

An Inexact Augmented Lagrangian Algorithm for Training Leaky ReLU Neural Network with Group Sparsity

Wei Liu

liuwei175@lsec.cc.ac.cn

Rensselaer Polytechnic Institute

Joint work with **Xin Liu** (AMSS, CAS), and **Xiaojun Chen** (PolyU)

Talk in SIAM OP23

June 2, 2023

Wei Liu, Xin Liu, and Xiaojun Chen, *Linearly Constrained Nonsmooth Optimization for Training Autoencoders*. SIAM Journal on Optimization, 2022, 32(3): 1931-1957

Wei Liu, Xin Liu, and Xiaojun Chen, *An Inexact Augmented Lagrangian algorithm for Training Leaky ReLU Neural Network with Group Sparsity*. Submitted to Journal of Machine Learning Research, under minor revision

1. The optimization problem for training deep neural networks

The Model for The Neural Network

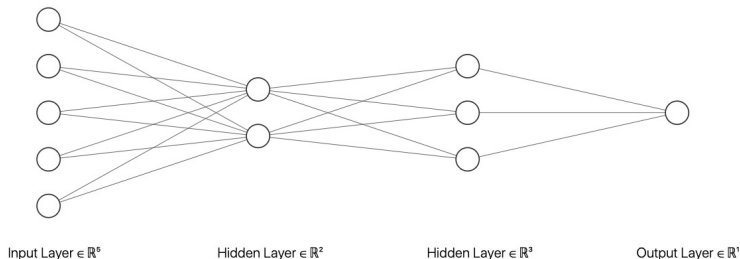
Neural Networks (NN)– deep learning

Given an input data $\{(x_n, y_n)\}_{n=1}^N$, where $x_n \in \mathbb{R}^{N_0}$, $y_n \in \mathbb{R}^{N_L}$

$$\min_{\substack{W_\ell, b_\ell, \\ \ell=1,2,\dots,L}} \frac{1}{N} \sum_{n=1}^N \|\sigma(W_L \sigma(\cdots \sigma(W_1 x_n + b_1) + b_2 \cdots) + b_L) - y_n\|^2$$

- N : the number of input data
- L : the number of layers
- N_ℓ : the number of neurons at the ℓ -th layer
- σ : the activation function.
Here we use leaky ReLU, i.e., $\sigma(z) = \max(z, \alpha z)$ with $1 > \alpha > 0$
- variables to be determined:
 - $W_\ell \in \mathbb{R}^{N_\ell \times N_{\ell-1}}$
 - $b_\ell \in \mathbb{R}^{N_\ell}$

An Example



Why leaky ReLU?

- ReLU and leaky ReLU reduce the vanishing gradient phenomenon by making the hidden layers sparse ([Sun 2019])
- the performance of the leaky ReLU network is reported to be slightly better than that of the ReLU network ([Mass et al. 2013, Pedamonti 2018])

Our Aimed Problem

DNN with group sparsity

$$(1) \quad \min_{w,b} \frac{1}{N} \sum_{n=1}^N \|\sigma(W_L \sigma(\cdots \sigma(W_1 x_n + b_1) + \cdots) + b_L) - y_n\|^2 + \mathcal{R}_1(w)$$

- \mathcal{R}_1 : group lasso regularizer with $\lambda_w > 0$, i.e.,

$$\mathcal{R}_1(w) := \lambda_w \sum_{\ell=1}^L \|W_\ell\|_{2,1} = \lambda_w \sum_{\ell=1}^L \sum_{j=1}^{N_{\ell-1}} \|(W_\ell)_{\cdot,j}\|$$

- $w = (\text{vec}(W_1)^T, \dots, \text{vec}(W_L)^T)^T \in \mathbb{R}^{\tilde{N}}$, $b = (b_1^T, \dots, b_L^T)^T \in \mathbb{R}^{\bar{N}}$
- $\tilde{N} := \sum_{\ell=1}^L N_\ell N_{\ell-1}$, $\bar{N} := \sum_{\ell=1}^L N_\ell$

Why group sparsity?

- pursuing the parameter sparsity and theoretical improvement in efficiency ([[Hoefler et al. 2021](#)])
- requiring less training time ([[Wen et al. 2016](#)])

Existing Approaches: SGD-based Methods

Properties

- calculating the gradient via the chain rule
- in cases where nonsmooth activation functions are used, a subgradient in a neighborhood is often used

Some SGD-based Methods

- Vanilla SGD ([Cramir 1946])
- Adagrad ([Duchi-Hazan-Singer 2011])
- AdagradDecay ([Duchi-Hazan-Singer 2011])
- Adadelata ([Zeiler 2012])
- Adam ([Kingma-Ba-Adam 2014])
- Adamax ([Kingma-Ba-Adam 2014])

Existing Approaches: SGD-based Methods (Cont'd)

Limitations

- neglecting the exactness in calculating the subgradient of the objective function \Rightarrow convergence?

chain rule does not hold

$$\min_{z \in \mathbb{R}} f(z) = \frac{1}{2} (-(zx_1)_+ + (zx_2)_+ + 1)^2 + \frac{1}{2} (-(zy_1)_+ + (zy_2)_+ + 1)^2$$

where $x = (-1, 1)$, $y = (-2, 0)$, $z_+ := \max\{0, z\}$.

