



# Global Convergence of Block Coordinate Descent in Deep Learning

Jinshan Zeng<sup>1,2,\*</sup>, Tim Tsz-Kit Lau<sup>3,\*</sup>, Shao-Bo Lin<sup>4</sup>, Yuan Yao<sup>2</sup>

<sup>1</sup> Jiangxi Normal Univ.

<sup>2</sup> HKUST

<sup>3</sup> Northwestern

<sup>4</sup> CityU HK

\* Equal contribution



Provide a general methodology to establish the global convergence of BCD methods in training deep neural networks to a critical point at a rate of  $\mathcal{O}(1/k)$ , where  $k$  is the number of iterations, without the *block multiconvexity* and *differentiability* assumptions

## INTRODUCTION

### Motivation of Block Coordinate Descent (BCD) in Deep Learning

- Gradient-based methods are commonly used in training deep neural networks (DNNs), but may suffer from various problems, especially when neural networks become deeper (i.e., more layers): Since the gradients of the loss function w.r.t. parameters of earlier layers involve those of later layers

- if one of the later layer gradients is zero, all gradients of previous layers are zero (known as **gradient vanishing**);
- if one of the later layer gradients is  $\pm\infty$  (in terms of computer precision), all gradients of previous layers are  $\pm\infty$  (known as **gradient exploding**)

- Gradient-free methods have recently been adapted to training DNNs:

- Block Coordinate Descent (BCD)
- Alternating Direction Method of Multipliers (ADMM)

- Advantages of Gradient-free Methods:

- Ability to deal with non-differentiable nonlinearities and potentially avoid **vanishing gradient**
- Can be easily implemented in a *distributed* and *parallel* fashion

## DNN TRAINING VIA BLOCK COORDINATE DESCENT

### Variable Splitting

- View parameters of hidden layers and the output layer as variable blocks

- Variable splitting**: Split the highly coupled network layer-wise to compose a surrogate loss function

- Notations:
  - $\mathcal{W} := \{\mathbf{W}_\ell\}_{\ell=1}^L$ : the set of layer parameters (bias vectors are absorbed)
  - $\mathcal{L} : \mathbb{R}^k \times \mathbb{R}^k \rightarrow \mathbb{R}_+ \cup \{0\}$ : some loss function
  - $\Phi(\mathbf{x}_i; \mathcal{W}) := \mathcal{L}_L(\mathbf{W}_L \sigma_{L-1}(\mathbf{W}_{L-1} \cdots \mathbf{W}_2 \sigma_1(\mathbf{W}_1 \mathbf{x}_i))$ : the neural network

- Empirical risk minimization**:

$$\min_{\mathcal{W}} \mathcal{R}_n(\Phi(\mathbf{X}; \mathcal{W}), \mathbf{Y}) := \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\Phi(\mathbf{x}_i; \mathcal{W}), \mathbf{y}_i)$$

### Two-Splitting Formulation

- Introduce one set of auxiliary variables  $\mathcal{V} := \{\mathbf{V}_\ell\}_{\ell=1}^L$ 

$$\min_{\mathcal{W}, \mathcal{V}} \mathcal{L}_0(\mathcal{W}, \mathcal{V}) := \mathcal{R}_n(\mathbf{V}_L; \mathbf{Y}) + \sum_{\ell=1}^L r_\ell(\mathbf{W}_\ell) + \sum_{\ell=1}^L s_\ell(\mathbf{V}_\ell)$$

subject to  $\mathbf{V}_\ell = \sigma_\ell(\mathbf{W}_\ell \mathbf{V}_{\ell-1}), \ell \in \{1, \dots, L\}$

