



ICML 2019  
Long Beach, CA

# Global Convergence of Block Coordinate Descent in Deep Learning

Jinshan Zeng<sup>1,2,\*</sup> Tim Tsz-Kit Lau<sup>3,\*</sup> Shaobo 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

Tim Tsz-Kit Lau

Department of Statistics **Northwestern University**

# INTRODUCTION

# MOTIVATION OF BLOCK COORDINATE DESCENT (BCD) IN DEEP LEARNING

- Gradient-based methods are commonly used in training deep neural networks
- But gradient-based methods may suffer from various problems for deep networks
- Gradients of the loss function w.r.t. parameters of earlier layers involve those of later layers
  - ⇒ **Gradient vanishing**
  - ⇒ **Gradient exploding**
- **First-order gradient-based** methods does not work well

# MOTIVATION OF BLOCK COORDINATE DESCENT (BCD) IN DEEP LEARNING

- **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:
  - Deal with non-differentiable nonlinearities
  - Potentially avoid **vanishing gradient**
  - Can be easily implemented in a *distributed* and *parallel* fashion

# BLOCK COORDINATE DESCENT IN DEEP LEARNING

# BLOCK COORDINATE DESCENT IN DEEP LEARNING

- 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
  - $\mathcal{L} : \mathbb{R}^k \times \mathbb{R}^k \rightarrow \mathbb{R}_+ \cup \{0\}$ : loss function
  - $\Phi(\mathbf{x}_i; \mathcal{W}) := \sigma_L(\mathbf{W}_L \sigma_{L-1}(\mathbf{W}_{L-1} \cdots \mathbf{W}_2 \sigma_1(\mathbf{W}_1 \mathbf{x}_i)))$ : 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 ways of variable splitting appear in the literature

## BCD IN DEEP LEARNING: TWO-SPLITTING FORMULATION

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

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

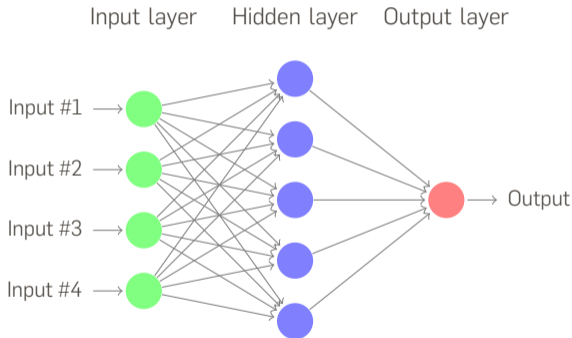
$$\text{subject to } \mathbf{V}_\ell = \sigma_\ell(\mathbf{W}_\ell \mathbf{V}_{\ell-1}), \ell \in \{1, \dots, L\}$$

- The functions  $r_\ell$  and  $s_\ell$  are regularizers
- Rewritten as unconstrained optimization:

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

- $\gamma > 0$  is a hyperparameter

## TWO-SPLITTING FORMULATION: GRAPHICAL ILLUSTRATION



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

- Jointly minimize the *distances* (in terms of **squared Frobenius norms**) between the input and the output of hidden layers
- E.g., define  $\mathbf{V}_0 := \mathbf{X}$ ,

$$\|\mathbf{V}_1 - \sigma_1(\mathbf{W}_1 \mathbf{V}_0)\|_F^2$$



## BCD IN DEEP LEARNING: THREE-SPLITTING FORMULATION

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

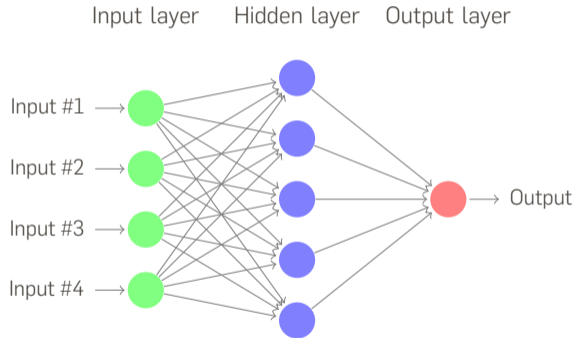
$$\min_{\mathcal{W}, \mathcal{V}, \mathcal{U}} \mathcal{L}_0(\mathcal{W}, \mathcal{V}) \quad \text{subject to} \quad \mathbf{U}_\ell = \mathbf{W}_\ell \mathbf{V}_{\ell-1}, \mathbf{V}_\ell = \sigma_\ell(\mathbf{U}_\ell), \ell \in \{1, \dots, L\}$$

