BLOCK COORDINATE DESCENT FOR NEURAL NET WORKS PROVABLY FINDS GLOBAL MINIMA

Anonymous authors

Paper under double-blind review

ABSTRACT

In this paper, we consider a block coordinate descent (BCD) algorithm for training deep neural networks and provide a new global convergence guarantee under strictly monotonically increasing activation functions. While existing works demonstrate convergence to stationary points for BCD in neural networks, our contribution is the first to prove convergence to global minima, ensuring arbitrarily small loss. We show that the loss with respect to the output layer decreases exponentially while the loss with respect to the hidden layers remains well-controlled. Additionally, we derive generalization bounds using the Rademacher complexity framework, demonstrating that BCD not only achieves strong optimization guarantees but also provides favorable generalization performance. Moreover, we propose a modified BCD algorithm with skip connections and non-negative projection, extending our convergence guarantees to ReLU activation, which are not strictly monotonic. Empirical experiments confirm our theoretical findings, showing that the BCD algorithm achieves a small loss for strictly monotonic and ReLU activations.

024 025 026

004

010 011

012

013

014

015

016

017

018

019

021

1 INTRODUCTION

027 028

Deep learning has led to significant advances across various domains, such as computer vision, natural language processing, and reinforcement learning, achieving unprecedented performance in numerous tasks. However, understanding the training dynamics and optimization behavior of deep neural networks remains an ongoing challenge due to the highly non-convex nature of their loss functions (Li et al. 2018). Proving convergence to global minima of gradient descent via backpropagation, particularly for deep networks with multiple layers, remains an open problem in the field. While the neural tangent kernel (NTK) regime (Jacot et al. 2018) addresses this problem by reducing the non-convex loss to the convex one in RKHS, it fails to fully explain the empirical success of deep learning because it often outperforms kernel methods, even if we employ NTK as the kernel.

Contrary to the backpropagation-based training, the block coordinate descent (BCD), which originated from the mathematical optimization field (see Tseng (2001), for example), is an optimization framework where we divide a variable into several blocks and optimize them alternately. BCD of-040 fers computational advantages by updating subsets of parameters iteratively, allowing for tractable 041 optimization of complex systems. The objective function appearing in the neural network training is 042 also highly non-convex, and to overcome this issue, BCD-based neural network optimization meth-043 ods have been proposed (Carreira-Perpinan & Wang, 2014; Askari et al., 2018; Lau et al., 2018; 044 Zhang & Brand, 2017; Patel et al., 2020; Zeng et al., 2019; Nakamura et al., 2021; Qiao et al., 2021; Zhang et al., 2022; Xu et al., 2024). When we apply BCD to neural network training, the most natural way is that we regard the weights of each layer as a block, and existing works adopt this 046 way. By the formulation of BCD, the loss function of the neural network can be divided into several 047 components, one of which coincides with a loss with respect to a layer. Compared to the original 048 loss, these divided ones have more accessible landscapes to optimize.

Based on such an advantage of BCD for neural networks, its theoretical perspective, mainly about its convergence guarantee, has been explored in recent years. However, existing theoretical works on BCD for neural networks (Zhang & Brand, 2017; Zeng et al., 2019; Zhang et al., 2022; Xu
et al., 2024) have only focused on the convergence to stationary points, points with zero gradients. Convergence to stationary points does not imply convergence to global minima, especially when

the objective function is highly non-convex, such as the loss that appears in the training of neural networks (Li et al.) 2018; Safran & Shamir, 2018).

How neural network training finds global minima has been one of the most significant topics in deep learning theory literature. However, existing guarantees on BCD remain in convergence to the stationary points. To bridge this gap, we aim to provide the convergence guarantee to the global minima of BCD for neural networks. To this end, we consider multi-layer neural networks and employ a BCD-type algorithm, updating the parameters using vanilla gradient descent. Our contribution can be summarized as follows:

- We prove the global convergence of a block coordinate descent (BCD) algorithm, where we train deep neural network models with strictly monotonically increasing activation. We ensure that the parameters attain arbitrarily small loss by proving that (i) the loss with respect to the output layer will decrease exponentially to zero and (ii) the loss with respect to the hidden layers remains small in every iteration. Through the analysis, we carefully evaluate the difference propagated from the output layer to the input layer. To the best of our knowledge, this is the first result that guarantees convergence to the global minima of neural networks with any number of layers beyond the NTK regime.
- We derive a generalization error bound of deep neural networks trained by BCD under settings with i.i.d. data. In the convergence analysis, we show that the norm of weight matrices of each layer can be bounded by a constant. Combining this and the Rademacher complexity argument from Bartlett et al. (2017) gives a upper bound on generalization error. Compared to the existing works on gradient descent, BCD enables us to provide the generalization gap bound for multi-layer neural networks with an optimization guarantee.
- A notable challenge in applying our approach to commonly used activation functions like ReLU is their non-monotonic nature. Since ReLU is not strictly monotonically increasing, our initial convergence result does not directly apply. To address this issue, we propose a modified BCD algorithm incorporating skip connections (He et al., 2016) and non-negative projection updates. This modification ensures that convergence guarantees extend to ReLU networks, thereby broadening the applicability of our method to real-world architectures that predominantly use ReLU activations.
 - We validate our theoretical findings through numerical experiments, showing that BCD for both strictly monotonic and ReLU activations achieves arbitrarily small loss values. These empirical results confirm the practical viability of our proposed methods, demonstrating their effectiveness in optimizing deep neural networks beyond theoretical guarantees.

1.1 OTHER RELATED WORKS

063

064

065

067

068

069

070

084

085

090 Convergence guarantee of GD/SGD for neural networks In recent years, theoretical works on 091 the convergence guarantee of (stochastic) gradient descent for neural networks have been inten-092 sively investigated. In the neural tangent kernel (NTK) regime (Jacot et al.) 2018; Allen-Zhu et al.) 2019b; Arora et al., 2019; Du et al., 2019; Zou et al., 2020), to name a few, the training dynamics of 094 deep neural networks can be approximated by the gradient descent in RKHS. While we can ensure 095 its global convergence by exploiting the convexity, the *feature learning ability* of neural networks, 096 which is considered one of the critical ingredients of the practical success of deep learning, is not reflected since the training dynamics are reduced to the kernel method. For example, the parameters 098 of networks trained by The NTK regime hardly move from their initial points as the number of pa-099 rameters increases. On the other hand, our analysis does not fall into such a situation. Moreover, our analysis does not require any *overparameterization* on hidden layers to ensure global convergence. 100

The mean-field (MF) regime (Nitanda & Suzuki) 2017; Chizat & Bach, 2018; Mei et al., 2019; Tzen & Raginsky, 2020; Pham & Nguyen, 2021; Nguyen & Pham, 2023) is another promising approach of investigating neural network training. It regards the training of parameters as that of (probability) measure over the parameters, by which we can convert the non-convex optimization with respect to the parameters to the convex one where the distribution of parameters itself is a variable to be optimized. While several studies ensure its global convergence by employing this convexity without loss of feature learning ability, most of their models only focus on two or three-layer networks, our analysis admits the any number of hidden layers.