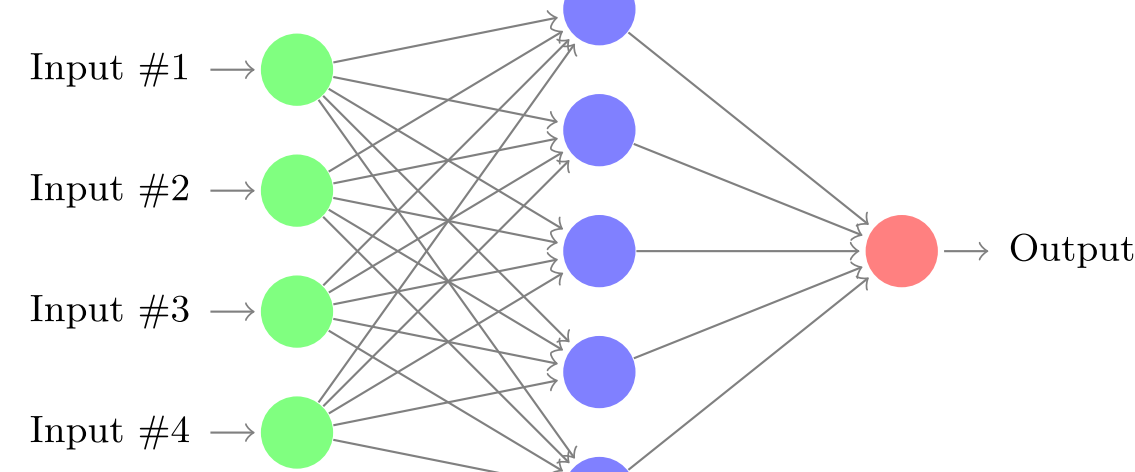
- The functions  $r_\ell$  and  $s_\ell$  are nonnegative functions revealing the priors of the weight variable  $\mathbf{W}_\ell$  and the state variable  $\mathbf{V}_\ell$  (i.e., regularizers)

- The constrained optimization is usually rewritten as an unconstrained one through:

$$\min_{\mathcal{W}, \mathcal{V}} \mathcal{L}(\mathcal{W}, \mathcal{V}) := \mathcal{L}_0(\mathcal{W}, \mathcal{V}) + \frac{\gamma}{2} \sum_{\ell=1}^L \|\mathbf{V}_\ell - \sigma_\ell(\mathbf{W}_\ell \mathbf{V}_{\ell-1})\|_F^2,$$

where  $\gamma > 0$  is a tuning parameter/hyperparameter

Input layer Hidden layer Output layer



$\mathbf{X} \in \mathbb{R}^{4 \times n}$   $\sigma_1(\mathbf{W}_1 \mathbf{X}) =: \mathbf{V}_1$   $\hat{\mathbf{Y}} = \mathbf{W}_2 \mathbf{V}_1$

- Jointly minimize the *squared distances* (in terms of **Frobenius norms**) between the input and the output of hidden layers

- E.g., define  $\mathbf{V}_0 := \mathbf{X}$ ,
 
$$\|\mathbf{V}_1 - \sigma_1(\mathbf{W}_1 \mathbf{V}_0)\|_F^2$$

### Three-Splitting Formulation

- Introduce two sets of auxiliary variables  $\mathcal{U} := \{\mathbf{U}_\ell\}_{\ell=1}^L$ ,  $\mathcal{V} := \{\mathbf{V}_\ell\}_{\ell=1}^L$ 

$$\min_{\mathcal{W}, \mathcal{V}, \mathcal{U}} \mathcal{L}_0(\mathcal{W}, \mathcal{V})$$

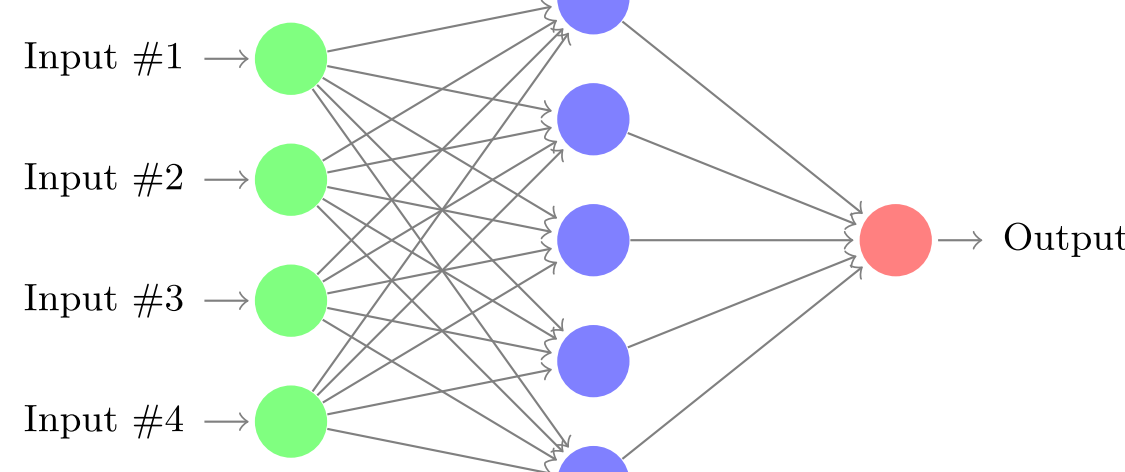
subject to  $\mathbf{U}_\ell = \mathbf{W}_\ell \mathbf{V}_{\ell-1}, \mathbf{V}_\ell = \sigma_\ell(\mathbf{U}_\ell), \ell \in \{1, \dots, L\}$