SGD calculate

$$\begin{aligned} g(z) := & (-(zx_1)_+ + (zx_2)_+ + 1)(-x_1 h(zx_1) + x_2 h(zx_2)) \\ & + (-(zy_1)_+ + (zy_2)_+ + 1)(-y_1 h(zy_1) + x_2 h(zy_2)) \end{aligned}$$

as one derivative of f at z , $h(z) := \text{sign}(z_+)$.

however, 0 is not a stationary point!

More examples: [April-July-Kummer 2011]

An equivalent Model

[Carreira Perpiñán-Wang 2012]

$$(2) \quad \min_{\substack{W_\ell, v_{n,\ell}, b_\ell, \\ \ell=1,2,\dots,L, n=1,2,\dots,N}} \frac{1}{N} \sum_{n=1}^N \|v_{n,L} - y_n\|_F^2 + \mathcal{R}_1(w)$$

s. t. $v_{n,\ell} = \sigma_\ell(W_\ell v_{n,\ell-1} + b_\ell)$
 $n = 1, \dots, N, \ell = 1, \dots, L$

- $v_{n,0} = x_n$
- $v_{n,\ell}$: the output of the ℓ -th layer with respect to the n -th input data

Solving (2) instead of (1) is able to alleviate the problems that SGD-based methods suffer

Existing Methods for Solving the Problem (2)

- ℓ_2 penalty method:
 - MAC [Carreira-Perpinan-Wang 2014]
 - proximal BCD (pBCD) [Lau-Zeng-Wu-Yao 2018]
 - pBCD [Zeng-Lau-Lin-Yao 2019]
- multi-block 'ADMM':
 - Taylor et. al. 2016
 - Zhang-Chen-Saligrama 2016
 - Evens-Latafat-Themelis 2020

The above methods enjoy no exact penalty results!

ℓ_1 penalty method: present the exact penalty result, and establish the subsequence convergence result ([Cui-He-Pang 2020])

- limitation in theory: excludes ReLU and leaky ReLU
- limitation in practice: no regularizer

⇒ Design some new ℓ_1 penalty method.

A New Formulation

$$(P) \quad \min_{w,b,v,u} \bar{O}(w,v) := \frac{1}{N} \sum_{n=1}^N \|v_{n,L} - y_n\|^2 + \mathcal{R}_1(w) + \mathcal{R}_2(v)$$
$$\text{s. t. } \sigma(u_{n,\ell}) - v_{n,\ell} = 0, u_{n,\ell} - (W_\ell v_{n,\ell-1} + b_\ell) = 0,$$
$$n \in [N], \ell \in [L].$$

- $v := (v_{1,1}^T, v_{2,1}^T, \dots, v_{1,L}^T, v_{2,L}^T, \dots, v_{N,L}^T)^T \in \mathbb{R}^m$
- $u = (u_{1,1}^T, u_{2,1}^T, \dots, u_{1,L}^T, u_{2,L}^T, \dots, u_{N,L}^T)^T \in \mathbb{R}^m$
- $\mathcal{R}_2(v) = \lambda_v \|v\|^2$: a regularization term
- Feasible set: $\Omega_1 := \{(w, b, v, u) : v - \sigma(u) = 0, u = \Psi(v)w + Ab\}$
- linear operator $\Psi(v) : \mathbb{R}^m \mapsto \mathbb{R}^{m \times \bar{N}}$, matrix $A \in \mathbb{R}^{m \times \bar{N}}$
- $m = N\bar{N}$

A New Penalty Approach

- we consider to have $v \geq \sigma(u)$ as a constraint and add a penalty term $\beta^T(v - \sigma(u))$ in the objective function
- $\beta = (\beta_1 e_{NN_1}^T, \dots, \beta_L e_{NN_L}^T)^T \in \mathbb{R}^m$
- we write $v \geq \sigma(u)$ by $v - u \geq 0$ and $v - \alpha u \geq 0$

$$(PP) \quad \begin{aligned} \min_{w, b, v, u} \quad & O(w, v, u) = \bar{O}(w, v) + \beta^T(v - \sigma(u)) \\ \text{s. t.} \quad & v - u \geq 0, v - \alpha u \geq 0, u = \Psi(v)w + Ab. \end{aligned}$$

- feasible set: Ω_2
- we focus on obtaining a **l(imiting)-stationary point of problem (PP)**

