

1

2 **Convergence Analysis of Block Newton Methods for** 3 **1D Shallow Neural Network Approximation**

4 Zhiqiang Cai^{1,*}, Anastassia Doktorova¹, Robert D. Falgout²
5 and César Herrera¹

6 ¹ Department of Mathematics, Purdue University, West Lafayette, IN 47907-2067,
7 USA

8 ² Livermore National Laboratory, Livermore, CA 94550-9698, USA

9 Received 13 February 2026; Accepted (in revised version) 18 March 2026

11

Abstract. This paper analyzes local convergence of the block Newton (BN) method introduced in [5, 6] for one-dimensional shallow neural network approximation to functions and diffusion-reaction problems. The BN method consists of the 2×2 block nonlinear Gauss-Seidel, linear Gauss-Seidel, or Jacobi method for outer iteration and the Newton method for inner iteration. The blocks are corresponding to the linear and the nonlinear parameters. Under some reasonable assumptions, we establish local convergence of the BN methods as well as the reduced BN (rBN) method for one-dimensional diffusion-reaction problems and least-squares function approximation. Unlike common optimization methods, the rBN allows for the reduction of the number of parameters during the optimization process when some neurons contribute little to the approximation or are at nearly optimal locations.

12 **AMS subject classifications:** to be provided by authors

13 **Key words:** Neural network, Elliptic problems, Least-Squares approximation, Newton's method,
14 Local convergence analysis.

15

16 **1 Introduction**

17 One dimensional ReLU shallow neural network (NN) with n neurons generates a set of
18 continuous piecewise linear functions. Specifically, the set with restriction of the biases
19 in the interval $I = (0,1)$ is given by

$$\mathcal{M}_n(I) = \left\{ \alpha + \sum_{i=0}^n c_i \sigma(x - b_i) : \alpha \in \mathbb{R}, c_i \in \mathbb{R}, 0 = b_0 < b_1 < \dots < b_n < b_{n+1} = 1 \right\},$$

*Corresponding author.

Email addresses: caiz@purdue.edu (Z. Cai), adoktoro@purdue.edu (A. Doktorova), rfalgout@llnl.gov (R. Falgout), herre125@purdue.edu (C. Herrera)

20 where $\sigma(t) = \max\{0, t\}$ is the ReLU activation function. Denote by $\mathbf{c} = (c_0, \dots, c_n)^T \in \mathbb{R}^{n+1}$
 21 and $\mathbf{b} = (b_1, \dots, b_n)^T \in \mathbb{R}^n$ the respective linear parameters and nonlinear parameters, and
 22 denote by $\boldsymbol{\theta} = (c_0, c_1, \dots, c_n, b_1, \dots, b_n)^T \in \mathbb{R}^{2n+1}$ all parameters. Notice that the weights of
 23 all hidden layer is chosen to be one by normalization (see [11]). Each function in $\mathcal{M}_n(I)$

$$v(x; \mathbf{c}, \mathbf{b}) = \alpha + \sum_{i=0}^n c_i \sigma(x - b_i)$$

24 depends on the parameters \mathbf{c} and \mathbf{b} and is a piecewise linear function of $x \in [0, 1]$ with
 25 respect to the partition by the breaking (mesh) points: $0 = b_0 < b_1 < \dots < b_n < b_{n+1} = 1$.

26 It is known (see, e.g., [11]) that the set $\mathcal{M}_n(I)$ is equivalent to the free-knot splines
 27 (FKS) (see [16]), where FKS utilizes local hat basis functions. FKS can substantially en-
 28 hance the approximation order and reduce the number of degrees of freedom for non-
 29 smooth functions (see [3] and the discussion in [5]). As an example of this, the order of
 30 the best approximation to $f(x) = x^\alpha$ ($0 < \alpha < 1$) on I is merely $\alpha < 1$ (i.e., $\mathcal{O}(n^{-\alpha})$) when
 31 using finite elements on a fixed uniform mesh, whereas for FKS, the order becomes one
 32 (i.e., $\mathcal{O}(n^{-1})$) no matter how small the exponent $\alpha > 0$ is (see, e.g., [8, 16]). This is a huge
 33 improvement in approximation.

34 Despite the remarkable approximation capability of FKS for non-smooth functions,
 35 there are two essential difficulties that have led numerical analysts moving away from
 36 FKS: (1) no successful extension of FKS to two or higher dimensions has been achieved,
 37 and (2) determining the optimal knot locations (the nonlinear parameters \mathbf{b}) results in a
 38 high-dimensional, non-convex optimization problem, that is computationally expensive
 39 and hence dismisses its benefit in approximation. To the best of our knowledge, there
 40 are no available efficient optimization schemes that enable the competitiveness of FKS.
 41 While the first issue may be addressed by employing neural networks due to their global
 42 supported basis functions, the second issue still remains a major challenge.

43 To address this challenge, a major advance on fast iterative solver was recently made
 44 in [5, 6] for solving non-convex optimization problems arising from shallow ReLU NN
 45 approximation to a given function or solutions of elliptic differential equations in one
 46 dimension. Specifically, a well-designed damped block Newton (dBN) method was de-
 47 veloped. First, the dBN adopts a classical outer-inner iterative strategy (see, e.g., [1, 2, 7,
 48 10, 12, 14, 15]), alternating between updates of the linear and nonlinear parameters. Sec-
 49 ond, to solve the resulting dense, ill-conditioned linear systems due to the global basis
 50 functions of NNs, the dBN uses the fact that the exact inversion of those matrices can
 51 be represented in terms of products of sparse matrices (see [6]). Third, the dBN deals di-
 52 rectly with singularities of the Hessian for the nonlinear parameters by removing neurons
 53 whose linear parameters are small or whose nonlinear parameters have reached nearly
 54 optimal locations. As a result, the computational cost per iteration of the dBN is $\mathcal{O}(n)$,
 55 and numerical experiments show that this method is capable of moving mesh points ef-
 56 fectively and efficiently. Beyond strong one dimensional results, the methodology of the
 57 dBN is conceptually promising for higher dimensions, as it demonstrates how to design

58 iterative solvers that exploit the problem structure and the approximation and geometric
59 meanings of the NN parameters.

60 The purpose of this paper is to provide a theoretical guarantee on why this sophisti-
61 cated dBN moves the mesh points efficiently. This will be done by analyzing local con-
62 vergence of block Newton (BN) methods. The BN method consists of the 2×2 block
63 nonlinear Gauss-Seidel, linear Gauss-Seidel, or Jacobi method for outer iteration and the
64 Newton method for inner iteration. The blocks are corresponding to the linear and the
65 nonlinear parameters. By following the machinery in [13] on local convergence of the
66 componentwise Gauss-Seidel method, we first develop a local convergence theory for
67 the block Newton methods applicable to shallow ReLU NNs in both one and multiple
68 dimensions, provided that the Hessian matrix at a critical point is symmetric positive
69 definite (SPD) and that the 2×2 block nonlinear Gauss-Seidel, linear Gauss-Seidel, or
70 Jacobi matrix is invertible.

71 By expressing the BN method as a fixed-point iteration, local convergence of the BN
72 method is established by showing that a norm of the corresponding Jacobian at the criti-
73 cal point is strictly less than one. Note that derivation of the Jacobian matrix is non-trivial
74 (see Lemma 3.1). To guarantee feasibility of each Newton step, the BN method is modi-
75 fied to the reduced BN (rBN) method that allows a reduction in the number of parameters
76 during the optimization process. Local convergence of the rBN is justified by showing
77 that some nonlinear parameters are at nearly optimal locations.

78 The paper is structured as follows. Section 2 introduces BN methods with three differ-
79 ent outer iteration methods. Local convergence analysis of these methods are presented
80 in Section 3. Sufficient conditions on SPD of the Hessian for one-dimensional problems
81 are derived in Section 4. Section 4.4 discusses the reduced Block Newton method and
82 some conclusions and remarks are presented in Section 5.

83 2 Block Newton methods

84 Let $\theta = \begin{pmatrix} \mathbf{c} \\ \mathbf{b} \end{pmatrix}$, where $\mathbf{c} \in \mathbb{R}^{n+1}$ and $\mathbf{b} \in \mathbb{R}^n$. Given an open set $\mathcal{D} \subseteq \mathbb{R}^{2n+1}$ and a twice
85 continuously differentiable function $F = F(\theta) = F(\mathbf{c}, \mathbf{b}) : \mathcal{D} \rightarrow \mathbb{R}$, we aim to find a minimizer
86 $\theta^* = \begin{pmatrix} \mathbf{c}^* \\ \mathbf{b}^* \end{pmatrix} \in \mathcal{D}$ such that

$$F(\theta^*) = \min_{\theta \in \mathcal{D}} F(\theta).$$

87 Optimality conditions imply that θ^* satisfies the system of nonlinear algebraic equations

$$\nabla_{\theta} F(\theta) = \begin{pmatrix} \nabla_{\mathbf{c}} F(\theta) \\ \nabla_{\mathbf{b}} F(\theta) \end{pmatrix} = \begin{pmatrix} \mathbf{0} \\ \mathbf{0} \end{pmatrix}, \quad (2.1)$$

88 where $\nabla_{\mathbf{c}}$ and $\nabla_{\mathbf{b}}$ denote the gradients of $F(\boldsymbol{\theta})$ with respect to \mathbf{c} and \mathbf{b} , respectively. The
89 Hessian matrix is given by

$$\nabla_{\boldsymbol{\theta}}^2 F(\boldsymbol{\theta}) = \begin{pmatrix} \nabla_{\mathbf{c}\mathbf{c}}^2 F(\boldsymbol{\theta}) & \nabla_{\mathbf{c}\mathbf{b}}^2 F(\boldsymbol{\theta}) \\ \nabla_{\mathbf{b}\mathbf{c}}^2 F(\boldsymbol{\theta}) & \nabla_{\mathbf{b}\mathbf{b}}^2 F(\boldsymbol{\theta}) \end{pmatrix} = \begin{pmatrix} \mathcal{H}_{11}(\boldsymbol{\theta}) & \mathcal{H}_{12}(\boldsymbol{\theta}) \\ \mathcal{H}_{21}(\boldsymbol{\theta}) & \mathcal{H}_{22}(\boldsymbol{\theta}) \end{pmatrix} \in \mathbb{R}^{(2n+1) \times (2n+1)},$$

90 where $\mathcal{H}_{ij}(\boldsymbol{\theta})$ for $i, j = 1, 2$ are given by

$$\mathcal{H}_{11}(\boldsymbol{\theta}) = \nabla_{\mathbf{c}\mathbf{c}}^2 F(\boldsymbol{\theta}), \quad \mathcal{H}_{12}(\boldsymbol{\theta}) = \nabla_{\mathbf{c}\mathbf{b}}^2 F(\boldsymbol{\theta}), \quad \mathcal{H}_{21}(\boldsymbol{\theta}) = \nabla_{\mathbf{b}\mathbf{c}}^2 F(\boldsymbol{\theta}), \quad \mathcal{H}_{22}(\boldsymbol{\theta}) = \nabla_{\mathbf{b}\mathbf{b}}^2 F(\boldsymbol{\theta}).$$

91 The nonlinear system in Eq. (2.1) can be solved using Newton's method, though this
92 may be computationally expensive. To reduce cost, we can use a block Newton (BN)
93 method by performing an outer-inner iteration that alternates between the variables \mathbf{c}
94 and \mathbf{b} . For the outer iteration, one may use a Gauss–Seidel or Jacobi scheme, and apply a
95 Newton iteration to each block during the inner solve.