- The constrained optimization is usually rewritten as an unconstrained one through:
 
$$\min_{\mathcal{W}, \mathcal{V}} \mathcal{L}(\mathcal{W}, \mathcal{V}) := \mathcal{L}_0(\mathcal{W}, \mathcal{V}) + \frac{\gamma}{2} \sum_{\ell=1}^L \|\mathbf{V}_\ell - \sigma_\ell(\mathbf{W}_\ell \mathbf{V}_{\ell-1})\|_F^2,$$

where  $\gamma > 0$  is a tuning parameter/hyperparameter

- Variables are more loosely coupled than those in the two-splitting formulation

Input layer Hidden layer Output layer



$\mathbf{X} \in \mathbb{R}^{4 \times n}$   $\mathbf{W}_1 \mathbf{X} =: \mathbf{U}_1$   $\sigma_1(\mathbf{U}_1) =: \mathbf{V}_1$   $\hat{\mathbf{Y}} = \mathbf{W}_2 \mathbf{V}_1$

- Jointly minimize the *squared distances* (in terms of **Frobenius norms**) between
  - the input and the *pre-activation* output of hidden layers
  - the *pre-activation* output and the *post-activation* output of hidden layers

- E.g., define  $\mathbf{V}_0 := \mathbf{X}$ ,
 
$$\|\mathbf{U}_1 - \mathbf{W}_1 \mathbf{V}_0\|_F^2 + \|\mathbf{V}_1 - \sigma_1(\mathbf{U}_1)\|_F^2$$

### Algorithms

#### Algorithm 1 Two-splitting BCD for DNN Training

Data:  $\mathbf{X} \in \mathbb{R}^{d \times n}$ ,  $\mathbf{Y} \in \mathbb{R}^{k \times n}$   
 Initialization:  $\{\mathbf{W}_\ell^{(0)}, \mathbf{V}_\ell^{(0)}\}_{\ell=1}^L$ ,  $\mathbf{V}_0^{(0)} \equiv \mathbf{V}_0 := \mathbf{X}$   
 Hyperparameters:  $\gamma > 0, \alpha > 0$   
 for  $t = 1, \dots$  do  
 $\mathbf{V}_L^{(t)} = \arg\min_{\mathbf{V}_L} \{s_L(\mathbf{V}_L) + \mathcal{R}_n(\mathbf{V}_L; \mathbf{Y}) + \frac{\gamma}{2} \|\mathbf{V}_L - \mathbf{W}_L^{(t-1)} \mathbf{V}_{L-1}^{(t-1)}\|_F^2 + \frac{\alpha}{2} \|\mathbf{V}_L - \mathbf{V}_L^{(t-1)}\|_F^2\}$   
 $\mathbf{W}_L^{(t)} = \arg\min_{\mathbf{W}_L} \{r_L(\mathbf{W}_L) + \frac{\gamma}{2} \|\mathbf{V}_L^{(t)} - \mathbf{W}_L \mathbf{V}_{L-1}^{(t-1)}\|_F^2 + \frac{\alpha}{2} \|\mathbf{W}_L - \mathbf{W}_L^{(t-1)}\|_F^2\}$   
 for  $\ell = L-1, \dots, 1$  do  
 $\mathbf{V}_\ell^{(t)} = \arg\min_{\mathbf{V}_\ell} \{s_\ell(\mathbf{V}_\ell) + \frac{\gamma}{2} \|\mathbf{V}_\ell - \sigma_\ell(\mathbf{W}_\ell^{(t-1)} \mathbf{V}_{\ell-1}^{(t-1)})\|_F^2 + \frac{\gamma}{2} \|\mathbf{V}_\ell^{(t)} - \sigma_{\ell+1}(\mathbf{W}_{\ell+1}^{(t)} \mathbf{V}_\ell)\|_F^2 + \frac{\alpha}{2} \|\mathbf{V}_\ell - \mathbf{V}_\ell^{(t-1)}\|_F^2\}$   
 $\mathbf{W}_\ell^{(t)} = \arg\min_{\mathbf{W}_\ell} \{r_\ell(\mathbf{W}_\ell) + \frac{\gamma}{2} \|\mathbf{V}_\ell^{(t)} - \sigma_\ell(\mathbf{W}_\ell \mathbf{V}_{\ell-1}^{(t-1)})\|_F^2 + \frac{\alpha}{2} \|\mathbf{W}_\ell - \mathbf{W}_\ell^{(t-1)}\|_F^2\}$   
 end for  
 end for

