An Improved Convergence Analysis of Cyclic Block Coordinate Descent-type Methods for Strongly Convex Minimization

Xingguo Li^{*}, Tuo Zhao[◇], Raman Arora[◇], Han Liu[†], Mingyi Hong[‡] *University of Minnesota *Johns Hopkins University [†]Princeton University [‡]Iowa State University

Background

A class of strongly convex minimization problems:

P1: $\min_{x \in \mathbb{R}^d} \mathcal{L}(x) + \mathcal{R}(x),$

- A partition of p blocks: $x = [x_1, ..., x_p]^\top$;
- $\mathcal{L}(x)$: differentiable and convex;
- $\mathcal{R}(x) = \sum_{j=1}^{p} \mathcal{R}_{j}(x_{j})$: strongly convex and possibly nonsmooth for each $\mathcal{R}_i(\cdot)$.

Popular examples:

- Elastic-net Penalized Regression:
 - $\min \frac{1}{2} ||b Ax||^2 + \lambda_1 ||x||^2 + \lambda_2 ||x||_1;$

Main Results

Table 1. Comparison between Our Results and Beck & Tetruashvili [1].

	Method	$\mathcal{L}(\cdot)$	$\mathcal{R}(\cdot)$	Improved Iteration Complexity	Beck & Tetruashvili [1]
[a]	CBCGD	Quadratic	Smooth	$\mathcal{O}\left(\mu^{-1} \log^2 p L^2 \log(1/\epsilon)\right)$	$\mathcal{O}\left(\mu^{-1}pL^2\log(1/\epsilon)\right)$
[b]	CBCGD	Quadratic	Nonsmooth	$\mathcal{O}\left(\mu^{-1}\log^2 pL^2\log(1/\epsilon)\right)$	N/A
[c]	CBCGD	General Convex	Smooth	$\mathcal{O}\left(\mu^{-1}p\cdot\min\{L^2,p\}\log(1/\epsilon)\right)$	$\mathcal{O}\left(\mu^{-1}pL^2\log(1/\epsilon)\right)$
[d]	CBCGD	General Convex	Nonsmooth	$\mathcal{O}\left(\mu^{-1}pL^2\log(1/\epsilon)\right)$	N/A
[e]	CBCM	Quadratic	Smooth	$\mathcal{O}\left(\mu^{-1}\log^2 pL^2\log(1/\epsilon)\right)$	N/A
[f]	CBCM	Quadratic	Nonsmooth	$\mathcal{O}\left(\mu^{-1}\log^2 pL^2\log(1/\epsilon)\right)$	N/A
[g]	CBCM	General Convex	Smooth	$\mathcal{O}\left(\mu^{-1}pL^2\log(1/\epsilon)\right)$	N/A
[h]	CBCM	General Convex	Nonsmooth	$\mathcal{O}\left(\mu^{-1}pL^2\log(1/\epsilon)\right)$	N/A

• Ridge Penalized Logistic Regression:

$$\min_{x} \frac{1}{n} \sum_{i=1}^{n} \log(1 + \exp[-b_{i}x \mathcal{A}_{i*}]) + \lambda ||x||^{2};$$
• Support Vector Machines:

$$\min_{x} \frac{1}{2n} \sum_{i=1}^{n} \max\{1 - b_{i}(x^{\top}A_{i*} - y), 0\}^{2} + \lambda ||x||^{2}.$$

Popular Solvers for model P1:

- Gradient decent methods;
- Alternating direction method of multipliers;
- Cyclic block coordinate descent methods.

Algorithms of Our Interests for model P1:

- (1) Cyclic Block Coordinate Minimization;
- (2) Cyclic Block Coordinate Gradient Descent.

Algorithm

Define two auxiliary variables:

Our main contributions:

Develop the iteration bounds of CBCM and CBCGD for different specifications on $\mathcal{L}(\cdot)$ and $\mathcal{R}(\cdot)$; (1)

Significantly improve the dependence on dimension p of the iteration bound of CBCGD for quadratic $\mathcal{L}(\cdot)$;

Improve the iteration bound of CBCGD for smooth $\mathcal{R}(\cdot)$. (3)

