}^{(k+1)} = (\mathbf{c}^{(k+1)}, \boldsymbol{\alpha}^{(k+1)})$ by solving the system of linear equations
 253

$$254 \quad (4.4) \quad \mathcal{A}(\mathbf{r}^{(k)}) \hat{\mathbf{c}}^{(k+1)} = \mathbf{f}(\mathbf{r}^{(k)}),$$

255 using either a direct or iterative solver.

256 (ii) In the case that every entry of $\mathbf{c}^{(k+1)}$ is nonzero, compute the search direction

$$257 \quad \mathbf{p}^{(k+1)} = \left(D^{-1}(\mathbf{c}^{(k+1)}) \otimes I_{d+1} \right) \mathbf{s}^{(k+1)},$$

258 where $\mathbf{s}^{(k)}$ is the solution of the Gauss-Newton system

$$259 \quad (4.5) \quad \mathcal{H}(\mathbf{r}^{(k)}) \mathbf{s}^{(k+1)} = -\mathbf{G}(\mathbf{c}^{(k+1)}, \boldsymbol{\alpha}^{(k+1)}, \mathbf{r}^{(k)}).$$

260 (iii) In the case that $\mathbf{c}^{(k+1)}$ has some entries with magnitudes less than a certain small threshold,
 261 the search direction $\mathbf{p}^{(k+1)}$ is obtained in a similar fashion as in (ii) by updating the entries of
 262 $\mathbf{p}^{(k)}$ that corresponds to the rest of entries of $\mathbf{c}^{(k+1)}$. The biases corresponding to those small
 263 entries will be assigned randomly.

264 (iv) Compute the nonlinear parameter

$$265 \quad \mathbf{r}^{(k+1)} = \mathbf{r}^{(k)} + \gamma_{k+1} \mathbf{p}^{(k+1)},$$

266 where the damping parameter γ_{k+1} is the solution of a one-dimensional optimization problem

$$267 \quad \gamma_{k+1} = \arg \min_{\gamma \in \mathbb{R}_0^+} \mathcal{J}_\mu \left(u_n(\cdot; \mathbf{c}^{(k+1)}, \mathbf{r}^{(k)} + \gamma \mathbf{p}^{(k+1)}) \right).$$

268 See Algorithm 4.1 for a pseudocode of the SgGN method.

269 The SgGN method needs the algebraic solutions of multiple linear systems in (4.4) and (4.5)
 270 during the iterations. Since the focus of this work is on investigating the SgGN method for shallow
 271 ReLU networks, here we just briefly mention some numerical issues considered in the method and
 272 leave the details to a forthcoming paper [7]. The linear systems may be solved with direct or iterative
 273 solvers. If $\mathbf{r}^{(k)}$ satisfies the assumption in Lemma 2.1, both the matrices $\mathcal{A}(\mathbf{r}^{(k)})$ and $\mathcal{H}(\mathbf{r}^{(k)})$ are
 274 symmetric and positive definite (see Lemmas 3.1 and 4.2). Nevertheless, both of them could be
 275 *nearly* singular in some applications, such as when the underlying target function has sharp layers.
 276 This is especially the case when the SgGN iterations start to converge so that some basis functions
 277 get close to each other. (A rigorous way to characterize how close they are will be given in [7].) In
 278 this work, we use direct solvers for the purpose of verifying the convergence of the SgGN algorithm.

279 **5. SgGN for Discrete Least-Squares Optimization Problems.** We then show how the
 280 SgGN method can be extended to discrete least-squares optimization problems. For a given data
 281 set $\{(\mathbf{x}^i, u^i)\}_{i=1}^m$ with $\mathbf{x}^i \in \Omega$ and $u^i \in \mathbb{R}$ and a given distribution function $0 \leq \mu(\mathbf{x}) \leq 1$, consider
 282 the following discrete least-squares minimization problem: finding $u_n(\mathbf{x}) \in \hat{\mathcal{M}}_n(\Omega)$ such that

$$283 \quad (5.1) \quad \mathcal{J}_{m,\mu}(u_n) = \min_{v \in \hat{\mathcal{M}}_n(\Omega)} \mathcal{J}_{m,\mu}(v),$$

284 where $\mathcal{J}_{m,\mu}(v)$ is the weighted discrete least-squares loss function given by

$$285 \quad \mathcal{J}_{m,\mu}(v) = \frac{1}{2} \sum_{i=1}^m \mu(\mathbf{x}^i) (v(\mathbf{x}^i) - u^i)^2 = \frac{1}{2} \|v - u\|_{m,\mu}^2.$$

Algorithm 4.1 A structure-guided Gauss-Newton (SgGN) method for (3.1)

Input: network parameters $\mathbf{r} = (\mathbf{r}_1, \dots, \mathbf{r}_n)$, data set $\{(\mathbf{x}^i, u^i)\}_{i=1}^N$, density function $\mu(\mathbf{x})$

Output: network parameters \mathbf{c}, \mathbf{r}

Initialize $\mathbf{r}^{(0)}$ in (2.4)

for $k = 0, 1, \dots$ **do**

▷ *Linear parameter* \mathbf{c}

Form $\mathcal{A}(\mathbf{r}^{(k)})$, $\mathbf{f}(\mathbf{r}^{(k)})$ in (3.4).