96 More specifically, let $(\mathbf{c}^{(k)}, \mathbf{b}^{(k)})$ denote the current iterate. Then the block nonlinear
97 Gauss-Seidel (see, e.g., [13]) method, as the outer iteration, computes the new iterate
98 $(\mathbf{c}^{(k+1)}, \mathbf{b}^{(k+1)})$ as follows:

99 1. Update the variable $\mathbf{c}^{(k+1)}$:

$$\mathbf{c}^{(k+1)} = \mathbf{c}^{(k)} - \left[\mathcal{H}_{11}(\mathbf{c}^{(k)}, \mathbf{b}^{(k)}) \right]^{-1} \nabla_{\mathbf{c}} F(\mathbf{c}^{(k)}, \mathbf{b}^{(k)}).$$

100 2. Update the variable $\mathbf{b}^{(k+1)}$:

$$\mathbf{b}^{(k+1)} = \mathbf{b}^{(k)} - \left[\mathcal{H}_{22}(\mathbf{c}^{(k+1)}, \mathbf{b}^{(k)}) \right]^{-1} \nabla_{\mathbf{b}} F(\mathbf{c}^{(k+1)}, \mathbf{b}^{(k)}).$$

101 In other words, the new iterate $(\mathbf{c}^{(k+1)}, \mathbf{b}^{(k+1)})$ is the solution of the block diagonal system
102 of nonlinear algebraic equations

$$\begin{pmatrix} \mathcal{H}_{11}(\mathbf{c}^{(k)}, \mathbf{b}^{(k)}) & \mathbf{0} \\ \mathbf{0} & \mathcal{H}_{22}(\mathbf{c}^{(k+1)}, \mathbf{b}^{(k)}) \end{pmatrix} \begin{pmatrix} \mathbf{c}^{(k+1)} - \mathbf{c}^{(k)} \\ \mathbf{b}^{(k+1)} - \mathbf{b}^{(k)} \end{pmatrix} = - \begin{pmatrix} \nabla_{\mathbf{c}} F(\mathbf{c}^{(k)}, \mathbf{b}^{(k)}) \\ \nabla_{\mathbf{b}} F(\mathbf{c}^{(k+1)}, \mathbf{b}^{(k)}) \end{pmatrix}, \quad (2.2)$$

103 which can be solved sequentially by computing two systems of linear algebraic equa-
104 tions.

105 This nonlinear Gauss-Seidel method differs from the classical linear Gauss–Seidel
106 method, where the new iterate is obtained by solving the block lower-triangular system
107 of linear equations

$$\begin{pmatrix} \mathcal{H}_{11}(\mathbf{c}^{(k)}, \mathbf{b}^{(k)}) & \mathbf{0} \\ \mathcal{H}_{21}(\mathbf{c}^{(k)}, \mathbf{b}^{(k)}) & \mathcal{H}_{22}(\mathbf{c}^{(k)}, \mathbf{b}^{(k)}) \end{pmatrix} \begin{pmatrix} \mathbf{c}^{(k+1)} - \mathbf{c}^{(k)} \\ \mathbf{b}^{(k+1)} - \mathbf{b}^{(k)} \end{pmatrix} = - \begin{pmatrix} \nabla_{\mathbf{c}} F(\mathbf{c}^{(k)}, \mathbf{b}^{(k)}) \\ \nabla_{\mathbf{b}} F(\mathbf{c}^{(k)}, \mathbf{b}^{(k)}) \end{pmatrix}. \quad (2.3)$$

108 On the other hand, in the Jacobi method, the new iterate $(\mathbf{c}^{(k+1)}, \mathbf{b}^{(k+1)})$ solves the block
109 diagonal system of linear algebraic equations

$$\begin{pmatrix} \mathcal{H}_{11}(\mathbf{c}^{(k)}, \mathbf{b}^{(k)}) & \mathbf{0} \\ \mathbf{0} & \mathcal{H}_{22}(\mathbf{c}^{(k)}, \mathbf{b}^{(k)}) \end{pmatrix} \begin{pmatrix} \mathbf{c}^{(k+1)} - \mathbf{c}^{(k)} \\ \mathbf{b}^{(k+1)} - \mathbf{b}^{(k)} \end{pmatrix} = - \begin{pmatrix} \nabla_{\mathbf{c}} F(\mathbf{c}^{(k)}, \mathbf{b}^{(k)}) \\ \nabla_{\mathbf{b}} F(\mathbf{c}^{(k)}, \mathbf{b}^{(k)}) \end{pmatrix}. \quad (2.4)$$

110 3 Convergence analysis

111 For brevity, we refer to the schemes defined in (2.2), (2.3), and (2.4) as NL-GS, L-GS, and
112 JB, respectively. This section presents analytic tools for deriving local convergence con-
113 ditions for the BN methods introduced in the previous section. To this end, we introduce
114 the following assumption:

115 **Invertibility Assumption:** there exists an open set $\mathcal{O} \subseteq \mathcal{D}$ such that $\mathcal{H}_{11}(\boldsymbol{\theta})$ is invertible;
116 moreover, $\mathcal{H}_{22}(G_1(\boldsymbol{\theta}), \mathbf{b})$ is invertible for NL-GS and $\mathcal{H}_{22}(\boldsymbol{\theta})$ is invertible for L-GS and JB.

117 Define the mapping

$$G(\boldsymbol{\theta}) = \begin{pmatrix} G_1(\boldsymbol{\theta}) \\ G_2(\boldsymbol{\theta}) \end{pmatrix},$$

118 where $G_1: \mathcal{D} \rightarrow \mathbb{R}^{n+1}$ and $G_2: \mathcal{D} \rightarrow \mathbb{R}^n$ are respectively given by

$$G_1(\boldsymbol{\theta}) = \mathbf{c} - \mathcal{H}_{11}^{-1}(\boldsymbol{\theta}) \nabla_{\mathbf{c}} F(\boldsymbol{\theta}), \quad (3.1)$$

119 and by

$$G_2(\boldsymbol{\theta}) = \begin{cases} \mathbf{b} - \mathcal{H}_{22}^{-1}(G_1(\boldsymbol{\theta}), \mathbf{b}) \nabla_{\mathbf{b}} F(G_1(\boldsymbol{\theta}), \mathbf{b}) & \text{for NL-GS,} \\ \mathbf{b} - \mathcal{H}_{22}^{-1}(\boldsymbol{\theta}) \left(\nabla_{\mathbf{b}} F(\boldsymbol{\theta}) + \mathcal{H}_{21}(\boldsymbol{\theta}) (G_1(\boldsymbol{\theta}) - \mathbf{c}) \right) & \text{for L-GS,} \\ \mathbf{b} - \mathcal{H}_{22}^{-1}(\boldsymbol{\theta}) \nabla_{\mathbf{b}} F(\boldsymbol{\theta}), & \text{for JB.} \end{cases} \quad (3.2)$$

120 Then the BN methods can be expressed as the fixed-point iteration of $G(\boldsymbol{\theta})$, i.e.,

$$\boldsymbol{\theta}^{k+1} = G(\boldsymbol{\theta}^k) \quad \text{for } k \in \mathbb{N} = \{0, 1, 2, \dots\}. \quad (3.3)$$

121 Denote by $\mathbf{J}_G(\boldsymbol{\theta}) \in \mathbb{R}^{(2n+1) \times (2n+1)}$ the Jacobian matrix of G at $\boldsymbol{\theta}$. A sufficient condition
122 for local convergence of the fixed-point iteration in (3.3) was given in Theorem 10.1.3
123 of [13]. For the convenience of readers, we state it below and provide its proof.

124 **Theorem 3.1** (Ostrowski). *Suppose that $G: \mathcal{O} \rightarrow \mathbb{R}^{2n+1}$ has a fixed point $\boldsymbol{\theta}^* \in \mathcal{O}$ and that the
125 mapping G is differentiable at $\boldsymbol{\theta}^*$. Denote by $\|\cdot\|$ a norm in \mathbb{R}^{2n+1} , if $\|\mathbf{J}_G(\boldsymbol{\theta}^*)\| = \sigma < 1$, then the
126 fixed-point iteration in (3.3) converges locally to $\boldsymbol{\theta}^*$.*

127 *Proof.* For any $0 < \epsilon < 1 - \sigma$, the assumption on the differentiability of G at θ^* implies that
 128 there exists a $\delta > 0$ neighbourhood centered at θ^* , $\mathcal{B}(\theta^*; \delta) = \{\theta \in \mathbb{R}^{2n+1} : \|\theta - \theta^*\| < \delta\} \subset \mathcal{O}$,
 129 such that

$$\|G(\theta) - [G(\theta^*) + \mathbf{J}_G(\theta^*)(\theta - \theta^*)]\| \leq \epsilon \|\theta - \theta^*\| \quad \text{for all } \theta \in \mathcal{B}(\theta^*; \delta). \quad (3.4)$$

130 For any $k \in \mathbb{N}$, the assumption that $\theta^* = G(\theta^*)$ and (3.3) give

$$\theta^{k+1} - \theta^* = G(\theta^k) - G(\theta^*) = \left(G(\theta^k) - [G(\theta^*) + \mathbf{J}_G(\theta^*)(\theta^k - \theta^*)] \right) + \mathbf{J}_G(\theta^*)(\theta^k - \theta^*),$$

131 which, together with the triangle inequality, the assumption that $\theta^k \in \mathcal{B}(\theta^*; \delta)$, and (3.4),
 132 implies

$$\|\theta^{k+1} - \theta^*\| \leq (\epsilon + \sigma) \|\theta^k - \theta^*\| < \|\theta^k - \theta^*\|.$$

If the initial θ^0 belongs to $\mathcal{B}(\theta^*; \delta)$, then the second inequality implies that $\theta^k \in \mathcal{B}(\theta^*; \delta)$ by induction. Hence, we have

$$\|\theta^{k+1} - \theta^*\| \leq (\epsilon + \sigma)^k \|\theta^0 - \theta^*\|,$$

133 which proves the theorem. □

134 Next, let

$$B(\theta) = \begin{pmatrix} \mathcal{H}_{11}(\theta) & \mathbf{0} \\ \mathcal{H}_{21}(\theta) & \mathcal{H}_{22}(\theta) \end{pmatrix} \quad \text{and} \quad B_1(\theta) = \begin{pmatrix} \mathcal{H}_{11}(\theta) & \mathbf{0} \\ \mathbf{0} & \mathcal{H}_{22}(G_1(\theta), \mathbf{b}) \end{pmatrix}.$$

135 Assume that θ^* is a minimizer of $F(\theta)$, then Taylor's expansion of $\nabla_{\theta} F(\theta)$ at θ^* gives

$$\nabla_{\theta} F(\theta) = \nabla_{\theta}^2 F(\theta^*)(\theta - \theta^*) + R(\theta; \theta^*), \quad (3.5)$$

where $R(\theta; \theta^*)$ is the remainder satisfying

$$\lim_{\theta \rightarrow \theta^*} \|R(\theta; \theta^*)\| / \|\theta - \theta^*\| = 0.$$

136 Similarly, expanding $\nabla_{\mathbf{b}} F(G_1(\theta), \mathbf{b})$ about θ and using (3.1), we have

$$\nabla_{\mathbf{b}} F(G_1(\theta), \mathbf{b}) = \nabla_{\mathbf{b}} F(\theta) - \mathcal{H}_{21}(\theta) \mathcal{H}_{11}^{-1}(\theta) \nabla_{\mathbf{c}} F(\theta) + \tilde{R}(\theta),$$

where the remainder $\tilde{R}(\theta)$ satisfies

$$\lim_{\theta \rightarrow \theta^*} \frac{\|\tilde{R}(\theta)\|}{\|\theta - \theta^*\|} = 0.$$