#### Algorithm 2 Three-splitting BCD for DNN training

Samples:  $\mathbf{X} \in \mathbb{R}^{d \times n}$ ,  $\mathbf{Y} \in \mathbb{R}^{k \times n}$   
 Initialization:  $\{\mathbf{W}_\ell^{(0)}, \mathbf{V}_\ell^{(0)}, \mathbf{U}_\ell^{(0)}\}_{\ell=1}^L$ ,  $\mathbf{V}_0^{(0)} \equiv \mathbf{V}_0 := \mathbf{X}$   
 Hyperparameters:  $\gamma > 0, \alpha > 0$   
 for  $t = 1, \dots$  do  
 $\mathbf{V}_L^{(t)} = \arg\min_{\mathbf{V}_L} \{s_L(\mathbf{V}_L) + \mathcal{R}_n(\mathbf{V}_L; \mathbf{Y}) + \frac{\gamma}{2} \|\mathbf{V}_L - \mathbf{U}_L^{(t-1)}\|_F^2 + \frac{\alpha}{2} \|\mathbf{V}_L - \mathbf{V}_L^{(t-1)}\|_F^2\}$   
 $\mathbf{U}_L^{(t)} = \arg\min_{\mathbf{U}_L} \{\frac{\gamma}{2} \|\mathbf{V}_L^{(t)} - \mathbf{U}_L\|_F^2 + \frac{\gamma}{2} \|\mathbf{U}_L - \mathbf{W}_L^{(t-1)} \mathbf{V}_{L-1}^{(t-1)}\|_F^2\}$   
 $\mathbf{W}_L^{(t)} = \arg\min_{\mathbf{W}_L} \{r_L(\mathbf{W}_L) + \frac{\gamma}{2} \|\mathbf{U}_L^{(t)} - \mathbf{W}_L \mathbf{V}_{L-1}^{(t-1)}\|_F^2 + \frac{\alpha}{2} \|\mathbf{W}_L - \mathbf{W}_L^{(t-1)}\|_F^2\}$   
 for  $\ell = L-1, \dots, 1$  do  
 $\mathbf{V}_\ell^{(t)} = \arg\min_{\mathbf{V}_\ell} \{s_\ell(\mathbf{V}_\ell) + \frac{\gamma}{2} \|\mathbf{V}_\ell - \sigma_\ell(\mathbf{U}_\ell^{(t-1)})\|_F^2 + \frac{\gamma}{2} \|\mathbf{U}_\ell^{(t)} - \mathbf{W}_{\ell+1}^{(t)} \mathbf{V}_\ell\|_F^2\}$   
 $\mathbf{U}_\ell^{(t)} = \arg\min_{\mathbf{U}_\ell} \{\frac{\gamma}{2} \|\mathbf{V}_\ell^{(t)} - \sigma_\ell(\mathbf{U}_\ell)\|_F^2 + \frac{\gamma}{2} \|\mathbf{U}_\ell - \mathbf{W}_\ell^{(t-1)} \mathbf{V}_{\ell-1}^{(t-1)}\|_F^2 + \frac{\alpha}{2} \|\mathbf{U}_\ell - \mathbf{U}_\ell^{(t-1)}\|_F^2\}$   
 $\mathbf{W}_\ell^{(t)} = \arg\min_{\mathbf{W}_\ell} \{r_\ell(\mathbf{W}_\ell) + \frac{\gamma}{2} \|\mathbf{U}_\ell^{(t)} - \mathbf{W}_\ell \mathbf{V}_{\ell-1}^{(t-1)}\|_F^2 + \frac{\alpha}{2} \|\mathbf{W}_\ell - \mathbf{W}_\ell^{(t-1)}\|_F^2\}$   
 end for  
 end for