$\mathbf{c}^{(k+1)} \leftarrow \mathcal{A}^{-1}(\mathbf{r}^{(k)})\mathbf{f}(\mathbf{r}^{(k)})$

▷ *Nonlinear parameter* \mathbf{r}

Form $\mathbf{G}(\mathbf{c}^{(k+1)}, \mathbf{r}^{(k)})$, $\mathcal{H}(\mathbf{r}^{(k)})$ in Equations (3.7) and (4.2)

$\mathbf{s}^{(k)} \leftarrow -\mathcal{H}^{-1}(\mathbf{r}^{(k)})\mathbf{G}(\mathbf{c}^{(k+1)}, \mathbf{r}^{(k)})$

$\mathbf{p}^{(k)} \leftarrow (D^{-1}(\mathbf{c}^{(k+1)}) \otimes I_{d+1})\mathbf{s}^{(k)}$

$\gamma_{k+1} \leftarrow \arg \min_{\gamma \in \mathbb{R}_0^+} \mathcal{J}_\mu(u_n(\cdot; \mathbf{c}^{(k+1)}, \mathbf{r}^{(k)} + \alpha\mathbf{p}^{(k)}))$

$\mathbf{r}^{(k+1)} \leftarrow \mathbf{r}^{(k)} + \gamma_{k+1}\mathbf{p}^{(k)}$

if a desired loss or a specified number of iterations is reached **then**

Return $\mathbf{c}^{(k+1)}, \mathbf{r}^{(k+1)}$

end if

end for

286

287 Here, $\|\cdot\|_{m,\mu}$ is the weighted discrete $L^2(\Omega)$ norm given by

$$288 \quad \|f\|_{m,\mu} = \left(\sum_{i=1}^m \mu(\mathbf{x}^i) f^2(\mathbf{x}^i) \right)^{\frac{1}{2}},$$

289 and the corresponding weighted discrete $L^2(\Omega)$ inner product is denoted by $\langle \cdot, \cdot \rangle_{m,\mu}$.

290 Problem (5.1) may be regarded as a discretization of (3.1) using a numerical integration of
 291 Monte-Carlo type. Hence, the SgGN method proposed in the previous section can be applied
 292 directly for iteratively solving the discrete least-squares minimization problem in (5.1) by simply
 293 replacing integrals by summations. For the reader's convenience, we describe the corresponding
 294 components of the SgGN method.

295 To this end, let

$$296 \quad u_n(\mathbf{x}) = \sum_{i=1}^n c_i \sigma_i(\mathbf{x}) + \alpha_0 = \sum_{i=1}^n c_i \sigma(\mathbf{w}_i \mathbf{x} + b_i) + \alpha_0$$

297 be a solution of (5.1), then

$$298 \quad \hat{\mathbf{c}} = (c_1, \dots, c_n, \alpha_0)^T \quad \text{and} \quad \mathbf{r}^T = (\mathbf{r}_1^T, \dots, \mathbf{r}_n^T) = (b_1, \mathbf{w}_1^T, \dots, b_n, \mathbf{w}_n^T)$$

299 are the linear and nonlinear parameters of $u_n(\mathbf{x})$, respectively. The blocks of $\mathcal{A}(\mathbf{r})$ and $\mathbf{f}(\mathbf{r})$ in (3.4)
 300 are given by

$$301 \quad \mathcal{A}_{11}(\mathbf{r}) = [a_{ij}(\mathbf{r})]_{n \times n}, \quad \mathcal{A}_{12}(\mathbf{r}) = [b_i(\mathbf{r})]_{n \times 1}, \quad \mathcal{A}_{22} = \langle 1, 1 \rangle_{m,\mu}, \quad \mathbf{f}_1(\mathbf{r}) = [f_i(\mathbf{r})]_{n \times 1}, \quad f_2 = \langle u, 1 \rangle_{m,\mu},$$

302 where entries of the block matrices are defined by

$$303 \quad \begin{cases} a_{ij}(\mathbf{r}) = \langle \sigma_j, \sigma_i \rangle_{m,\mu}, & f_i(\mathbf{r}) = \langle u, \sigma_i \rangle_{m,\mu}, & \text{for } i, j = 1, \dots, n, \\ b_i(\mathbf{r}) = \langle \sigma_i, 1 \rangle_{m,\mu}, & & \text{for } i = 1, \dots, n, \end{cases}$$

304 respectively. The gradient vector of $\mathcal{J}_{m,\mu}(u_n)$ with respect to \mathbf{r} is

$$305 \quad \mathbf{G}(\mathbf{c}, \boldsymbol{\alpha}, \mathbf{r}) = \sum_{i=1}^m \mu(\mathbf{x}^i) (u_n(\mathbf{x}^i) - u(\mathbf{x}^i)) \mathbf{H}(\mathbf{x}^i) \otimes \mathbf{y}^i,$$

306 and the layer Gauss-Newton matrix is

$$307 \quad \mathcal{H}(\mathbf{r}) = \sum_{i=1}^m \mu(\mathbf{x}^i) (\mathbf{H}(\mathbf{x}^i) \mathbf{H}^T(\mathbf{x}^i)) \otimes (\mathbf{y}^i (\mathbf{y}^i)^T).$$

308 The damping parameter γ_{k+1} is the solution of a one-dimensional optimization problem

$$309 \quad \gamma_{k+1} = \arg \min_{\gamma \in \mathbb{R}_0^+} \mathcal{J}_{m,\mu} \left(u_n(\cdot; \mathbf{c}^{(k+1)}, \mathbf{r}^{(k)} + \gamma \mathbf{p}^{(k+1)}) \right).$$

310 **6. Numerical Experiments.** In this section, a series of numerical experiments is conducted
 311 to test the effectiveness and accuracy of the proposed SgGN algorithm. We also compare it with
 312 several existing optimization algorithms commonly used for neural networks: Adam [13], BFGS
 313 [4, 10, 12, 22], and the Gauss-Newton-based KFRA method introduced in [2] which is considered
 314 more applicable than the original Gauss-Newton-based KFAC method [17].

315 Our test problems encompass various types of target functions, including step functions in one
 316 and two dimensions, a delta-like function in 1D, and a continuous piece-wise linear function in 2D.
 317 These functions are well-suited for accurate approximation using shallow neural networks. However,
 318 they pose significant challenges for optimization algorithms due to the presence of discontinuities
 319 or sharp transition layers. As indicated in [5, 15], the nonlinear parameters $\mathbf{r}_i = (b_i, \mathbf{w}_i^T)^T \in \mathbb{R}^{d+1}$
 320 correspond to the breaking points/lines of the neurons which form a physical partition of the domain.
 321 Therefore, the efficiency of an optimization algorithm can be measured as the efficiency of moving
 322 those breaking points/lines from a uniform distribution to the nearly “optimal” locations according
 323 to the underlying target function.

