

Convergence of a Block Coordinate Descent Method for Nondifferentiable Minimization

Paul Tseng

Presenter: Lei Tang

Department of CSE
Arizona State University

Nov. 7th, 2008

- Popular method for minimizing a real-valued continuously differentiable function f of n variables, subject to bound constraint, is (block) coordinate descent (BCD).
- In this work, coordinate descent actually refers to alternating optimization(AO). Each step find the exact minimizer.
- Popular for its efficiency, simplicity and scalability.
- Applied to large-scale SVM, Lasso etc.
- Unfortunately, the convergence of coordinate descent is not clear. Not like steepest descent method.
- In this work, it is shown that if the function satisfy some mild conditions, BCD converges to the stationary point.

- Popular method for minimizing a real-valued continuously differentiable function f of n variables, subject to bound constraint, is (block) coordinate descent (BCD).
- In this work, coordinate descent actually refers to alternating optimization(AO). Each step find the exact minimizer.
- Popular for its efficiency, simplicity and scalability.
- Applied to large-scale SVM, Lasso etc.
- Unfortunately, the convergence of coordinate descent is not clear. Not like steepest descent method.
- In this work, it is shown that if the function satisfy some mild conditions, BCD converges to the stationary point.

Questions?

- ❶ Does BCD Converge?
- ❷ Does BCD Converge to the local minimizer?
- ❸ When does BCD converge to the stationary point?
- ❹ What's the convergence rate?

Existing works

- Convergence of coordinate descent method requires typically that f be strictly convex (or quasiconvex and hemivariate) differentiable
- the strict convexity is relaxed to pseudoconvexity, which allows f to have non-unique minimum along coordinate directions.
- If f is not differentiable, the coordinate descent method may get stuck at a nonstationary point even when f is convex.
- However, this method still works when the nondifferentiable part of f is separable.

$$f(x_1, \dots, x_N) = f_0(x_1, \dots, x_N) + \sum_{k=1}^N f_k(x_k)$$

where f_k is non-differentiable, each x_k represents one block.

- This work shows that BCD converges to a stationary point if f_0 has certain smoothness property.

- Convergence of coordinate descent method requires typically that f be strictly convex (or quasiconvex and hemivariate) differentiable
- the strict convexity is relaxed to pseudoconvexity, which allows f to have non-unique minimum along coordinate directions.
- If f is not differentiable, the coordinate descent method may get stuck at a nonstationary point even when f is convex.
- However, this method still works **when the nondifferentiable part of f is seperable**.

$$f(x_1, \dots, x_N) = f_0(x_1, \dots, x_N) + \sum_{k=1}^N f_k(x_k)$$

where f_k is non-differentiable, each x_k represents one block.

- This work shows that BCD converges to a stationary point if f_0 has certain smoothness property.

An Example of Alternating Optimization

$$\phi_1(x, y, z) = -xy - yz - zx + (x - 1)_+^2 + (-x - 1)_+^2 + (y - 1)_+^2 + (-y - 1)_+^2 + (z - 1)_+^2 + (-z - 1)_+^2$$

Note that the optimal x given fixed y and z is

$$x = \text{sign}(y + z) \left(1 + \frac{1}{2}|y + z| \right)$$

Suppose you start from $(-1 - \epsilon, 1 + \frac{1}{2}\epsilon, -1 - \frac{1}{4}\epsilon)$:

$$\begin{array}{ll} (1 + \frac{1}{8}\epsilon, 1 + \frac{1}{2}\epsilon, -1 - \frac{1}{4}\epsilon) & (-1 - \frac{1}{64}\epsilon, -1 - \frac{1}{16}\epsilon, 1 + \frac{1}{32}\epsilon) \\ (1 + \frac{1}{8}\epsilon, -1 - \frac{1}{16}\epsilon, -1 - \frac{1}{4}\epsilon) & (-1 - \frac{1}{64}\epsilon, 1 + \frac{1}{128}\epsilon, 1 + \frac{1}{32}\epsilon) \\ (1 + \frac{1}{8}\epsilon, -1 - \frac{1}{16}\epsilon, 1 + \frac{1}{32}\epsilon) & (-1 - \frac{1}{64}\epsilon, 1 + \frac{1}{128}\epsilon, -1 - \frac{1}{256}\epsilon) \end{array}$$

Cycle around 6 edges of the cube $(\pm 1, \pm 1, \pm 1)!!$

An Example of Alternating Optimization

$$\phi_1(x, y, z) = -xy - yz - zx + (x - 1)_+^2 + (-x - 1)_+^2 + (y - 1)_+^2 + (-y - 1)_+^2 + (z - 1)_+^2 + (-z - 1)_+^2$$