Proof Sketch for Quadratic $\mathcal{L}(\cdot)$ via CBCGD

Let $\mathcal{F}(x) = \mathcal{L}(x) + \mathcal{R}(x)$, $L_{\min}^{\mu} = \min_{j} L_{j} + \mu_{j}$, $x^{*} = \operatorname{argmin}_{x} \mathcal{F}(x)$, $d_{\max} = \max_{j} \{d_{j} : x_{j} \in \mathbb{R}^{d_{j}}\}$, ϵ be a pre-specified accuracy of the objective value:

(1) Characterize the successive descent after each CBCGD iteration:

 $\mathcal{F}(x^{(t)}) - \mathcal{F}(x^{(t+1)}) \ge \frac{L_{\min}^{\mu}}{2} ||x^{(t)} - x^{(t+1)}||^2;$

(2) Characterize the gap towards the optimal objective value after each CBCGD iteration:

$$\mathcal{F}(x^{(t+1)}) - \mathcal{F}(x^*) \le \frac{L^2 \log^2(2p \cdot d_{\max})}{2\mu} ||x^{(t+1)} - x^{(t)}||^2$$

A new proof technique: Symmetrification.

(3) Combine (1) and (2): To guarantee $\mathcal{F}(x^{(t)}) - \mathcal{F}(x^*) \leq \epsilon$,

$$\begin{aligned} x_{1:(j-1)}^{(t+1)} &= [x_1^{(t+1)\top}, \dots, x_{j-1}^{(t+1)\top}]^\top, \\ x_{(j+1):p}^{(t)} &= [x_{j+1}^{(t)\top}, \dots, x_p^{(t)\top}]^\top. \end{aligned}$$

At t + 1-th iteration, we take (1) CBCM: For all j = 1, ..., p, $x_{j}^{(t+1)} = \underset{x_{j}}{\operatorname{argmin}} \mathcal{L}\left(x_{1:(j-1)}^{(t+1)}, x_{j}, x_{(j+1):p}^{(t)}\right) + \mathcal{R}_{j}(x_{j});$

(2) CBCGD: For all
$$j = 1, ..., p$$
,
 $x_j^{(t+1)} = \operatorname*{argmin}_{x_j} \nabla_j \mathcal{L} \Big(x_{1:(j-1)}^{(t+1)}, x_j, x_{(j+1):p}^{(t)} \Big)^{\top} (x_j - x_j^{(t)})$
 $+ \frac{\eta_j}{2} ||x_j - x_j^{(t)}||^2 + \mathcal{R}_j(x_j),$

 $\eta_i > 0$: the step-size parameter for updating x_i .

Assumptions

Assumption 1. $\nabla \mathcal{L}(\cdot)$ is **Lipschitz continuous** and **blockwise Lipschitz continuous**, i.e., $\exists L$ and L_i 's such that $\forall x, x' \in \mathbb{R}^d$ and $\forall j = 1, \dots, p$, we have

 $||\nabla \mathcal{L}(x') - \nabla \mathcal{L}(x)|| \le L||x - x'||,$



Tightness of the Iteration Complexity for Quadratic $\mathcal{L}(\cdot)$

Consider: $\min_x \mathcal{H}(x) := ||Bx||^2$, $x^* = \operatorname{argmin} \mathcal{H}(x) = [0, 0, \dots, 0]^\top$. Let



 $||\nabla_j \mathcal{L}\left(x_{1:(j-1)}, x'_j, x_{(j+1):p}\right) - \nabla_j \mathcal{L}(x)|| \le L_j ||x_j - x'_j||.$

Assumption 2. $\mathcal{R}(\cdot)$ is strongly convex and blockwise **strongly convex**, i.e., $\exists \mu$ and μ_i 's such that $\forall x, x' \in \mathbb{R}^d$ and $\forall j = 1, \dots, p$, we have

 $\mathcal{R}(x) \ge \mathcal{R}(x') + (x - x')^{\top} \xi' + \frac{\mu}{2} ||x - x'||^2,$ $\mathcal{R}_j(x_j) \ge \mathcal{R}_j(x'_j) + (x_j - x'_j)^\top \xi'_j + \frac{\mu_j}{2} ||x_j - x'_j||^2.$

Lower Bound

Iteration Complexity is independent of p and cannot be further improved when λ_{max} and λ_{min} of Hessian do not scale with p.

References

[1] A. Beck and L. Tetruashvili. On the convergence of block coordinate descent type methods. SIAM Journal on Optimization, 23(4):2037–2060, 2013.







