

Introduction

- Training deep neural networks (DNNs) efficiently is challenging due to the associated **highly nonconvex** optimization
- The backpropagation (**backprop**) algorithm has long been the most widely used algorithm for gradient computation of parameters of DNNs
- Backprop suffers from various issues, e.g., vanishing gradients
- Various methods to alleviate this issue, e.g., ReLUs and LSTM, but unable to completely tackle this inherent problem to backprop
- One viable alternative is to adopt gradient-free methods, including (but not limited to) the alternating direction method of multipliers (**ADMM**) and the block coordinate descent (**BCD**)
- Mainly to decompose the highly coupled and composite DNN training objective into several *loosely coupled* and *almost separable* simple subproblems

Problem Formulation

Notations:

- A feedforward neural network with L hidden layers
- d_ℓ : the number of nodes of the ℓ -th layer
- N : the number of training samples
- K : the number of classes, i.e., $d_{L+1} = K$
- $\mathbf{x}_j \in \mathbb{R}^{d_0}$: the j -th training data, $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_N) \in \mathbb{R}^{d_0 \times N}$
- $\mathbf{y}_j \in \mathbb{R}^{d_{L+1}}$: the one-hot vector of its corresponding label, \mathbf{y}_j : the i -th entry of the column vector \mathbf{y}_j , $\mathbf{Y} = (\mathbf{y}_1, \dots, \mathbf{y}_N) \in \mathbb{R}^{d_{L+1} \times N}$
- $\mathbf{W}_\ell \in \mathbb{R}^{d_\ell \times d_{\ell-1}}$: the weight matrix between the ℓ -th and $(\ell-1)$ -th hidden layers
- $\mathbf{b}_\ell \in \mathbb{R}^{d_\ell}$: the bias vector of the ℓ -th hidden layer
- $\mathbf{W}_{L+1} \in \mathbb{R}^{d_{L+1} \times d_L}$: the weight matrix between the last hidden layer and the output layer
- $\mathbf{b}_{L+1} \in \mathbb{R}^{d_{L+1}}$: the bias vector of the output layer
- h : general activation function
- $\mathcal{W} := \{\mathbf{W}_\ell\}_{\ell=1}^{L+1}$ and $\mathbf{b} := \{\mathbf{b}_\ell\}_{\ell=1}^{L+1}$
- Optimization problem of regularized DNNs:

$$\min_{\mathbf{a}, \mathbf{z}, \mathcal{W}, \mathbf{b}} F(\mathbf{a}, \mathcal{W}, \mathbf{b}) \equiv \gamma_{L+1} \sum_{j=1}^N \mathcal{L}(\mathbf{W}_{L+1} \mathbf{a}_{L,j} + \mathbf{b}_{L+1}; \mathbf{y}_j) + \sum_{\ell=1}^{L+1} r_\ell(\mathbf{W}_\ell)$$
 subject to $\mathbf{z}_{\ell,j} = \mathbf{W}_\ell \mathbf{a}_{\ell-1,j} + \mathbf{b}_\ell$, $\mathbf{a}_{\ell,j} = h(\mathbf{z}_{\ell,j})$ for all $\ell \in \{1, \dots, L\}, j \in \{1, \dots, N\}$
 - \mathcal{L} : generic loss function
 - r_ℓ : convex but possibly nonsmooth regularizers
 - $\mathbf{a} := \{\mathbf{a}_{\ell,j}\}_{\ell=1}^L$: the set of all activation vectors
 - $\mathbf{z} := \{\mathbf{z}_{\ell,j}\}_{\ell=1}^L$, $\mathbf{a}_{0,j} := \mathbf{x}_j$ for all $j \in \{1, \dots, N\}$, and $\gamma_{L+1} > 0$