324 In the comparison study, we used BFGS as a baseline. For each test, BFGS was first repeated
 325 30 times, and we report the median loss and the corresponding approximation result. We then
 326 run the other two methods, KFAC and our SgGN, using the same number of iterations. For the
 327 first-order Adam optimizer, we run all test problems with a relatively larger number of iterations
 328 until the corresponding loss functions plateau.

329 It is important to note that different methods may have varying computation complexities
 330 per iteration. For instance, computational cost of BFGS [20] and the KFRA [2] is $O(n^2)$ per
 331 iteration, while cost of Adam is only $O(n)$. Our SgGN involves solutions of two dense linear
 332 system with coefficient matrices $\mathcal{A}(\mathbf{r}^{(k)})$ and $\mathcal{H}(\mathbf{r}^{(k)})$, respectively, as in Algorithm 4.1. These
 333 matrices are typically very ill conditioned for the test problems considered here. In our current
 334 implementation, truncated SVDs are used for the solution and its cost is $O(rn^2)$ with r depending
 335 on the accuracy. This suffices our purpose of comparing the convergence of the SgGN with the other
 336 methods. Although we use the number of iterations as one of the reference metrics for the efficiency
 337 performance of the optimization, our major focus for comparison in this study is the quality of the
 338 solution. As shown for the test problems, the SgGN can often converge to much more accurate
 339 approximation than the other methods.

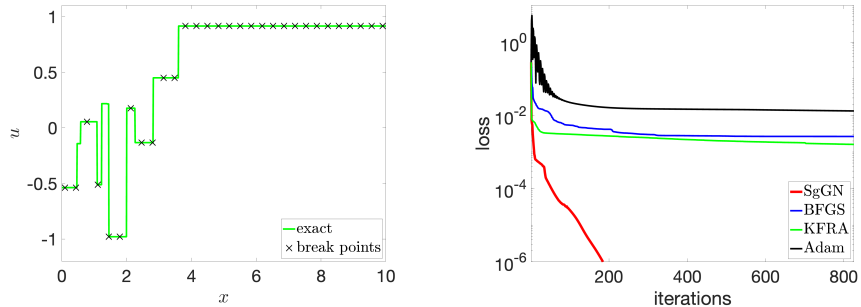
340 The detailed parameter setting for each method is listed in Table 6.1. All the methods start
 341 with the same initial as listed in the Initialization section in Table 6.1. The integration of the loss
 342 function $\mathcal{J}_\mu(u_n)$ is computed by the composite mid-point rule over a uniform partition with the
 343 mesh size $h = 0.01$.

344 **6.1. One-dimensional piece-wise constant function.** The first test problem is a one-
 345 dimensional piece-wise constant function defined in the interval $[0, 10]$ with ten pieces from a skewed

Table 6.1: List of parameters used in the methods, where the parameters for BFGS and Adam are referred to the MATLAB deep learning toolbox [18].

BFGS	
<code>net.trainParam.min_grad</code>	minimum performance gradient with value 0
<code>net.trainParam.max_fail</code>	maximum validation failures with value 10^4
<code>net.trainParam.epochs</code>	maximum number of epochs to train with value 10^4
KFRA	
γ	A damping parameter for the approximated Gauss-Newton matrix induced by the full Gauss-Newton
Adam	
<code>InitialLearnRate</code>	initial learning rate α_1
<code>DropRateFactor</code>	multiplicative factor α_f by which the learning rate drops
<code>DropPeriod</code>	number of epochs that passes between adjustments to the learning rate, denoted by T
Initialization	
Linear coefficient \mathbf{c}	initialized by a narrow normal distribution $\mathcal{N}(0, 0.01)$
nonlinear parameter \mathbf{r}_i	the corresponding breaking hyper-planes uniformly partition the domain

346 distribution (see Figure 6.1(a)). We use these skewed pieces to test whether an optimizer can
 347 move the uniformly initialized breaking points to catch those discontinuities in a target function.
 348 Theoretically, 20 neurons are enough to approximate a ten-piece step function with a given accuracy
 349 $\epsilon > 0$ [5, 6]. Due to the uncertainty of solving a non-convex optimization problem, 30 neurons are
 350 used in the test. For Adam method, we adjusted the learning rate for the best performance and the
 351 reported results are obtained using $\alpha_1 = 0.1$, $\alpha_f = 0.5$, and $T = 1000$. For the KFRA, $\gamma = 0.01$.



(a) Target function u and initial breaking points (b) Loss curves of the optimization methods

Fig. 6.1: One-dimensional piece-wise constant function approximation: target function, initial breaking points, and optimization loss curves.

352 The loss decay curves are depicted in Figure 6.1(b). While the loss curve for SgGN continues
 353 to decline even after 200 iterations, the curves for the other three methods decay slowly, reaching
 354 values close to their final training loss. Table 6.2 compares the least squares loss values of the

Table 6.2: Comparison for one-dimensional piece-wise constant function.

Method	SgGN		BFGS		KFRA	Adam
Iteration	9	825	207	825	825	10,000
$\mathcal{J}_{m,\mu}$	8.76E-4	6.56E-9	4.03E-3	2.65E-3	1.61E-3	8.14E-3