Let

$$B_2(\boldsymbol{\theta}) = \begin{pmatrix} I & \mathbf{0} \\ -\mathcal{H}_{21}(\boldsymbol{\theta})\mathcal{H}_{11}^{-1}(\boldsymbol{\theta}) & I \end{pmatrix},$$

137 then combining with (3.5) yields

$$\begin{pmatrix} \nabla_{\mathbf{c}}F(\boldsymbol{\theta}) \\ \nabla_{\mathbf{b}}F(G_1(\boldsymbol{\theta}), \mathbf{b}) \end{pmatrix} = B_2(\boldsymbol{\theta})\nabla_{\boldsymbol{\theta}}F(\boldsymbol{\theta}) + \begin{pmatrix} \mathbf{0} \\ \tilde{R}(\boldsymbol{\theta}) \end{pmatrix} = B_2(\boldsymbol{\theta})\nabla_{\boldsymbol{\theta}}^2F(\boldsymbol{\theta}^*)(\boldsymbol{\theta} - \boldsymbol{\theta}^*) + \hat{R}(\boldsymbol{\theta}; \boldsymbol{\theta}^*), \quad (3.6)$$

where

$$\hat{R}(\boldsymbol{\theta}; \boldsymbol{\theta}^*) = B_2(\boldsymbol{\theta})R(\boldsymbol{\theta}; \boldsymbol{\theta}^*) + \begin{pmatrix} \mathbf{0} \\ \tilde{R}(\boldsymbol{\theta}) \end{pmatrix}.$$

138 The next lemma provides a formula for $\mathbf{J}_G(\boldsymbol{\theta}^*)$ when $\nabla_{\boldsymbol{\theta}}^2F(\boldsymbol{\theta}^*)$ is symmetric and positive definite (SPD) for Gauss–Seidel schemes. For the componentwise Gauss–Seidel method, this result was proved in Theorem 10.3.3 of [13]. Here, we use block Gauss–Seidel methods.

142 **Lemma 3.1.** *Let $G : \mathcal{O} \rightarrow \mathbb{R}^{2n+1}$ be the mapping defined in (3.2) for the L-GS or NL-GS, and let $\boldsymbol{\theta}^* \in \mathcal{O}$ be a fixed point of G . Assume that $\nabla_{\boldsymbol{\theta}}^2F(\boldsymbol{\theta}^*)$ is SPD. Then*

$$\mathbf{J}_G(\boldsymbol{\theta}^*) = I_{2n+1} - B^{-1}(\boldsymbol{\theta}^*)\nabla_{\boldsymbol{\theta}}^2F(\boldsymbol{\theta}^*), \quad (3.7)$$

144 where I_{2n+1} is the order- $(2n+1)$ identity matrix.

145 *Proof.* By definition of differentiability, to show the validity of (3.7), it suffices to show that if $\mathbf{J}_G(\boldsymbol{\theta}^*)$ is given in (3.7), then

$$\lim_{\boldsymbol{\theta} \rightarrow \boldsymbol{\theta}^*} \frac{\|G(\boldsymbol{\theta}) - [G(\boldsymbol{\theta}^*) + \mathbf{J}_G(\boldsymbol{\theta}^*)(\boldsymbol{\theta} - \boldsymbol{\theta}^*)]\|}{\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|} = 0. \quad (3.8)$$

147 To this end, for $\mathbf{J}_G(\boldsymbol{\theta}^*)$ given in (3.7), we have

$$B(\boldsymbol{\theta}^*)\mathbf{J}_G(\boldsymbol{\theta}^*) = \begin{pmatrix} \mathbf{0} & -\mathcal{H}_{12}(\boldsymbol{\theta}^*) \\ \mathbf{0} & \mathbf{0} \end{pmatrix} \quad \text{and} \quad B(\boldsymbol{\theta}^*) - B(\boldsymbol{\theta}^*)\mathbf{J}_G(\boldsymbol{\theta}^*) = \nabla_{\boldsymbol{\theta}}^2F(\boldsymbol{\theta}^*),$$

148 which, together with the assumption that $\boldsymbol{\theta}^*$ is a fixed-point of $G(\boldsymbol{\theta})$, implies

$$\begin{aligned} \mathbf{a}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) &\equiv B(\boldsymbol{\theta}^*) \left(G(\boldsymbol{\theta}) - [G(\boldsymbol{\theta}^*) + \mathbf{J}_G(\boldsymbol{\theta}^*)(\boldsymbol{\theta} - \boldsymbol{\theta}^*)] \right) \\ &= B(\boldsymbol{\theta}^*) \left(G(\boldsymbol{\theta}) - \boldsymbol{\theta} \right) + \nabla_{\boldsymbol{\theta}}^2F(\boldsymbol{\theta}^*)(\boldsymbol{\theta} - \boldsymbol{\theta}^*). \end{aligned}$$

149 For the L-GS, (3.5) gives

$$G(\boldsymbol{\theta}) - \boldsymbol{\theta} = -B^{-1}(\boldsymbol{\theta})\nabla_{\boldsymbol{\theta}}F(\boldsymbol{\theta}) = -B^{-1}(\boldsymbol{\theta})\nabla_{\boldsymbol{\theta}}^2F(\boldsymbol{\theta}^*)(\boldsymbol{\theta} - \boldsymbol{\theta}^*) - B^{-1}(\boldsymbol{\theta})R(\boldsymbol{\theta}; \boldsymbol{\theta}^*),$$

150 which implies

$$\mathbf{a}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \left(B(\boldsymbol{\theta}) - B(\boldsymbol{\theta}^*) \right) B^{-1}(\boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}}^2 F(\boldsymbol{\theta}^*) (\boldsymbol{\theta} - \boldsymbol{\theta}^*) - B(\boldsymbol{\theta}^*) B^{-1}(\boldsymbol{\theta}) R(\boldsymbol{\theta}; \boldsymbol{\theta}^*). \quad (3.9)$$

151 For the NL-GS, (3.6) leads to

$$G(\boldsymbol{\theta}) - \boldsymbol{\theta} = -B_1^{-1}(\boldsymbol{\theta}) B_2(\boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}}^2 F(\boldsymbol{\theta}^*) (\boldsymbol{\theta} - \boldsymbol{\theta}^*) - B_1^{-1}(\boldsymbol{\theta}) \hat{R}(\boldsymbol{\theta}; \boldsymbol{\theta}^*),$$

which, together with

$$B_3(\boldsymbol{\theta}) \equiv B_2^{-1}(\boldsymbol{\theta}) B_1(\boldsymbol{\theta}) = \begin{pmatrix} \mathcal{H}_{11}(\boldsymbol{\theta}) & \mathbf{0} \\ \mathcal{H}_{21}(\boldsymbol{\theta}) & \mathcal{H}(G_1(\boldsymbol{\theta}), \mathbf{b}) \end{pmatrix},$$

152 yields

$$\mathbf{a}(\boldsymbol{\theta}, \boldsymbol{\theta}^*) = \left(B_3(\boldsymbol{\theta}) - B(\boldsymbol{\theta}^*) \right) B_1^{-1}(\boldsymbol{\theta}) B_2(\boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}}^2 F(\boldsymbol{\theta}^*) (\boldsymbol{\theta} - \boldsymbol{\theta}^*) - B(\boldsymbol{\theta}^*) B_1^{-1}(\boldsymbol{\theta}) \hat{R}(\boldsymbol{\theta}; \boldsymbol{\theta}^*). \quad (3.10)$$

153 By the assumptions that $F(\boldsymbol{\theta})$ is twice continuously differentiable and that $\nabla_{\boldsymbol{\theta}}^2 F(\boldsymbol{\theta}^*)$ is
 154 SPD, there exists a $\delta > 0$ neighbourhood centered at $\boldsymbol{\theta}^*$, $\mathcal{B}(\boldsymbol{\theta}^*; \delta) = \{\boldsymbol{\theta} \in \mathbb{R}^{2n+1} : \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| <$
 155 $\delta\} \subset \mathcal{O}$, such that $\|B^{-1}(\boldsymbol{\theta})\|$ and $\|B_1^{-1}(\boldsymbol{\theta})\|$ are bounded for all $\boldsymbol{\theta} \in \mathcal{B}(\boldsymbol{\theta}^*; \delta)$ and that

$$\lim_{\boldsymbol{\theta} \rightarrow \boldsymbol{\theta}^*} \left(B(\boldsymbol{\theta}) - B(\boldsymbol{\theta}^*) \right) = \mathbf{0} \quad \text{and} \quad \lim_{\boldsymbol{\theta} \rightarrow \boldsymbol{\theta}^*} \left(B_3(\boldsymbol{\theta}) - B(\boldsymbol{\theta}^*) \right) = \mathbf{0}.$$

156 Now, (3.8) is a direct consequence of (3.9), (3.10), the triangle inequality, and the facts that

$$\lim_{\boldsymbol{\theta} \rightarrow \boldsymbol{\theta}^*} \|R(\boldsymbol{\theta}; \boldsymbol{\theta}^*)\| / \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| = 0 \quad \text{and} \quad \lim_{\boldsymbol{\theta} \rightarrow \boldsymbol{\theta}^*} \|\hat{R}(\boldsymbol{\theta}; \boldsymbol{\theta}^*)\| / \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\| = 0.$$

157 This completes the proof of the lemma. \square

158 The following lemma was stated and proved in [9]. For convenience of readers, a
 159 brief proof is provided.

160 **Lemma 3.2.** Assume that $A \in \mathbb{R}^{m \times m}$ is SPD and that $M \in \mathbb{R}^{m \times m}$ is invertible. Then the matrix
 161 $M + M^T - A$ is SPD if and only if

$$\|I_{2n+1} - M^{-1}A\|_A < 1,$$

162 where $\|\mathbf{v}\|_A = \sqrt{\langle A\mathbf{v}, \mathbf{v} \rangle}$ and $\langle \cdot, \cdot \rangle$ is the standard inner product in \mathbb{R}^m .

163 *Proof.* For any $\mathbf{v} \in \mathbb{R}^m$, it is easy to check that

$$\|(I - M^{-1}A)\mathbf{v}\|_A^2 = \langle A\mathbf{v}, \mathbf{v} \rangle - \left\langle (M + M^T - A)(M^{-1}A)\mathbf{v}, (M^{-1}A)\mathbf{v} \right\rangle,$$

164 which implies the validity of the lemma. \square

165 Now, we are ready to present a sufficient condition for local convergence of the fixed-
 166 point iteration in (3.3) for the NL-GS and L-GS.