Definitions

- Clarke subdifferential: $\partial^c f(\bar{z}) = \text{co} \{ \lim_{z \rightarrow \bar{z}} \nabla f(z) : f \text{ is smooth at } z \}$
- limiting subdifferential: $\partial f(\bar{z}) := \left\{ v : \exists z^k \xrightarrow{f} \bar{z}, v^k \rightarrow v \text{ such that } \liminf_{z \rightarrow z^k} \frac{f(z) - f(z^k) - \langle v^k, z - z^k \rangle}{\|z - z^k\|} \geq 0, \forall k \right\}$
- a point $\bar{z} \in \mathcal{Z}$ is a l-stationary point, a C(larke)-stationary point of $\min_{z \in \mathcal{Z}} f(z)$ if $0 \in \partial f(\bar{z}) + \mathcal{N}_{\mathcal{Z}}(\bar{z})$, $0 \in \partial^c f(\bar{z}) + \mathcal{N}_{\mathcal{Z}}^c(\bar{z})$, respectively

A New Penalty Approach

- we consider to have $v \geq \sigma(u)$ as a constraint and add a penalty term $\beta^T(v - \sigma(u))$ in the objective function
- $\beta = (\beta_1 e_{NN_1}^T, \dots, \beta_L e_{NN_L}^T)^T \in \mathbb{R}^m$
- we write $v \geq \sigma(u)$ by $v - u \geq 0$ and $v - \alpha u \geq 0$

$$(PP) \quad \min_{w, b, v, u} O(w, v, u) = \bar{O}(w, v) + \beta^T(v - \sigma(u))$$
$$\text{s. t. } v - u \geq 0, v - \alpha u \geq 0, u = \Psi(v)w + Ab.$$

- feasible set: Ω_2
- we focus on obtaining a **limiting-stationary point of problem (PP)**

Definitions

- Clarke subdifferential: $\partial^c f(\bar{z}) = \text{co} \{ \lim_{z \rightarrow \bar{z}} \nabla f(z) : f \text{ is smooth at } z \}$
- limiting subdifferential: $\partial f(\bar{z}) := \left\{ v : \exists z^k \xrightarrow{f} \bar{z}, v^k \rightarrow v \text{ such that } \liminf_{z \rightarrow z^k} \frac{f(z) - f(z^k) - \langle v^k, z - z^k \rangle}{\|z - z^k\|} \geq 0, \forall k \right\}$
- a point $\bar{z} \in \mathcal{Z}$ is a l-stationary point, a C(larke)-stationary point of $\min_{z \in \mathcal{Z}} f(z)$ if $0 \in \partial f(\bar{z}) + \mathcal{N}_{\mathcal{Z}}(\bar{z})$, $0 \in \partial^c f(\bar{z}) + \mathcal{N}_{\mathcal{Z}}^c(\bar{z})$, respectively

Why l-stationary point?

- *l*-stationary is stronger than C-stationary

An example

Consider
$$\min_{w_1 \in \mathbb{R}, w_2 \in \mathbb{R}, b_1 \in \mathbb{R}, b_2 \in \mathbb{R}} f(w_1, w_2, b_1, b_2) :=$$

$$((w_2 \sigma(w_1 + b_1) + b_2) + 1)^2 + ((w_2 \sigma(2w_1 + b_1) + b_2) - 1)^2. \quad (3)$$

let $w_2^* = 1, b_1^* = 0, w_1^* = 0, b_2^* = 0,$

$$\partial^c f(w_1^*, w_2^*, b_1^*, b_2^*) = \left\{ (t, 0, s, 0)^T : t \in [2\alpha - 4, 2 - 4\alpha], s \in [-2 + 2\alpha, 2 - 2\alpha] \right\},$$

$$\partial(f(w_1^*, w_2^*, b_1^*, b_2^*))$$

$$= \left\{ (-2\alpha, 0, 0, 0)^T, (2\alpha - 4, 0, 2\alpha - 2, 0)^T, (2 - 4\alpha, 0, 2 - 2\alpha, 0)^T, (-2, 0, 0, 0)^T \right\},$$

$$f(w_1^* + \epsilon, w_2^*, b_1^*, b_2^*) = 5\epsilon^2 - 2\epsilon + 2 < 2 = f(w_1^*, w_2^*, b_1^*, b_2^*), \quad \epsilon : \text{ a small positive scalar}$$

for any $0 < \alpha < \frac{1}{2}$, $(w_1^*, w_2^*, b_1^*, b_2^*)$ is a C-stationary point, but is not a l-stationary point and local minimizer.