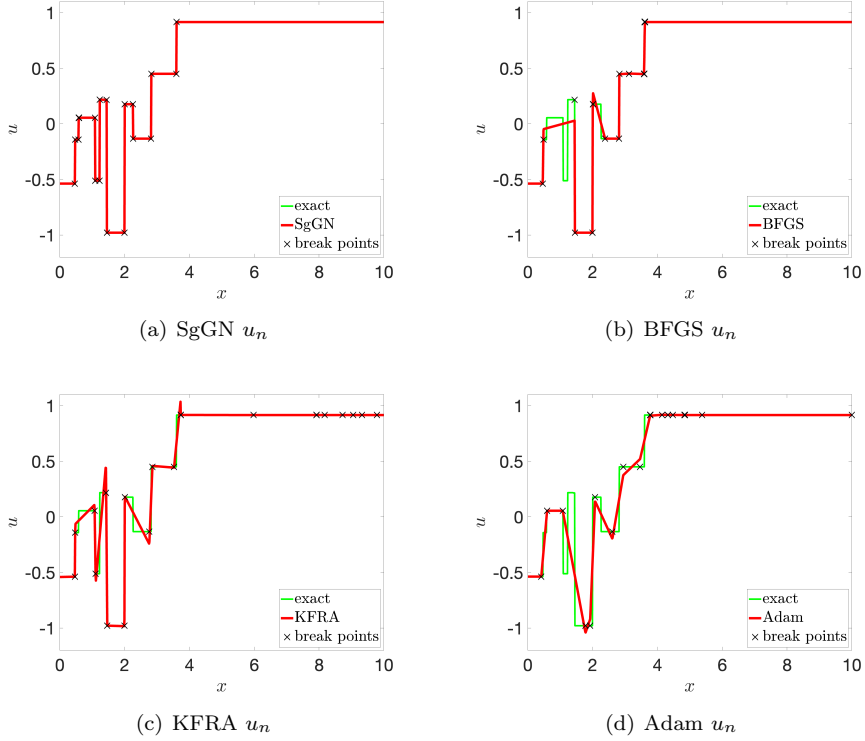


Fig. 6.2: One-dimensional piece-wise constant function approximation results.

355 four methods. Notably, SgGN achieves a loss value of magnitude 10^{-9} , significantly lower than the
 356 10^{-3} magnitude obtained by the other methods. Further insight into these loss comparisons can
 357 be derived from the approximation results shown in Figure 6.2. Specifically, only SgGN accurately
 358 captures all steps (Figure 6.2(a)), by precisely aligning the breaking points to the discontinuities.
 359 In contrast, the other methods (Figures 6.2(b) to 6.2(d)) either overlook certain discontinuous steps
 360 or induce overshooting.

361 **6.2. One-dimensional delta-like function.** In the second experiment, we consider a smooth
 362 but sharp delta-like function [23]

363
$$u(x) = \sum_{i=1}^k \frac{1}{d_i(x - x_i)^2 + 1}, \quad x \in [-1.5, 1.5],$$

364 where k is the number of centers, x_i is the center position, and d_i is to control center width so that
 365 the larger d_i is, the narrower width is.

366 In the experiment, we set $k = 3$ with centers $\{x_1, x_2, x_3\} = \left\{-\frac{\pi^2}{10}, -\left(\pi - \frac{5}{2}\right), \frac{\sqrt{85}}{10}\right\}$, and width
 367 parameters $\{d_1, d_2, d_3\} = \{10^4, 10^3, 5 \times 10^3\}$. We employed 15 neurons in the first hidden layer for
 368 our tests. To evaluate the optimization methods, we initialized the neural network with uniformly
 369 distributed breaking points, as illustrated in Figure 6.3(a). Our primary objectives are twofold:
 370 first, to assess each optimization method’s ability to accurately capture all three irrational peak
 371 centers from an rationally initialized breaking points and second, to determine whether the method
 372 can adaptively allocate the 15 breaking points to account for the differing widths of the centered
 373 delta-like peaks. For Adam, $\alpha_1 = 0.02$, $\alpha_f = 0.6$ and $T = 2000$, while in KFRA, $\gamma = 0.0001$.

374 As illustrated in Figure 6.3(b), the loss curve for SgGN not only decays more rapidly than
 375 those of the other three methods but also converges to a better solution within about 50 iterations.
 376 The SgGN method converges at a loss magnitude of 10^{-4} , which is substantially lower than the
 377 $10^{-3} \sim 10^{-2}$ magnitudes exhibited by the other methods (see Table 6.3). As further shown in
 378 Figure 6.4, SgGN successfully moves the breaking points to align with all center positions and
 379 adaptively distributes the remaining neuron breaking points to accurately approximate sharp peaks
 380 of varying widths. However, the other methods (see Figures 6.4(b) to 6.4(d)) fail to capture these
 381 nuances; they either fail to capture all peaks or do not distribute the breaking points proportionally.

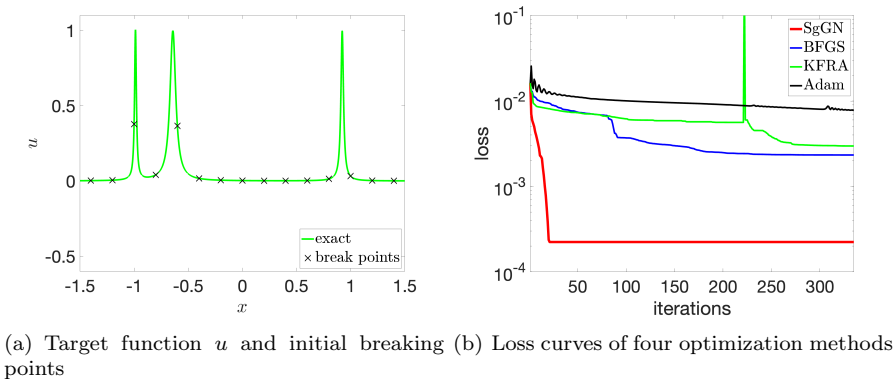


Fig. 6.3: One-dimensional delta-like function approximation: target function, initial breaking points, and optimization loss curves.

Table 6.3: Comparison for one-dimensional delta-like function.

Method	SgGN		BFGS		KFRA	Adam
Iteration	12	334	91	334	334	10,000
$\mathcal{J}_{m,\mu}$	2.24E-3	2.23E-4	3.74E-3	2.33E-3	2.97E-3	3.94E-3

382 **6.3. Two-dimensional piece-wise constant function.** Next, we consider a 2D piece-wise
 383 constant function defined in the domain $[-1, 1]^2$:

$$384 \quad u(x) = \begin{cases} 1, & -0.5 \leq x + y \leq 0.5, \\ -1, & \text{otherwise.} \end{cases}$$

