



# ReHLine: Regularized Composite ReLU-ReHU Loss Minimization with Linear Computation and Linear Convergence

Ben Dai\* (CUHK)

Yixuan Qiu\* (SUFE)



## Introduction

In this paper, we consider a general regularized ERM based on a convex (but possible **nonsmooth**) PLQ loss with linear **constraints**:

$$\min_{\beta \in \mathbb{R}^d} \sum_{i=1}^n L_i(\mathbf{x}_i^\top \beta) + \frac{1}{2} \|\beta\|_2^2, \quad \text{s.t. } \mathbf{A}\beta + \mathbf{b} \geq \mathbf{0}, \quad (1)$$

where  $\mathbf{x}_i \in \mathbb{R}^d$  is the covariate vector for the  $i$ -th observation,  $\beta \in \mathbb{R}^d$  is an unknown coefficient vector,  $\mathbf{A} \in \mathbb{R}^{K \times d}$  and  $\mathbf{b} \in \mathbb{R}^K$  are posed as linear inequality constraints for  $\beta$ , and  $L_i(\cdot) \geq 0$  is a convex piecewise linear-quadratic loss (PLQ) loss function.

**Table 1.** Overview of existing algorithms in solving (1).

ALGORITHM	COMPLEXITY (PER ITERATION)	#ITERATION	COMPLEXITY (TOTAL)
P-GD	$\mathcal{O}(n)$	$\mathcal{O}(\varepsilon^{-1})$ [6]	$\mathcal{O}(n\varepsilon^{-1})$
CD	$\mathcal{O}(n^2)$	$\mathcal{O}(\log(\varepsilon^{-1}))$ [31]	$\mathcal{O}(n^2 \log(\varepsilon^{-1}))$
IPM	$\mathcal{O}(n^2)$	$\mathcal{O}(\log(\varepsilon^{-1}))$ [18]	$\mathcal{O}(n^2 \log(\varepsilon^{-1}))$
ADMM	$\mathcal{O}(n^2)$	$o(\varepsilon^{-1})$ [9, 20]	$o(n^2 \varepsilon^{-1})$
SDCA	$\mathcal{O}(n)$	$\mathcal{O}(\varepsilon^{-1})$ [39]	$\mathcal{O}(n\varepsilon^{-1})$
ReHLine (ours)	$\mathcal{O}(n)$	$\mathcal{O}(\log(\varepsilon^{-1}))$	$\mathcal{O}(n \log(\varepsilon^{-1}))$

**Contribution.** Compared with existing algorithms, the proposed ReHLine solver has four appealing "*linear properties*":

- It applies to any convex piecewise linear-quadratic loss function (potential for non-smoothness included).
- In addition, it supports linear constraints on the parameter vector.
- The optimization algorithm has a provable linear convergence rate.
- The per-iteration computation is linear in the sample size.

## The ReHLine decomposition

**Definition 1.** A function  $L(z)$  is composite ReLU-ReHU, if there exist  $\mathbf{u}, \mathbf{v} \in \mathbb{R}^L$  and  $\tau, \mathbf{s}, \mathbf{t} \in \mathbb{R}^H$  such that

$$L(z) = \sum_{l=1}^L \text{ReLU}(u_l z + v_l) + \sum_{h=1}^H \text{ReHU}_{\tau_h}(s_h z + t_h), \quad (3)$$

where  $\text{ReLU}(z) = z_+$ , and  $\text{ReHU}_{\tau_h}(z)$  defined in (2).

**Theorem 1.** A loss function  $L$  is convex PLQ iff it is composite ReLU-ReHU.

Table 2. Some widely used composite ReHLine losses as in (3).

PROBLEM	LOSS ( $L_i(z_i)$ )	COMPOSITE ReLU-ReHU PARAMETERS
SVM	$c_i(1 - y_i z_i)_+$	$u_{1i} = -c_i y_i, v_{1i} = c_i$
s-SVM	$c_i \text{ReHU}_1(-(y_i z_i - 1))$	$s_{1i} = -\sqrt{c_i} y_i, t_{1i} = \sqrt{c_i}, \tau = \sqrt{c_i}$
SVM <sup>2</sup>	$c_i((1 - y_i z_i)_+)^2$	$s_{1i} = -\sqrt{2c_i} y_i, t_{1i} = \sqrt{2c_i}, \tau = \infty$
LAD	$c_i  y_i - z_i $	$u_{1i} = c_i, v_{1i} = -c_i y_i, u_{2i} = -c_i, v_{2i} = c_i y_i$
SVR	$c_i( y_i - z_i  - \varepsilon)_+$	$u_{1i} = c_i, v_{1i} = -(y_i + \varepsilon), u_{2i} = -c_i, v_{2i} = y_i - \varepsilon$
QR	$c_i \rho_\kappa(y_i - z_i)$	$u_{1i} = -c_i \kappa, v_{1i} = \kappa c_i y_i, u_{2i} = c_i(1 - \kappa), v_{2i} = -c_i(1 - \kappa) y_i$