2. Theoretical Results

Stationary Point of (P)

$v - \sigma(u) = 0$ can be rewritten as complementary constraints:

$v - u \geq 0, (v - u)(v - \alpha u) = 0, v - \alpha u \geq 0$, We call $(w^*, b^*, v^*, u^*) \in \Omega_1$ a

MPCC W-stationary point [Scheel-Scholtes 2000; Guo-Chen 2021] of (P), if there exist $\mu^1 \in \mathbb{R}^m, \mu^2 \in \mathbb{R}^m$ and $\xi \in \mathbb{R}^m$ such that

$$0 = \nabla_w \bar{O}(w^*, v^*) + \Psi(v^*)^T \xi, \quad 0 = A^T \xi$$

$$0 = \nabla_v \bar{O}(w^*, v^*) - \mu^1 - \mu^2 + \nabla_v \xi^T (u^* - \Psi(v^*) w^*)$$

$$0 = \mu^1 + \alpha \mu^2 + \xi$$

$$(\mu^1)^T (v^* - u^*) = 0, \quad (\mu^2)^T (v^* - \alpha u^*) = 0$$

MPCC W-stationary point + $\mu_i^1 \mu_i^2 \geq 0, \forall i : u_i^* = v_i^* = 0 \Rightarrow$ MPCC
C-stationary point

Lemma 1

*NNAMCQ^a holds for the constraints set of problem (P) \Rightarrow Any local minimizer of (P) is its **MPCC C-stationary point***

^a[Ye-Zhang 2013]

Stationary Point of (P)

$v - \sigma(u) = 0$ can be rewritten as complementary constraints:

$v - u \geq 0, (v - u)(v - \alpha u) = 0, v - \alpha u \geq 0$, We call $(w^*, b^*, v^*, u^*) \in \Omega_1$ a

MPCC W-stationary point [Scheel-Scholtes 2000; Guo-Chen 2021] of (P), if there exist $\mu^1 \in \mathbb{R}^m, \mu^2 \in \mathbb{R}^m$ and $\xi \in \mathbb{R}^m$ such that

$$0 = \nabla_w \bar{O}(w^*, v^*) + \Psi(v^*)^T \xi, \quad 0 = A^T \xi$$

$$0 = \nabla_v \bar{O}(w^*, v^*) - \mu^1 - \mu^2 + \nabla_v \xi^T (u^* - \Psi(v^*) w^*)$$

$$0 = \mu^1 + \alpha \mu^2 + \xi$$

$$(\mu^1)^T (v^* - u^*) = 0, \quad (\mu^2)^T (v^* - \alpha u^*) = 0$$

MPCC W-stationary point + $\mu_i^1 \mu_i^2 \geq 0, \forall i : u_i^* = v_i^* = 0 \Rightarrow$ MPCC
C-stationary point

Lemma 1

*NNAMCQ^a holds for the constraints set of problem (P) \Rightarrow Any local minimizer of (P) is its **MPCC C-stationary point***

^a[Ye-Zhang 2013]

Stationary Point of (PP)

We call (w^*, b^*, v^*, u^*) a **KKT point** of (PP), if there exists $\mu^1 \in \mathbb{R}_+^m$, $\mu^2 \in \mathbb{R}_+^m$ and $\xi \in \mathbb{R}^m$ such that

$$0 = \nabla_w \bar{O}(w^*, v^*) + \Psi(v^*)^T \xi, \quad 0 = A^T \xi$$

$$0 = \nabla_v \bar{O}(w^*, v^*) + \beta - \mu^1 - \mu^2 + \nabla_v \xi^T (u^* - \Psi(v^*) w^*)$$

$$0 \in \partial_u (-\beta^T \sigma(u^*)) + \mu^1 + \alpha \mu^2 + \xi$$

$$(\mu^1)^T (v^* - u^*) = 0, \quad (\mu^2)^T (v^* - \alpha u^*) = 0$$

Lemma 2

MFCQ^a holds for the constraints set of problem (PP).

^a[Mangasarian 1994]

- $\Rightarrow (w^*, b^*, v^*, u^*)$ is a **l-stationary point** of (PP) if and only if (w^*, b^*, v^*, u^*) is a **KKT point** of (PP)
- \Rightarrow any local minimizer of (PP) is its l-stationary point

Stationary Point of (PP)

We call (w^*, b^*, v^*, u^*) a **KKT point** of (PP), if there exists $\mu^1 \in \mathbb{R}_+^m$, $\mu^2 \in \mathbb{R}_+^m$ and $\xi \in \mathbb{R}^m$ such that

$$0 = \nabla_w \bar{O}(w^*, v^*) + \Psi(v^*)^T \xi, \quad 0 = A^T \xi$$

$$0 = \nabla_v \bar{O}(w^*, v^*) + \beta - \mu^1 - \mu^2 + \nabla_v \xi^T (u^* - \Psi(v^*) w^*)$$

$$0 \in \partial_u (-\beta^T \sigma(u^*)) + \mu^1 + \alpha \mu^2 + \xi$$

$$(\mu^1)^T (v^* - u^*) = 0, \quad (\mu^2)^T (v^* - \alpha u^*) = 0$$

Lemma 2

MFCQ^a holds for the constraints set of problem (PP).