385 In the previous two examples, we used more neurons than the minimum required to mitigate the

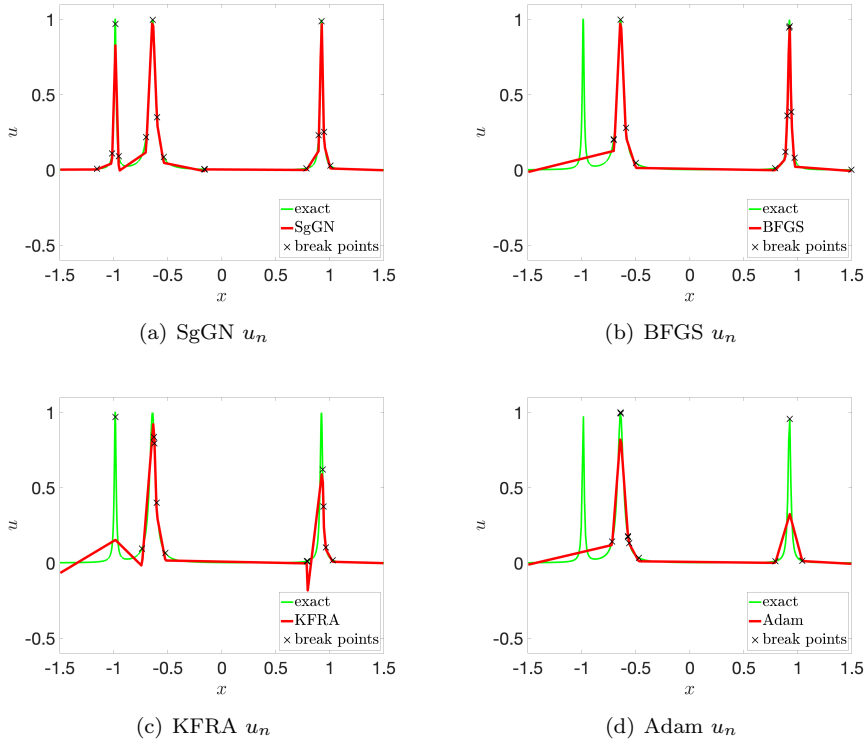


Fig. 6.4: One-dimensional delta-like function approximation: optimization loss curves and approximation results using four optimization methods.

386 uncertainties arising from the non-convex optimization problems. In contrast, this test focuses on
 387 utilizing only the minimum number of neurons necessary to compare the performances of different
 388 methods. As shown in Figure 6.5(a), each discontinuous segment can, in theory, be approximated
 389 using just two neurons. An effective approximation would place a pair of neurons for each discontinuous
 390 line, with the proximity of the corresponding breaking lines serving as a measure of the
 391 quality of the approximation. Consequently, we opt for using only 4 neurons in this example. For
 392 Adam method, $\alpha_1 = 0.01$, $\alpha_f = 0.8$ and $T = 2000$ and in KFRA, $\gamma = 0.005$.

393 As shown in Figure 6.5(c), the SgGN loss curve not only decays at a faster rate than those of the
 394 other three methods but also reaches the final training loss in approximately 20 iterations. Table 6.4
 395 compares the least squares loss values after 142 iterations for the second-order methods and 10,000
 396 iterations for Adam. Given the integration mesh size $h = 0.01$, there exists a theoretical lower bound
 397 on the proximity of the breaking lines, thus imposing a limit on the minimal achievable loss value.
 398 Even so, SgGN attains a loss magnitude of 10^{-3} , lower than the 10^{-2} magnitudes demonstrated by
 399 the other methods. Furthermore, as displayed in Figure 6.5. SgGN (see Figure 6.5(h)) accurately
 400 positioned all the four breaking lines to capture the discontinuities. In contrast, the other three
 401 methods (see Figures 6.5(i) and 6.5(k)) captured only one side of the discontinuity lines.

402 **6.4. Two-dimensional function in $\hat{\mathcal{M}}_n(\Omega)$.** For the previous three problems, the target
 403 functions do not reside within the defined network function space $\hat{\mathcal{M}}_n(\Omega)$. Consequently, the re-
 404 sultant loss function does not converge to zero. This non-zero convergence precludes a definitive