## GLOBAL CONVERGENCE ANALYSIS

### Main Assumptions

**Assumption 1** Suppose that

- the loss function  $\mathcal{L}$  is a proper lower semicontinuous<sup>1</sup> and nonnegative function,
- the activation functions  $\sigma_\ell$  ( $\ell = 1, \dots, L-1$ ) are Lipschitz continuous on any bounded set,
- the regularizers  $r_\ell$  and  $s_\ell$  ( $\ell = 1, \dots, L-1$ ) are nonnegative lower semicontinuous convex functions, and
- all these functions  $\mathcal{L}$ ,  $\sigma_\ell$ ,  $r_\ell$  and  $s_\ell$  ( $\ell = 1, \dots, L-1$ ) are either real analytic or semialgebraic, and continuous on their domains.

**Proposition 1** Examples satisfying Assumption 1 include:

- $\mathfrak{L}$  is the squared, logistic, hinge, or cross-entropy losses;
- $\sigma_\ell$  is ReLU, leaky ReLU, sigmoid, hyperbolic tangent, linear, polynomial, or softplus activations;
- $r_\ell$  and  $s_\ell$  are the squared  $\ell_2$  norm, the  $\ell_1$  norm, the elastic net, the indicator function of some nonempty closed convex set (such as the nonnegative closed half space, box set or a closed interval  $[0, 1]$ ), or 0 if no regularization.

### Main Theorem

**Theorem 1** Let  $\{\mathcal{Q}^t := (\{\mathbf{W}_\ell^t\}_{\ell=1}^L, \{\mathbf{V}_\ell^t\}_{\ell=1}^L)\}_{t \in \mathbb{N}}$  and  $\{\mathcal{P}^t := (\{\mathbf{W}_\ell^t\}_{\ell=1}^L, \{\mathbf{V}_\ell^t\}_{\ell=1}^L, \{\mathbf{U}_\ell^t\}_{\ell=1}^L)\}_{t \in \mathbb{N}}$  be the sequences generated by Algorithms 1 and 2, respectively. Suppose that Assumption 1 holds, and that one of the following conditions holds: (i) there exists a convergent subsequence  $\{\mathcal{Q}^{t_j}\}_{j \in \mathbb{N}}$  (resp.  $\{\mathcal{P}^{t_j}\}_{j \in \mathbb{N}}$ ); (ii)  $r_\ell$  is coercive for any  $\ell = 1, \dots, L$ ; (iii)  $\mathcal{L}$  (resp.  $\bar{\mathcal{L}}$ ) is coercive. Then for any  $\alpha > 0$ ,  $\gamma > 0$  and any finite initialization  $\mathcal{Q}^0$  (resp.  $\mathcal{P}^0$ ), the following hold

- $\{\mathcal{L}(\mathcal{Q}^t)\}_{t \in \mathbb{N}}$  (resp.  $\{\bar{\mathcal{L}}(\mathcal{P}^t)\}_{t \in \mathbb{N}}$ ) converges to some  $\mathcal{L}^*$  (resp.  $\bar{\mathcal{L}}^*$ ).
- $\{\mathcal{Q}^t\}_{t \in \mathbb{N}}$  (resp.  $\{\mathcal{P}^t\}_{t \in \mathbb{N}}$ ) converges to a critical point of  $\mathcal{L}$  (resp.  $\bar{\mathcal{L}}$ ).
- $\frac{1}{T} \sum_{t=1}^T \|\mathbf{g}^t\|_F^2 \rightarrow 0$  at the rate  $\mathcal{O}(1/T)$  where  $\mathbf{g}^t \in \partial \mathcal{L}(\mathcal{Q}^t)$ . Similarly,  $\frac{1}{T} \sum_{t=1}^T \|\bar{\mathbf{g}}^t\|_F^2 \rightarrow 0$  at the rate  $\mathcal{O}(1/T)$  where  $\bar{\mathbf{g}}^t \in \partial \bar{\mathcal{L}}(\mathcal{P}^t)$ .