108 More recently, Banerjee et al. (2023) proposed restricted strong convexity (RSC) to analyze neural 109 network training, which derives the global convergence guarantee by assuming that the gradient and 110 output of neural networks correlate with each other during the training. However, Banerjee et al. 111 (2023) still requires an analysis of this correlation assumption and does not fully explain the nature 112 of global convergence in the training.

Generalization error bound of multi-layer neural networks Investigation of generalization er-114 ror analysis for multi-layer neural networks has been explored in recent years (Neyshabur et al.) 115 2015; Wei & Ma, 2019; Bartlett et al., 2017; Neyshabur et al., 2017; Golowich et al., 2018; Bartlett 116 et al., 2019; Arora et al., 2018; Suzuki et al., 2020). These works give a generalization error by 117 evaluating the complexity of neural networks from various perspective, such as the VC-dimension, the norm of parameters of networks, and so on. On the other hand, most of these results do not consider the optimization, but we also demonstrate the global convergence guarantee. Moreover, several works on generalization error analysis go beyond two-layer networks. However, most focus only on three-layer networks (Allen-Zhu & Li, 2019; Allen-Zhu et al., 2019a).

122 123 124

145

113

118

119

120

121

PRELIMINARIES 2

125 2.1 NOTATIONS 126

For an integer n, we define $[n] := \{1, \ldots, n\}$. For $x \in \mathbb{R}^d$, ||x|| denotes its Euclidean norm. We 127 denote the *d*-dimensional identity matrix by I_d . For $A \in \mathbb{R}^{n \times m}$, $||A||_F \coloneqq \sqrt{\sum_{i,j} A_{ij}^2}$ denotes its 128 129 Frobenius norm, and $||A||_{op} \coloneqq \max_{||x|| \le 1} ||Ax||$ denotes its operator norm. For two symmetric matrices 130 A and B, we denote $A \prec B$ $(A \preceq B)$ if and only if the matrix B - A is positive (non-negative) definite. For $x = (x_1, \ldots, x_d)^\top \in \mathbb{R}^d$, $\operatorname{diag}(x) \in \mathbb{R}^{d \times d}$ denotes a diagonal matrix whose *j*-th 131 132 diagonal component is x_i . 133

134 2.2 PROBLEM SETTINGS

135 Here, we introduce problem settings we consider in this paper. We observe n training examples 136 $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$, where $x_i \in \mathbb{R}^{d_{in}}$ is a feature vector and $y_i \in \mathbb{R}^{d_{out}}$ is a label. Let X =137 $(x_1 \dots x_n)^{\top} \in \mathbb{R}^{n \times d_{in}}$. Throughout the analysis, we consider high-dimensional settings $n \leq d_{in}$. 138 Moreover, we make an assumption about the matrix X as follows: 139

Assumption 1 (Data matrix is full row rank). rank(X) = n. 140

141 This assumption is required to show the global convergence. As we will see in the proof of the main 142 result, we cannot ensure the existence of global minima without Assumption 143

A multi-layer neural network is defined by 144

$$f_{NN}(x) \coloneqq W_L \sigma(W_{L-1} \sigma(\dots W_2 \sigma(W_1 x)) \dots)$$

146 where σ is element-wise activation and $W_1 \in \mathbb{R}^{r \times d_{in}}, W_j \in \mathbb{R}^{r \times r}$ for $j \in \{2, \ldots, L-1\}$, and 147 $W_L \in \mathbb{R}^{d_{out} \times r}$. We consider that all the hidden layers have the same width r. 148

Then, we make the following assumption on the activation function. 149

150 **Assumption 2** (Activation). $\sigma : \mathbb{R} \to \mathbb{R}$ is monotonically increasing and satisfies $\sigma(0) = 0$. Especially, there exists a constant $0 < \alpha < 2$ such that $\inf_{x \in \mathbb{R}} \sigma'(x) \ge \alpha$ holds Moreover, σ is ℓ -Lipschitz, i.e., for any $u_1, u_2 \in \mathbb{R}$, $|\sigma(u_1) - \sigma(u_2)| \le \ell |u_1 - u_2|$ holds. 151 152 153

154 A typical example of activation function satisfying Assumption 2 is LeakyReLU activation $x \mapsto$ $\max\{x, ax\}$ (a < 1): which satisfies Assumption 2 with $\alpha = a$ and $\ell = 1$. We note that other activation, such as ReLU $x \mapsto \max\{x, 0\}$, does not satisfy Assumption 2. We also provide the 156 global convergence algorithm when we use the ReLU activation in Section 5157

158 Under this formulation of neural networks, we formalize the regression problem 159

160
$$\min_{\mathbf{W}} \sum_{i=1}^{n} (f_{NN}(x_i) - y_i)^2,$$
(1)

¹If σ is not differentiable, we assume that $\sigma(x_1) - \sigma(x_2) \ge \alpha(x_1 - x_2)$ for any $x_1, x_2 \in \mathbb{R}$.

where $\mathbf{W} = (W_1, \dots, W_L)$. One of the most straightforward approaches to solve(1) is (stochastic) gradient method, in which the parameters are updated using the loss gradient. Conversely, we employ a layer-wise optimization method called *block coordinate descent*, as we introduce in the following section.

BLOCK COORDINATE DESCENT

In this section, after we introduce the basic notion of the *block coordinate descent* (BCD), we provide the algorithm we consider in this paper. BCD, which originated from the mathematical optimization field (see Tseng (2001), for example), is an optimization framework where we divide a variable into several blocks and optimize them alternately.

In BCD, instead of directly utilizing the loss (1), we introduce auxiliary parameters $V_{1,i} \dots V_{L,i}$. $V_{j,i}$ aims to approximate the output of *j*-th layer for the *i*-th sample x_i . By construction, we have $V_{j,i} \in \mathbb{R}^r$ for $j = 1, \dots, L-1$. By using these auxiliary parameters, we reformulate (1) as follows:

$$\min_{\mathbf{W},\mathbf{V}} F(\mathbf{W},\mathbf{V}) \coloneqq \sum_{i=1}^{n} \left[\|W_L V_{L-1,i} - y_i\|^2 + \gamma \sum_{j=1}^{L-1} \|\sigma(W_j V_{j-1,i}) - V_{j,i}\|^2 \right],$$
(2)

where $\gamma > 0$ is a hyperparameter and we denote $V_{0,i} \coloneqq x_i$, $\mathbf{W} = (W_1, \dots, W_L)$, and $\mathbf{V} = (V_{1,1}, \dots, V_{L-1,n})$. In the reformulated problem (2), the second term represents the loss at the *j*-th layer, indicating how $V_{j,i}$ approximates the output of the layer given the input x_i . The first term represents the loss at the output layer, showing how close the outputs of the network with the approximated (hidden layer) output $V_{j,i}$ are to the training labels y_1, \dots, y_n . By the construction, if ($\mathbf{W}^*, \mathbf{V}^*$) satisfies F = 0 in (2), \mathbf{W}^* is the optimal solution of (1).