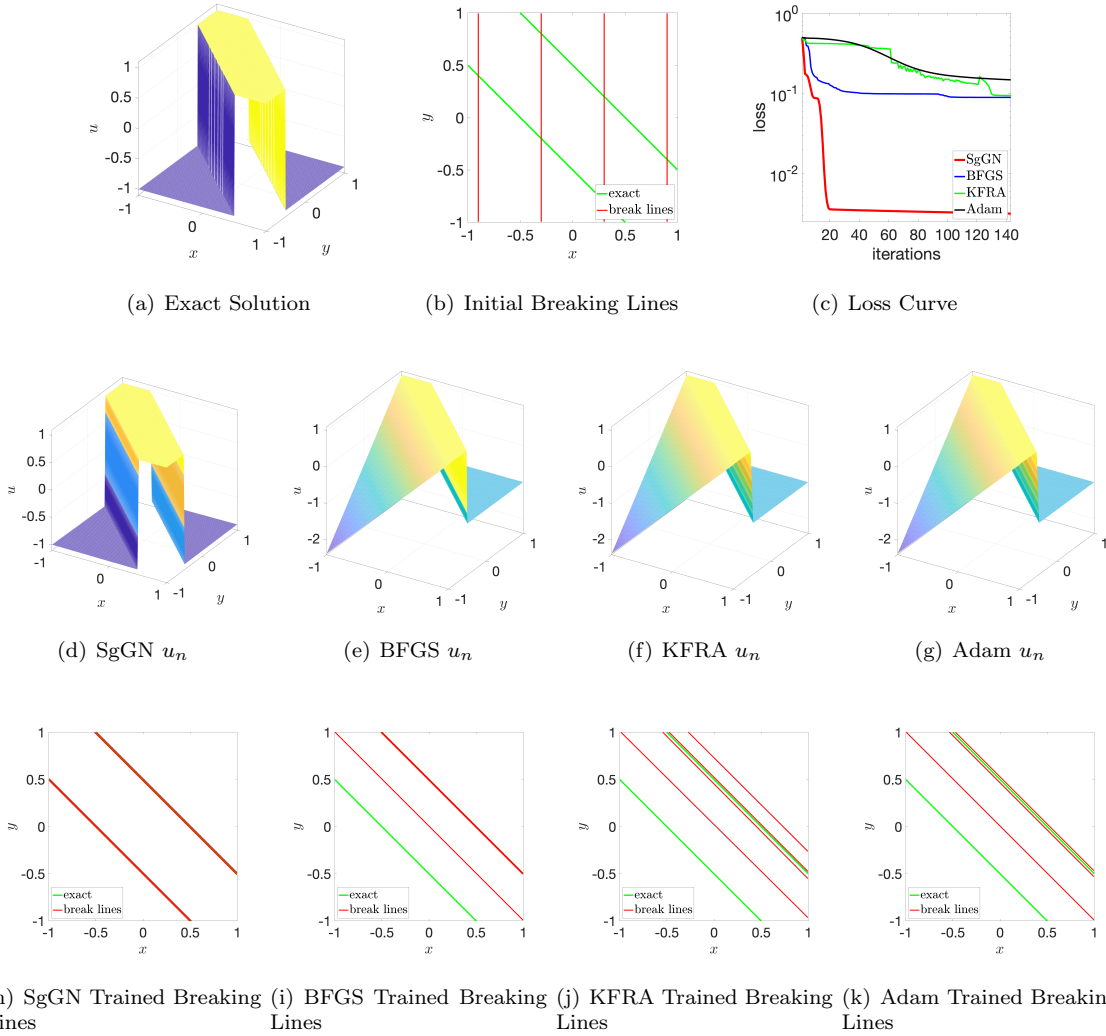


Fig. 6.5: Two-dimensional piece-wise constant function approximation: target function, initial breaking lines, optimization loss curves and approximation results using four optimization methods.

Table 6.4: Comparison for a two-dimensional piece-wise constant function.

Method	SgGN		BFGS		KFRA	Adam
Iteration	9	142	100	142	142	10,000
$\mathcal{J}_{m,\mu}$	8.82E-2	3.16E-3	9.20E-2	8.92E-2	9.40E-2	9.23E-2

405 determination of the optimal approximation for each target function within the specified function
 406 space. To better assess the performance of parameter optimization, we introduce an artificial func-

407 tion from $\hat{\mathcal{M}}_n(\Omega)$ with randomly selected optimal parameters \mathbf{c}^* and \mathbf{r}^* that

408 (6.1)
$$u(\mathbf{x}) = \sum_{i=1}^N c_i^* \phi_i(\mathbf{x}; \mathbf{r}_i^*) + \alpha_0^*,$$

409 where N is number of neurons. In this case, the known optimal parameters \mathbf{c}^* and \mathbf{r}^* allow us to
 410 directly evaluate the performance of different optimization methods by tracking the movement of
 411 the breaking lines toward these optimal values. When too many neurons are provided, any initial
 412 positions of the breaking lines would naturally be close to some of the optimal ones. Moreover, with
 413 uniform initialization, some of the initial positions will be inherently close to the optimal solutions.
 414 To focus on the movement of the breaking lines, we opt for initialization along only horizontal and
 415 vertical axes, and we limit the number of neurons to 5. For Adam method, we have $\alpha_1 = 0.1$,
 416 $\alpha_f = 0.5$, $T = 2000$ for horizontal initialization, and $\alpha_1 = 0.1$, $\alpha_f = 0.8$, $T = 3000$ for vertical
 417 initialization respectively. In KFRA, $\gamma = 0.005$.

Table 6.5: Comparison for a two-dimensional piece-wise linear function with horizontal initial break-
 ing lines (VI) and vertical initial breaking lines (VI).

Method	SgGN		BFGS		KFRA	Adam
Iteration	99	207	204	207	207	10,000
$\mathcal{J}_{m,\mu}$ (HI)	6.28E-22	6.68E-27	7.50E-22	7.50E-22	6.12E-2	1.17E-5
Iteration	4	105	30	105	105	10,000
$\mathcal{J}_{m,\mu}$ (VI)	2.35E-4	4.34E-26	5.21E-4	2.71E-4	5.56E-2	2.15E-4

418 The loss decay curves are depicted in [Figures 6.6\(b\) and 6.6\(d\)](#). SgGN reaches the magnitude
 419 10^{-10} within just 50 iterations while the other three methods decay more slowly, requiring a greater
 420 number of iterations to reach their final training loss. [Table 6.5](#) compares the least squares loss
 421 values when initialized horizontally and vertically, respectively. Notably, SgGN excels by achieving
 422 a near-zero loss value in both scenarios, a level of accuracy significantly higher than what the other
 423 methods managed to accomplish. Further examination of the approximation results, as shown in
 424 [Figure 6.6](#), reveals more nuanced insights. SgGN accurately moves all five breaking lines to their
 425 optimal positions under both horizontal and vertical initializations (see [Figures 6.6\(e\) and 6.6\(i\)](#)).
 426 On the other hand, while the other methods ([Figures 6.6\(f\), 6.6\(h\), 6.6\(j\) and 6.6\(l\)](#)) either perform
 427 well only with horizontal initializations or fail altogether to correctly position all the breaking lines.

428 **6.5. Data science application.** Lastly, we test an application of shallow network in data
 429 science. The task is to predict the age of abalones using physical measurements, drawing data from
 430 the UCI dataset [\[19\]](#). Since there is no prior knowledge about network structure setting, we test
 431 two shallow networks where the numbers of neurons are 40 and 80 respectively. The parameters for
 432 Adam in the two tests are $\alpha_1 = 0.1$, $\alpha_f = 0.8$, $T = 1500$, and in KFRA, $\gamma = 1$.

Table 6.6: Comparison for the data science problem.