- Rewritten as unconstrained optimization:

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

- Variables **more loosely coupled** than those in two-splitting

# THREE-SPLITTING FORMULATION: GRAPHICAL ILLUSTRATION



$$\mathbf{X} \in \mathbb{R}^{4 \times n}$$

$$\begin{aligned} \mathbf{W}_1 \mathbf{X} &=: \mathbf{U}_1 \\ \sigma_1(\mathbf{U}_1) &=: \mathbf{V}_1 \end{aligned}$$

$$\hat{\mathbf{Y}} = \mathbf{W}_2 \mathbf{V}_1$$

- Jointly minimize the *distances* (in terms of **squared Frobenius norms**) between
  1. the input and the *pre-activation* output of hidden layers
  2. the *pre-activation* output and the *post-activation* output of hidden layers
- E.g., define  $\mathbf{V}_0 := \mathbf{X}$ ,

$$\begin{aligned} &\|\mathbf{U}_1 - \mathbf{W}_1 \mathbf{V}_0\|_F^2 \\ &+ \|\mathbf{V}_1 - \sigma_1(\mathbf{U}_1)\|_F^2 \end{aligned}$$

# BLOCK COORDINATE DESCENT (BCD) ALGORITHMS

# BLOCK COORDINATE DESCENT (BCD) ALGORITHMS

- Devise algorithms for training DNNs based on the two formulations
- Update all the variables cyclically while fixing the remaining blocks
- Update in a backward order as in **backpropagation**
- Adopt the *proximal update strategies*

## BCD ALGORITHM (TWO-SPLITTING)

## 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:**  $\{\mathbf{W}_\ell^{(0)}, \mathbf{V}_\ell^{(0)}\}_{\ell=1}^L$ ,  $\mathbf{V}_0^{(t)} \equiv \mathbf{V}_0 := \mathbf{X}$ **Hyperparameters:**  $\gamma > 0$ ,  $\alpha > 0$ for  $t = 1, \dots$  do

$$\mathbf{V}_L^{(t)} = \operatorname{argmin}_{\mathbf{V}_L} \left\{ s_L(\mathbf{V}_L) + \mathcal{R}_n(\mathbf{V}_L; \mathbf{Y}) + \frac{\gamma}{2} \|\mathbf{V}_L - \mathbf{W}_L^{(t-1)} \mathbf{V}_{L-1}^{(t-1)}\|_F^2 + \frac{\alpha}{2} \|\mathbf{V}_L - \mathbf{V}_L^{(t-1)}\|_F^2 \right\}$$

$$\mathbf{W}_L^{(t)} = \operatorname{argmin}_{\mathbf{W}_L} \left\{ r_L(\mathbf{W}_L) + \frac{\gamma}{2} \|\mathbf{V}_L^{(t)} - \mathbf{W}_L \mathbf{V}_{L-1}^{(t-1)}\|_F^2 + \frac{\alpha}{2} \|\mathbf{W}_L - \mathbf{W}_L^{(t-1)}\|_F^2 \right\}$$

for  $\ell = L - 1, \dots, 1$  do

$$\mathbf{V}_\ell^{(t)} = \operatorname{argmin}_{\mathbf{V}_\ell} \left\{ s_\ell(\mathbf{V}_\ell) + \frac{\gamma}{2} \|\mathbf{V}_\ell - \sigma_\ell(\mathbf{W}_\ell^{(t-1)} \mathbf{V}_{\ell-1}^{(t-1)})\|_F^2 + \frac{\gamma}{2} \|\mathbf{V}_{\ell+1}^{(t)} - \sigma_{\ell+1}(\mathbf{W}_{\ell+1}^{(t)} \mathbf{V}_\ell)\|_F^2 + \frac{\alpha}{2} \|\mathbf{V}_\ell - \mathbf{V}_\ell^{(t-1)}\|_F^2 \right\}$$

$$\mathbf{W}_\ell^{(t)} = \operatorname{argmin}_{\mathbf{W}_\ell} \left\{ r_\ell(\mathbf{W}_\ell) + \frac{\gamma}{2} \|\mathbf{V}_\ell^{(t)} - \sigma_\ell(\mathbf{W}_\ell \mathbf{V}_{\ell-1}^{(t-1)})\|_F^2 + \frac{\alpha}{2} \|\mathbf{W}_\ell - \mathbf{W}_\ell^{(t-1)}\|_F^2 \right\}$$