167 **Theorem 3.2.** *If $\nabla_{\theta}^2 F(\theta^*)$ is SPD, then the fixed point iteration (3.3) for the L-GS or NL-GS
 168 converges locally to θ^* in the norm induced by $\nabla_{\theta}^2 F(\theta^*)$.*

169 *Proof.* Symmetry and positive definiteness of the Hessian $\nabla_{\theta}^2 F(\theta^*)$ implies that

$$B(\theta^*)^T + B(\theta^*) - \nabla_{\theta}^2 F(\theta^*) = \begin{pmatrix} \mathcal{H}_{11}(\theta^*) & \mathbf{0} \\ \mathbf{0} & \mathcal{H}_{22}(\theta^*) \end{pmatrix}$$

170 is SPD. Hence, by Lemma 3.2, we have

$$\|I_{2n+1} - B(\theta^*)^{-1} \nabla_{\theta}^2 F(\theta^*)\|_{\nabla_{\theta}^2 F(\theta^*)} = \sigma < 1,$$

171 which, together with Theorem 3.1 and Lemma 3.1, implies the validity of the theorem.
 172 This completes the proof. □

173 **Remark 3.1.** For the Jacobi method described in (2.4), if $\nabla^2 F(\theta^*)$ is SPD, then a similar
 174 argument as that of Lemma 3.1 gives the following Jacobian matrix $J_G(\theta^*)$ of $G(\theta)$ at the
 175 fixed point θ^*

$$J_G(\theta^*) = I_{2n+1} - \tilde{B}(\theta^*)^{-1} \nabla_{\theta}^2 F(\theta^*), \quad \text{where } \tilde{B}(\theta) = \begin{pmatrix} \mathcal{H}_{11}(\theta) & \mathbf{0} \\ \mathbf{0} & \mathcal{H}_{22}(\theta) \end{pmatrix}.$$

176 Moreover, the iteration is locally convergent to θ^* if and only if

$$\tilde{B}(\theta^*)^T + \tilde{B}(\theta^*) - \nabla_{\theta}^2 F(\theta^*) = \begin{pmatrix} \mathcal{H}_{11}(\theta^*) & -\mathcal{H}_{12}(\theta^*) \\ -\mathcal{H}_{21}(\theta^*) & \mathcal{H}_{22}(\theta^*) \end{pmatrix}$$

177 is SPD.

178 4 Applications

179 This section applies the local convergence theory developed in the previous section to
 180 one-dimensional least-squares (LS) function approximation and diffusion-reaction (DR)
 181 problems studied in [6].

182 For the diffusion-reaction problem, let u be the exact solution of the differential equa-
 183 tion

$$\begin{cases} -(a(x)u'(x))' + r(x)u(x) = f(x) & \text{in } I = (0,1), \\ u(0) = \alpha, \quad u(1) = \beta, \end{cases} \quad (4.1)$$

184 where the diffusion coefficient $a(x)$, the reaction coefficient $r(x)$, and the right-hand
 185 side $f(x)$ are given real-valued functions defined on I . Assume that $a(x)$ and $r(x)$ are

186 bounded below by the respective positive constant $\mu > 0$ and non-negative constant $r_0 \geq 0$
 187 almost everywhere on I . Assume that $f(x) \in C(I)$ and that $a(x) \in C^1(I)$. For the least-
 188 squares approximation problem, we assume that the weight function $r(x) \geq r_0 > 0$. For
 189 these problems, the NN approximation is to seek $u_n^*(x) = u_n(x; \theta^*) \in \mathcal{M}_n(I)$ such that

$$u_n^*(x) = u_n(x; \theta^*) = \underset{u_n \in \mathcal{M}_n(I)}{\operatorname{argmin}} J(u_n(\cdot; \theta)), \quad (4.2)$$

190 where the functional $J(v)$ is given by

$$J(v) = \begin{cases} \frac{1}{2} \int_0^1 [a(x)(v'(x))^2 + r(x)(v(x))^2] dx - \int_0^1 f(x)v(x) dx + \frac{\gamma}{2}(v(1) - \beta)^2, & \text{DR,} \\ \frac{1}{2} \int_0^1 r(x)[(v(x) - u(x))^2] dx, & \text{LS.} \end{cases} \quad (4.3)$$

191 4.1 Hessian

192 First, we calculate the Hessian matrices. To this end, denote by

$$H(t) = \begin{cases} 1, & t > 0, \\ \frac{1}{2}, & t = 0, \\ 0, & t < 0, \end{cases} \quad \text{and} \quad \delta(t) = \begin{cases} +\infty, & t = 0, \\ 0, & t \neq 0, \end{cases}$$

193 the Heaviside (unit) step function and the Dirac delta function, respectively. Clearly,
 194 $H(t) = \sigma'(t)$ and $\delta(t) = \sigma''(t)$ everywhere except at $t = 0$. For each $i = 0, 1, \dots, n$, let

$$\sigma_i(x) = \sigma(x - b_i), \quad H_i(x) = H(x - b_i), \quad \delta_i(x) = \delta(x - b_i),$$

195 and let

$$\begin{aligned} \mathbf{\Sigma}_{n+1}(x) &= \begin{pmatrix} \sigma_0(x) \\ \sigma_1(x) \\ \vdots \\ \sigma_n(x) \end{pmatrix}, & \mathbf{H}_n(x) &= \begin{pmatrix} H_1(x) \\ \vdots \\ H_n(x) \end{pmatrix}, \\ \mathbf{H}_{n+1}(x) &= \begin{pmatrix} H_0(x) \\ H_1(x) \\ \vdots \\ H_n(x) \end{pmatrix}, & \mathbf{\Lambda}_n(x) &= \begin{pmatrix} \delta_1(x) \\ \vdots \\ \delta_n(x) \end{pmatrix}. \end{aligned}$$

196 For $i = 1, \dots, n$, let

$$g_i = g_i(\theta) = \begin{cases} r(b_i)u_n(b_i) - f(b_i) - a'(b_i)u_n'(b_i), & \text{DR,} \\ r(b_i)(u_n(b_i) - u(b_i)), & \text{LS,} \end{cases} \quad (4.4)$$

197 where $u'_n(b_i) := \sum_{j=0}^{i-1} c_j + \frac{c_i}{2}$, and set

$$\mathbf{g}(\boldsymbol{\theta}) = (g_1, \dots, g_n)^T, \quad \mathbf{d} = (1 - b_1, \dots, 1 - b_n)^T, \quad \text{and} \quad \hat{\mathbf{c}} = (c_1, \dots, c_n)^T.$$

198 For any $u_n(x; \boldsymbol{\theta}) \in \mathcal{M}_n(I)$, the value of the functional $J(u_n(\cdot; \boldsymbol{\theta}))$ at $u_n(x; \boldsymbol{\theta})$ is a function of
 199 parameters $\boldsymbol{\theta}$. For simplicity of notation, denote this function by F , i.e., $F(\boldsymbol{\theta}) = J(u_n(\cdot; \boldsymbol{\theta}))$.
 200 In [6], we derived the principle blocks of the Hessian matrix $\nabla_{\boldsymbol{\theta}}^2 F(\boldsymbol{\theta})$ as follows

$$\mathcal{H}_{11}(\boldsymbol{\theta}) = \begin{cases} \int_0^1 a(x) \mathbf{H}_{n+1} \mathbf{H}_{n+1}^T dx + \int_0^1 r(x) \boldsymbol{\Sigma}_{n+1} \boldsymbol{\Sigma}_{n+1}^T dx + \gamma \mathbf{d} \mathbf{d}^T, & \text{DR,} \\ \int_0^1 r(x) \boldsymbol{\Sigma}_{n+1} \boldsymbol{\Sigma}_{n+1}^T dx, & \text{LS,} \end{cases} \quad (4.5)$$

$$\mathcal{H}_{22}(\boldsymbol{\theta}) = \begin{cases} \mathbf{D}(\hat{\mathbf{c}}) \mathbf{D}(\mathbf{g}) + \mathbf{D}(\hat{\mathbf{c}}) \left(\int_0^1 r(x) \mathbf{H}_n \mathbf{H}_n^T dx \right) \mathbf{D}(\hat{\mathbf{c}}) + \gamma \hat{\mathbf{c}} \hat{\mathbf{c}}^T, & \text{DR,} \\ \mathbf{D}(\hat{\mathbf{c}}) \mathbf{D}(\mathbf{g}) + \mathbf{D}(\hat{\mathbf{c}}) \left(\int_0^1 r(x) \mathbf{H}_n \mathbf{H}_n^T dx \right) \mathbf{D}(\hat{\mathbf{c}}), & \text{LS,} \end{cases} \quad (4.6)$$

201 where $\mathbf{D}(\mathbf{g}) = \text{diag}(g_1, \dots, g_n)$ and $\mathbf{D}(\hat{\mathbf{c}}) = \text{diag}(c_1, \dots, c_n)$ are diagonal matrices.

202 To compute the off-diagonal block $\mathcal{H}_{12}(\boldsymbol{\theta}) = \left(\frac{\partial^2 F(\boldsymbol{\theta})}{\partial c_i \partial b_j} \right)_{(n+1) \times n}$ of the Hessian matrix
 203 $\nabla_{\boldsymbol{\theta}}^2 F(\boldsymbol{\theta})$, let

$$F_j(\boldsymbol{\theta}) = \begin{cases} \int_0^1 \left(f(x) H_j(x) - a(x) u'_n(x) \delta_j(x) - r(x) u_n(x) H_j(x) \right) dx - \gamma (u_n(1) - \beta), & \text{DR,} \\ \int_0^1 r(x) (u(x) - u_n(x)) H_j(x) dx, & \text{LS,} \end{cases}$$

204 for $j = 1, \dots, n$. From (3.5) and (6.3) in [5], we have

$$\frac{\partial}{\partial b_j} F(\boldsymbol{\theta}) = c_j F_j(\boldsymbol{\theta}), \quad (4.7)$$

205 then together with the fact that $\frac{\partial c_j}{\partial c_i} = \delta_{ij}$ (the Kronecker delta), we obtain

$$\frac{\partial^2 F}{\partial c_i \partial b_j}(\boldsymbol{\theta}) = \delta_{ij} F_j(\boldsymbol{\theta}) + c_j \frac{\partial F_j}{\partial c_i}(\boldsymbol{\theta}). \quad (4.8)$$

206 **Lemma 4.1.** For the functional in (4.3), the off-diagonal block $\mathcal{H}_{12}(\boldsymbol{\theta})$ has the form

$$\mathcal{H}_{12}(\boldsymbol{\theta}) = \begin{cases} \mathcal{R}(\boldsymbol{\theta}) - \left(\int_0^1 a(x) \mathbf{H}_{n+1} \boldsymbol{\Lambda}_n^T dx + \int_0^1 r(x) \boldsymbol{\Sigma}_{n+1} \mathbf{H}_n^T dx \right) \mathbf{D}(\hat{\mathbf{c}}) - \gamma \mathbf{d} \hat{\mathbf{c}}^T, & \text{DR,} \\ \mathcal{R}(\boldsymbol{\theta}) - \left(\int_0^1 r(x) \boldsymbol{\Sigma}_{n+1} \mathbf{H}_n^T dx \right) \mathbf{D}(\hat{\mathbf{c}}), & \text{LS,} \end{cases}$$

207 where $F_0(\boldsymbol{\theta}) = 0$ and $\mathcal{R}(\boldsymbol{\theta}) = (\mathcal{R}_{ij}(\boldsymbol{\theta})) = (\delta_{ij} F_i(\boldsymbol{\theta}))_{(n+1) \times n}$.