Taken together, (1) can be rewritten as:

$$\min_{\beta \in \mathbb{R}^d} \sum_{i=1}^n \sum_{l=1}^L \text{ReLU}(u_{li} \mathbf{x}_i^\top \beta + v_{li}) + \sum_{i=1}^n \sum_{h=1}^H \text{ReHU}_{\tau_{hi}}(s_{hi} \mathbf{x}_i^\top \beta + t_{hi}) + \frac{1}{2} \|\beta\|_2^2, \quad \text{s.t. } \mathbf{A}\beta + \mathbf{b} \geq \mathbf{0}, \quad (4)$$

where  $\mathbf{U} = (u_{li}), \mathbf{V} = (v_{li}) \in \mathbb{R}^{L \times n}$  and  $\mathbf{S} = (s_{hi}), \mathbf{T} = (t_{hi}), \mathbf{T} = (\tau_{hi}) \in \mathbb{R}^{H \times n}$  are the ReLU-ReHU loss parameters, as illustrated in Table 2.

The **Lagrangian dual**, which is a box-QP, is presented in **Theorem 2** of our paper, which is derived using the Karush-Kuhn-Tucker (KKT) condition:

$$\hat{\beta} = \sum_{k=1}^K \hat{\xi}_k \mathbf{a}_k - \sum_{i=1}^n \mathbf{x}_i \left( \sum_{l=1}^L \hat{\lambda}_{li} u_{li} + \sum_{h=1}^H \hat{\gamma}_{hi} s_{hi} \right) = \mathbf{A}^\top \hat{\xi} - \bar{\mathbf{U}}_{(3)} \text{vec}(\hat{\Lambda}) - \bar{\mathbf{S}}_{(3)} \text{vec}(\hat{\Gamma}).$$

## Algorithm and Results

This proposed **ReHLine** is based on the coordinate descent (CD), drawing inspiration from *Liblinear*. Its motivation is to utilize the **linear structure** in the KKT conditions, and simultaneously update primal and dual variables, considerably reducing the computational complexity for CD updates. For illustration, we only demo one dual variable, see the full details in the paper.

**Canonical CD updates.** By excluding the terms unrelated to  $\lambda_{li}$ :

$$\lambda_{li}^{\text{new}} = \mathcal{P}_{[0,1]} \left( \frac{u_{li} \mathbf{x}_i^\top \left( \sum_{k=1}^K \xi_k \mathbf{a}_k - \sum_{(l',i') \neq (l,i)} \lambda_{l'i'} u_{l'i'} \mathbf{x}_{i'} - \sum_{h',i'} \gamma_{h'i'} s_{h'i'} \mathbf{x}_{i'} \right) + v_{li}}{u_{li}^2 \|\mathbf{x}_i\|_2^2} \right), \quad (10)$$

updating one  $\lambda_{li}$  value requires  $\mathcal{O}(K + nd + nL + nH)$  of computation. Adding all variables together, the canonical CD update rule for one full cycle has a computational complexity of  $\mathcal{O}((K + nd + nL + nH)(K + nL + nH))$ .

**ReHLine updates** significantly reduces the computational complexity of canonical CD by updating  $\beta$  according to the KKT condition (9) after each update of a dual variable.

$$\lambda_{li}^{\text{new}} = \mathcal{P}_{[0,1]} \left( \lambda_{li}^{\text{old}} - \frac{\nabla_{\lambda_{li}} \mathcal{L}(\lambda^{\text{old}})}{u_{li}^2 \|\mathbf{x}_i\|_2^2} \right) = \mathcal{P}_{[0,1]} \left( \lambda_{li}^{\text{old}} + \frac{u_{li} \mathbf{x}_i^\top \beta^{\text{old}} + v_{li}}{u_{li}^2 \|\mathbf{x}_i\|_2^2} \right).$$

Accordingly, the primal variable  $\beta$  is updated as

$$\beta^{\text{new}} = \beta^{\text{old}} - (\lambda_{li}^{\text{new}} - \lambda_{li}^{\text{old}}) u_{li} \mathbf{x}_i,$$

updating one  $\lambda_{li}$  value only requires  $\mathcal{O}(d)$  of computation. Adding all variables together, the ReHLine update rule for one full cycle has a computational complexity of  $\mathcal{O}((K + nL + nH)d)$ .

**Theorem 2.** Let  $\mu^{(q)}$  be a sequence of iterates generated by ReHLine. Then the dual objective **converges at least linearly** to that of  $\mu^*$ .