One of the benefits of the reformulation (2) is that we can treat the objective function with respect to the weights of each layer (W_1, \ldots, W_L) separately. Such a simplification results not only in a faster implementation (e.g., parallelization) but also a favorable loss landscape, including theoretical tractability. While various methods for optimizing (2) have been explored, we consider a relatively simple one, updating weights W_i and auxiliary variables $V_{j,i}$ sequentially from the output layer. Specifically, we update the variables in order $W_L \to V_{L-1,i} \to W_{L-1} \to \dots V_{1,i} \to W_1$ by using the objective function (2). We summarize the algorithm considered in this paper in Algorithm 1 From now on, we explain its detailed procedure.

Algorithm 1: Block coordinate descent

input : $(W_1)_{ab} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1/d_{in}), (W_j)_{ab} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1/r)$ for all j = 2, ..., L, $V_{0,i} = x_i$. 1 K: max outer iteration, K_{V, K_W} : max inner iteration, $\eta_V, \eta_W^{(1)}, \eta_W^{(2)}$: step size; $W_j \leftarrow$ output of Algorithm 2 with inputs s_1, s_2 , and W_j for j = 2, ..., L; 3 $V_{j,i} \leftarrow \sigma(W_j V_{j-1,i})$ for all j = 1, ..., L - 1 and i = 1, ..., n.; 4 for $k \leftarrow 1$ to K do
$$\begin{split} & W_L \leftarrow W_L - \eta_W^{(1)} \nabla_{W_L} \sum_{i=1}^n \|W_L V_{L-1,i} - y_i\|^2; \\ & \text{for } i \leftarrow 1 \text{ to } n \text{ do} \end{split}$$
 $| V_{L-1,i} \leftarrow V_{L-1,i} - \eta_V \nabla_{V_{L-1,i}} || W_L V_{L-1,i} - y_i ||^2;$ for $j \leftarrow L - 1$ to 2 do $W_j \leftarrow W_j - \gamma \eta_W^{(1)} \sum_{i=1}^n \nabla_{W_j} \|\sigma(W_j V_{j,i}) - V_{j+1,i}\|^2;$ for $i \leftarrow 1$ to n do for $k_{inner} \leftarrow 1$ to K_V do $\| V_{j-1,i} \leftarrow V_{j-1,i} - \gamma \eta_V \nabla_{V_{j-1,i}} \| \sigma(W_j V_{j-1,i}) - V_{j,i} \|^2;$ for $k_{inner} \leftarrow 1$ to K_W do $W_1 \leftarrow W_1 - \gamma \eta_W^{(2)} \sum_{i=1}^n \nabla_{W_1} \| \sigma(W_1 V_{0,i}) - V_{1,i} \|^2;$

Initialization We consider Gaussian initialization for W_j ; that is, each element of W_1 is sampled from $\mathcal{N}(0, d_{in}^{-1})$, and each element of W_j (j = 2, ..., L) is sampled from $\mathcal{N}(0, r^{-1})$. After that, we apply singular value bounding (SVB) (Jia et al., 2017) to W_j (j = 2, ..., L). In SVB, we

217	Algorithm 2: Singular Value Bounding
218	input : W_i : matrix, (s_1, s_2) : lower and upper bounds on the singular values
219	1 $(U, \Sigma, V) \leftarrow$: Singular value decomposition of $W_i = U\Sigma V$;
220	2 for $s \leftarrow diagonal components of \Sigma do$
221	$s \leftarrow \max\{s_1, \min\{s_2, s\}\};$
222	

224 conduct the singular value decomposition of W_j as $W_j = U\Sigma V^{\top}$, where U and V are orthogonal 225 matrices, and Σ is a non-negative diagonal matrix. Since W_j is full-rank with probability 1 over 226 the initialization, we also have $\Sigma \in \mathbb{R}^{r \times r}$ with probability 1. After SVB, we adjust each diagonal 227 component of Σ to be within the interval $[s_1, s_2]$. Then, we utilize $W_j = U\Sigma'V^{\top}$ as the initial 228 parameter of W_j , where Σ' be the matrix obtained by the adjustment. We summarize this procedure 229 in Algorithm 2.

230 In Jia et al. (2017), SVB is conducted at every epoch to enhance the stability of the training and pre-231 diction performance of stochastic gradient descent. The upper and lower bounding of the singular 232 value prevents the amplifying or vanishing of a gradient in the backpropagation. Applying SVB also 233 has several advantages in BCD, not only for practical reasons but also from a theoretical perspec-234 tive. First, the regularity of W_j results in a preferable condition number of the objective function 235 $\|\sigma(W_j V_{j-1,i}) - V_{j,i}\|^2$ in F, the loss at the j-th layer. Moreover, the upper bound on the singular value prevents V_j from becoming extremely large at the initialization. 236

237 **Remark 3.1.** While Jia et al. (2017) applies SVB at every epoch, we use it only at the initialization. 238 By setting the step size not too large, we can ensure that all the singular values of W_i remain in 239 a bounded interval, as we show in the proof, with which we enjoy the same benefit throughout the training. 240

After initializing W_j , we initialize V_j in an exact manner, i.e., $V_{j,i} = \sigma(W_j V_{j-1,i})$ for all j =242 $1, \ldots, L-1$ and $i = 1, \ldots, n$. While we can employ any initialization scheme for V_j , the exact 243 manner results in $\|\sigma(W_i V_{i-1,i}) - V_{i,i}\|^2 = 0$ at the initialization, leading to faster convergence. 244

Update of V For optimizing W and V, we utilize vanilla gradient descent. We employ a common step size η_V for each $V_{j,i}$ and perform multiple updates using the loss $\|\sigma(W_i V_{j-1,i}) - V_{j,i}\|^2$ (line 6, 12), given by

$$V_{L-1,i} \leftarrow V_{L-1,i} - \eta_V \nabla_{V_{L-1,i}} \| W_L V_{L-1,i} - y_i \|^2$$
(3)

250 251

and

241

245

246

247

248

249

252

257

259 260 261

266 267

269

$$V_{j-1,i} \leftarrow V_{j-1,i} - \gamma \eta_V \nabla_{V_{j-1,i}} \| \sigma(W_j V_{j-1,i}) - V_{j,i} \|^2$$

The first update (3) can be interpreted as solving the linear equation $W_L V_{L-1,i} = y_i$, which has a 253 solution if the matrix W_L is full row rank. We assume that the activation satisfies Assumption 2. In 254 this case, since the mapping $\sigma: \mathbb{R} \to \mathbb{R}$ is a bijection, there exists an inverse map σ^{-1} , and training 255 V_j can be viewed as equivalent to solving the linear equation $W_{j+1}V_{j,i} = \sigma^{-1}(V_{j+1,i})$. Therefore, 256 it is expected that $V_{j,i}$ converges to the solution via gradient descent with a suitable choice of η_V as long as the matrix $W_j \in \mathbb{R}^{r \times r}$ is regular. 258

Update of W For the update of
$$W_j$$
 $(j = 1, ..., L)$, we use the loss function at *j*-th layer, that is,
 $\sum_{i=1}^{n} ||W_L V_{L-1,i} - y_i||^2$ for W_L and $\sum_{i=1}^{n} ||\sigma(W_j V_{j-1,i}) - V_{j,i}||^2$ for W_j $(j = 1, ..., L-1)$.