end for

end for

## BCD ALGORITHM (THREE-SPLITTING)

## 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:**  $\{\mathbf{W}_\ell^{(0)}, \mathbf{V}_\ell^{(0)}, \mathbf{U}_\ell^{(0)}\}_{\ell=1}^L$ ,  $\mathbf{V}_0^{(t)} \equiv \mathbf{V}_0 := \mathbf{X}$ **Hyperparameters:**  $\gamma > 0, \alpha > 0$ for  $t = 1, \dots$  do

$$\mathbf{V}_L^{(t)} = \operatorname{argmin}_{\mathbf{V}_L} \left\{ s_L(\mathbf{V}_L) + \mathcal{R}_n(\mathbf{V}_L; \mathbf{Y}) + \frac{\gamma}{2} \|\mathbf{V}_L - \mathbf{U}_L^{(t-1)}\|_F^2 + \frac{\alpha}{2} \|\mathbf{V}_L - \mathbf{V}_L^{(t-1)}\|_F^2 \right\}$$

$$\mathbf{U}_L^{(t)} = \operatorname{argmin}_{\mathbf{U}_L} \left\{ \frac{\gamma}{2} \|\mathbf{V}_L^{(t)} - \mathbf{U}_L\|_F^2 + \frac{\gamma}{2} \|\mathbf{U}_L - \mathbf{W}_L^{(t-1)} \mathbf{V}_{L-1}^{(t-1)}\|_F^2 \right\}$$

$$\mathbf{W}_L^{(t)} = \operatorname{argmin}_{\mathbf{W}_L} \left\{ r_L(\mathbf{W}_L) + \frac{\gamma}{2} \|\mathbf{U}_L^{(t)} - \mathbf{W}_L \mathbf{V}_{L-1}^{(t-1)}\|_F^2 + \frac{\alpha}{2} \|\mathbf{W}_L - \mathbf{W}_L^{(t-1)}\|_F^2 \right\}$$

for  $\ell = L - 1, \dots, 1$  do

$$\mathbf{V}_\ell^{(t)} = \operatorname{argmin}_{\mathbf{V}_\ell} \left\{ s_\ell(\mathbf{V}_\ell) + \frac{\gamma}{2} \|\mathbf{V}_\ell - \sigma_\ell(\mathbf{U}_\ell^{(t-1)})\|_F^2 + \frac{\gamma}{2} \|\mathbf{U}_{\ell+1}^{(t)} - \mathbf{W}_{\ell+1}^{(t)} \mathbf{V}_\ell\|_F^2 \right\}$$

$$\mathbf{U}_\ell^{(t)} = \operatorname{argmin}_{\mathbf{U}_\ell} \left\{ \frac{\gamma}{2} \|\mathbf{V}_\ell^{(t)} - \sigma_\ell(\mathbf{U}_\ell)\|_F^2 + \frac{\gamma}{2} \|\mathbf{U}_\ell - \mathbf{W}_\ell^{(t-1)} \mathbf{V}_{\ell-1}^{(t-1)}\|_F^2 + \frac{\alpha}{2} \|\mathbf{U}_\ell - \mathbf{U}_\ell^{(t-1)}\|_F^2 \right\}$$

$$\mathbf{W}_\ell^{(t)} = \operatorname{argmin}_{\mathbf{W}_\ell} \left\{ r_\ell(\mathbf{W}_\ell) + \frac{\gamma}{2} \|\mathbf{U}_\ell^{(t)} - \mathbf{W}_\ell \mathbf{V}_{\ell-1}^{(t-1)}\|_F^2 + \frac{\alpha}{2} \|\mathbf{W}_\ell - \mathbf{W}_\ell^{(t-1)}\|_F^2 \right\}$$

end for

end for

# GLOBAL CONVERGENCE ANALYSIS

# ASSUMPTIONS OF THE FUNCTIONS FOR CONVERGENCE GUARANTEES

## Assumption

Suppose that

- (a) the loss function  $\mathcal{L}$  is a proper lower semicontinuous<sup>1</sup> and nonnegative function,
- (b) the activation functions  $\sigma_\ell$  ( $\ell = 1 \dots, L - 1$ ) are Lipschitz continuous on any bounded set,
- (c) the regularizers  $r_\ell$  and  $s_\ell$  ( $\ell = 1 \dots, L - 1$ ) are nonnegative lower semicontinuous convex functions, and
- (d) 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.