208 *Proof.* The lemma is a direct consequence of (4.8) and the fact that

$$\frac{\partial F_j}{\partial c_i}(\boldsymbol{\theta}) = \begin{cases} -\int_0^1 a(x)H_i(x)\delta_j(x)dx - \int_0^1 r(x)\sigma_i(x)H_j(x)dx - \gamma(1-b_i), & \text{DR,} \\ -\int_0^1 r(x)\sigma_i(x)H_j(x)dx, & \text{LS,} \end{cases}$$

209 for all $i=0,1,\dots,n$. □

210 For a given function $w: I \rightarrow \mathbb{R}$, define the matrices

$$\mathbb{H}_\Sigma(w; \boldsymbol{\theta}) = \begin{pmatrix} \int_0^1 w(x)\boldsymbol{\Sigma}_{n+1}\boldsymbol{\Sigma}_{n+1}^T dx & -\int_0^1 w(x)\boldsymbol{\Sigma}_{n+1}\mathbf{H}_n^T dx \\ -\int_0^1 w(x)\mathbf{H}_n\boldsymbol{\Sigma}_{n+1}^T dx & \int_0^1 w(x)\mathbf{H}_n\mathbf{H}_n^T dx \end{pmatrix}$$

211 and

$$\mathbb{H}_\Lambda(w; \boldsymbol{\theta}) = \begin{pmatrix} \int_0^1 w(x)\mathbf{H}_{n+1}\mathbf{H}_{n+1}^T dx & -\int_0^1 w(x)\mathbf{H}_{n+1}\boldsymbol{\Lambda}_n^T dx \\ -\int_0^1 w(x)\boldsymbol{\Lambda}_n\mathbf{H}_{n+1}^T dx & \mathbf{0} \end{pmatrix}.$$

212 Let $\mathbb{H}_{\Sigma+\Lambda}(\boldsymbol{\theta}) = \mathbb{H}_\Sigma(r; \boldsymbol{\theta}) + \mathbb{H}_\Lambda(a; \boldsymbol{\theta})$. Combining (4.5), (4.6), and Lemma 4.1, the Hessian
213 matrix at $\boldsymbol{\theta}$ has the form

$$\nabla_{\boldsymbol{\theta}}^2 F(\boldsymbol{\theta}) = \mathbb{R}(\boldsymbol{\theta}) + \begin{cases} \mathcal{D}(\hat{\mathbf{c}})\mathbb{H}_{\Sigma+\Lambda}(\boldsymbol{\theta})\mathcal{D}(\hat{\mathbf{c}}) + \mathcal{D}_0(\boldsymbol{\theta}) + \gamma \begin{pmatrix} \mathbf{d} \\ -\hat{\mathbf{c}} \end{pmatrix} (\mathbf{d}^T, -\hat{\mathbf{c}}^T), & \text{DR,} \\ \mathcal{D}(\hat{\mathbf{c}})\mathbb{H}_\Sigma(r; \boldsymbol{\theta})\mathcal{D}(\hat{\mathbf{c}}) + \mathcal{D}_0(\boldsymbol{\theta}), & \text{LS,} \end{cases} \quad (4.9)$$

where

$$\mathcal{D}(\hat{\mathbf{c}}) = \begin{pmatrix} I_{n+1} & \mathbf{0} \\ \mathbf{0} & \mathbf{D}(\hat{\mathbf{c}}) \end{pmatrix}, \quad \mathcal{D}_0(\boldsymbol{\theta}) = \begin{pmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{D}(\hat{\mathbf{c}})\mathbf{D}(\mathbf{g}) \end{pmatrix} \quad \text{and} \quad \mathbb{R}(\boldsymbol{\theta}) = \begin{pmatrix} \mathbf{0} & \mathcal{R}(\boldsymbol{\theta}) \\ \mathcal{R}^T(\boldsymbol{\theta}) & \mathbf{0} \end{pmatrix}.$$

214 4.2 Local convergence

215 In this section, we present sufficient conditions for the local convergence of the BN method
216 for both the diffusion-reaction and the least-squares approximation problems.

217 First, we establish lower bounds of the smallest eigenvalues of the matrices $\mathbb{H}_\Sigma(w; \boldsymbol{\theta})$
218 and $\mathbb{H}_\Lambda(w; \boldsymbol{\theta})$. To this end, let

$$h_i = b_{i+1} - b_i, \quad i=0, \dots, n; \quad h_{\min} = \min_{0 \leq i \leq n} \{h_i\}, \quad \tilde{h}_i = \min\{h_{i-1}, h_i\}, \quad i=1, \dots, n.$$

219 **Lemma 4.2.** *Suppose that $0 \leq w_0 \leq w(x) \in C(I)$. Then, for any $\alpha = (\alpha_0, \alpha_1, \dots, \alpha_n)^T \in \mathbb{R}^{n+1}$,*
 220 *$\beta = (\beta_1, \dots, \beta_n)^T \in \mathbb{R}^n$, and $\tau \in (0, 1]$, we have*

$$\begin{cases} (\alpha^T, \beta^T) \mathbb{H}_\Sigma(w; \theta) \begin{pmatrix} \alpha \\ \beta \end{pmatrix} \geq \frac{w_0 h_{\min}^3}{96} |\alpha|^2 + \frac{w_0}{24} \sum_{i=1}^n \tilde{h}_i \beta_i^2 \\ (\alpha^T, \beta^T) \mathbb{H}_\Lambda(w; \theta) \begin{pmatrix} \alpha \\ \beta \end{pmatrix} \geq \frac{w_0(1-\tau)h_{\min}}{4} |\alpha|^2 - \frac{1}{2\tau w_0} \sum_{i=1}^n w(b_i)^2 (h_{i-1}^{-1} + h_i^{-1}) \beta_i^2, \end{cases} \quad (4.10)$$

221 where $|\alpha|$ is the magnitude of α .

222 *Proof.* For any $\alpha = (\alpha_0, \alpha_1, \dots, \alpha_n)^T \in \mathbb{R}^{n+1}$ and any $\beta = (\beta_1, \dots, \beta_n)^T \in \mathbb{R}^n$, let $\beta_0 = 0$ and
 223 define the following piecewise linear and constant functions

$$\alpha_i(x) = \sum_{j=0}^i \alpha_j \sigma_j(x) = \sum_{j=0}^i \alpha_j \sigma(x - b_j) \quad \text{and} \quad \beta_i(x) = \sum_{j=0}^i \beta_j H_j(x) \quad \text{for } i = 0, 1, \dots, n,$$

224 respectively. For $x \in (b_i, b_{i+1})$, $\alpha_n(x) = \alpha_i(x)$ and $\beta_n(x) = \beta_i(x) = \sum_{j=0}^i \beta_j := a_i$. Clearly, we
 225 have

$$(\alpha^T, \beta^T) \mathbb{H}_\Sigma(w; \theta) \begin{pmatrix} \alpha \\ \beta \end{pmatrix} = \int_0^1 w(\alpha_n(x) - \beta_n(x))^2 dx \geq w_0 \sum_{i=0}^n \int_{b_i}^{b_{i+1}} (\alpha_i(x) - a_i)^2 dx. \quad (4.11)$$

By the Simpson rule,

$$\alpha_i(b_i) = \alpha_{i-1}(b_i) \quad \text{and} \quad \sum_{i=0}^n (\alpha_i(b_{i+1}) - a_i)^2 \geq \sum_{i=1}^n (\alpha_{i-1}(b_i) - a_{i-1})^2,$$

226 we obtain

$$\begin{aligned} & 12 \sum_{i=0}^n \int_{b_i}^{b_{i+1}} (\alpha_i(x) - a_i)^2 dx \geq 2 \sum_{i=0}^n h_i [(\alpha_i(b_i) - a_i)^2 + (\alpha_i(b_{i+1}) - a_i)^2] \\ & \geq \sum_{i=0}^n h_i [(\alpha_i(b_i) - a_i)^2 + (\alpha_i(b_{i+1}) - a_i)^2] + \sum_{i=1}^n h_i (\alpha_{i-1}(b_i) - a_i)^2 + \sum_{i=1}^n h_{i-1} (\alpha_{i-1}(b_i) - a_{i-1})^2 \\ & \geq \frac{1}{2} \sum_{i=0}^n h_i (\alpha_i(b_i) - \alpha_i(b_{i+1}))^2 + \frac{1}{2} \sum_{i=0}^n \tilde{h}_i (a_i - a_{i-1})^2 = \frac{1}{2} \sum_{i=0}^n h_i^3 \left(\sum_{j=0}^i \alpha_j \right)^2 + \frac{1}{2} \sum_{i=0}^n \tilde{h}_i \beta_i^2. \end{aligned}$$

227 Now, the first inequality in (4.10) is a direct consequence of (4.11) and the following in-
 228 equality

$$|\alpha|^2 = \sum_{i=0}^n \left(\sum_{j=0}^i \alpha_j - \sum_{j=0}^{i-1} \alpha_j \right)^2 \leq 2 \sum_{i=0}^n \left(\sum_{j=0}^i \alpha_j \right)^2 + 2 \sum_{i=0}^n \left(\sum_{j=0}^{i-1} \alpha_j \right)^2 \leq 4 \sum_{i=0}^n \left(\sum_{j=0}^i \alpha_j \right)^2. \quad (4.12)$$

229 To prove the validity of the second inequality in (4.10), let $\delta^{\epsilon^+}(t)$ and $\delta^{\epsilon^-}(t)$ be ap-
 230 proximations to the Dirac delta function such that

$$\delta^{\epsilon^+}(t) = \begin{cases} \frac{1}{\epsilon + \rho}, & t \in [-\rho/2, \epsilon + \rho/2], \\ 0, & \text{otherwise,} \end{cases} \quad \text{and} \quad \delta^{\epsilon^-}(t) = \begin{cases} \frac{1}{\epsilon + \rho}, & t \in [-\rho/2 - \epsilon, \rho/2], \\ 0, & \text{otherwise,} \end{cases}$$

231 for any $\epsilon > 0$ and $\rho \in (0, h_{\min})$.

232 For each $i = 1, \dots, n$, let $\epsilon_i = \tau h_i / 2$, $\delta_i^+(x) = \delta^{\epsilon_i^+}(x - b_i)$, and $\delta_i^-(x) = \delta^{\epsilon_i^-}(x - b_i)$. Set

$$\begin{aligned} \hat{\alpha}_n(x) &= \sum_{i=0}^n \alpha_i H_i(x), & \hat{\beta}_n(x) &= \sum_{i=1}^n \beta_i \delta_i(x), \\ \beta_+(x) &= \sum_{i=1}^n w(b_i) \beta_i \delta_i^+(x), & \beta_-(x) &= \sum_{i=1}^n w(b_i) \beta_i \delta_i^-(x). \end{aligned}$$

233 For all $i = 0, \dots, n$ and $j = 1, \dots, n$, it is easy to check that

$$\int_0^1 w(x) H_i(x) \delta_j(x) dx = \frac{1}{2} \left(\int_0^1 w(b_j) H_i(x) \delta_j^+(x) dx + \int_0^1 w(b_j) H_i(x) \delta_j^-(x) dx \right),$$

234 which, together with multiplying by $\alpha_i \beta_j$ and summing over i and j , implies

$$\int_0^1 w(x) \hat{\alpha}_n(x) \hat{\beta}_n(x) dx = \frac{1}{2} \left(\int_0^1 \hat{\alpha}_n(x) \beta_+(x) dx + \int_0^1 \hat{\alpha}_n(x) \beta_-(x) dx \right).$$