For W_2, \ldots, W_L , we use a common step size $\eta_W^{(1)}$ and conduct the gradient descent update:

$$W_L \leftarrow W_L - \eta_W^{(1)} \nabla_{W_L} \sum_{i=1}^n \|W_L V_{L-1,i} - y_i\|^2$$

and

$$W_{j} \leftarrow W_{j} - \gamma \eta_{W}^{(1)} \sum_{i=1}^{n} \nabla_{W_{j}} \| \sigma(W_{j}V_{j,i}) - V_{j+1,i} \|^{2}$$

for each iteration (line 5, 9). For W_1 , we employ a different step size $\eta_W^{(2)}$ and apply 268

$$W_1 \leftarrow W_1 - \gamma \eta_W^{(2)} \sum_{i=1}^n \nabla_{W_1} \| \sigma(W_1 V_{0,i}) - V_{1,i} \|^2$$

270 multiple (K_W) times for each iteration (line 14). These update manners are required to attain 271 the global convergence. With respect to the loss of the second to L-th layer, we update both 272 W_j and $V_{j-1,i}$. In particular, by applying multiple updates to $V_{j-1,i}$, we can ensure linear convergence of the loss $\sum_{i=1}^{n} \|\sigma(W_j V_{j-1,i}) - V_{j,i}\|^2$ for each iteration while the singular values of matrix W_j are upper and lower bounded. On the other hand, the existence of W^* satisfying 273 274 275 $\sum_{i=1}^{n} \|\sigma(W^*V_{j-1,i}) - V_{j,i}\|^2 = 0$ is not ensured, particularly in the case where n > r. Hence, it is not necessary to update W_j for multiple times. Furthermore, as the number of iterations increases, it 276 277 becomes less likely to maintain the regularity of the matrix W_j . This is why we only apply gradient 278 descent once to W_j (j = 2, ..., L). On the other hand, in the first layer, the input $V_{0,i} = x_i$ is fixed, and we need to demonstrate linear convergence of the loss $\sum_{i=1}^{n} \|\sigma(W_1 V_{0,i}) - V_{1,i}\|^2$ through the update of W_1 . In the overparameterized setting $d_{in} \ge n$, if the data matrix satisfies rank (U) = n, 279 280 we can ensure the existence of a global minima W^* satisfying $\sum_{i=1}^n \|\sigma(W^*V_{0,i}) - V_{1,i}\|^2 = 0$, and 281 282 hence linear convergence under a suitable choice of $\eta_{W}^{(2)}$. 283

Remark 3.2. Concerning the recent progress of the block coordinate descent algorithms applied 284 to deep learning, as represented by (Jia et al. 2017; Zhang & Brand 2017; Lau et al. 2018; 285 Patel et al., 2020), among others, we employ a relatively simple approach using vanilla gradient 286 descent without any regularization, focusing on devising the loss function and the order in which the 287 parameters are updated. While our convergence proof is based on this specific setup, our analysis 288 can be extended to encompass more complex scenarios. Our algorithm is adaptable to different 289 settings, including potential applications to other loss functions and problems, such as classification 290 problems, and the inclusion of regularization terms. We discuss possible extensions in Appendix \overline{A}

291 292

293

294

295

296

297

298 299

300

4 GLOBAL CONVERGENCE OF BLOCK COORDINATE DESCENT

In this section, we show that BCD for neural networks with an activation satisfying Assumption 2 finds global minima, in other words, the objective value F converges to an arbitrarily small value. In this section, we consider the case with single output $(d_{out} = 1)$. We discuss its extension to the multi-output case in Appendix B Moreover, for the single output case, we provide a bound on the generalization error under the i.i.d. setting by utilizing the Rademacher complexity argument.

4.1 GLOBAL CONVERGENCE WITH MONOTONICALLY INCREASING ACTIVATION

301 Here, we consider the case of single outputs $d_{out} = 1$. In this case, the objective function is 302 described by

303 304 305

306

307

311

314

$$\min_{\mathbf{W},\mathbf{V}} F(\mathbf{W},\mathbf{V}) \coloneqq \sum_{i=1}^{n} \left[\left(W_L V_{L-1,i} - y_i \right)^2 + \gamma \sum_{j=1}^{L-1} \| \sigma(W_j V_{j-1,i}) - V_{j,i} \|^2 \right].$$
(4)

We now state the first main result, the global convergence of BCD with activation satisfying Assumption 2

308 **Theorem 4.1** (BCD finds global minima of neural networks). We assume that activation σ satisfies 309 Assumption 2 and there exists a constant $C_V > 0$ such that $\lambda_{\max}(V_j V_j^{\top}) \leq C_V$ for $j = 1, \dots, L-1$ 310 during training. We denote $s \coloneqq \sigma_{\min}(X) > 0$. Let $R_i = |W_L V_{L-1,i} - y_i|$ at the initial value of the objective function with respect to the output layer, and define $R \coloneqq \sum_{i=1}^n R_i^2$, $R_{\max} \coloneqq \max_i R_i$, 312 313

and
$$C_K \coloneqq \left(\frac{2}{\alpha}\right)^L \left(4R_{\max}\eta_V + \frac{2}{2-\alpha}\sqrt{\epsilon}\right)$$

Then, under $(s_1, s_2) = (\frac{3}{4}, \frac{5}{4}), \eta_V \leq \frac{1}{8\alpha\ell^2}, \eta_W^{(1)} \leq \frac{\eta_V^{-1}}{8\sqrt{r}C_V K} \left(\frac{\alpha}{2}\right)^L, \eta_W^{(2)} \leq \frac{1}{2\ell^2 \cdot \max\|x_i\|}, and$

$$K = \left\lceil \frac{2}{\eta_V} \log\left(\frac{3R}{\epsilon}\right) \right\rceil, K_V = \left\lceil \frac{1}{\gamma \alpha \ell \eta_V} \log\left(\frac{3\gamma (L-2)rnC_K^2}{\epsilon}\right) \right\rceil, K_W = \left\lceil \frac{1}{4\gamma s \alpha^2 \eta_W^{(2)}} \log\left(\frac{3rnC_K^2}{\epsilon}\right) \right\rceil,$$

320 it holds $F(\mathbf{W}, \mathbf{V}) \leq \epsilon$, where $\mathbf{W} = (W_1, \ldots, W_L)$ and $\mathbf{V} = (V_{1,1}, \ldots, V_{L-1,n})$ are the parame-321 ters obtained by the output of Algorithm 1322

The proof can be see in Appendix C. Theorem 4.1 exhibits that BCD provably finds a global min-323 imum under a suitable choice of hyperparameters. While the definitions of K, K_V and K_W are somewhat complex, the total number of gradient computation to achieve ϵ error is bounded by $\tilde{O}(K(LK_V + K_W)) = \tilde{O}(\log^2(\frac{1}{\epsilon})).$

The proof consists of two parts: (i) the loss with respect to the output layer is monotonically decreasing in the outer loop, and (ii) the loss with respect to the hidden layer remains sufficiently small at the end of each iteration. We provide more detail to Appendix C due to page limitations.