Note that the optimal x given fixed y and z is

$$x = \text{sign}(y + z) \left(1 + \frac{1}{2}|y + z| \right)$$

Suppose you start from $(-1 - \epsilon, 1 + \frac{1}{2}\epsilon, -1 - \frac{1}{4}\epsilon)$:

$$\begin{array}{ll} (1 + \frac{1}{8}\epsilon, 1 + \frac{1}{2}\epsilon, -1 - \frac{1}{4}\epsilon) & (-1 - \frac{1}{64}\epsilon, -1 - \frac{1}{16}\epsilon, 1 + \frac{1}{32}\epsilon) \\ (1 + \frac{1}{8}\epsilon, -1 - \frac{1}{16}\epsilon, -1 - \frac{1}{4}\epsilon) & (-1 - \frac{1}{64}\epsilon, 1 + \frac{1}{128}\epsilon, 1 + \frac{1}{32}\epsilon) \\ (1 + \frac{1}{8}\epsilon, -1 - \frac{1}{16}\epsilon, 1 + \frac{1}{32}\epsilon) & (-1 - \frac{1}{64}\epsilon, 1 + \frac{1}{128}\epsilon, -1 - \frac{1}{256}\epsilon) \end{array}$$

Cycle around 6 edges of the cube $(\pm 1, \pm 1, \pm 1)!!$

An Example of Alternating Optimization

$$\phi_1(x, y, z) = -xy - yz - zx + (x - 1)_+^2 + (-x - 1)_+^2 + (y - 1)_+^2 + (-y - 1)_+^2 + (z - 1)_+^2 + (-z - 1)_+^2$$

Note that the optimal x given fixed y and z is

$$x = \text{sign}(y + z) \left(1 + \frac{1}{2}|y + z| \right)$$

Suppose you start from $(-1 - \epsilon, 1 + \frac{1}{2}\epsilon, -1 - \frac{1}{4}\epsilon)$:

$$\begin{array}{ll} (1 + \frac{1}{8}\epsilon, 1 + \frac{1}{2}\epsilon, -1 - \frac{1}{4}\epsilon) & (-1 - \frac{1}{64}\epsilon, -1 - \frac{1}{16}\epsilon, 1 + \frac{1}{32}\epsilon) \\ (1 + \frac{1}{8}\epsilon, -1 - \frac{1}{16}\epsilon, -1 - \frac{1}{4}\epsilon) & (-1 - \frac{1}{64}\epsilon, 1 + \frac{1}{128}\epsilon, 1 + \frac{1}{32}\epsilon) \\ (1 + \frac{1}{8}\epsilon, -1 - \frac{1}{16}\epsilon, 1 + \frac{1}{32}\epsilon) & (-1 - \frac{1}{64}\epsilon, 1 + \frac{1}{128}\epsilon, -1 - \frac{1}{256}\epsilon) \end{array}$$

Cycle around 6 edges of the cube $(\pm 1, \pm 1, \pm 1)!!$

An Example of Alternating Optimization

$$\phi_1(x, y, z) = -xy - yz - zx + (x - 1)_+^2 + (-x - 1)_+^2 + (y - 1)_+^2 + (-y - 1)_+^2 + (z - 1)_+^2 + (-z - 1)_+^2$$

Note that the optimal x given fixed y and z is

$$x = \text{sign}(y + z) \left(1 + \frac{1}{2}|y + z| \right)$$

Suppose you start from $(-1 - \epsilon, 1 + \frac{1}{2}\epsilon, -1 - \frac{1}{4}\epsilon)$:

$$\begin{array}{ll} (1 + \frac{1}{8}\epsilon, 1 + \frac{1}{2}\epsilon, -1 - \frac{1}{4}\epsilon) & (-1 - \frac{1}{64}\epsilon, -1 - \frac{1}{16}\epsilon, 1 + \frac{1}{32}\epsilon) \\ (1 + \frac{1}{8}\epsilon, -1 - \frac{1}{16}\epsilon, -1 - \frac{1}{4}\epsilon) & (-1 - \frac{1}{64}\epsilon, 1 + \frac{1}{128}\epsilon, 1 + \frac{1}{32}\epsilon) \\ (1 + \frac{1}{8}\epsilon, -1 - \frac{1}{16}\epsilon, 1 + \frac{1}{32}\epsilon) & (-1 - \frac{1}{64}\epsilon, 1 + \frac{1}{128}\epsilon, -1 - \frac{1}{256}\epsilon) \end{array}$$

Cycle around 6 edges of the cube $(\pm 1, \pm 1, \pm 1)!!$

An Example of Alternating Optimization

$$\phi_1(x, y, z) = -xy - yz - zx + (x - 1)_+^2 + (-x - 1)_+^2 + (y - 1)_+^2 + (-y - 1)_+^2 + (z - 1)_+^2 + (-z - 1)_+^2$$