The Proximal BCD Algorithm

- **Input:** training data $\{(\mathbf{x}_j, \mathbf{y}_j)\}_{j=1}^N$ and regularization parameters $\{\gamma_\ell\}_{\ell=1}^{L+1}$, $\mathbf{a}_{L+1,j}^{(k)} := \mathbf{y}_j$ for all $k \in \mathbb{N}$ and $j \in \{1, \dots, N\}$
- Initialize $\mathcal{W}^{(0)}$, $\mathbf{b}^{(0)}$, $\alpha_i^{(0)} \in (0, \infty)$, $\beta_i^{(0)} \in (0, 1)$ for all i , and $s, t \in (0, 1)$
- $\mathbf{a}^{(0)}$ initialized by *forward propagation*
- Iterative update each block using a *proximal step*
- *Adaptive momentums* used in implementations
- Update of the parameters of the last layer:

$$\mathbf{W}_{L+1}^* = \operatorname{argmin}_{\mathbf{W}_{L+1}} \gamma_{L+1} \mathcal{L}(\mathbf{W}_{L+1} \mathbf{a}_{L,j}^{(k-1)} + \mathbf{b}_{L+1}^{(k-1)}; \mathbf{y}_j) + \frac{\alpha_1^{(k-1)}}{2} \|\mathbf{W}_{L+1} - \mathbf{W}_{L+1}^{(k-1)}\|_F^2 + r_{L+1}(\mathbf{W}_{L+1})$$

$$\mathbf{W}_{L+1}^{(k)} = \mathbf{W}_{L+1}^{(k-1)} + \beta_1^{(k-1)} (\mathbf{W}_{L+1}^* - \mathbf{W}_{L+1}^{(k-1)})$$

$$\mathbf{b}_{L+1}^* = \operatorname{argmin}_{\mathbf{b}_{L+1}} \gamma_{L+1} \mathcal{L}(\mathbf{W}_{L+1}^{(k-1)} \mathbf{a}_{L,j}^{(k-1)} + \mathbf{b}_{L+1}; \mathbf{y}_j) + \frac{\alpha_1^{(k-1)}}{2} \|\mathbf{b}_{L+1} - \mathbf{b}_{L+1}^{(k-1)}\|^2$$

$$\mathbf{b}_{L+1}^{(k)} = \mathbf{b}_{L+1}^{(k-1)} + \beta_1^{(k-1)} (\mathbf{b}_{L+1}^* - \mathbf{b}_{L+1}^{(k-1)})$$
- Update of the hidden layers, $\ell = L, \dots, 1$:

1. Activation vectors:

$$\mathbf{a}_{\ell,j}^* = \operatorname{argmin}_{\mathbf{a}_{\ell,j}} \frac{\gamma_{\ell+1}}{2} \|\mathbf{W}_{\ell+1}^{(k)} \mathbf{a}_{\ell,j} + \mathbf{b}_{\ell+1}^{(k)} - \mathbf{a}_{\ell+1,j}^{(k)} + \mathbf{u}_{\ell+1,j}^{(k)}\|^2 + \frac{\gamma_\ell}{2} \|\mathbf{W}_\ell^{(k-1)} \mathbf{a}_{\ell-1,j}^{(k-1)} + \mathbf{b}_\ell^{(k-1)} - \mathbf{a}_{\ell,j} + \mathbf{u}_{\ell,j}^{(k)}\|^2 + \frac{\alpha_{2(L-\ell+1)}^{(k-1)}}{2} \|\mathbf{a}_{\ell,j} - \mathbf{a}_{\ell,j}^{(k-1)}\|^2 + \iota_{\mathcal{S}_\ell}(\mathbf{a}_{\ell,j}) \text{ for all } j \in \{1, \dots, N\}$$

$$\mathbf{a}_{\ell,j}^{(k)} = \mathbf{a}_{\ell,j}^{(k-1)} + \beta_{2(L-\ell+1)}^{(k-1)} (\mathbf{a}_{\ell,j}^* - \mathbf{a}_{\ell,j}^{(k-1)}) \text{ for all } j \in \{1, \dots, N\}$$

2. Auxiliary variables:

$$\mathbf{u}_{\ell,j}^{(k)} = \operatorname{argmin}_{\mathbf{u}_{\ell,j}} \frac{\gamma_\ell}{2} \|\mathbf{W}_\ell^{(k-1)} \mathbf{a}_{\ell-1,j}^{(k-1)} + \mathbf{b}_\ell^{(k-1)} - \mathbf{a}_{\ell,j}^{(k)} + \mathbf{u}_{\ell,j}^{(k)}\|^2$$