### Extensions

#### Extension to Prox-Linear

- In the  $\mathbf{V}_L$ -update of both Algorithms 1 and 2, the empirical risk is involved in the optimization problems
- Generally hard to obtain its closed-form solution except for the square loss
- Use *prox-linear* update strategies for other smooth losses such as the logistic, cross-entropy, and exponential losses
- For some parameter  $\alpha > 0$ , the  $\mathbf{V}_L$ -update in Algorithm 1 is  $\mathbf{V}_L^{(t)} = \arg\min_{\mathbf{V}_L} \{s_L(\mathbf{V}_L) + \langle \nabla \mathcal{R}_n(\mathbf{V}_L^{(t-1)}; \mathbf{Y}), \mathbf{V}_L - \mathbf{V}_L^{(t-1)} \rangle + \frac{\gamma}{2} \|\mathbf{V}_L - \mathbf{V}_L^{(t-1)}\|_F^2 + \frac{\alpha}{2} \|\mathbf{V}_L - \mathbf{W}_L^{(t-1)} \mathbf{V}_{L-1}^{(t-1)}\|_F^2\}$ ,
 
$$\mathbf{V}_L^{(t)} = \arg\min_{\mathbf{V}_L} \{s_L(\mathbf{V}_L) + \langle \nabla \mathcal{R}_n(\mathbf{V}_L^{(t-1)}; \mathbf{Y}), \mathbf{V}_L - \mathbf{V}_L^{(t-1)} \rangle + \frac{\gamma}{2} \|\mathbf{V}_L - \mathbf{V}_L^{(t-1)}\|_F^2 + \frac{\alpha}{2} \|\mathbf{V}_L - \mathbf{U}_L^{(t-1)}\|_F^2\}. \quad (1)$$
- The  $\mathbf{V}_L$ -update in Algorithm 2 is  $\mathbf{V}_L^{(t)} = \arg\min_{\mathbf{V}_L} \{s_L(\mathbf{V}_L) + \langle \nabla \mathcal{R}_n(\mathbf{V}_L^{(t-1)}; \mathbf{Y}), \mathbf{V}_L - \mathbf{V}_L^{(t-1)} \rangle + \frac{\gamma}{2} \|\mathbf{V}_L - \mathbf{V}_L^{(t-1)}\|_F^2 + \frac{\gamma}{2} \|\mathbf{V}_L - \mathbf{U}_L^{(t-1)}\|_F^2\}$ .
 
$$\mathbf{V}_L^{(t)} = \arg\min_{\mathbf{V}_L} \{s_L(\mathbf{V}_L) + \langle \nabla \mathcal{R}_n(\mathbf{V}_L^{(t-1)}; \mathbf{Y}), \mathbf{V}_L - \mathbf{V}_L^{(t-1)} \rangle + \frac{\gamma}{2} \|\mathbf{V}_L - \mathbf{V}_L^{(t-1)}\|_F^2 + \frac{\gamma}{2} \|\mathbf{V}_L - \mathbf{U}_L^{(t-1)}\|_F^2\}. \quad (2)$$

**Theorem 2 (Global convergence for prox-linear update)** Consider adopting the prox-linear updates (1), (2) for the  $\mathbf{V}_L$ -subproblems in Algorithms 1 and 2, respectively. Under the conditions of Theorem 1, if further  $\nabla \mathcal{R}_n$  is Lipschitz continuous with a Lipschitz constant  $L_R$  and  $\alpha > \max\{0, \frac{L_R - 1}{2}\}$ , then all claims in Theorem Theorem 1 still hold for both algorithms.