Note that the optimal x given fixed y and z is

$$x = \text{sign}(y + z) \left(1 + \frac{1}{2}|y + z| \right)$$

Suppose you start from $(-1 - \epsilon, 1 + \frac{1}{2}\epsilon, -1 - \frac{1}{4}\epsilon)$:

$$\begin{array}{ll} (1 + \frac{1}{8}\epsilon, 1 + \frac{1}{2}\epsilon, -1 - \frac{1}{4}\epsilon) & (-1 - \frac{1}{64}\epsilon, -1 - \frac{1}{16}\epsilon, 1 + \frac{1}{32}\epsilon) \\ (1 + \frac{1}{8}\epsilon, -1 - \frac{1}{16}\epsilon, -1 - \frac{1}{4}\epsilon) & (-1 - \frac{1}{64}\epsilon, 1 + \frac{1}{128}\epsilon, 1 + \frac{1}{32}\epsilon) \\ (1 + \frac{1}{8}\epsilon, -1 - \frac{1}{16}\epsilon, 1 + \frac{1}{32}\epsilon) & (-1 - \frac{1}{64}\epsilon, 1 + \frac{1}{128}\epsilon, -1 - \frac{1}{256}\epsilon) \end{array}$$

Cycle around 6 edges of the cube $(\pm 1, \pm 1, \pm 1)!!$

Some Examples

- The gradient in the example is not zero at any $(\pm 1, \pm 1, \pm 1)$.
- The example we show is unstable to perturbations.
- The example has non-smooth 2nd derivatives.
- More complicated examples could be constructed to show that even if the function is infinitely differentiable, stable cyclic behavior still occurs, whose gradient is bounded away from zero in the limiting path.
- Please see *On Search Directions for Minimization Algorithms*, Mathematical Programming, 1974.

Some Examples

- The gradient in the example is not zero at any $(\pm 1, \pm 1, \pm 1)$.
- The example we show is unstable to perturbations.
- The example has non-smooth 2nd derivatives.
- More complicated examples could be constructed to show that even if the function is infinitely differentiable, stable cyclic behavior still occurs, whose gradient is bounded away from zero in the limiting path.
- Please see *On Search Directions for Minimization Algorithms*, Mathematical Programming, 1974.

Some Examples

- The gradient in the example is not zero at any $(\pm 1, \pm 1, \pm 1)$.
- The example we show is unstable to perturbations.
- The example has non-smooth 2nd derivatives.
- More complicated examples could be constructed to show that **even if the function is infinitely differentiable, stable cyclic behavior still occurs, whose gradient is bounded away from zero in the limiting path.**
- Please see *On Search Directions for Minimization Algorithms*, Mathematical Programming, 1974.

Alternating Optimization Algorithm

AO-1 Let $\Psi_i \subseteq \mathfrak{R}^{p_i}$, for $i = 1, \dots, t$, and let $\Psi = \Psi_1 \times \dots \times \Psi_t$. Partition $x \in \mathfrak{R}^s$ as

$$x = (X_1, X_2, \dots, X_t)^T, \quad \text{with} \quad X_i \in \mathfrak{R}^{p_i} \quad \text{for} \quad i = 1, \dots, t,$$

$\bigcup_{i=1}^t X_i = X$; $X_i \cap X_j = \emptyset$ for $i \neq j$; and $s = \sum_{i=1}^t p_i$. Pick an initial iterate $x^{(0)} = (X_1^{(0)}, X_2^{(0)}, \dots, X_t^{(0)})^T \in \Psi = \Psi_1 \times \dots \times \Psi_t$, a vector norm $\|\cdot\|$, termination threshold ε , and iteration limit L . Set $r = 0$.

AO-2 For $i = 1, \dots, t$, compute the restricted minimizer

$$X_i^{(r+1)} = \arg \min_{X_i \in \Psi_i \subset \mathfrak{R}^{p_i}} \left\{ f(\mathbf{x}_1^{(r+1)}, \dots, \mathbf{x}_{t-1}^{(r+1)}, X_i, \mathbf{x}_{t+1}^{(r)}, \dots, \mathbf{x}_t^{(r)}) \right\} \quad (1)$$

AO-3 If $\|x^{(r+1)} - x^{(r)}\| \leq \varepsilon$ or $r > L$, then quit; otherwise, set $r = r+1$ and go to AO-2.

Figure: Alternating Optimization Algorithm

Before we go into the proof details, I would like to introduce some convergence properties of AO that might be useful.

Typically, we have this EU assumption:

Existence and Uniqueness (EU) Assumption. Let $\Psi_i \subseteq \mathbb{R}^{P_i}$, $i = 1, \dots, t$; and let $\Psi = \Psi_1 \times \dots \times \Psi_t$. Partition $x = (X_1, \dots, X_t)^T$, $X_i \in \mathbb{R}^{P_i}$, and let $g(X_i) = f(\bar{x}_1, \dots, \bar{x}_{i-1}, X_i, \bar{x}_{i+1}, \dots, \bar{x}_t)$, $i = 1, \dots, t$. If $x \in \Psi$, then $g(X_i)$ has a unique (global) minimizer for $X_i \in \Psi_i$. (EU)