We should note that the claims presented in Theorem 4.1 lie outside the framework of the so-called NTK regime (Jacot et al.) 2018), among others. Specifically, while the NTK regime assumes that the parameters of neural networks remain almost unchanged during training, our analysis allows for scenarios where the parameters undergo changes of $\Omega(1)$.

Remark 4.2. The assumption in Theorem 4.1, $\lambda_{\max}(V_j V_j^{\top}) \leq C_V$, ensures that the auxiliary parameters $V_{j,i}$ are bounded during training. While we assume the existence of C_V in Theorem 4.1we can provide a quantitative bound on the C_V as $C_V = O((\gamma \eta_V \ell n K K_V)^2)$ (note that this bound may not be tight). We provide a detailed derivation of this bound in Appendix D.

338 339 340

4.2 GENERALIZATION ERROR BOUND

The objective of this subsection is to show that BCD Algorithm 1 does not only have a strong convergence guarantee, but also attains favorable generalization performance. To this end, we need to make an assumption about the data distribution.

Assumption 3. The training sample $\{(x_i, y_i)\}_{i=1}^n$ is independently sampled from a distribution (x, y) ~ P. Under the distribution P, it holds that $||x|| \le B_X$ and $|y| \le B_Y$ almost surely.

The first statement defines the data distribution, which is essential and standard requirement for describing the generalization error. The one requires that inputs and labels should be bounded with probability one, which is also standard.

350 We then provide the following result on the generalization error bound.

Theorem 4.3 (Generalization error bound). Let \hat{f}_{NN} be the output of Algorithm 1 under the same condition as Theorem 4.1 Then, if Assumption 3 holds,

357 358

359

351

352

$$\mathbb{E}_{(x,y)\sim P}\left[\left(\hat{f}_{NN}(x) - y\right)^{2}\right] \leq \frac{1}{n} \sum_{i=1}^{n} \left(\hat{f}_{NN}(x_{i}) - y_{i}\right)^{2} \\
+ \tilde{O}\left(\frac{\|X\|}{n} (B_{Y} + 2^{L}\ell^{L-1}B_{X}) d_{in}^{\frac{1}{2}} L^{\frac{3}{2}} (2r)^{\frac{L}{2}} \log r + (B_{Y} + 2^{L}\ell^{L-1}B_{X})^{2} \sqrt{\frac{\log(1/\delta)}{n}}\right).$$

with probability at least $1 - \delta$ over the training sample $\{(x_i, y_i)\}_{i=1}^n$.

The proof can be seen in Appendix E. Notably, Theorem 4.3 provides a bound on the generalization error for multi-layer neural networks with optimization guarantees, beyond the NTK regime. To obtain Theorem 4.3, we utilize a result from Bartlett et al. (2017), which evaluates the generalization gap using the spectral norms of the weight matrix of each layer. As mentioned in the previous section, we can show that the spectral norm (equal to the maximum singular value) of W_j is upper bounded. Combining this with the result from Bartlett et al. (2017), we can derive the generalization gap of BCD (see Appendix E for details).

5 RELU ACTIVATION

In this section, we propose a BCD algorithm specifically for the ReLU activation $\sigma(x) := \max\{x, 0\}$, which has been excluded in Theorem 4.1 due to Assumption 2. The difficulty in handling the ReLU activation is that it only takes non-negative values. For attaining zero loss for a hidden layer $\|\sigma(W_j V_{j-1}) - V_j\|^2$, we need to prevent V_j from taking negative value due to this non-negativity. Therefore, we must exclude such situations by modifying Algorithm 1.

374 375

376

367

368

5.1 BCD FOR NEURAL NETWORKS WITH SKIP CONNECTION

As a solution to overcome the difficulty of ReLU activation, we consider ResNet (He et al.) 2016) type networks, where the neural networks includes skip connection. With skip connection, the

objective function treated in BCD is given by

$$\min_{\mathbf{W},\mathbf{V}} F(\mathbf{W},\mathbf{V}) \coloneqq \sum_{i=1}^{n} \left[(W_L V_{L-1,i} - y_i)^2 + \gamma \sum_{j=1}^{L-1} \|\sigma(W_j V_{j-1,i}) + V_{j-1,i} - V_{j,i}\|^2 \right]$$

where the loss of the hidden layer, $\gamma \sum_{j=1}^{L-1} \|\sigma(W_j V_{j-1,i}) + V_{j-1,i} - V_{j,i}\|^2$ differs from (4). We describe the modified algorithm in Algorithm 3. We use the notation $V_{j,i}^+ = \max\{V_{j,i}, 0\}$.

Algorithm 3: Block coordinate descent: ReLU

387 **Input:** $(W_1)_{ab} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1/d_{in}), (W_j)_{ab} \stackrel{i.i.d.}{\sim} \mathcal{N}(0, 1/r)$ for all $j = 2, \ldots, L, V_{0,i} = x_i$ 388 1 K: max iteration, K_{in} : max inner iteration, η_V , $\eta_W^{(1)}$, $\eta_W^{(2)}$: step size; 2 $W_j \leftarrow \text{SVB}(W_j)$ with inputs s_1, s_2 , and W_j for j = 2, ..., L-1; 389 390 $V_{j,i} = \sigma(W_j V_{j-1,i}) + V_{j-1,i}$ for all $j = 1, \dots, L-1$ and $i = 1, \dots, n$; 391 4 for $k \leftarrow 1$ to K do 392 for $i \leftarrow 1$ to n do 5 393 $V_{L-1,i} \leftarrow (V_{L-1,i} - \eta_V \nabla_{V_{L-1,i}} \| W_L V_{L-1,i} - y_i \|^2)^+;$ 6 394 $W_{L-1} \leftarrow W_{L-1} - \gamma \eta_W^{(1)} \sum_{i=1}^n \nabla_{W_{L-1}} \| \sigma(W_{L-1}V_{L-2,i}) + V_{L-2,i} - V_{L-1,i} \|^2;$ 7 for $j \leftarrow L - 1$ to 2 do 8 396 $W_{j} \leftarrow W_{j} - \gamma \eta_{W}^{(1)} \sum_{i=1}^{n} \nabla_{W_{j}} \| \sigma(W_{j} V_{j-1,i}) + V_{j-1,i} - V_{j,i} \|^{2};$ 397 9 for $i \leftarrow 1$ to n do 10 for $k_{inner} \leftarrow 1$ to K_V do 399 11 $| V_{j-1,i} \leftarrow V_{j-1,i} - \gamma \eta_V \nabla_{V_{j-1,i}} \| \sigma(W_j V_{j-1,i}) + V_{j-1,i} - V_{j,i} \|^2;$ 12 400 $V_{j-1,i} \leftarrow (V_{j-1,i})^+;$ 401 13 for $k_{inner} \leftarrow 1$ to K_W do 14 402 $W_1 \leftarrow W_1 - \gamma \eta_W^{(2)} \sum_{i=1}^n \nabla_{W_1} \| W_1 V_{0,i} - V_{1,i} \|^2;$ 403 15