^a[Mangasarian 1994]

- $\Rightarrow (w^*, b^*, v^*, u^*)$ is a **l-stationary point** of (PP) if and only if (w^*, b^*, v^*, u^*) is a **KKT point** of (PP)
- \Rightarrow any local minimizer of (PP) is its l-stationary point

Main Theoretical Results

Let $\theta > \frac{1}{N}\|X\|_F^2$ and

$$\Omega_\theta = \{(w, b, v, u) : v - u \geq 0, v - \alpha u \geq 0, u = \Psi(v)w + Ab, O(w, v, u) \leq \theta\}$$

Theorem 3

The set Ω_θ is bounded and the solution set to problem (PP) is nonempty and bounded.

(PP) : global(local) minimizers



(P) : global(local) minimizers

I-stationary point \Leftrightarrow KKT point

some conditions \Downarrow

MPCC C-stationary point



MPCC W-stationary point

- $\beta_\ell > LL_{\bar{O}} \max\{\theta_w, 1\}^L + 2 \sum_{j=\ell+1}^L \beta_j \theta_w \max\{\theta_w, 1\}^{j-\ell-1}$, $L_{\bar{O}}$ is the Lipschitz modulus of \bar{O} over Ω_θ
- $O(w, v, u) < \theta$

Main Theoretical Results

Let $\theta > \frac{1}{N}\|X\|_F^2$ and

$$\Omega_\theta = \{(w, b, v, u) : v - u \geq 0, v - \alpha u \geq 0, u = \Psi(v)w + Ab, O(w, v, u) \leq \theta\}$$

Theorem 3

The set Ω_θ is bounded and the solution set to problem (PP) is nonempty and bounded.

(PP) : global(local) minimizers



(P) : global(local) minimizers

I-stationary point \Leftrightarrow KKT point

some conditions \Downarrow

MPCC C-stationary point



MPCC W-stationary point

- $\beta_\ell > LL_{\bar{O}} \max\{\theta_w, 1\}^L + 2 \sum_{j=\ell+1}^L \beta_j \theta_w \max\{\theta_w, 1\}^{j-\ell-1}$, $L_{\bar{O}}$ is the Lipschitz modulus of \bar{O} over Ω_θ
- $O(w, v, u) < \theta$

Extensions to ReLU Networks

$\alpha = 0 \Rightarrow$ the solution set of (PP) may be **unbounded**



solve (PP) over a constructed set

$$\min_{w,b,v,u} O(w, v, u)$$

$$\text{s. t. } v - u \geq 0, v - \alpha u \geq 0, u = \Psi(v)w + Ab \quad (\text{PP}_b)$$

$$b \geq -e_{\overline{N}} \overline{N}(\theta_w + \theta_v)$$

- the solution set to (PP_b) is nonempty and bounded
- any local minimizer of (PP_b) is its l-stationary point, and a MPCC-W stationary point of problem (P)
- the designed algorithm can be used to solve ReLU networks

3. Algorithm Framework and Convergence Analysis

An Inexact Augmented Lagrangian Method

[Lu-Zhang 2012; Chen et al. 2017]

Augmented Lagrangian function

$$\mathcal{L}_\rho(w, b, v, u; \xi) := O(w, v, u) + \langle \xi, u - \Psi(v)w - Ab \rangle + \frac{\rho}{2} \|u - \Psi(v)w - Ab\|^2$$

- $\rho \in \mathbb{R}_+$
- $\xi \in \mathbb{R}^m$

Subproblem

$$\min_{(w, b, v, u): v \geq u, v \geq \alpha u} \mathcal{L}_\rho(w, b, v, u; \xi) \quad (4)$$

- **IALAM**: IALM framework with subproblem solved by an Alternating Minimization method

Algorithm Framework

The inexact augmented Lagrangian method for solving problem (PP) (IALM framework)

- ➊ Input: $(w^{(0)}, b^{(0)}, v^{(0)}, u^{(0)}) \in \Omega_\theta$, $\rho^{(0)} > 0$, $\eta_1 \in (0, 1)$, $\eta_2, \eta_3 > 1$, $\xi^{(0)} \in \mathbb{R}^m$, $\gamma \in \mathbb{N}_+$. $k := 1$
- ➋ let $(\xi, \rho) = (\xi^{(k-1)}, \rho^{(k-1)})$, solve (4) inexactly
- ➌ Calculate $\xi^{(k)} := \xi^{(k-1)} + \rho^{(k-1)}(u^{(k)} - \Psi(v^{(k)})w^{(k)} - Ab^{(k)})$
- ➍ if $k \leq \gamma$, let $\rho^{(k)} = \rho^{(k-1)}$. If $k > \gamma$ and

$$\|u^{(k)} - \Psi(v^{(k)})w^{(k)} - Ab^{(k)}\| \leq \eta_1 \max_{t=k-\gamma, \dots, k-1} \|u^{(t)} - \Psi(v^{(t)})w^{(t)} - Ab^{(t)}\|$$

let $\rho^{(k)} = \rho^{(k-1)}$. otherwise, let

$$\rho^{(k)} = \max \left\{ \rho^{(k-1)} / \eta_2, \|\xi^{(k)}\|^{1+\eta_3} \right\}$$