Theorem 2 [10]. Suppose that (EU) holds for $f: \mathfrak{R}^s \mapsto \mathfrak{R}$. Let $x = (X_1, \dots, X_t)^T$, and $\Psi = \Psi_1 \times \dots \times \Psi_t$, where Ψ_i is a compact subset of \mathfrak{R}^{P_i} , $i = 1, \dots, t$. Let $\{x^{(r+1)} = T(x^{(r)})\}$ denote the AO iterate sequence begun at $x^{(0)} \in \Psi$, and denote the fixed points of T as $\Omega = \{x \in \Psi : x = T(x)\}$. Then:

(i) if $x^* \in \Omega$, then $x^* = (X_1^*, \dots, X_t^*)^T$ satisfies, for $i = 1, \dots, t$,

$$X_i^* = \arg \min_{X_i \in \Psi_i \subset \mathfrak{R}^{P_i}} \left\{ f(\mathbf{x}_1^*, \dots, \mathbf{x}_{i-1}^*, X_i, \mathbf{x}_{i+1}^*, \dots, \mathbf{x}_t^*) \right\};$$

(ii) $f(x^{(r+1)}) \leq f(x^{(r)})$, equality if and only if $x^{(r)} \in \Omega$;

(iii) either: (a) $\exists x^* \in \Omega$ and $r_0 \in \mathfrak{R}$ so that $x^{(r)} = x^*$ for all $r \geq r_0$;

or (b) the limit of every convergence subsequence of $\{x^{(r)}\}$ is in Ω .

- Under certain conditions, all limit points of an AO sequence are either **saddle points** of a special type of minimizers.
- However, not all saddle point can be captured by AO. Only those which looks like a minimizer along the grouped coordinate (X_1, X_2 , etc) can be captured.
- The potential for convergence to a saddle point is a “price” need to pay.
- What if strict convex functions? **Converge to the global optimal q-linearly**

- Under certain conditions, all limit points of an AO sequence are either **saddle points** of a special type of minimizers.
- However, not all saddle point can be captured by AO. Only those which looks like a minimizer along the grouped coordinate (X_1, X_2 , etc) can be captured.
- The potential for convergence to a saddle point is a “price” need to pay.
- What if strict convex functions? **Converge to the global optimal q-linearly**

Theorem 3 [10]. Let x^* be a local minimizer of $f: \mathfrak{X}^S \mapsto \mathfrak{R}$ for which $\nabla^2 f(x^*)$ is positive definite, and let f be C^2 in a neighborhood $N(x^*, \delta)$. Let $0 < \varepsilon \leq \delta$ be chosen so that f is strictly convex on $N(x^*, \varepsilon)$. Finally, assume that if $y = (\mathfrak{X}_1, \dots, \mathfrak{X}_{i-1}, Y_i, \mathfrak{X}_{i+1}, \dots, \mathfrak{X}_t)^T \in N(x^*, \varepsilon)$, and Y_i^* locally minimizes $g_1(Y_i) = f(\mathfrak{X}_1, \dots, \mathfrak{X}_{i-1}, Y_i, \mathfrak{X}_{i+1}, \dots, \mathfrak{X}_t)$, then Y_i^* is also the unique global minimizer of g_1 :

Then for any $x^{(0)} \in N(x^*, \varepsilon)$, the corresponding AO iterate sequence $\{x^{(r+1)} = T(x^{(r)})\} \rightarrow x^*$ q-linearly.

- The previous two results are making strong assumptions:
 - Each restricted minimization problem has a unique solution.
 - Strict convexity near the optimal.
- Here, study the functions with relaxed assumptions:
 - Minimize a nondifferentiable (nonconvex) function $f(x_1, \dots, x_N)$ with certain separability and regularity properties.
 - Converge to a stationary point if f is
 - pseudoconvex in every pair of coordinate blocks from among $N - 1$ coordinate blocks; or
 - f has at most one minimum in each of $N - 2$ coordinate blocks
 - If f is quasiconvex and hemivariate in every coordinate block, the assumption could be relaxed further.

- The previous two results are making strong assumptions:
 - Each restricted minimization problem has a unique solution.
 - Strict convexity near the optimal.
- Here, study the functions with relaxed assumptions:
 - Minimize a nondifferentiable (nonconvex) function $f(x_1, \dots, x_N)$ with certain separability and regularity properties.
 - Converge to a stationary point if f is
 - pseudoconvex in every pair of coordinate blocks from among $N - 1$ coordinate blocks; or
 - f has at most one minimum in each of $N - 2$ coordinate blocks
 - If f is quasiconvex and hemivariate in every coordinate block, the assumption could be relaxed further.