404 405

407

378

379 380

382

384

386

The initialization and update of W_1, \ldots, W_{L-1} are common in Algorithm 1 and Algorithm 3 How-406 ever, there are several differences between the two algorithms in their update procedures. First, in Algorithm 3, we apply the non-negative projection $V \mapsto V^+$ for each $V_{j,i}$ after the inner loop 408 finishes. This is required for the non-negativity of ReLU: to ensure the solvability of the equation 409 $\|\sigma(W_j V_{j-1}) - V_j\|^2 = 0$. Next, we do not update W_L in Algorithm 3. This is required to ensure 410 the existence of $V_{L-1,i}$ satisfying $W_L V_{L-1,i} = y_i$ under the condition $V_{L-1,i} \ge 0$. To verify this, 411 we first provide the following lemma.

412 **Lemma 5.1.** Suppose that the vector W_L has both positive and negative entries. Then, for any y_i , 413 there exists a non-negative vector $V_{L-1,i}$ satisfying $W_L V_{L-1,i} = y_i$. 414

This lemma implies that, to ensure the global convergence for arbitrary training label y_i , it is suf-415 ficient to check that W_L has both positive and negative components. Clearly, such a situation will 416 occur frequently as the with of the hidden layer r increases. Indeed, by the symmetry of the Gaus-417 sian distribution, this probability is calculated as $1-2 \cdot \left(\frac{1}{2}\right)^r = 1-2^{-r+1}$. Additionally, we provide 418 a high probability bound on the norm of the positive and negative components of W_L , which deter-419 mines the convergence speed of the gradient descent. 420

Lemma 5.2. Let
$$W_L^{\top} \sim \mathcal{N}(0, r^{-1}I_r)$$
, $w_+ \coloneqq \max\{W_L, \mathbf{0}^{\top}\}$, and $w_- \coloneqq \min\{W_L, \mathbf{0}^{\top}\}$. Then,
for any $\delta > 0$, with probability at least $1 - 2\delta$, $\min\{\|w_+\|^2, \|w_-\|^2\} \ge \frac{1}{2} - \sqrt{\frac{8\log(2/\delta)}{r}}$ holds.

424 Since it is not trivial that the similar inequality holds for each iteration when considering the update 425 of W_L , we assume that W_L is fixed during training for simplicity. 426

Similarly to the problem (\mathbf{A}) considered in the previous section, we consider 1-dimensional outputs 427 here. We then formally state the convergence result of Algorithm 3 applied to networks with ReLU 428 activation and skip connections. 429

Theorem 5.3 (Global convergence of BCD with ReLU activation). We assume that there exists 430 a constant $C_V > 0$ such that $\lambda_{\max}(V_j V_j^{\top}) \leq C_V$ for $j = 1, \ldots, L-1$ during training. We 431 denote $s \coloneqq \sigma_{\min}(X)$. Let $R_i = |W_L V_{L-1,i} - y_i|$ at the initial value of the objective function with

$$\begin{aligned} & \text{respect to the output layer, and define } R \coloneqq \sum_{i=1}^{n} R_{i}^{2}, R_{\max} \coloneqq \max_{i} R_{i}, \text{ and } C_{K} \coloneqq (4R_{\max}\eta_{V} + 5\sqrt{\epsilon}) \left(\frac{3}{2}\right)^{L}. \text{ Then, under } (s_{1}, s_{2}) = (0, \frac{1}{4}), \eta_{V} \le \frac{1}{2\min\{\|w_{+}\|^{2}, \|w_{-}\|^{2}\}}, \eta_{W}^{(1)} \le \frac{\eta_{V}^{-1}}{24\sqrt{r}C_{V}K} \left(\frac{2}{3}\right)^{L}, \\ & \eta_{W}^{(2)} \le \frac{1}{2\cdot\max_{i}\|x_{i}\|}, \text{ and} \\ & \text{A36} \quad K = \left[\frac{1}{4\eta_{V}\min\{\|w_{+}\|^{2}, \|w_{-}\|^{2}\}}\log\left(\frac{3R}{\epsilon}\right)\right], K_{V} = \left[\frac{3}{4\gamma\eta_{V}}\log\left(\frac{49(L-2)rnC_{K}^{2}}{3\epsilon}\right)\right], K_{W} = \left[\frac{1}{4\gamma s\eta_{W}^{(2)}}\log\left(\frac{C_{K}^{2}}{\epsilon}\right)\right], \\ & \text{A40} \quad \text{it holds } E(\mathbf{W}, \mathbf{W}) \le \operatorname{conders} \mathbf{W}, \quad (\mathbf{W}, \mathbf{W}) = \operatorname{conders} \mathbf{W}, \quad (\mathbf{W}, \mathbf{W}) = \operatorname{conders} \mathbf{W}, \quad (\mathbf{W}, \mathbf{W}) = \operatorname{conders} \mathbf{W}, \\ & \text{A41} \quad \mathbf{W} = \left[\frac{1}{4\eta_{V}}\min\{\|w_{+}\|^{2}, \|w_{-}\|^{2}\}}\log\left(\frac{3R}{\epsilon}\right)\right], \\ & \text{A42} \quad \text{it holds } E(\mathbf{W}, \mathbf{W}) \le \operatorname{conders} \mathbf{W}, \quad (\mathbf{W}, \mathbf{W}) \le \operatorname{cond} \mathbf{W}, \quad (\mathbf{W}, \mathbf{W}) \le \operatorname{conders} \mathbf{W}, \\ & \text{A43} \quad \mathbf{W} = \left[\frac{1}{4\eta_{V}}\min\{\|w_{+}\|^{2}, \|w_{-}\|^{2}\}}\log\left(\frac{3R}{\epsilon}\right)\right], \\ & \text{A440} \quad \text{it holds } E(\mathbf{W}, \mathbf{W}) \le \operatorname{conders} \mathbf{W}, \quad (\mathbf{W}, \mathbf{W}) \le \operatorname{cond} \mathbf{W}, \quad (\mathbf{W}, \mathbf{W}) \le \operatorname{conders} \mathbf{W}, \\ & \text{A440} \quad \mathbf{W} = \left[\frac{1}{4\eta_{V}}\min\{\|w_{+}\|^{2}, \|w_{-}\|^{2}\}}\log\left(\frac{3R}{\epsilon}\right)\right], \\ & \text{A440} \quad \text{A440} \quad \mathbf{W} = \left[\frac{1}{4\eta_{V}}\exp\left(\frac{3}{2}\right) + \left[\frac{1}{2\eta_{V}}\exp\left(\frac{3}{2}\right)\right], \\ & \text{A440} \quad \mathbf{W} = \left[\frac{1}{4\eta_{V}}\exp\left(\frac{3}{2}\right)\right], \\ &$$