**Table 5.** The running times of SOTA solvers on ML tasks using the *Benchopt*. “**X**” indicates cases where the solver produced an invalid solution or exceeded the allotted time limit (“objective” for failure on objective function, and “both” for both objective and feasibility). **Speed-up** refers to the speed-up in running time achieved by **ReHLine**.

TASK	DATASET	ECOS	MOSEK	SCS	DCCP	REHLINE
FairSVM	SPF ( $\times 1e-4$ )	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	4.25( $\pm 0.5$ )
	Philippine ( $\times 1e-2$ )	1550( $\pm 0.6$ )	87.4( $\pm 0.2$ )	130( $\pm 42$ )	1137( $\pm 9.2$ )	1.03( $\pm 0.2$ )
	Sylva-prior ( $\times 1e-2$ )	<b>X</b>	<b>X</b>	<b>X</b>	<b>X</b>	0.47( $\pm 0.1$ )
	Creditcard ( $\times 1e-1$ )	175( $\pm 0.2$ )	64.2( $\pm 0.1$ )	161( $\pm 405$ )	<b>X</b>	0.64( $\pm 0.2$ )
	Fail/Succeed	2/2	2/2	2/2	3/1	0/4
	Speed-up (on Creditcard)	273x	100x	252x	$\infty$	—

TASK	DATASET	ECOS	MOSEK	SCS	REHLINE
ElasticQR	LD ( $\times 1e-4$ )	<b>X</b>	106( $\pm 7$ )	34.9( $\pm 25.0$ )	2.60( $\pm 0.30$ )
	Kin8nm ( $\times 1e-3$ )	<b>X</b>	92.0( $\pm 1.0$ )	63.1( $\pm 58.5$ )	4.12( $\pm 0.95$ )
	House-8L ( $\times 1e-3$ )	887( $\pm 161$ )	277( $\pm 34$ )	<b>X</b>	7.21( $\pm 1.99$ )
	Topo-2-1 ( $\times 1e-2$ )	4752( $\pm 2015$ )	<b>X</b>	<b>X</b>	3.04( $\pm 0.49$ )
	BT ( $\times 1e-0$ )	7079( $\pm 2517$ )	<b>X</b>	<b>X</b>	2.49( $\pm 0.56$ )
	Fail/Succeed	3/2	2/3	3/2	0/5
	Speed-up (on BT)	2843x	$\infty$	$\infty$	—

TASK	DATASET	ECOS	MOSEK	SCS	HQREG	REHLINE
RidgeHuber	Liver-disorders ( $\times 1e-4$ )	<b>X</b>	<b>X</b>	<b>X</b>	4.90( $\pm 0.00$ )	1.40( $\pm 0.20$ )
	Kin8nm ( $\times 1e-3$ )	<b>X</b>	<b>X</b>	<b>X</b>	1.58( $\pm 0.21$ )	2.04( $\pm 0.30$ )
	House-8L ( $\times 1e-3$ )	<b>X</b>	925( $\pm 2$ )	<b>X</b>	2.42( $\pm 0.34$ )	0.80( $\pm 0.21$ )
	Topo-2-1 ( $\times 1e-2$ )	2620( $\pm 1040$ )	267( $\pm 1$ )	213( $\pm 2$ )	3.53( $\pm 0.67$ )	1.78( $\pm 0.32$ )
	BT ( $\times 1e-1$ )	<b>X</b>	2384( $\pm 433$ )	<b>X</b>	12.5( $\pm 1.8$ )	5.28( $\pm 1.31$ )
	Fail/Succeed	4/1	2/3	4/1	0/5	0/5
	Speed-up (on BT)	$\infty$	452x	$\infty$	2.37x	—

TASK	DATASET	ECOS	MOSEK	SCS	LIBLINEAR	REHLINE
SVM	SPF ( $\times 1e-4$ )	<b>X</b>	372( $\pm 1$ )	237( $\pm 27$ )	12.7( $\pm 0.1$ )	3.90( $\pm 0.10$ )
	Philippine ( $\times 1e-2$ )	1653( $\pm 41$ )	86.5( $\pm 0.2$ )	153( $\pm 146$ )	1.80( $\pm 0.02$ )	0.82( $\pm 0.02$ )
	Sylva-prior ( $\times 1e-3$ )	<b>X</b>	731( $\pm 2$ )	843( $\pm 1006$ )	16.0( $\pm 0.6$ )	4.08( $\pm 0.84$ )
	Creditcard ( $\times 1e-2$ )	2111( $\pm 804$ )	<b>X</b>	1731( $\pm 4510$ )	23.1( $\pm 2.5$ )	5.08( $\pm 1.45$ )
	Fail/Succeed	2/2	1/3	0/4	0/4	0/4
	Speed-up (on Creditcard)	415x	$\infty$	340x	4.5x	—