---

<sup>1</sup>A function  $f : \mathcal{X} \rightarrow \mathbb{R}$  is called *lower semicontinuous* if  $\liminf_{\mathbf{x} \rightarrow \mathbf{x}_0} f(\mathbf{x}) \geq f(\mathbf{x}_0)$  for any  $\mathbf{x}_0 \in \mathcal{X}$ .



# EXAMPLES OF THE FUNCTIONS

## Proposition

Examples satisfying Assumption 1 include:

- (a)  $\mathcal{L}$  is the squared, logistic, hinge, or cross-entropy losses;
- (b)  $\sigma_\ell$  is ReLU, leaky ReLU, sigmoid, hyperbolic tangent, linear, polynomial, or softplus activations;
- (c)  $r_\ell$  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

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<sup>2</sup> 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

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

<sup>2</sup>An extended-real-valued function  $h : \mathbb{R}^p \rightarrow \mathbb{R} \cup \{+\infty\}$  is called *coercive* if and only if  $h(\mathbf{x}) \rightarrow +\infty$  as  $\|\mathbf{x}\| \rightarrow +\infty$ .

# EXTENSIONS

## Extensions

1. *Prox-linear updates* instead of proximal update strategies
2. *Residual Networks (ResNets)* with skip connections

**Global convergence** of both extensions are also proved

# PROOF IDEAS

# 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:

- (a) **Multiconvexity** and **differentiability** assumptions, and
- (b) **(Blockwise) Lipschitz differentiability** assumption on the unregularized part of objective function

## PROOF IDEAS

- 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

# PROOF IDEAS

To overcome these challenges:

- (i) Exploit the **proximal strategies** for all the *non-strongly convex* subproblems to cheaply obtain the desired *sufficient descent* property
- (ii) Take advantage of the **Lipschitz continuity** of the activations as well as the specific splitting formulations to yield the desired *relative error* property

# SUMMARY OF THEORETICAL RESULTS OF THIS PAPER

## Theoretical Results

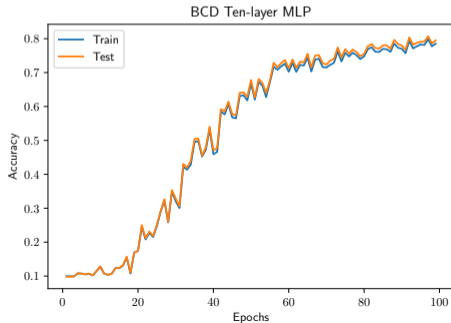
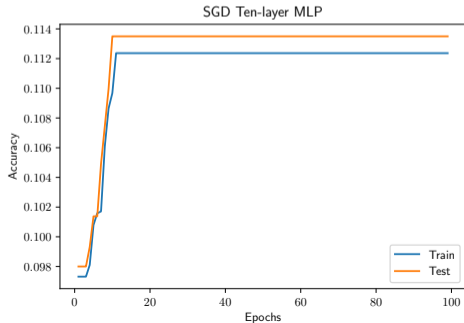
1. Global convergence to a critical point at a rate of  $\mathcal{O}(1/T)$ , where  $T$  is the number of iterations
2. Further, if the initialization is sufficiently close to some global minimum of  $\mathcal{L}$  or  $\bar{\mathcal{L}}$ , then both the sequences generated by Algorithms 1 and 2 converges to their corresponding global minima
3. **Comparison with the convergence of SGD/stochastic subgradient method:**
  - *BCD*: **Global (whole sequence)** convergence
  - *SGD* (Davis et al., 2019): **Subsequence** convergence



DEMONSTRATION

# DEMONSTRATION

- 10-class classification for the MNIST handwritten digit (0–9) dataset (with 60K training samples; 10K test samples)
- Fully-connected neural network (MLP)
- 10 hidden layers
- Comparison of training and test accuracies (after 100 epochs)



# Poster #78

Paper: <http://proceedings.mlr.press/v97/zeng19a.html>

GitHub: <https://github.com/timlautk/BCD-for-DNNs-PyTorch>

The End

Thank you!