3. Parameters:

$$\mathbf{W}_\ell^* = \operatorname{argmin}_{\mathbf{W}_\ell} \sum_{j=1}^N \frac{\gamma_\ell}{2} \|\mathbf{W}_\ell \mathbf{a}_{\ell-1,j}^{(k-1)} + \mathbf{b}_\ell^{(k-1)} - \mathbf{a}_{\ell,j}^{(k)} + \mathbf{u}_{\ell,j}^{(k)}\|^2 + \frac{\alpha_{2(L-\ell+1)}^{(k-1)}}{2} \|\mathbf{W}_\ell - \mathbf{W}_\ell^{(k-1)}\|_F^2 + r_\ell(\mathbf{W}_\ell)$$

$$\mathbf{W}_\ell^{(k)} = \mathbf{W}_\ell^{(k-1)} + \beta_{2(L-\ell+1)}^{(k-1)} (\mathbf{W}_\ell^* - \mathbf{W}_\ell^{(k-1)})$$

$$\mathbf{b}_\ell^* = \operatorname{argmin}_{\mathbf{b}_\ell} \sum_{j=1}^N \frac{\gamma_\ell}{2} \|\mathbf{W}_\ell^{(k-1)} \mathbf{a}_{\ell-1,j}^{(k-1)} + \mathbf{b}_\ell - \mathbf{a}_{\ell,j}^{(k)} + \mathbf{u}_{\ell,j}^{(k)}\|^2 + \frac{\alpha_{2(L-\ell+1)}^{(k-1)}}{2} \|\mathbf{b}_\ell - \mathbf{b}_\ell^{(k-1)}\|^2$$

$$\mathbf{b}_\ell^{(k)} = \mathbf{b}_\ell^{(k-1)} + \beta_{2(L-\ell+1)}^{(k-1)} (\mathbf{b}_\ell^* - \mathbf{b}_\ell^{(k-1)})$$

Convergence Results

Theorem 1 (Global convergence) Under certain assumptions and the fact that the sequence generated by the proposed algorithm has a finite limit point where \tilde{F} satisfies the **Kurdyka-Łojasiewicz (KL) property**, this sequence converges to such a limit point, which is a critical point of (3).

Contributions

- Propose a novel algorithm based on **BCD of Gauss-Seidel type**
- Define the loss function using the quadratic penalty method by unrolling the nested structure of DNNs into separate “single-layer” training tasks
- Simplifications of commonly used activation functions as **projections** onto *nonempty closed convex sets*
- The overall loss function is **block multiconvex**
- Obtain global convergence guarantees under the framework of the **Kurdyka-Łojasiewicz (KL) property**

Reformulation (quadratic penalty method):

$$\min_{\mathbf{a}, \mathbf{z}, \mathcal{W}, \mathbf{b}} F(\mathbf{a}, \mathcal{W}, \mathbf{b}) + \sum_{j=1}^N \sum_{\ell=1}^L \frac{\rho_\ell}{2} \|h(\mathbf{z}_{\ell,j}) - \mathbf{a}_{\ell,j}\|^2 + \sum_{j=1}^N \sum_{\ell=1}^L \frac{\gamma_\ell}{2} \|\mathbf{W}_\ell \mathbf{a}_{\ell-1,j} + \mathbf{b}_\ell - \mathbf{z}_{\ell,j}\|^2 \quad (2)$$