- ➎ Let $k := k + 1$, if stop criterion is not met, return to step 2.
- ➏ Output: $(w^{(k)}, b^{(k)}, v^{(k)}, u^{(k)})$

- the feasibility violation decreases non-monotonically.

Alternating Minimization Method for Solving (4)

Main Idea: solve (4) by splitting it into (w, b) and (v, u) blocks

why?

- the variable dimension of (4) is large
- w, b are variables of (1), meanwhile v, u are the auxiliary variables
- restricted to both of these two blocks are easy to solve
 - the (w, b) subproblem is strongly convex
 - the (v, u) subproblem has a closed-form unique solution

Alternating Minimization Method for Solving (4) (Cont'd)

Step 1: $(w^{(j)}, b^{(j)}) \rightarrow (w^{(j+1)}, b^{(j+1)})$

$$(w^{(j+1)}, b^{(j+1)}) := \arg \min_{w, b} \mathcal{L}_\rho(w, b, v^{(j)}, u^{(j)}; \xi)^1 \quad (5)$$

Step 2: $(v^{(j)}, u^{(j)}) \rightarrow (v^{(j+1)}, u^{(j+1)})$

$$(v^{(j+1)}, u^{(j+1)}) := \arg \min_{(v, u): v \geq u, v \geq \alpha u} \mathcal{L}_\rho(w^{(j+1)}, b^{(j+1)}, v, u; \xi) + \mathcal{P}(u, v; u^{(j)}, v^{(j)}, \tau^{(j)}) \quad (6)$$

$$\mathcal{P}(u, v; u^{(j)}, v^{(j)}, \tau^{(j)}) := \frac{1}{2} \sum_{n=1}^N \sum_{\ell=2}^L \left\| \begin{pmatrix} v_{n, \ell-1} \\ u_{n, \ell} \end{pmatrix} - \begin{pmatrix} v_{n, \ell-1}^{(j)} \\ u_{n, \ell}^{(j)} \end{pmatrix} \right\|_{S_\ell^{(j)}}^2 + \frac{\tau_1}{2} \sum_{n=1}^N \|u_{n, 1} - u_{n, 1}^{(j)}\|^2$$

- $\tau_1 > 0$, $\tau^{(j)} := (\tau_2^{(j)}, \dots, \tau_L^{(j)})^T \in \mathbb{R}^{L-1}$
- $S_\ell^{(j)} := \tau_\ell^{(j)} I_{N_\ell + N_{\ell-1}} - \rho \begin{bmatrix} -W_\ell^{(j+1)} & I_{N_\ell} \end{bmatrix}^T \begin{bmatrix} -W_\ell^{(j+1)} & I_{N_\ell} \end{bmatrix} \geq \tau_1 I_{N_\ell + N_{\ell-1}}$
- $\tau_\ell^{(j)} := \rho \left\| \begin{bmatrix} -W_\ell^{(j+1)} & I_{N_\ell} \end{bmatrix} \right\|^2 + \tau_1$

¹proximal gradient method [Dai-Fletcher 2005]

Alternating Minimization Method for Solving (4) (Cont'd)

Alternating Minimization Algorithm for Solving (4)

- ➊ Input: $A, \xi, \rho > 0, (w^{(0)}, b^{(0)}, v^{(0)}, u^{(0)})$. let $\tau_1 > 0, J = 0$
- ➋ Update $(w^{(J+1)}, b^{(J+1)})$ by solving problem (6)
- ➌ Update $(u^{(J+1)}, v^{(J+1)})$ by solving problem (6)
- ➍ Set $J := J + 1$. If the stop criterion is not met, return to Step 2
- ➎ Output: $(w^{(J)}, b^{(J)}, v^{(J)}, u^{(J)})$

Theorem 4

Let $\{(w^{(j)}, b^{(j)}, v^{(j)}, u^{(j)})\}$ be the sequence generated by the Alternating Minimization Algorithm. Then any accumulation point (w^, b^*, v^*, u^*) of $\{(w^{(j)}, b^{(j)}, v^{(j)}, u^{(j)})\}$ is a KKT point of (4).*

Theorem 5

Let $\{(w^{(k)}, b^{(k)}, v^{(k)}, u^{(k)})\}$ be the sequence generated by IALAM with $\eta_3 > 1$. Then the following statements hold.