#### Extension to ResNet Training

- Consider the simplified ResNet training problem (**two-splitting formulation**):

$$\min_{\mathcal{W}, \mathcal{V}} \mathcal{R}_n(\mathbf{V}_L; \mathbf{Y}) + \sum_{\ell=1}^L r_\ell(\mathbf{W}_\ell) + \sum_{\ell=1}^L s_\ell(\mathbf{V}_\ell)$$

subject to  $\mathbf{V}_\ell = \mathbf{V}_{\ell-1} = \sigma_\ell(\mathbf{W}_\ell \mathbf{V}_{\ell-1}), \ell \in \{1, \dots, L\} \quad (3)$

- Three-splitting formulation**:

$$\min_{\mathcal{W}, \mathcal{V}, \mathcal{U}} \mathcal{R}_n(\mathbf{V}_L; \mathbf{Y}) + \sum_{\ell=1}^L r_\ell(\mathbf{W}_\ell) + \sum_{\ell=1}^L s_\ell(\mathbf{V}_\ell)$$

subject to  $\mathbf{U}_\ell = \mathbf{W}_\ell \mathbf{V}_{\ell-1}, \mathbf{V}_\ell - \mathbf{V}_{\ell-1} = \sigma_\ell(\mathbf{U}_\ell), \ell \in \{1, \dots, L\}$

$$\bar{\mathcal{L}}_{\text{res}}(\mathcal{W}, \mathcal{V}, \mathcal{U}) := \mathcal{R}_n(\mathbf{V}_L; \mathbf{Y}) + \sum_{\ell=1}^L r_\ell(\mathbf{W}_\ell) + \sum_{\ell=1}^L s_\ell(\mathbf{V}_\ell) + \frac{\gamma}{2} \sum_{\ell=1}^L [\|\mathbf{V}_\ell - \mathbf{V}_{\ell-1} - \sigma_\ell(\mathbf{U}_\ell)\|_F^2 + \|\mathbf{U}_\ell - \mathbf{W}_\ell \mathbf{V}_{\ell-1}\|_F^2]$$

#### Algorithm 3 BCD for DNN Training with ResNets

Samples:  $\mathbf{X} \in \mathbb{R}^{d_0 \times n}$ ,  $\mathbf{Y} \in \mathbb{R}^{d_N \times n}$ ,  $\mathbf{V}_0^{(0)} \equiv \mathbf{V}_0 := \mathbf{X}$   
 Initialization:  $\{\mathbf{W}_\ell^0, \mathbf{V}_\ell^0, \mathbf{U}_\ell^0\}_{\ell=1}^L$   
 Parameters:  $\gamma > 0, \alpha > 0$   
 for  $t = 1, \dots$  do  
 $\mathbf{V}_L^{(t)} = \arg\min_{\mathbf{V}_L} \{s_L(\mathbf{V}_L) + \mathcal{R}_n(\mathbf{V}_L; \mathbf{Y}) + \frac{\gamma}{2} \|\mathbf{V}_L - \mathbf{V}_{L-1}^{(t-1)} - \mathbf{U}_L^{(t-1)}\|_F^2 + \frac{\alpha}{2} \|\mathbf{V}_L - \mathbf{V}_L^{(t-1)}\|_F^2\}$   
 $\mathbf{U}_L^{(t)} = \arg\min_{\mathbf{U}_L} \{\frac{\gamma}{2} \|\mathbf{V}_L^{(t)} - \mathbf{V}_{L-1}^{(t-1)} - \mathbf{U}_L\|_F^2 + \|\mathbf{U}_L - \mathbf{W}_L^{(t-1)} \mathbf{V}_{L-1}^{(t-1)}\|_F^2\}$   
 $\mathbf{W}_L^{(t)} = \arg\min_{\mathbf{W}_L} \{r_L(\mathbf{W}_L) + \frac{\gamma}{2} \|\mathbf{U}_L^{(t)} - \mathbf{W}_L \mathbf{V}_{L-1}^{(t-1)}\|_F^2 + \frac{\alpha}{2} \|\mathbf{W}_L - \mathbf{W}_L^{(t-1)}\|_F^2\}$   
 for  $\ell = L-1, \dots, 1$  do  
 $\mathbf{V}_\ell^{(t)} = \arg\min_{\mathbf{V}_\ell} \{s_\ell(\mathbf{V}_\ell) + \frac{\gamma}{2} \|\mathbf{V}_\ell - \mathbf{V}_{\ell-1}^{(t-1)} - \sigma_\ell(\mathbf{U}_\ell^{(t-1)})\|_F^2 + \|\mathbf{V}_{\ell+1}^{(t)} - \mathbf{V}_\ell - \sigma_{\ell+1}(\mathbf{U}_{\ell+1}^{(t)})\|_F^2 + \|\mathbf{U}_{\ell+1}^{(t)} - \mathbf{W}_{\ell+1}^{(t)} \mathbf{V}_\ell\|_F^2\}$   
 $\mathbf{U}_\ell^{(t)} = \arg\min_{\mathbf{U}_\ell} \{\frac{\gamma}{2} \|\mathbf{V}_\ell^{(t)} - \mathbf{V}_{\ell-1}^{(t-1)} - \sigma_\ell(\mathbf{U}_\ell)\|_F^2 + \|\mathbf{U}_\ell - \mathbf{W}_\ell^{(t-1)} \mathbf{V}_{\ell-1}^{(t-1)}\|_F^2 + \frac{\alpha}{2} \|\mathbf{U}_\ell - \mathbf{U}_\ell^{(t-1)}\|_F^2\}$   
 $\mathbf{W}_\ell^{(t)} = \arg\min_{\mathbf{W}_\ell} \{r_\ell(\mathbf{W}_\ell) + \frac{\gamma}{2} \|\mathbf{U}_\ell^{(t)} - \mathbf{W}_\ell \mathbf{V}_{\ell-1}^{(t-1)}\|_F^2 + \frac{\alpha}{2} \|\mathbf{W}_\ell - \mathbf{W}_\ell^{(t-1)}\|_F^2\}$   
 end for  
 end for