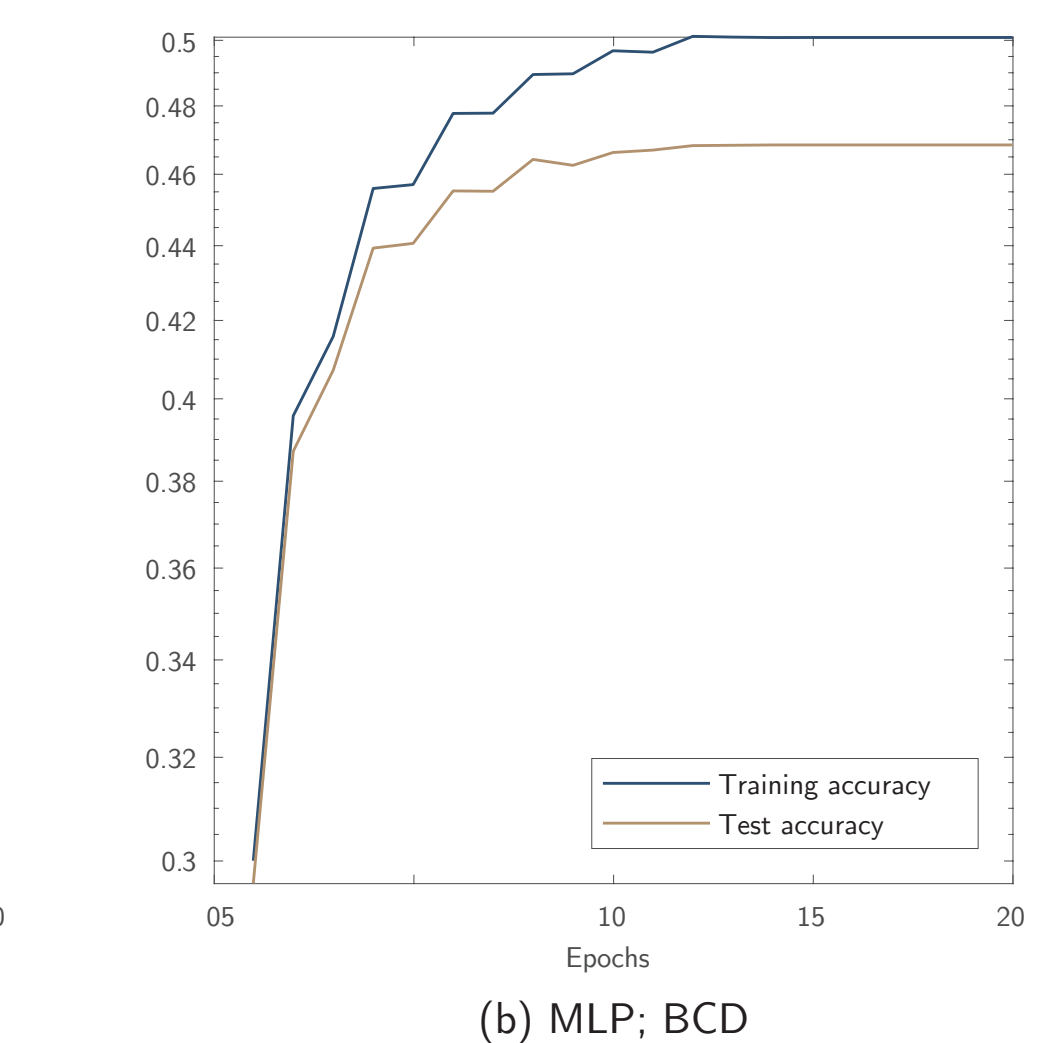
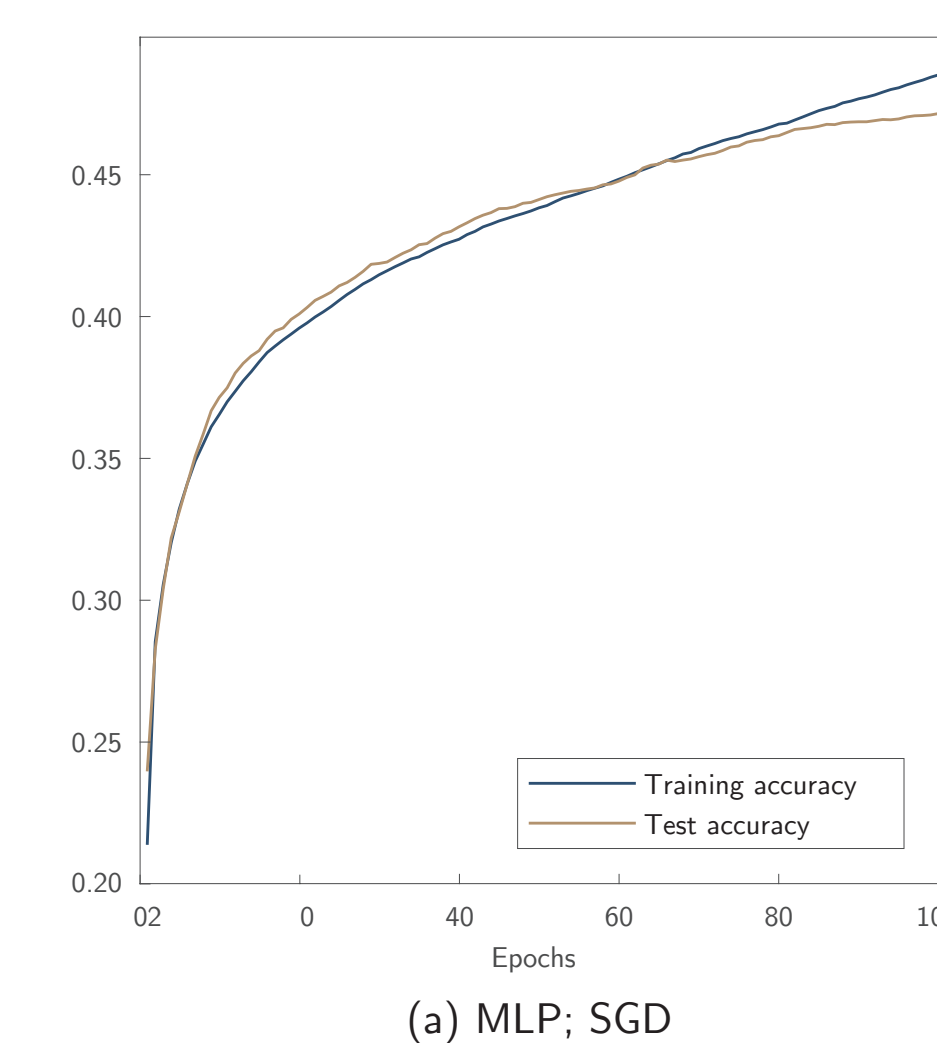
- $\|\cdot\|$ is the standard Euclidean norm, $\rho_\ell > 0$ and $\gamma_\ell > 0$ for all $\ell \in \{1, \dots, L\}$
- Allows for any general activation functions such as hyperbolic tangent, sigmoid and ReLU
- The formulation (2) is generally hard to solve explicitly
- Can be simplified if we consider the constraint $\mathbf{a}_{\ell,j} = h(\mathbf{z}_{\ell,j})$ as a **projection** onto a *convex set*
- E.g., ReLU can be thought of as a projection onto the closed upper half-space
- Introduction of a set of auxiliary variable $\{\mathbf{u}_{\ell,j}\}_{\ell=1}^L$:

$$\min_{\mathbf{a}, \mathcal{W}, \mathbf{b}, \mathbf{u}} \tilde{F}(\mathbf{a}, \mathcal{W}, \mathbf{b}, \mathbf{u}) \equiv F(\mathbf{a}, \mathcal{W}, \mathbf{b}) + \sum_{j=1}^N \sum_{\ell=1}^L \left[\frac{\gamma_\ell}{2} \|\mathbf{W}_\ell \mathbf{a}_{\ell-1,j} + \mathbf{b}_\ell - \mathbf{a}_{\ell,j} + \mathbf{u}_{\ell,j}\|^2 + \iota_{\mathcal{S}_\ell}(\mathbf{a}_{\ell,j}) \right] \quad (3)$$
 - \mathcal{S}_ℓ : a nonempty closed convex set
 - $\iota_{\mathcal{S}_\ell}$: the indicator function of a nonempty closed convex set \mathcal{S}_ℓ such that

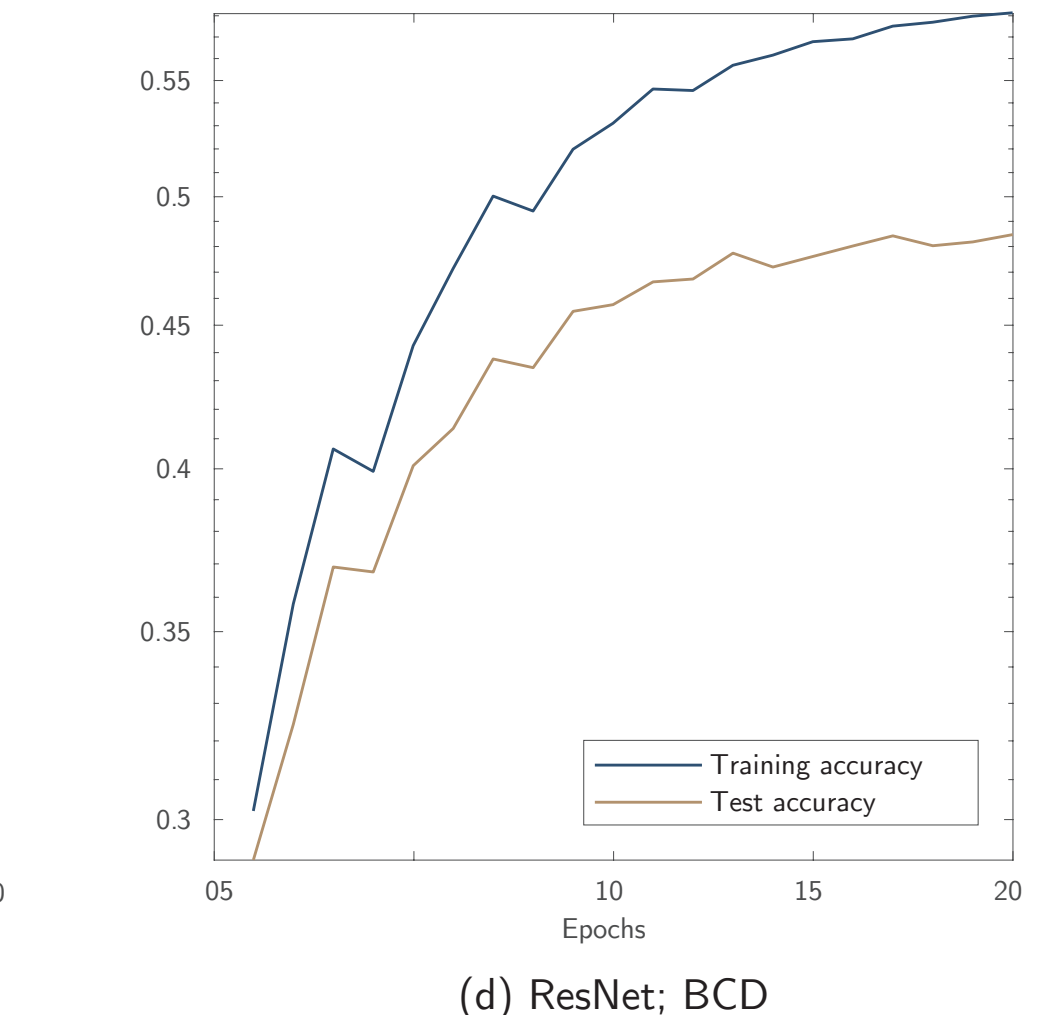
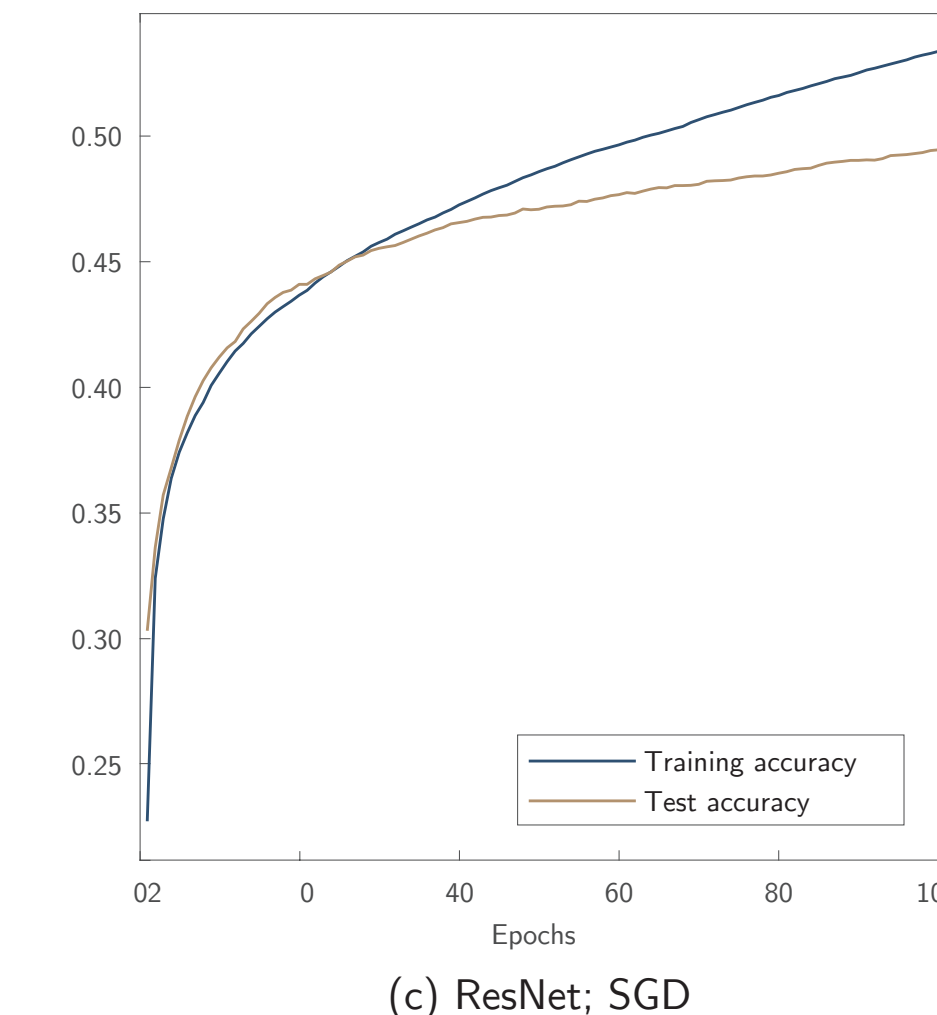
$$\iota_{\mathcal{S}_\ell}(\mathbf{u}) = \begin{cases} 0 & \text{for } \mathbf{u} \in \mathcal{S}_\ell \\ +\infty & \text{otherwise} \end{cases}$$
 - For the case of the ReLU activation function, $\mathcal{S}_\ell = \mathbb{R}_+^{d_\ell}$
- The objective function (3) is **block multiconvex** which allows for established convergence guarantees in existing literature using the proposed algorithm

Experimental Results

- CIFAR-10 with 50,000 training and 10,000 test samples
- 3,072-4,000-4,000-10 MLP
- 3,072-4,000-3,072-4,000-10 DNN with a residual connection in the second hidden layer (ResNet)



- The proximal BCD algorithm (20 epochs) and Backprop (SGD; 100 epochs)
- Squared losses; ReLU activations; no regularizations
- Final test accuracies (medians of 5 runs):



1. MLP with SGD: 0.4765
2. MLP with BCD: 0.4682
3. ResNet with SGD: 0.494
4. ResNet with BCD: 0.4843

Conclusion

- Proposed an efficient BCD algorithm and established its convergence guarantees according to our block multiconvex formulation
- Three major advantages of BCD include:
 1. **High per epoch efficiency at early stages** (observed in the figures), i.e., the training and test accuracies of BCD grow much faster than SGD in terms of epoch at the early stage
 2. **Good scalability**, i.e., BCD can be implemented in a distributed and parallel manner via data parallelism on multi-core CPUs
 3. **Gradient free**, i.e., gradient computations are unnecessary used for the updates
- They generally require more memory than SGD method
- A future direction is to study *the feasibility of the stochastic and parallel block coordinate descent methods*