- (a) $\liminf_{k \rightarrow \infty} \|u^{(k)} - \Psi(v^{(k)})w^{(k)} - Ab^{(k)}\| = 0$ and the sequence $\{(w^{(k)}, b^{(k)}, v^{(k)}, u^{(k)})\}$ has at least one accumulation point.
- (b) $\liminf_{k \rightarrow \infty} \text{dist}((w^{(k)}, b^{(k)}, v^{(k)}, u^{(k)}), \mathcal{Z}^*) = 0$, where \mathcal{Z}^* is the set of KKT points of (PP).
- (c) If in addition that $\gamma = 1$, then $\lim_{k \rightarrow \infty} \|u^{(k)} - \Psi(v^{(k)})w^{(k)} - Ab^{(k)}\| = 0$. Furthermore, any accumulation point (w^*, b^*, v^*, u^*) of $\{(w^{(k)}, b^{(k)}, v^{(k)}, u^{(k)})\}$ is a KKT point of problem (PP).

- Theoretical Contribution: Different from the existing methods [Lu-Zhang 2012; Chen et al. 2017], we prove the existence of the accumulation point. Moreover, our designed algorithm support $\gamma > 1$
- Extensions: IALAM can be applied to many kind of networks

Theorem 5

Let $\{(w^{(k)}, b^{(k)}, v^{(k)}, u^{(k)})\}$ be the sequence generated by IALAM with $\eta_3 > 1$. Then the following statements hold.

- (a) $\liminf_{k \rightarrow \infty} \|u^{(k)} - \Psi(v^{(k)})w^{(k)} - Ab^{(k)}\| = 0$ and the sequence $\{(w^{(k)}, b^{(k)}, v^{(k)}, u^{(k)})\}$ has at least one accumulation point.
- (b) $\liminf_{k \rightarrow \infty} \text{dist}((w^{(k)}, b^{(k)}, v^{(k)}, u^{(k)}), \mathcal{Z}^*) = 0$, where \mathcal{Z}^* is the set of KKT points of (PP).
- (c) If in addition that $\gamma = 1$, then $\lim_{k \rightarrow \infty} \|u^{(k)} - \Psi(v^{(k)})w^{(k)} - Ab^{(k)}\| = 0$. Furthermore, any accumulation point (w^*, b^*, v^*, u^*) of $\{(w^{(k)}, b^{(k)}, v^{(k)}, u^{(k)})\}$ is a KKT point of problem (PP).

- **Theoretical Contribution:** Different from the existing methods [Lu-Zhang 2012; Chen et al. 2017], we prove the existence of **the accumulation point**. Moreover, our designed algorithm support $\gamma > 1$
- **Extensions:** IALAM can be applied to many kind of networks

4. Numerical Experiments

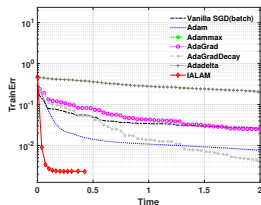
Default settings

- Stop criterion: $\rho^{(k)} > 10^3 \rho^{(0)}$
- Initialization: $W_\ell^{(0)} = \text{randn}(N_\ell, N_\ell - 1)/N$, $b^{(0)} = 0$, $u_{n,\ell}^{(0)} = W_\ell^{(0)} v_{n,\ell-1}^{(0)}$, $v_{n,\ell}^{(0)} = \sigma(u_{n,\ell}^{(0)})$
- Test problem: randomly generated synthetic dataset and MNIST dataset
- $N_{\text{test}} = \lceil N/5 \rceil$

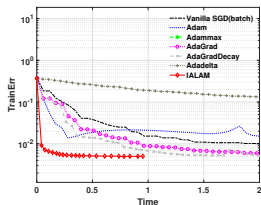
Output evaluation.

- **TrainErr** = $\frac{1}{N} \sum_{n=1}^N \|\sigma(W_L \sigma(\cdots \sigma(W_1 x_n + b_1) + b_2 \cdots) + b_L) - y_n\|^2$
- **TestErr** = $\frac{1}{N} \sum_{n=N+1}^{N+N_{\text{test}}} \|\sigma(W_L \sigma(\cdots \sigma(W_1 x_n + b_1) + b_2 \cdots) + b_L) - y_n\|^2$
- **FeasVi** = $\frac{1}{N} \sum_{n=1}^N \sum_{\ell=1}^L \|v_{n,\ell} - \sigma(u_{n,\ell})\|^2 + \frac{1}{N} \sum_{n=1}^N \sum_{\ell=1}^L \|u_{n,\ell} - (W_\ell v_{n,\ell-1} + b_\ell)\|^2$
- **Column Sparse Ratio**: the ratio of columns in all the matrices W_ℓ , totaling $\sum_{\ell=0}^{L-1} N_\ell$ columns, where the l_2 norm values are below the tolerance ϵ .
- **Accuracy** (for the training set) and **TestAcc**