- **Effective domain:** $dom\ h = \{x \in R^m \mid h(x) < \infty\}$
- A function f is **proper** if $f \neq \infty$.
- A space is **compact** if it is closed and bounded.
- **Lower Directional derivative:**

$$h'(x; d) = \liminf_{\lambda \rightarrow 0} \frac{h(x + \lambda d) - h(x)}{\lambda}$$

- **Gateaux-Differentiable:**

$$h'(x; d) = \lim_{\lambda \rightarrow 0} \frac{h(x + \lambda d) - h(x)}{\lambda} = \frac{d}{d\lambda} h(x + \lambda d)|_{\lambda=0}$$

If the transformation $H(d) : d \rightarrow h'(x; d)$ is **continuous** and **linear**, then F is said to be Gateaux differentiable at u .

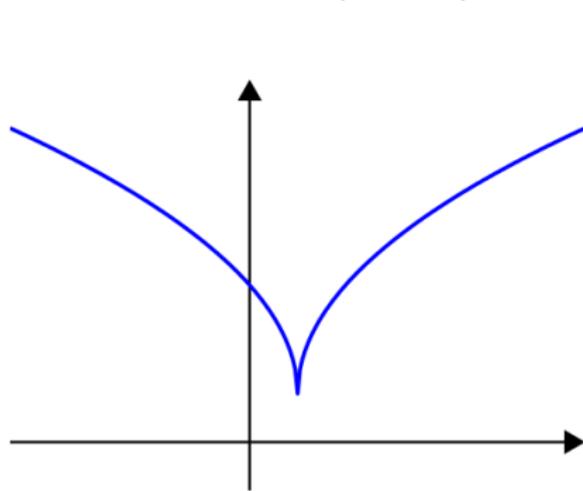
In other words,

$$h'(x; \alpha d) = \alpha h'(x; d);$$

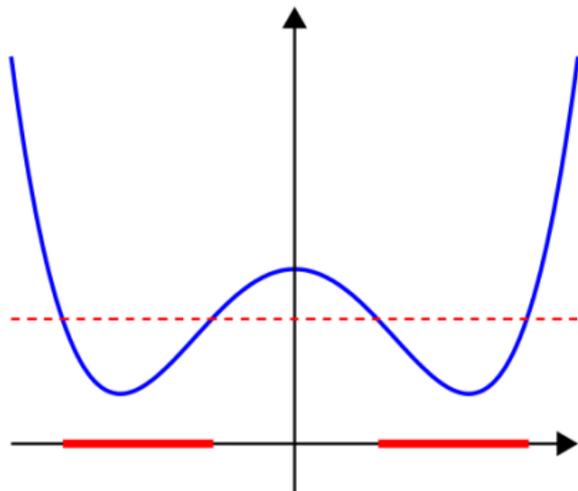
$$h'(x; (d_1 + d_2)) = h'(x; d_1) + h'(x; d_2)$$

QuasiConvex

- **Quasiconvex**: a real-valued function defined on an interval or on a convex subset or a real vector space such that the inverse image of any set of the form $(-\infty, a)$ is a convex set.



Quasiconvex but not convex



Not Quasiconvex

$$h(\lambda x + (1 - \lambda)y) \leq \max(h(x), h(y)) \quad \forall \lambda \in [0, 1]$$

or

$$h(x + \lambda d) \leq \max(h(x), h(x + d))$$

- **Pseudoconvex**: a function satisfying the following property:

$$h(x + d) \geq h(x), \quad \text{whenever } x \in \text{dom } h \text{ and } h'(x; , d) \geq 0$$

- $\arctan(x)$ is pseudo convex, but not convex. Its derivative is

$$\frac{1}{1 + x^2}$$

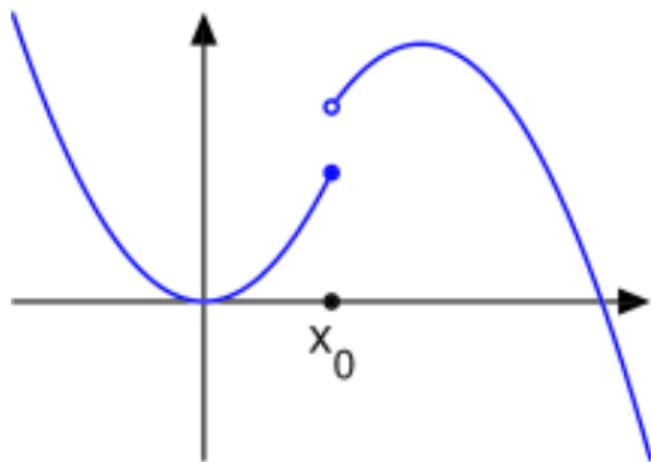
which is always positive. But it's not convex function.

- **hemivariate**: h is not constant on any line segment belonging to $\text{dom } h$. Used to guarantee the unique minimizer for each restricted minimization problem.

Lower Semi-continuous

- lower semi-continuous:

$$\liminf_{x \rightarrow x_0} f(x) \geq f(x_0)$$



- A Lower Semi-Continuous Function indicates that the limit point x_0 (if in the effective domain), the function value f is always smaller than the limiting value of f .