it holds $F(\mathbf{W}, \mathbf{V}) \leq \epsilon$, where $\mathbf{W} = (W_1, \ldots, W_L)$ and $\mathbf{V} = (V_{1,1}, \ldots, V_{L-1,n})$ are the parameters obtained by the output of Algorithm 3

The proof can be seen in Appendix F Thus, we obtain a global convergence guarantee of BCD for networks with ReLU activation.

NUMERICAL EXPERIMENT

In this section, we conduct numerical experiments to verify our theoretical findings. Particularly, we numerically confirm that BCD converges to a global minimum for monotonically increasing activation (Algorithm 1) and ReLU (Algorithm 3) using an artificial dataset.





Figure 1: Loss of Algorithm 1 with LeakyReLU

Figure 2: Loss of Algorithm 3 with ReLU

6.1 MONOTONICALLY INCREASING ACTIVATION

First, we conduct a numerical experiment for a monotonicall y increasing activation. We apply Algorithm 1 to a neural network with four hidden layers, each with r = 30 nodes, and LeakyReLU activation $\sigma(x) = \max\{x, 0.5x\}$, which satisfies Assumption 2 with $\alpha = 0.5$ and $\ell = 1$. We prepare n = 500 training samples from a teacher network with a single hidden layer and the same activation. We set $d_{in} = 600$, sample x_i from the normal distribution, and define y_i as the output of the teacher network. For hyperparameters, we employ $K_V = K_V = 100$ and $\eta_V = \eta_W^{(1)} = \eta_W^{(2)} = 1$.

Figure 1 shows the result. The black line means the training error, i.e., $\frac{1}{n} \sum_{i=1}^{n} (f_{NN}(x_i) - y_i)^2$. Other lines represent the loss of j-th layer, i.e, $\sum_{i=1}^{n} \|\sigma(W_j V_{j-1,i}) - V_{j,i}\|^2$ for $j \in \{1, 2, 3, 4\}$. We can observe that the training error monotonically decreases while the losses for each layer remain small, which reflects our theoretical findings.

6.2 RELU ACTIVATION

Next, we experimentally examine BCD for ReLU activation using Algorithm 3. We apply Algo-rithm 3 to a neural network with for hidden layers, r = 30, ReLU activation and skip connec-tion. Similarly to the monotinically increasing activation, we prepare a dataset with n = 500 and $d_{in} = 600$ using a teacher network. For hyperparameters, we employ $K_V = K_V = 100$ and $\eta_V = \eta_W^{(1)} = \eta_W^{(2)} = 1.$

Figure 2 shows the result. Like Figure 1, the black line means the training error. Other lines represent the loss of *j*-th layer, i.e, $\sum_{i=1}^{n} \|\sigma(W_j V_{j-1,i}) + V_{j-1,i} - V_{j,i}\|^2$ for $j \in \{1, 2, 3, 4\}$. We can observe the same convergence procedure here: the training error monotonically decreases and the losses for each layer remain small.

Additionally, we plot the training loss without using the skip connection as the dashed black line.
While the training loss for BCD without skip connections does not decrease due to the difficulty of maintaining non-negativity, the skip connection drastically improves BCD training.

7 CONCLUSION

490

491

500 501

502

503

504

531

492 In this paper, we proposed a block coordinate descent (BCD) algorithm for training deep neural 493 networks and ensured the convergence to global minima for networks with strictly monotonically 494 increasing activation functions. We also derived a generalization bound using Rademacher complexity, ensuring both strong optimization and generalization performance. For ReLU activations, we 495 introduced a modified BCD algorithm with skip connections and non-negative projection updates to 496 ensure convergence. Empirical validation demonstrated the practical effectiveness of our algorithms 497 for both monotonic and ReLU activations. Overall, this work advances the understanding of BCD 498 in neural networks, offering provable convergence and generalization guarantees. 499

References

- Zeyuan Allen-Zhu and Yuanzhi Li. What can resnet learn efficiently, going beyond kernels? Advances in Neural Information Processing Systems, 32, 2019.
- Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in overparameterized neural networks, going beyond two layers. *Advances in neural information processing systems*, 32, 2019a.
- Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over parameterization. In *International conference on machine learning*, pp. 242–252. PMLR, 2019b.
- Sanjeev Arora, Rong Ge, Behnam Neyshabur, and Yi Zhang. Stronger generalization bounds for deep nets via a compression approach. In *International conference on machine learning*, pp. 254–263. PMLR, 2018.
- Sanjeev Arora, Simon Du, Wei Hu, Zhiyuan Li, and Ruosong Wang. Fine-grained analysis of optimization and generalization for overparameterized two-layer neural networks. In *International Conference on Machine Learning*, pp. 322–332. PMLR, 2019.
- Armin Askari, Geoffrey Negiar, Rajiv Sambharya, and Laurent El Ghaoui. Lifted neural networks.
 arXiv preprint arXiv:1805.01532, May 2018.
- Arindam Banerjee, Pedro Cisneros-Velarde, Libin Zhu, and Misha Belkin. Restricted strong convexity of deep learning models with smooth activations. In *The Eleventh International Conference on Learning Representations*, 2023. URL https://openreview.net/forum?id=
 PINRbk7h01.
- Peter L Bartlett, Dylan J Foster, and Matus J Telgarsky. Spectrally-normalized margin bounds for neural networks. *Advances in neural information processing systems*, 30, 2017.
- Peter L Bartlett, Nick Harvey, Christopher Liaw, and Abbas Mehrabian. Nearly-tight vc-dimension and pseudodimension bounds for piecewise linear neural networks. *Journal of Machine Learning Research*, 20(63):1–17, 2019.
- Miguel Carreira-Perpinan and Weiran Wang. Distributed optimization of deeply nested systems. In
 Artificial Intelligence and Statistics, pp. 10–19. PMLR, 2014.
- Lenaic Chizat and Francis Bach. On the global convergence of gradient descent for overparameterized models using optimal transport. Advances in neural information processing systems, 31, 2018.
- Simon Du, Jason Lee, Haochuan Li, Liwei Wang, and Xiyu Zhai. Gradient descent finds global
 minima of deep neural networks. In *International conference on machine learning*, pp. 1675–1685. PMLR, 2019.