Method		SgGN		BFGS		KFRA	Adam
40 neurons	Iteration	35	200	50	200	200	10,000
	$\mathcal{J}_{m,\mu}$	1.90	1.90	2.09	1.89	2.59	2.02
80 neurons	Iteration	9	200	40	200	200	10,000
	$\mathcal{J}_{m,\mu}$	2.00	1.70	2.14	2.02	2.59	1.93

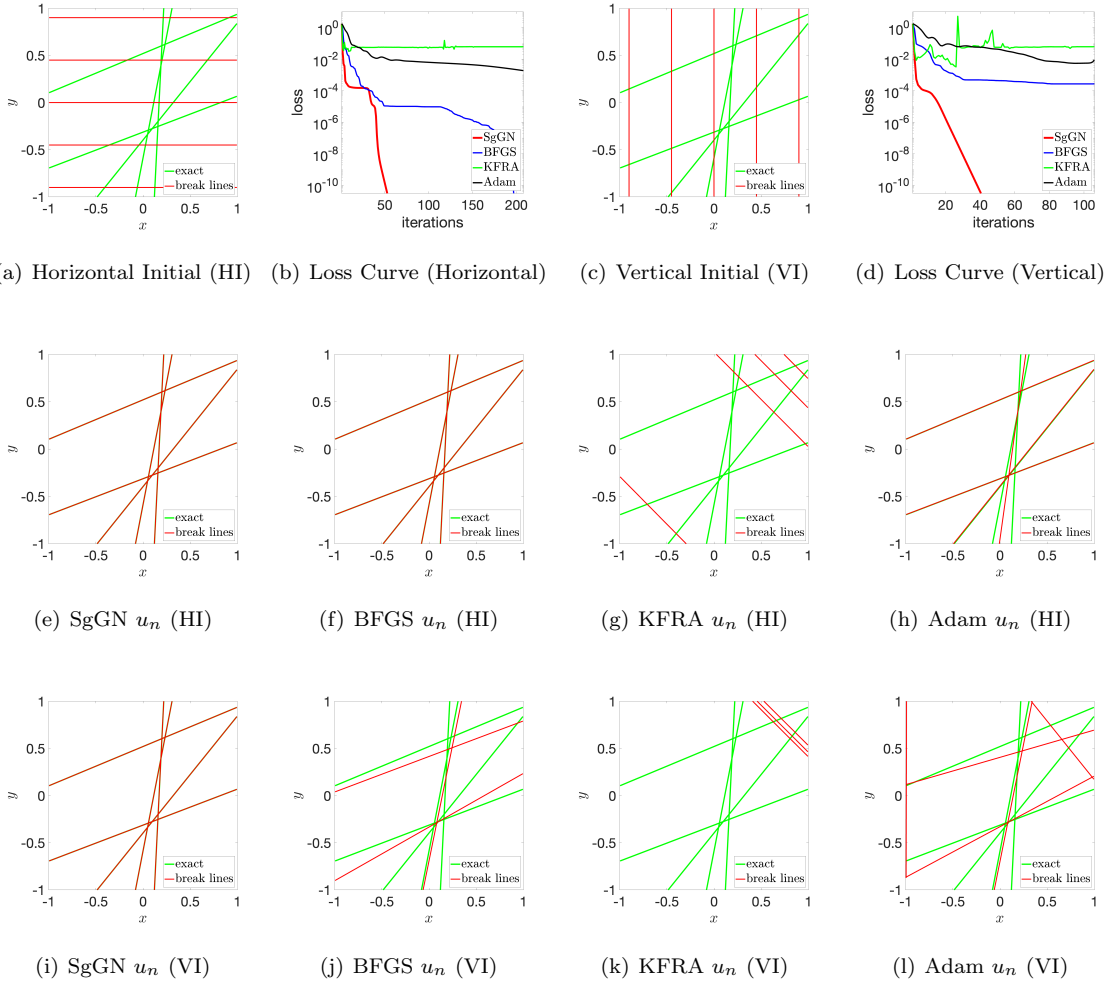


Fig. 6.6: Two-dimensional piece-wise linear function approximation: target function, initial breaking lines, optimization loss curves and approximation results using the optimization methods with horizontal initial (HI) and vertical initial (VI) breaking lines.

433 The loss decay curve for SgGN reaches a magnitude close to the final loss just within just 35
 434 iterations(see Figure 6.7). In contrast, the loss curves for the other three methods decay more slowly,
 435 requiring a greater number of iterations to to reach their final training loss. Table 6.6 compares
 436 the least squares loss values after 200 iterations for the second-order methods and 10,000 iterations
 437 for Adam, with number of neurons 40 and 80 respectively. Both SgGN and BFGS have similar
 438 performance in the scenario with 40 neurons. However, SgGN demonstrates a better least squares
 439 loss than the other three methods with 80 neurons. This unequivocally attests to SgGN’s superiority
 440 in terms of both the accuracy and effectiveness for this data science application.

441 **7. Conclusions and Discussions.** Newton’s method in optimization is a second-order iter-
 442 ative method for numerically solving optimization problems with general objective functions. One

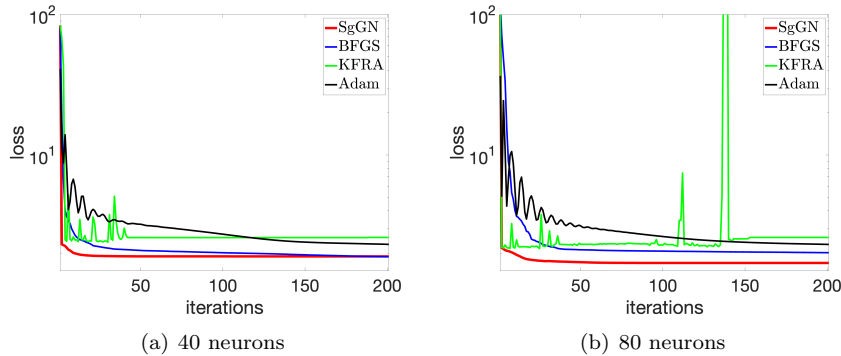


Fig. 6.7: Data science application: optimization loss curves using the optimization methods.

443 of its variants, BFGS has been successfully applied to NN-based machine learning applications.
 444 For nonlinear least-squares problems, one may use methods of Gauss-Newton type that exploit the
 445 quadratic form of the objective function. The structured-guided Gauss-Newton (SgGN) method
 446 introduced in this paper is an iterative method for solving nonlinear least-squares problems using
 447 shallow ReLU NN as a model. In addition to the least squares structure aspect, SgGN method
 448 effectively makes use of the structure of the network. Guided by both structure types, the method
 449 has some attractive features. One feature is the guarantee of the positive definiteness of the mass
 450 and layer Gauss-Newton matrices without the need of extra shifting like in usual Gauss-Newton
 451 methods.