235 Then the quadratic form of $\mathbb{H}_\Lambda(w; \theta)$ is bounded below by

$$\begin{aligned} & \begin{pmatrix} \boldsymbol{\alpha}^T & \boldsymbol{\beta}^T \end{pmatrix} \mathbb{H}_\Lambda(w; \theta) \begin{pmatrix} \boldsymbol{\alpha} \\ \boldsymbol{\beta} \end{pmatrix} \\ &= \int_0^1 w(x) \hat{\alpha}_n^2(x) dx - 2 \int_0^1 w(x) \hat{\alpha}_n(x) \hat{\beta}_n(x) dx \\ &\geq w_0 \int_0^1 \hat{\alpha}_n^2(x) dx - \int_0^1 \hat{\alpha}_n(x) \beta_+(x) dx - \int_0^1 \hat{\alpha}_n(x) \beta_-(x) dx \\ &= \frac{1}{2} \int_0^1 \left(\sqrt{w_0} \hat{\alpha}_n(x) - \frac{1}{\sqrt{w_0}} \beta_+(x) \right)^2 dx + \frac{1}{2} \int_0^1 \left(\sqrt{w_0} \hat{\alpha}_n(x) - \frac{1}{\sqrt{w_0}} \beta_-(x) \right)^2 dx \\ &\quad - \frac{1}{2w_0} \left(\int_0^1 \beta_+^2(x) + \beta_-^2(x) dx \right). \end{aligned} \tag{4.13}$$

236 To bound the first term below, let $\beta_{n+1} = c_{n+1} = \beta_0 = 0$ and $\epsilon_0 = \frac{h_0}{2}$. It follows from the facts

237 that $\hat{\alpha}_n(x)$ and $\beta_+(x)$ are piecewise constants, $h_i - \epsilon_i = \epsilon_i + (1 - \tau)h_i$ and (4.12) that

$$\begin{aligned} & \int_0^1 \left(\sqrt{w_0} \hat{\alpha}_n(x) - \frac{1}{\sqrt{w_0}} \beta_+(x) \right)^2 dx \\ & \geq \sum_{i=0}^n \left(\int_{b_i}^{b_i + \epsilon_i + \rho/2} + \int_{b_i + \epsilon_i + \rho/2}^{b_{i+1} - \rho/2} \right) \left(\sqrt{w_0} \hat{\alpha}_n(x) - \frac{1}{\sqrt{w_0}} \beta_+(x) \right)^2 dx \\ & = \sum_{i=0}^n \left(\epsilon_i + \frac{\rho}{2} \right) \left(\sqrt{w_0} \sum_{j=0}^i \alpha_j - \frac{w(b_i) \beta_i}{(\epsilon_i + \rho) \sqrt{w_0}} \right)^2 + \sum_{i=0}^n (h_i - \epsilon_i - \rho) \left(\sqrt{w_0} \sum_{j=0}^i \alpha_j \right)^2 \\ & \geq \sum_{i=0}^n \frac{\epsilon_i}{2} \left(\frac{w(b_i) \beta_i}{(\epsilon_i + \rho) \sqrt{w_0}} \right)^2 + \sum_{i=0}^n [(1 - \tau)h_i - \rho] \left(\sqrt{w_0} \sum_{j=0}^i \alpha_j \right)^2 \\ & \geq \frac{1}{2w_0} \sum_{i=1}^n \frac{\epsilon_i w^2(b_i) \beta_i^2}{(\epsilon_i + \rho)^2} + w_0 \frac{(1 - \tau)h_{\min} - \rho}{4} |\alpha|^2. \end{aligned} \tag{4.14}$$

238 Similarly, for the second term in (4.13), following an argument analogous to the one used
239 for inequality (4.14), we obtain

$$\begin{aligned} & \int_0^1 \left(\sqrt{w_0} \hat{\alpha}_n(x) - \frac{1}{\sqrt{w_0}} \beta_-(x) \right)^2 dx \\ & \geq \sum_{i=1}^{n+1} \left(\int_{b_{i-1} + \rho/2}^{b_{i-1} - \epsilon_{i-1} - \rho/2} + \int_{b_{i-1} - \epsilon_{i-1} - \rho/2}^{b_i} \right) \left(\sqrt{w_0} \hat{\alpha}_n(x) - \frac{1}{\sqrt{w_0}} \beta_-(x) \right)^2 dx \\ & \geq \frac{1}{2w_0} \sum_{i=1}^n \frac{\epsilon_{i-1} w^2(b_i) \beta_i^2}{(\epsilon_{i-1} + \rho)^2} + w_0 \frac{(1 - \tau)h_{\min} - \rho}{4} |\alpha|^2. \end{aligned} \tag{4.15}$$

Since

$$\int_0^1 \beta_+^2(x) dx = \sum_{i=1}^n (\epsilon_i + \rho)^{-1} w^2(b_i) \beta_i^2 \quad \text{and} \quad \int_0^1 \beta_-^2(x) dx = \sum_{i=1}^n (\epsilon_{i-1} + \rho)^{-1} w^2(b_i) \beta_i^2,$$

240 after letting $\rho \rightarrow 0$, it follows from (4.13), (4.14) and (4.15) that

$$\left(\alpha^T, \beta^T \right) \mathbb{H}_\Lambda(w; \theta) \begin{pmatrix} \alpha \\ \beta \end{pmatrix} \geq \frac{w_0(1 - \tau)h_{\min}}{4} |\alpha|^2 - \frac{1}{2\tau w_0} \sum_{i=1}^n w(b_i)^2 \left(h_{i-1}^{-1} + h_i^{-1} \right) \beta_i^2,$$

241 and the lemma is proved. □

242 **Lemma 4.3.** For all $i \in \{1, \dots, n\}$, assume that $c_i \neq 0$ and that

$$F_i^2(\theta) < \begin{cases} c_i^2 \left(\frac{r_0 h_{\min}^3}{96} + \frac{\mu(1 - \tau)h_{\min}}{4} \right) \left(\frac{g_i}{c_i} + \frac{r_0 \tilde{h}_i}{24} - \frac{a^2(b_i)}{2\tau\mu} \left(h_{i-1}^{-1} + h_i^{-1} \right) \right), & \text{DR,} \\ c_i^2 \left(\frac{r_0 h_{\min}^3}{96} \right) \left(\frac{g_i}{c_i} + \frac{r_0 \tilde{h}_i}{24} \right), & \text{LS,} \end{cases} \tag{4.16}$$

243 for some $0 < \tau < 1$, then $\nabla_\theta^2 F(\theta) = \nabla_\theta^2 J(u_n(\cdot; \theta))$ is SPD.

244 *Proof.* For any $\alpha = (\alpha_0, \alpha_1, \dots, \alpha_n)^T \in \mathbb{R}^{n+1}$ and $\beta = (\beta_1, \dots, \beta_n)^T \in \mathbb{R}^n$, let $\hat{\beta} = \mathbf{D}(\hat{\mathbf{c}})\beta$. For the
 245 diffusion-reaction problem, it follows from (4.9), the fact that $(\alpha^T \mathbf{d} - \beta^T \hat{\mathbf{c}})^2 \geq 0$, (4.10), and
 246 the assumptions on the lower bounds of $a(x)$ and $r(x)$ that

$$\begin{aligned} & (\alpha^T, \beta^T) \nabla_{\theta}^2 F(\theta) \begin{pmatrix} \alpha \\ \beta \end{pmatrix} \\ &= (\alpha^T, \beta^T) \mathbb{R}(\theta) \begin{pmatrix} \alpha \\ \beta \end{pmatrix} + (\alpha^T, \hat{\beta}^T) \mathbb{H}_{\Sigma+\Lambda}(\theta) \begin{pmatrix} \alpha \\ \hat{\beta} \end{pmatrix} + \hat{\beta}^T \mathbf{D}(\mathbf{g}) \mathbf{D}^{-1}(\hat{\mathbf{c}}) \hat{\beta} + \gamma (\alpha^T \mathbf{d} - \beta^T \hat{\mathbf{c}})^2 \\ &\geq (\alpha^T, \beta^T) \mathbb{R}(\theta) \begin{pmatrix} \alpha \\ \beta \end{pmatrix} + (\alpha^T, \hat{\beta}^T) \mathbb{H}_{\Sigma+\Lambda}(\theta) \begin{pmatrix} \alpha \\ \hat{\beta} \end{pmatrix} + \hat{\beta}^T \mathbf{D}(\mathbf{g}) \mathbf{D}^{-1}(\hat{\mathbf{c}}) \hat{\beta} \\ &\geq 2 \sum_{i=1}^n \alpha_i \beta_i F_i + \frac{r_0 (h_{\min})^3}{96} |\alpha|^2 + \frac{r_0}{24} \sum_{i=1}^n \tilde{h}_i (c_i \beta_i)^2 + \frac{\mu(1-\tau)h_{\min}}{4} |\alpha|^2 \\ &\quad + \sum_{i=1}^n \left(\frac{g_i}{c_i} - \frac{a^2(b_i)}{2\mu} (h_{i-1}^{-1} + h_i^{-1}) \right) (c_i \beta_i)^2 \\ &= \sum_{i=1}^n \left[\left(\frac{r_0 h_{\min}^3}{96} + \frac{\mu(1-\tau)h_{\min}}{4} \right) \alpha_i^2 + 2\alpha_i \beta_i F_i + \left(\frac{g_i}{c_i} + \frac{r_0 \tilde{h}_i}{24} - \frac{a^2(b_i)}{2\tau\mu} (h_{i-1}^{-1} + h_i^{-1}) \right) c_i^2 \beta_i^2 \right]. \end{aligned}$$

247 In a similar fashion, we have

$$(\alpha^T, \beta^T) \nabla_{\theta}^2 F(\theta) \begin{pmatrix} \alpha \\ \beta \end{pmatrix} \geq \sum_{i=1}^n \left[\left(\frac{r_0 h_{\min}^3}{96} \right) \alpha_i^2 + 2\alpha_i \beta_i F_i + \left(\frac{g_i}{c_i} + \frac{r_0 \tilde{h}_i}{24} \right) c_i^2 \beta_i^2 \right]$$

248 for the least-squares approximation problem. Notice that $a_1 x^2 + 2a_2 xy + a_3 y^2 > 0$ for all
 249 $(x, y) \neq (0, 0)$ if $a_1 > 0$, $a_3 > 0$, and $a_2^2 < a_1 a_3$. Now, the positive definiteness of $\nabla_{\theta}^2 F(\theta)$ is a
 250 direct consequence of (4.16). \square

251 If $c_j \neq 0$, by (4.7), $F_j(\theta)$ vanishes at a critical point θ^* . Hence, $\mathcal{R}(\theta^*) = \mathbf{0}$, which gives

$$\nabla_{\theta}^2 F(\theta^*) = \begin{cases} \mathcal{D}(\hat{\mathbf{c}}^*) \mathbb{H}_{\Sigma+\Lambda}(\theta^*) \mathcal{D}(\hat{\mathbf{c}}^*) + \mathcal{D}_0(\theta^*) + \gamma \begin{pmatrix} \mathbf{d} \\ -\hat{\mathbf{c}}^* \end{pmatrix} (\mathbf{d}^T, -(\hat{\mathbf{c}}^*)^T), & \text{DR,} \\ \mathcal{D}(\hat{\mathbf{c}}^*) \mathbb{H}_{\Sigma}(r; \theta^*) \mathcal{D}(\hat{\mathbf{c}}^*) + \mathcal{D}_0(\theta^*), & \text{LS.} \end{cases}$$

252 **Theorem 4.1.** For an open set $\mathcal{O} \subseteq \mathbb{R}^{2n+1}$ satisfying the invertibility assumption, let $\theta^* = \begin{pmatrix} \mathbf{c}^* \\ \mathbf{b}^* \end{pmatrix} \in$
 253 \mathcal{O} be a minimizer of problem (4.2). For all $i \in \{1, \dots, n\}$, assume that $c_i^* \neq 0$ and that

$$\frac{g_i^*}{c_i^*} + \frac{r_0 \tilde{h}_i^*}{24} > \begin{cases} \frac{a^2(b_i^*)}{2\mu} \left(\frac{1}{h_{i-1}^*} + \frac{1}{h_i^*} \right), & \text{DR,} \\ 0, & \text{LS.} \end{cases} \quad (4.17)$$

254 Then the Hessian matrix at the critical point θ^* , $\nabla_{\theta}^2 F(\theta^*) = \nabla_{\theta}^2 J(u_n(\cdot; \theta^*))$, is SPD. Moreover,
 255 the BN method using either NL-GS or L-GS converges locally to θ^* in the norm induced by
 256 $\nabla_{\theta}^2 F(\theta^*)$.