540 541 542	Spencer Frei, Yuan Cao, and Quanquan Gu. Agnostic learning of a single neuron with gradient descent. Advances in Neural Information Processing Systems, 33:5417–5428, 2020.
543 544	Noah Golowich, Alexander Rakhlin, and Ohad Shamir. Size-independent sample complexity of neural networks. In <i>Conference On Learning Theory</i> , pp. 297–299. PMLR, 2018.
545 546 547	Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog- nition. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 770–778, 2016.
549 550	Arthur Jacot, Franck Gabriel, and Clément Hongler. Neural tangent kernel: Convergence and gen- eralization in neural networks. <i>Advances in neural information processing systems</i> , 31, 2018.
551 552 553	Kui Jia, Dacheng Tao, Shenghua Gao, and Xiangmin Xu. Improving training of deep neural net- works via singular value bounding. In <i>Proceedings of the IEEE Conference on Computer Vision</i> <i>and Pattern Recognition</i> , pp. 4344–4352, 2017.
554 555 556	Tim Tsz-Kit Lau, Jinshan Zeng, Baoyuan Wu, and Yuan Yao. A proximal block coordinate descent algorithm for deep neural network training. <i>arXiv preprint arXiv:1803.09082</i> , 2018.
557 558	Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. Visualizing the loss land- scape of neural nets. <i>Advances in neural information processing systems</i> , 31, 2018.
560 561 562	Song Mei, Theodor Misiakiewicz, and Andrea Montanari. Mean-field theory of two-layers neural networks: dimension-free bounds and kernel limit. In <i>Conference on learning theory</i> , pp. 2388–2464. PMLR, 2019.
563 564	Mehryar Mohri, Afshin Rostamizadeh, and Ameet Talwalkar. <i>Foundations of machine learning</i> . MIT press, 2018.
565 566 567	Kensuke Nakamura, Stefano Soatto, and Byung-Woo Hong. Block-cyclic stochastic coordinate descent for deep neural networks. <i>Neural Networks</i> , 139:348–357, 2021.
568 569	Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. Norm-based capacity control in neural networks. In <i>Conference on learning theory</i> , pp. 1376–1401. PMLR, 2015.
570 571 572 573	Behnam Neyshabur, Srinadh Bhojanapalli, and Nathan Srebro. A pac-bayesian approach to spectrally-normalized margin bounds for neural networks. <i>arXiv preprint arXiv:1707.09564</i> , 2017.
574 575	Phan-Minh Nguyen and Huy Tuan Pham. A rigorous framework for the mean field limit of multi- layer neural networks. <i>Mathematical Statistics and Learning</i> , 6(3):201–357, 2023.
576 577 578	Atsushi Nitanda and Taiji Suzuki. Stochastic particle gradient descent for infinite ensembles. <i>arXiv</i> preprint arXiv:1712.05438, 2017.
579 580 581	Ravi G Patel, Nathaniel A Trask, Mamikon A Gulian, and Eric C Cyr. A block coordinate descent optimizer for classification problems exploiting convexity. <i>arXiv preprint arXiv:2006.10123</i> , 2020.
582 583 584	Huy Tuan Pham and Phan-Minh Nguyen. Global convergence of three-layer neural networks in the mean field regime. <i>arXiv preprint arXiv:2105.05228</i> , 2021.
585 586 587	Linbo Qiao, Tao Sun, Hengyue Pan, and Dongsheng Li. Inertial proximal deep learning alternating minimization for efficient neutral network training. In <i>ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)</i> , pp. 3895–3899. IEEE, 2021.
588 589 590	Itay Safran and Ohad Shamir. Spurious local minima are common in two-layer relu neural networks. In <i>International conference on machine learning</i> , pp. 4433–4441. PMLR, 2018.
591 592 593	Taiji Suzuki, Hiroshi Abe, and Tomoaki Nishimura. Compression based bound for non-compressed network: unified generalization error analysis of large compressible deep neural network. In <i>International Conference on Learning Representations</i> , 2020. URL https://openreview.net/forum?id=ByeGzlrKwH .

594 595 596	Paul Tseng. Convergence of a block coordinate descent method for nondifferentiable minimization. <i>Journal of optimization theory and applications</i> , 109:475–494, 2001.
597 598 599	Belinda Tzen and Maxim Raginsky. A mean-field theory of lazy training in two-layer neural nets: entropic regularization and controlled mckean-vlasov dynamics. <i>arXiv preprint arXiv:2002.01987</i> , 2020.
600 601	Martin J Wainwright. <i>High-dimensional statistics: A non-asymptotic viewpoint</i> , volume 48. Cambridge university press, 2019.
603 604	Colin Wei and Tengyu Ma. Data-dependent sample complexity of deep neural networks via lipschitz augmentation. <i>Advances in Neural Information Processing Systems</i> , 32, 2019.
605 606	Jintao Xu, Chenglong Bao, and Wenxun Xing. Convergence rates of training deep neural networks via alternating minimization methods. <i>Optimization Letters</i> , 18(4):909–923, 2024.
608 609	Tian Ye and Simon S Du. Global convergence of gradient descent for asymmetric low-rank matrix factorization. <i>Advances in Neural Information Processing Systems</i> , 34:1429–1439, 2021.
610 611 612	Gilad Yehudai and Shamir Ohad. Learning a single neuron with gradient methods. In Jacob Aber- nethy and Shivani Agarwal (eds.), <i>Proceedings of Thirty Third Conference on Learning Theory</i> , volume 125 of <i>Proceedings of Machine Learning Research</i> , pp. 3756–3786. PMLR, 2020.
613 614 615 616	Jinshan Zeng, Tim Tsz-Kit Lau, Shaobo Lin, and Yuan Yao. Global convergence of block coordi- nate descent in deep learning. In <i>International conference on machine learning</i> , pp. 7313–7323. PMLR, 2019.
617 618	Hui Zhang, Shenglong Zhou, Geoffrey Ye Li, and Naihua Xiu. 0/1 deep neural networks via block coordinate descent. <i>arXiv preprint arXiv:2206.09379</i> , 2022.
619 620 621	Ziming Zhang and Matthew Brand. Convergent block coordinate descent for training tikhonov regularized deep neural networks. <i>Advances in Neural Information Processing Systems</i> , 30, 2017.
622 623	Difan Zou, Yuan Cao, Dongruo Zhou, and Quanquan Gu. Gradient descent optimizes over- parameterized deep relu networks. <i>Machine learning</i> , 109:467–492, 2020.
624 625	
626	
627	
628	
629	
630	
631	
632	
633	
634	
635	
636	
637	
638	
639	
640	
641	
642	
643	
644	
645	
646	
647	