452 Another feature is the rapid convergence in practice. the SgGN method was tested for several
 453 one and two dimensional least-squares problems which are difficult for commonly used training
 454 algorithms in machine learning such as BFGS and Adam. The loss curves for all the test problems
 455 clearly show the superior convergence of SgGN. SgGN often out-performs those methods by a large
 456 margin. This conclusion is further strengthened by examining the ability and effectiveness of the
 457 methods in moving the breaking hyper-planes (breaking points for one dimension and breaking lines
 458 for two dimensions).

459 Each iteration of the SgGN requires linear solvers to approximately invert the mass matrix
 460 $\mathcal{A}(\mathbf{r}^{(k)})$ and the layer Gauss-Newton matrix $\mathcal{H}(\mathbf{r}^{(k)})$ for the linear and nonlinear parameters, re-
 461 spectively. While both the matrices are symmetric and positive definite, they are nevertheless
 462 ill-conditioned. In the numerical experiments reported in this paper, the truncated SVD is used
 463 as the linear solver, albeit at a significant computational cost. More efficient linear solvers will be
 464 investigated in an upcoming paper [7].

465 REFERENCES

466 [1] M. AINSWORTH AND Y. SHIN, *Active neuron least squares: a training method for multivariate rectified neural*
 467 *networks*, SIAM Journal on Scientific Computing, 44 (2022), pp. A1807–C366.
 468 [2] A. BOTEV, H. RITTER, AND D. BARBER, *Practical Gauss-Newton optimisation for deep learning*, in Int'l. Conf.
 469 Mach. Learning, PMLR, 2017, pp. 557–565.
 470 [3] L. BOTTOU, F. E. CURTIS, AND J. NOCEDAL, *Optimization methods for large-scale machine learning*, SIAM
 471 Review, 60 (2018), p. 223–311.
 472 [4] C. G. BROYDEN, *The convergence of a class of double-rank minimization algorithms 1. general considerations*,
 473 IMA J. Appl. Math., 6 (1970), pp. 76–90.
 474 [5] Z. CAI, J. CHEN, AND M. LIU, *Least-squares ReLU neural network (LSNN) method for linear advection-reaction*
 475 *equation*, J. Comput. Phys., (443 (2021) 110514).

- 476 [6] Z. CAI, J. CHOI, AND M. LIU, *Least-squares neural network (LSNN) method for linear advection-reaction*
477 *equation: general discontinuous interface*, arXiv:2301.06156v3[math.NA], (2023).
- 478 [7] Z. CAI, T. DING, M. LIU, X. LIU, AND J. XIA, *Mass and Gauss-Newton matrices for shallow ReLU neural*
479 *network*, manuscript.
- 480 [8] Z. CAI AND M. LIU, *Self-adaptive ReLU neural network method in least-squares data fitting*, Handbook of
481 *Research on Adaptive Artificial Intelligence*, (submitted).
- 482 [9] J. E. DENNIS JR AND R. B. SCHNABEL, *Numerical methods for unconstrained optimization and nonlinear*
483 *equations*, SIAM, 1996.
- 484 [10] R. FLETCHER, *A new approach to variable metric algorithms*, Computer J., 13 (1970), pp. 317–322.
- 485 [11] C. GAMBELLA, B. GHADDAR, AND J. NAOUM-SAWAYA, *Optimization problems for machine learning: a survey*,
486 arXiv:1901.05331 [math.OC], (2019).
- 487 [12] D. GOLDFARB, *A family of variable-metric methods derived by variational means*, Math. Comp., 24 (1970),
488 pp. 23–26.
- 489 [13] D. P. KINGMA AND J. BA, *ADAM: A method for stochastic optimization*, in Int’l Conf. Representation Learning,
490 San Diego, 2015; arXiv preprint arXiv:1412.6980.
- 491 [14] K. LEVENBERG, *A method for the solution of certain non-linear problems in least squares*, Quart. Appl. Math.,
492 2 (1944), pp. 164–168.
- 493 [15] M. LIU, Z. CAI, AND J. CHEN, *Adaptive two-layer ReLU neural network: I. best least-squares approximation*,
494 Comput. Math. Appl., 113 (2022), pp. 34–44.
- 495 [16] D. W. MARQUARDT, *An algorithm for least-squares estimation of nonlinear parameters*, J. SIAM, 11 (1963),
496 pp. 431–441.
- 497 [17] J. MARTENS AND R. GROSSE, *Optimizing neural networks with Kronecker-factored approximate curvature*, in
498 Int’l. Conf. Mach. Learning, PMLR, 2015, pp. 2408–2417.
- 499 [18] THE MATHWORKS INC., *Deep Learning Toolbox*, 2023.
- 500 [19] W. NASH, T. SELLERS, S. TALBOT, A. CAWTHORN, AND W. FORD, *Abalone*. UCI Machine Learning Repository,
501 1995. DOI: <https://doi.org/10.24432/C55C7W>.
- 502 [20] J. NOCEDAL AND S. J. WRIGHT, *Numerical Optimization*, Springer, 2006.
- 503 [21] J. M. ORTEGA AND W. C. RHEINBOLDT, *Iterative Solution of Nonlinear Equations in Several Variables*, SIAM,
504 2000.
- 505 [22] D. F. SHANNO, *Conditioning of quasi-Newton methods for function minimization*, Math. Comp., 24 (1970),
506 pp. 647–656.
- 507 [23] J. SHEN, Y. WANG, AND J. XIA, *Fast structured direct spectral methods for differential equations with variable*
508 *coefficients, i. the one-dimensional case*, SIAM J. Sci. Comput., 38 (2016), pp. A28–A54.
- 509 [24] S. SUN, Z. CAO, H. ZHU, AND J. ZHAO, *A survey of optimization methods from a machine learning perspective*,
510 IEEE Trans. on Cybernetics, 50 (2020), pp. 3668 – 3681.