257 *Proof.* This is a direct consequence of Lemma 4.3, (4.17), and the fact that $F_i(\theta^*) = 0$ for all
 258 $i = 1, \dots, n$. \square

259 4.3 Feasibility of BN

260 This section discusses feasibility of the BN with either NL-GS or L-GS, or equivalently,
 261 invertibility of $\mathcal{H}_{11}(\theta)$ and $\mathcal{H}_{22}(\theta)$. Invertibility of $\mathcal{H}_{11}(\theta)$ was shown in [6]; moreover,
 262 $\mathcal{H}_{11}(\theta)$ is SPD.

263 Regarding $\mathcal{H}_{22}(\theta)$, the formula in (4.6) uses derivative of the diffusion coefficient $a(x)$.
 264 As discussed in Section 5 of [5], if $a(x)$ is not differential at b_i for some $i \in \{1, \dots, n\}$, then
 265 b_i lies at the physical interface. This means that the (mesh) breakpoint b_i is already at a
 266 good location and hence should be fixed without further update.

267 Below, we discuss invertibility of $\mathcal{H}_{22}(\theta)$ at the k^{th} iteration. For simplicity of notation,
 268 the parameters are still denoted by $\hat{\mathbf{c}}$ and \mathbf{b} without the superscript (k). The assumption
 269 $c_i \neq 0$ for all $i \in \{1, \dots, n\}$ is a necessary condition for invertibility of $\mathcal{H}_{22}(\theta)$. During the
 270 iterative process, this condition is not always satisfied. In the case that $c_l = 0$ for some
 271 $l \in \{1, \dots, n\}$, then $\mathcal{H}_{22}(\theta)$ is singular. To deal with such singularity, similar to [4–6], the
 272 nonlinear parameter b_l is not updated at the current step, because the corresponding neu-
 273 ron has no contribution to the current approximation. When this happens several times
 274 to a neuron, then we can either remove or redistribute this neuron. Based on the above
 275 discussion, we introduce the following set of indices for the non-contributing neurons

$$S_1 = \{i \in \{1, \dots, n\} : |c_i| < \tau_1 \text{ or } b_i \notin I\}, \quad (4.18)$$

276 where $0 \leq \tau_1 < 1$ is a small parameter.

277 Next, at each step, to update the nonlinear parameters of neurons whose indices are
 278 not in S_1 , let us denote by $\hat{\mathcal{H}}_{22}(\theta)$ the reduced Hessian after removing neurons with
 279 indices in S_1 . For the diffusion problem, $\hat{\mathcal{H}}_{22}(\theta)$ is invertible if and only if $g_i \neq 0$ for all
 280 $i \notin S_1$. Under the assumption that $u_n(x) \approx u(x)$, i.e., $u_n(x)$ is a good approximation to
 281 $u(x)$, it was shown in [5] that

$$g_i \approx a(b_i)u''(b_i), \quad (4.19)$$

282 where g_i is defined in (4.4). As discussed in [5], $g_i \approx 0$ implies that b_i is either at a nearly
 283 optimal location or no need for update. In both cases, b_i is not updated. Hence, we
 284 introduce the following set of indices for non-updating neurons

$$S_2 = \left\{ i \in \{1, \dots, n\} \setminus S_1 : \frac{|g_i|}{a(b_i)} \leq \tau_2 \text{ or } a'(b_i) \text{ DNE} \right\}, \quad (4.20)$$

285 where $0 \leq \tau_2 < 1$ is a small parameter.

286 The above strategy does apply to the diffusion-reaction problem because (4.4) and
287 (4.1) imply the validity of (4.19):

$$g_i = r(b_i)u_n(b_i) - f(b_i) - a'(b_i)u'_n(b_i) \approx r(b_i)u(b_i) - f(b_i) - a'(b_i)u'(b_i) = a(b_i)u''(b_i).$$

288 Nevertheless, the sufficient condition on the positive definiteness of $\mathcal{H}_{22}(\theta)$ was given in
289 Lemma 5.2 of [6]:

$$\frac{g_i}{c_i} + \frac{r_0 \tilde{h}_i}{4} > 0 \quad \text{for all } i \in \{1, \dots, n\} \setminus S_1, \quad (4.21)$$

290 where $r_0 \geq 0$ for DR and $r_0 > 0$ for LS. Notice that the linear parameter $c_i = u'_n(b_i^+) - u'_n(b_i^-)$
291 means the change of the slope of $u_n(x)$ at b_i (see Section 5 of [5]). Hence, $c_i \approx u'(b_i^+) -$
292 $u'(b_i^-)$ and $u''(b_i)$ have the same sign, which implies

$$\frac{g_i}{c_i} \approx \frac{a(b_i)u''(b_i)}{c_i} = \frac{a(b_i)|u''(b_i)|}{|c_i|},$$

293 which, together with positivity of $a(b_i)$, implies (4.21) is almost always valid.

294 For the LS problem, (4.21) is sufficient but not necessary for the positive definite-
295 ness of $\mathcal{H}_{22}(\theta)$. Moreover, positive definiteness differs from invertibility. Since $g_i =$
296 $r(b_i)(u_n(b_i) - u(b_i))$ is assumed to be small, for implementation purposes, we do not up-
297 date the i^{th} neuron if $g_i/c_i < -1/\tau_3$ is negatively large, where parameter $\tau_3 \in (0, 1)$. In this
298 case, $|c_i| < \tau_3 |g_i|$ is small if g_i is small. Hence, let

$$S_2 = \left\{ i \in \{1, \dots, n\} \setminus S_1 : \frac{g_i}{c_i} < 0 \text{ and } |c_i| < \tau_3 |g_i| \right\}. \quad (4.22)$$

299 Notice that the reduced Hessian $\hat{\mathcal{H}}_{22}(\theta)$ is not guaranteed to be nonsingular after remov-
300 ing neurons with indices in $S_1 \cup S_2$. In cases where the reduced Hessian $\hat{\mathcal{H}}_{22}(\theta)$ is not
301 invertible, one might consider the Gauss-Newton matrix

$$\mathcal{G}_{22}(\theta) = \mathcal{D}(\hat{\mathbf{c}})\mathbb{H}_{\Sigma}(r; \theta)\mathcal{D}(\hat{\mathbf{c}}),$$

302 which is SPD under the assumption that $c_i \neq 0$ for all $i \in \{1, \dots, n\}$.

303 4.4 Implementation

304 In this section, we outline the implementation of the BN method using a reduced non-
305 linear system constructed based on the sets (4.18), (4.20), and (4.22). The NL-GS scheme,
306 implemented in earlier works [5,6], is used to illustrate the implementation, which is also
307 applicable to the other two schemes described in Section 2.

308 To this end, let

$$S = \{1, \dots, n\} \setminus (S_1 \cup S_2)$$

309 be the set of indices for the neurons which remain in the system, and denote by $S^c = S_1 \cup S_2$
310 the complement of S , the set of indices to be removed. Next, for a given vector $\mathbf{v} \in \mathbb{R}^n$,

311 denote by $\mathbf{v}_S \in \mathbb{R}^{n-|S^c|}$ as the vector obtained by removing entries whose indices belong
 312 to S^c , where $|S^c|$ denotes the number of indices in S^c . Similarly, for a matrix $\mathbf{B} \in \mathbb{R}^{n \times n}$,
 313 define $\mathbf{B}_S \in \mathbb{R}^{(n-|S^c|) \times (n-|S^c|)}$ as the matrix obtained by removing the rows of columns of
 314 \mathbf{B} corresponding to the indices in S^c . Then the *reduced* search direction vector for the
 315 Newton step is defined as

$$\mathbf{d}_R(\mathbf{c}, \mathbf{b}) = \mathbf{d}_R(\boldsymbol{\theta}) = - \left(\mathcal{H}_{22}(\boldsymbol{\theta})_S \right)^{-1} \left(\nabla_{\mathbf{b}} F(\boldsymbol{\theta}) \right)_S. \quad (4.23)$$

316 We are now ready to present the reduced block Newton (rBN) algorithm with the NL-
 317 GS scheme (see Algorithm 4.1 for pseudocode). Given the previous iterate $\mathbf{b}^{(k)}$. Then the
 318 current iterate $\boldsymbol{\theta}^{(k+1)} = \begin{pmatrix} \mathbf{c}^{(k+1)} \\ \mathbf{b}^{(k+1)} \end{pmatrix}$ is computed as follows:

319 (i) Compute the linear parameters

$$\mathbf{c}^{(k+1)} = \mathbf{c}^{(k)} - \mathcal{H}_{11}(\mathbf{c}^{(k)}, \mathbf{b}^{(k)})^{-1} \nabla_{\mathbf{c}} F(\mathbf{c}^{(k)}, \mathbf{b}^{(k)}).$$

320 (ii) Compute the search direction $\mathbf{p}^{(k)} = (p_1^{(k)}, \dots, p_n^{(k)})^T$ by

$$(p_i^{(k)})_{i \in S} = \mathbf{d}_R(\mathbf{c}^{(k+1)}, \mathbf{b}^{(k)}) \quad \text{and} \quad (p_i^{(k)})_{i \notin S} = \mathbf{0}, \quad (4.24)$$

321 where $\mathbf{d}_R(\mathbf{c}^{(k+1)}, \mathbf{b}^{(k)})$ is the *reduced* Newton's direction vector defined in (4.23).

322 (iii) Compute the nonlinear parameters

$$\mathbf{b}^{(k+1)} = \mathbf{b}^{(k)} + \mathbf{p}^{(k)}.$$

323 (iv) Redistribute non-contributing breakpoints $b_i^{(k+1)}$ for all $i \in S_1$ and sort $\mathbf{b}^{(k+1)}$.

324 **Remark 4.1.** The redistribution implemented for the numerical results shown in [5, 6] in
 325 step (iv) was carried out as follows: For a neuron $b_l^{(k+1)}$ satisfying $l \in S_1$, we set

$$b_l^{(k+1)} \leftarrow \frac{b_{m-1}^{(k+1)} + b_m^{(k+1)}}{2},$$

326 where $m \in \{1, \dots, n+1\}$ is an integer chosen uniformly at random.

Algorithm 4.1 Reduced block Newton (rBN) method for (4.2).