Stationary Point & Regular Function

- z is a **stationary point** if

$$f'(z; d) \geq 0, \quad \forall d$$

- f is **regular** if $\forall d = (d_1, \dots, d_N)$ which satisfy

$$f'(z; (0, \dots, d_k, \dots, 0)) \geq 0 \implies f'(z; d) \geq 0$$

- **coordinatewise minimum point**:

$$f(z + (0, \dots, d_k, \dots, 0)) \geq f(z), \quad \forall d_k$$

- This is less strong than the following condition:

$$f'(z; d) = \sum_{k=1}^N f'(z; (0, \dots, d_k, \dots, 0)), \quad \text{for all } d = (d_1, \dots, d_N)$$

Stationary Point & Regular Function

- z is a **stationary point** if

$$f'(z; d) \geq 0, \quad \forall d$$

- f is **regular** if $\forall d = (d_1, \dots, d_N)$ which satisfy

$$f'(z; (0, \dots, d_k, \dots, 0)) \geq 0 \implies f'(z; d) \geq 0$$

- **coordinatewise minimum point:**

$$f(z + (0, \dots, d_k, \dots, 0)) \geq f(z), \quad \forall d_k$$

- This is less strong than the following condition:

$$f'(z; d) = \sum_{k=1}^N f'(z; (0, \dots, d_k, \dots, 0)), \quad \text{for all } d = (d_1, \dots, d_N)$$

Stationary Point & Regular Function

- z is a **stationary point** if

$$f'(z; d) \geq 0, \quad \forall d$$

- f is **regular** if $\forall d = (d_1, \dots, d_N)$ which satisfy

$$f'(z; (0, \dots, d_k, \dots, 0)) \geq 0 \implies f'(z; d) \geq 0$$

- **coordinatewise minimum point**:

$$f(z + (0, \dots, d_k, \dots, 0)) \geq f(z), \quad \forall d_k$$

- This is less strong than the following condition:

$$f'(z; d) = \sum_{k=1}^N f'(z; (0, \dots, d_k, \dots, 0)), \quad \text{for all } d = (d_1, \dots, d_N)$$

An example of Regular Function with no additive property

$$f(x_1, x_2) = \phi(x_1, x_2) + \phi(-x_1, x_2) + \phi(x_1, -x_2) + \phi(-x_1, -x_2)$$

where $\phi(a, b) = \max\{0, a + b - \sqrt{a^2 + b^2}\}$

It's easy to verify that

$$f'(\mathbf{0}; (d_1, 0)) = 0, \quad f'(\mathbf{z}; (0, d_2)) = 0;$$

$$f'(\mathbf{0}; d) = |d_1| + |d_2| - \sqrt{d_1^2 + d_2^2} \neq f'(\mathbf{0}; (d_1, 0)) + f'(\mathbf{0}; (0, d_2))$$

Stationary Point = Coordinate-wise Minimum?

- z is a **stationary point** if

$$f'(z; d) \geq 0, \quad \forall d$$

- f is **regular** if $\forall d = (d_1, \dots, d_N)$ which satisfy

$$f'(z; (0, \dots, d_k, \dots, 0)) \geq 0 \implies f'(z; d) \geq 0$$

- **coordinatewise minimum point:**

$$f(z + (0, \dots, d_k, \dots, 0)) \geq f(z), \quad \forall d_k$$

- A coordinatewise minimum point z is a stationary point whenever f is regular at z .
- When is a function regular?

Stationary Point = Coordinate-wise Minimum?

- z is a **stationary point** if

$$f'(z; d) \geq 0, \quad \forall d$$

- f is **regular** if $\forall d = (d_1, \dots, d_N)$ which satisfy

$$f'(z; (0, \dots, d_k, \dots, 0)) \geq 0 \implies f'(z; d) \geq 0$$

- **coordinatewise minimum point:**

$$f(z + (0, \dots, d_k, \dots, 0)) \geq f(z), \quad \forall d_k$$

- A coordinatewise minimum point z is a stationary point whenever f is regular at z .
- When is a function regular?

Smoothness Assumptions

$$f(x_1, \dots, x_N) = f_0(x_1, \dots, x_N) + \sum_{k=1}^N f_k(x_k)$$

- A1** $\text{dom } f_0$ is open and f_0 is Gateaux-differentiable on $\text{dom } f_0$.
- A2** f_0 is Gateaux-differentiable on $\text{int}(\text{dom } f_0)$ and for every $z \in \text{dom } f \cap \text{bdry}(\text{dom } f_0)$, there exist

$$f(z + (0, \dots, d_k, \dots, 0)) < f(z)$$

Essentially the minimizer never occurs at the boundary point.

Lemma 3.1 Under A1, f is regular at each $z \in \text{dom } f$; Under A2, f is regular at each coordinatewise minimum point z of f .