**Theorem 3 (Convergence of BCD for ResNets)** Let  $\{\{\mathbf{W}_\ell^t, \mathbf{V}_\ell^t, \mathbf{U}_\ell^t\}_{\ell=1}^L\}_{t \in \mathbb{N}}$  be a sequence generated by BCD for the DNN training model with ResNets (i.e., Algorithm 3). Let assumptions of Theorem 1 hold. Then all claims in Theorem 1 still hold for BCD with ResNets by replacing  $\bar{\mathcal{L}}$  with  $\bar{\mathcal{L}}_{\text{res}}$ .

Moreover, consider adopting the prox-linear update for the  $\mathbf{V}_L$ -subproblem in Algorithm 3, then under the assumptions of Theorem 2, all claims of Theorem 2 still hold for Algorithm 3.

## PROOF IDEAS

Four key ingredients:

- The *sufficient descent* condition
- The *relative error* condition
- The *continuity condition* of the objective function
- The *Kurdyka-Łojasiewicz* property of the objective function

Establishing the *sufficient descent* and the *relative error* conditions require two kinds of assumptions:

- multiconvexity and differentiability assumptions, and
- (blockwise) Lipschitz differentiability assumption on the unregularized part of objective function

- In our cases, the unregularized part of  $\mathcal{L}$  in **two-splitting formulation**,

$$\mathcal{R}_n(\mathbf{V}_L; \mathbf{Y}) + \frac{\gamma}{2} \sum_{\ell=1}^L \|\mathbf{V}_\ell - \sigma_\ell(\mathbf{W}_\ell \mathbf{V}_{\ell-1})\|_F^2,$$

and that of  $\bar{\mathcal{L}}$  in **three-splitting formulation**,

$$\mathcal{R}_n(\mathbf{V}_L; \mathbf{Y}) + \frac{\gamma}{2} \sum_{\ell=1}^L [\|\mathbf{V}_\ell - \sigma_\ell(\mathbf{U}_\ell)\|_F^2 + \|\mathbf{U}_\ell - \mathbf{W}_\ell \mathbf{V}_{\ell-1}\|_F^2]$$

usually do not satisfy any of **assumption (a)** and **assumption (b)**

- E.g., when  $\sigma_\ell$  is ReLU or leaky ReLU, the functions  $\|\mathbf{V}_\ell - \sigma_\ell(\mathbf{W}_\ell \mathbf{V}_{\ell-1})\|_F^2$  and  $\|\mathbf{V}_\ell - \sigma_\ell(\mathbf{U}_\ell)\|_F^2$  are non-differentiable and nonconvex with respect to  $\mathbf{W}_\ell$ -block and  $\mathbf{U}_\ell$ -block, respectively

- To overcome these challenges:
  - Exploit the proximal strategies for all the non-strongly convex subproblems (see Algorithm 2) to cheaply obtain the desired *sufficient descent* property (see Lemma 1)
  - Take advantage