Require: Initial network parameters $\mathbf{b}^{(0)}$

Ensure: Network parameters \mathbf{c}, \mathbf{b}

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

▷ Linear parameters

$\mathbf{c}^{(k+1)} \leftarrow \mathbf{c}^{(k)} - \mathcal{H}_{11}(\mathbf{c}^{(k)}, \mathbf{b}^{(k)})^{-1} \nabla_{\mathbf{c}} F(\mathbf{c}^{(k)}, \mathbf{b}^{(k)})$

▷ nonlinear parameters

Compute the search direction $\mathbf{p}^{(k)}$ as in (4.24)

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

▷ Redistribute non-contributing neurons and sort $\mathbf{b}^{(k+1)}$

end for

327 4.4.1 Numerical experiment

328 To illustrate the performance of the reduced BN method outlined above, consider a sin-
329 gularly perturbed reaction-diffusion equation:

$$\begin{cases} -\varepsilon^2 u''(x) + u(x) = f(x), & x \in (-1, 1), \\ u(-1) = u(1) = 0. \end{cases} \quad (4.25)$$

330 For

$$\begin{aligned} f(x) = & -2 \left(\varepsilon - 4x^2 \tanh \left(\frac{1}{\varepsilon} \left(x^2 - \frac{1}{4} \right) \right) \right) \left(1 / \cosh \left(\frac{1}{\varepsilon} \left(x^2 - \frac{1}{4} \right) \right) \right)^2 \\ & + \tanh \left(\frac{1}{\varepsilon} \left(x^2 - \frac{1}{4} \right) \right) - \tanh \left(\frac{3}{4\varepsilon} \right), \end{aligned}$$

331 problem (4.25) has the following exact solution:

$$u(x) = \tanh \left(\frac{1}{\varepsilon} \left(x^2 - \frac{1}{4} \right) \right) - \tanh \left(\frac{3}{4\varepsilon} \right). \quad (4.26)$$

332 Denote by $|\cdot|_1$ the H^1 seminorm on $(-1, 1)$. For $\nu = \varepsilon^2 = 10^{-6}$, the solution of (4.26) exhibits
333 sharp interior layers. It is well-known that overshooting and oscillations occur when us-
334 ing continuous piecewise linear approximation on a uniform partition (see Fig. 1(a) for
335 $n = 16$). Using these uniform mesh points as an initial for the nonlinear parameters, 100
336 iterations of the BN method moves them efficiently toward the interior layers, and the re-
337 sulting approximation is greatly improved (see Fig. 1(b)). This example demonstrates
338 importance of non-uniform mesh in approximation and efficiency of the BN method
339 in non-convex optimization. For more numerical experiments, we refer readers to [5]
340 and [6].

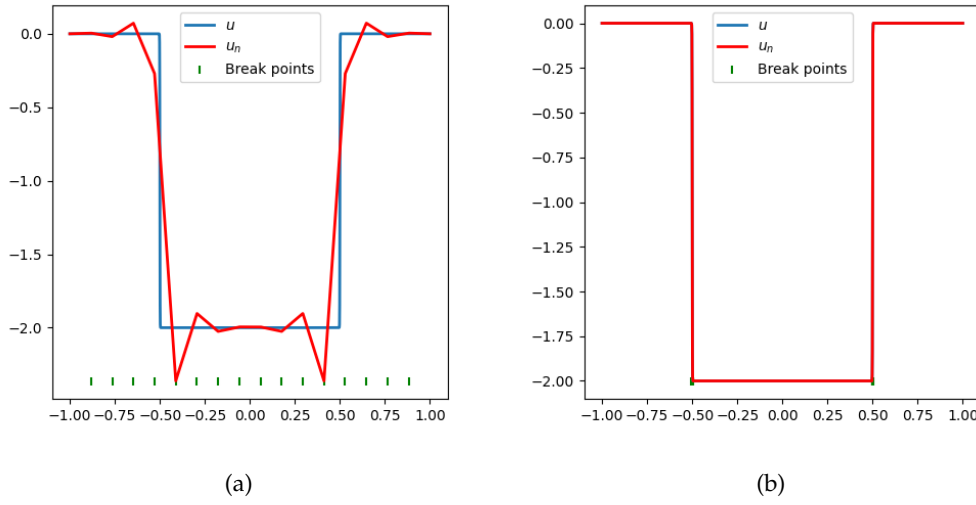


Figure 1: For $\nu = \varepsilon^2 = 10^{-6}$: (a) initial NN model with 16 uniform breakpoints; $\frac{|u - u_n|_1}{|u|_1} = 0.988$, (b) optimized NN model with 16 breakpoints, 100 iterations, $\frac{|u - u_n|_1}{|u|_1} = 0.173$.

341 **4.5 Reduced nonlinear system**

342 This section studies the rBN method introduced in Section 4.4. As explained in Section
 343 4.4, the iterative process may fix the location of certain breakpoints. We analyze the local
 344 convergence of the method under the assumption that no additional breakpoints will be
 345 fixed after a finite number of iterations.

346 To this end, for a given $\alpha \in \mathbb{R}$, assume that there are $k < n$ fixed breakpoints denoted
 347 by

$$0 = \tilde{b}_0 < \tilde{b}_1 < \dots < \tilde{b}_k < 1.$$

348 Define the set of shallow ReLU neural networks with k fixed and $n - k$ moving break-
 349 points as follows:

$$\mathcal{M}_{n,k}(I) = \left\{ \alpha + \sum_{i=0}^k c_i \sigma(x - \tilde{b}_i) + \sum_{j=1}^{n-k} c_{k+j} \sigma(x - b_j) : c_l \in \mathbb{R}, 0 < b_1 < \dots < b_{n-k} < 1 \right\}.$$

350 For any $u_{n,k}(x) \in \mathcal{M}_{n,k}(I)$, denote the corresponding parameters by

$$\theta_R = \begin{pmatrix} \mathbf{c} \\ \mathbf{b}_R \end{pmatrix} \in \mathbb{R}^{2n+1-k},$$

351 where $\mathbf{c}=(c_0,\dots,c_n)^T\in\mathbb{R}^{n+1}$ are again the linear parameters and $\mathbf{b}_R=(b_1,\dots,b_{n-k})^T\in\mathbb{R}^{n-k}$
 352 are the nonlinear parameters. After a certain number of the BN iterations, the optimiza-
 353 tion becomes to seek $u_{n,k}^*(x)=u_{n,k}(x;\boldsymbol{\theta}_R^*)\in\mathcal{M}_{n,k}(I)$ such that

$$u_{n,k}^*(x)=u_{n,k}^*(x;\boldsymbol{\theta}_R^*)=\underset{u_{n,k}\in\mathcal{M}_{n,k}(I)}{\operatorname{argmin}} J(u_{n,k}(\cdot;\boldsymbol{\theta}_R)).$$

354 The analysis presented in the preceding sections extends to the rBN method. Further-
 355 more, following a similar reasoning as in Theorem 4.1, if $c_{k+i}^*\neq 0$ and

$$\frac{g_i^*}{c_{i+k}^*}+\frac{r_0\tilde{h}_i^*}{24}>\begin{cases} \frac{a^2(b_i^*)}{2\mu}\left(\frac{1}{h_{i-1}^*}+\frac{1}{h_i^*}\right), & \text{DR,} \\ 0, & \text{LS,} \end{cases}$$

356 for all $i\in\{1,\dots,n-k\}$, then the BN method applied to this reduced parameter set con-
 357 verges locally to $\boldsymbol{\theta}_R^*=\begin{pmatrix} \mathbf{c}^* \\ \mathbf{b}_R^* \end{pmatrix}$.

358 **Remark 4.2.** In Section 3, we presented the local convergence analysis assuming the di-
 359 mensions of the linear and nonlinear parameters are $n+1$ and n , respectively. For the
 360 rBN method discussed here, the dimensions are $n+1$ and $n-k$. However, the proofs of
 361 Lemma 3.2 and the other lemmas in that section do not depend on the specific dimension
 362 of the parameter space. Therefore, the local convergence analysis extends directly to the
 363 rBN method.

364 5 Conclusions

365 A local convergence analysis for the BN methods including the reduced BN (rBN) method
 366 is presented in this paper. The rBN differs from the common optimization methods in re-
 367 ducing the number of parameters during the optimization process. We established a local
 368 convergence of the BN methods under some reasonable conditions for one-dimensional
 369 diffusion-reaction problems and least-squares function approximation. The analysis may
 370 be extended to higher dimensions for shallow neural network approximation problems.

371 Acknowledgements

372 This work was performed under the auspices of the U.S. Department of Energy by
 373 Lawrence Livermore National Laboratory under Contract DE-AC52-07NA27344. This
 374 work was supported by the NNSA Advanced Simulation and Computing (ASC) Pro-
 375 gram LLNL-JRNL-2015820.

376 **References**

- 377 [1] M. Ainsworth and Y. Shin, Plateau phenomenon in gradient descent training of ReLU net-
378 works: Explanation, quantification and avoidance, *SIAM J. Sci. Comput.*, 43 (2020), A3438–
379 A3468.
- 380 [2] M. Ainsworth and Y. Shin, Active neuron least squares: A training method for multivariate
381 rectified neural networks, *SIAM J. Sci. Comput.*, 44(4) (2022), A2253–A2275.
- 382 [3] H. G. Burchard, Splines (with optimal knots) are better, *Appl. Anal.*, 3 (1974), 309–319.
- 383 [4] Z. Cai, T. Ding, M. Liu, X. Liu, and J. Xia, A structure-guided gauss-newton method for
384 shallow ReLU neural network, arXiv:2404.05064v1 [cs.LG], 2024.
- 385 [5] Z. Cai, A. Doktorova, R. D. Falgout, and C. Herrera, Efficient shallow Ritz method for 1D
386 diffusion problems, *Comput. Math. Appl.*, 200 (2025), 349–363.
- 387 [6] Z. Cai, A. Doktorova, R. D. Falgout, and C. Herrera, Efficient shallow Ritz method for 1D
388 diffusion-reaction problems, *SIAM J. Sci. Comput.*, (2025), S414–S435.
- 389 [7] E. C. Cyr, M. A. Gulian, R. G. Patel, M. Perego, and N. A. Trask, Robust training and ini-
390 tialization of deep neural networks: An adaptive basis viewpoint, *Proceedings of Machine*
391 *Learning Research*, 107 (2020), 512–536.
- 392 [8] I. Daubechies, R. DeVore, S. Foucart, B. Hanin, and G. Petrova, Nonlinear approximation
393 and (deep) ReLU networks, *Constr. Approx.*, 55(1) (2022), 127–172.
- 394 [9] R. Falgout and P. Vassilevski, On generalizing the algebraic multigrid framework, *SIAM J.*
395 *Numer. Anal.*, 42 (2004), 1669–1693.
- 396 [10] G. H. Golub and V. Pereyra, The differentiation of pseudo-inverses and nonlinear least
397 squares problems whose variables separate, *SIAM J. Numer. Anal.*, 10(2) (1973), 413–432.
- 398 [11] M. Liu, Z. Cai, and J. Chen. Adaptive two-layer ReLU neural network: I. best least-squares
399 approximation. *Comput. Math. Appl.*, 113 (2022), 34–44.
- 400 [12] X. Liu and Y. Yuan, On the separable nonlinear least squares problems, *J. Comput. Math.*,
401 26(3) (2008), 390–403.
- 402 [13] J. M. Ortega and W. C. Rheinboldt, *Iterative Solution of Nonlinear Equations in Several Vari-*
403 *ables*, *Classics in Applied Mathematics*, Society for Industrial and Applied Mathematics,
404 1970.
- 405 [14] J. Park, J. Xu, and X. Xu, A neuron-wise subspace correction method for the finite neuron
406 method, arXiv:2211.12031 [math.NA], 2022.
- 407 [15] A. Ruhe and P. A. Wedin. Algorithms for separable nonlinear least squares problems. *SIAM*
408 *Rev.*, 22(3) (1980), 318–337.
- 409 [16] L. Schumaker, *Spline Functions: Basic Theory*, Wiley, New York, 1981.