Smoothness Assumptions

$$f(x_1, \dots, x_N) = f_0(x_1, \dots, x_N) + \sum_{k=1}^N f_k(x_k)$$

- A1** $\text{dom } f_0$ is open and f_0 is Gateaux-differentiable on $\text{dom } f_0$.
- A2** f_0 is Gateaux-differentiable on $\text{int}(\text{dom } f_0)$ and for every $z \in \text{dom } f \cap \text{bdry}(\text{dom } f_0)$, there exist

$$f(z + (0, \dots, d_k, \dots, 0)) < f(z)$$

Essentially the minimizer never occurs at the boundary point.

Lemma 3.1 Under A1, f is regular at each $z \in \text{dom } f$; Under A2, f is regular at each coordinatewise minimum point z of f .

Proof for Lemma 3.1

Lemma 3.1 Under A1, f is regular at each $z \in \text{dom} f$; Under A2, f is regular at each coordinatewise minimum point z of f .

Under A1, if $z \in \text{dom} f \implies z \in \text{dom} f_0$; Under A2, $z \in \text{int}(\text{dom} f_0)$ for any d such that $f'(z; (0, \dots, d_k, \dots, 0)) \geq 0 \quad k = 1, \dots, N$

We need to prove $f'(z; d) \geq 0$.

$$f'(z; d) = \underbrace{\langle \nabla_0(z), d \rangle}_{\text{Gateaux-differentiable}} + \liminf_{\lambda \downarrow 0} \sum_{k=1}^N [f_k(x_k + \lambda d_k) - f_k(x_k)] / \lambda$$

$$\geq \langle \nabla_0(z), d \rangle + \sum_{k=1}^N \liminf_{\lambda \downarrow 0} [f_k(x_k + \lambda d_k) - f_k(x_k)] / \lambda \quad (1)$$

$$= \langle \nabla f_0(z), d \rangle + \sum_{k=1}^N f'_k(z_k; d_k) \quad (2)$$

$$= \sum_{k=1}^N f'(z; (0, \dots, d_k, \dots, 0)) \geq 0 \quad (3)$$

Proof for Lemma 3.1

Lemma 3.1 Under A1, f is regular at each $z \in \text{dom} f$; Under A2, f is regular at each coordinatewise minimum point z of f .

Under A1, if $z \in \text{dom} f \implies z \in \text{dom} f_0$; Under A2, $z \in \text{int}(\text{dom} f_0)$ for any d such that $f'(z; (0, \dots, d_k, \dots, 0)) \geq 0 \quad k = 1, \dots, N$

We need to prove $f'(z; d) \geq 0$.

$$f'(z; d) = \underbrace{\langle \nabla_0(z), d \rangle}_{\text{Gateaux-differentiable}} + \liminf_{\lambda \downarrow 0} \sum_{k=1}^N [f_k(x_k + \lambda d_k) - f_k(x_k)]/\lambda$$

$$\geq \langle \nabla_0(z), d \rangle + \sum_{k=1}^N \liminf_{\lambda \downarrow 0} [f_k(x_k + \lambda d_k) - f_k(x_k)]/\lambda \quad (1)$$

$$= \langle \nabla f_0(z), d \rangle + \sum_{k=1}^N f'_k(z_k; d_k) \quad (2)$$

$$= \sum_{k=1}^N f'(z; (0, \dots, d_k, \dots, 0)) \geq 0 \quad (3)$$

- This work makes the assumption of A1 or A2.
- Under such assumptions, a coordinate-wise minimum is a stationary point.
- So the following convergence analysis just need to show that **the algorithm converges to a coordinate-wise minimum point.**
- **A1 & A2 only care about the smoothness of f_0 .** Even if f_1, \dots, f_N are not smooth, the claim here is still valid.
- Need additional properties to guarantee the convergence.

BCD Method.

Initialization. Choose any $x^0 = (x_1^0, \dots, x_N^0) \in \text{dom } f$.

Iteration $r+1$, $r \geq 0$. Given $x^r = (x_1^r, \dots, x_N^r) \in \text{dom } f$, choose an index $s \in \{1, \dots, N\}$ and compute a new iterate

$$x^{r+1} = (x_1^{r+1}, \dots, x_N^{r+1}) \in \text{dom } f$$

satisfying

$$x_s^{r+1} \in \arg \min_{x_s} f(x_1^r, \dots, x_{s-1}^r, x_s, x_{s+1}^r, \dots, x_N^r), \quad (2)$$

$$x_j^{r+1} = x_j^r, \quad \forall j \neq s. \quad (3)$$

Essentially Cyclic Rule. There exists a constant $T \geq N$ such that every index $s \in \{1, \dots, N\}$ is chosen at least once between the r th iteration and the $(r + T - 1)$ th iteration, for all r .

A well-known special case of this rule, for which $T = N$, is given below.

Cyclic Rule. Choose $s = k$ at iterations $k, k + N, k + 2N, \dots$, for $k = 1, \dots, N$.