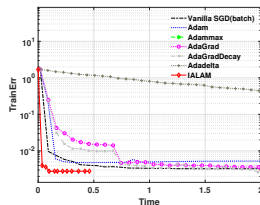
Comparisons Among IALAM and SGD-based Approaches on the Synthetic Dataset



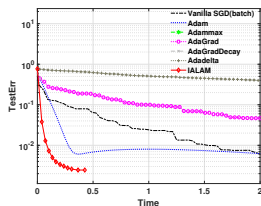
(a) $L = 2, N_1 = 10$



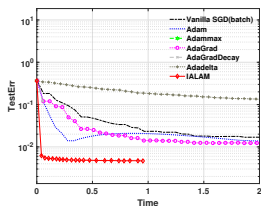
(b) $L = 3, N_1 = N_2 = 5$



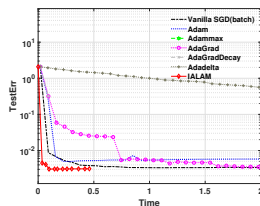
(c) $L = 4, N_1 = 4,$
 $N_2 = N_3 = 3$



(d) $L = 2, N_1 = 10$

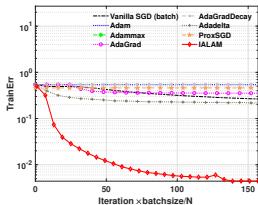


(e) $L = 3, N_1 = N_2 = 5$

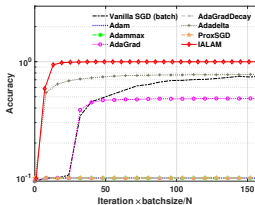


(f) $L = 4, N_1 = 4,$
 $N_2 = N_3 = 3$

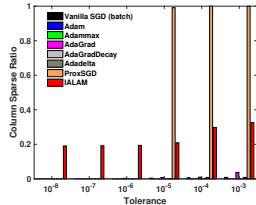
Comparisons Among IALAM and SGD-based Approaches on MNIST



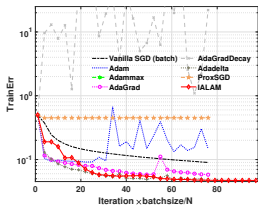
(a) TrainErr



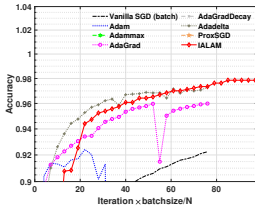
(b) Accuracy



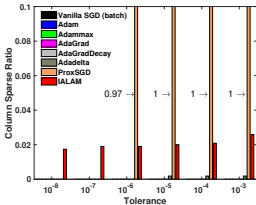
(c) Column Sparse Ratio



(d) TrainErr



(e) Accuracy



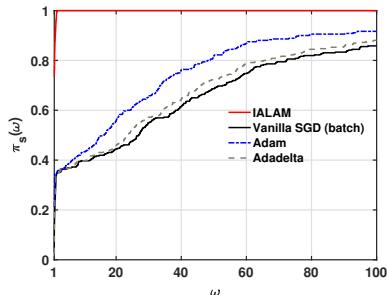
(f) Column Sparse Ratio

(a)–(c): $N = 1000$, $N_1 = 100$, $N_2 = 50$, $L = 3$

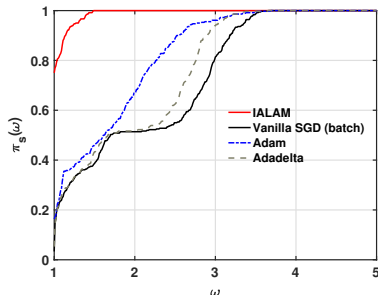
(d)–(f): $N = 60000$, $N_1 = 200$, $N_2 = 100$, $L = 3$

Overall Performance Profile on MNIST

[Dolan and More, 2002]



(a) TrainErr



(b) TestErr

FIG 1: Performance profile for IALAM, Valinna SGD, Adadelata and Adam on TrainErr and TestErr.

We select 720 test problems based on MNIST data set with different network parameter combinations

For the optimization problem (P) toward a nonconvex nonsmooth neuron network

- design a new model (PP) with bounded solution set
- present **exact penalization**:
 - local minimizers
 - l-stationary points
- design a **IALAM algorithm** for nonconvex problems.
 - converges to a **KKT point/ l-stationary point** of (PP)
 - the feasible violation decreases non-monotonic
 - the existence of approximation point is not required

Thanks for your listening!

Email: liuwei175@lsec.cc.ac.cn