Assuming f continuous, without using the Special Structure

Theorem 4.1 Assume the level set $X^0 = \{x : f(x) \leq f(x^0)\}$ is compact and that f is continuous on X^0 . Then, the sequence generated by BCD is defined and bounded. Moreover,

- (a) If $f(x_1, \dots, x_N)$ is pseudoconvex in (x_k, x_i) for every $i, k \in \{1, \dots, N\}$, and if f is regular at every $x \in X^0$, then every cluster point of $\{x^r\}$ is a stationary point of f .
- (b) If $f(x_1, \dots, x_N)$ is pseudoconvex in (x_k, x_i) for every $i, k \in \{1, \dots, N-1\}$, if f is regular at every $x \in X^0$, and if the cyclic rule is used, then every cluster point of $\{x^r\}_{r \equiv (N-1) \pmod N}$ is a stationary point of f .
- (c) If $f(x_1, \dots, x_N)$ has at most one minimum in x_k for $k = 2, \dots, N-1$, and if the cyclic rule is used, then every cluster point z of $\{x^r\}_{r \equiv (N-1) \pmod N}$ is a coordinatewise minimum point of f . In addition, if f is regular at z , then z is a stationary point of f .

- **Goal:** To show that the BCD algorithm converges to z such that

$$f(z + (0, \dots, d_k, \dots, 0)) \geq f(z); \quad \forall d_k, k = 1, \dots, N$$

- The stationary point property is obtained if the function is regular.
- The key process is to show the following by induction:
for $j = 1, \dots, T - 1$,

$$f(z^j) \leq f(z^j + (0, \dots, d_k, \dots, 0)), \quad \forall d_k, \forall k = s^1, \dots, s^j.$$

- $X^0 = \{x : f(x) \leq f(x^0)\}$ is compact
- $\Rightarrow f(x^{r+1}) \leq f(x^r)$ and $x^{r+1} \in X^0$ for all $r = 0, 1, \dots$
- $\Rightarrow \{x^r\}$ is bounded.
- \Rightarrow Consider any subsequence $\{x^r\}_{r \in R}$, converging to z , where $R \subseteq \{0, 1, \dots\}$,
 $\{x^{r-T+1+j}\}_{r \in R}$ is bounded.
- By passing r to a subsubsequence, we have
- $\Rightarrow \{x^{r-T+1+j}\}_{r \in R} \rightarrow z^j, j = 1, \dots, T$

Note that $z^{T-1} = z;$

$$\Rightarrow \underbrace{f(x^0) \geq \lim_{r \rightarrow \infty} f(x^r) = f(z^1) = \dots = f(z^T)}_{\text{f decreasing monotonically, and f is continuous}}$$

f decreasing monotonically, and f is continuous

Assume that the index s chosen at iteration $r - T + 1 + j$, $j \in \{1, \dots, T\}$, is the same for all $r \in R$ (denoted as s^j), then

$$f(x^{r-T+1+j}) \leq f(x^{r-T+1+j} + (0, \dots, d_{s^j}, \dots, 0)), \quad \forall d^{s^j}, j = 1, \dots, T$$

$$x_k^{r-T+1+j} = x_K^{r-T+j} \quad \forall k \neq s^j, j = 2, \dots, T$$

Based on the continuity of f on X^0 , we have

$$f(z^j) \leq f(z^j + (0, \dots, d_{s^j}, \dots, 0)), \quad \forall d_{s^j}, j = 1, \dots, T$$

$$z_k^j = z_k^{j-1} \quad \forall k \neq s^j, j = 2, \dots, T$$

$$\Rightarrow f(z^{j-1}) = f(z^j) \leq \underbrace{f(z^{j-1} + (0, \dots, d_{s^j}, \dots, 0))}_{z^j \text{ and } z^{j-1} \text{ only differ at index } s^j}$$

$$\forall d_{s^j}, j = 2, \dots, T$$

The limit point z^{j-1} is also the directional minimizer for d_{s^j} .

if f is pseudoconvex in (x_k, x_i) , $\forall i, k \in s^1 \cup \dots, \cup s^{T-1}$

We have

$$f(z^{j-1}) \leq f(z^{j-1} + (0, \dots, d_{sj}, \dots, 0)), \quad j = 2, \dots, T$$

- (a). f is pseudoconvex in (x_k, x_i) for every i, k in $\{1, \dots, N\}$
- (b). f is pseudoconvex in (x_k, x_i) for every i, k in $\{1, \dots, N-1\}$
 \Rightarrow if f is pseudoconvex in (x_k, x_i) , $\forall i, k \in s^1 \cup \dots, \cup s^{T-1}$

Claim for $j = 1, \dots, T-1$,

$$f(z^j) \leq f(z^j + (0, \dots, d_k, \dots, 0)), \quad \forall d_k, \forall k = s^1, \dots, s^j. \quad (4)$$

Note that

$$f(z) = f(z^{T-1}) \leq f(z^{T-1} + (0, \dots, d_{sT}, \dots, 0))$$

Then we have z is a coordinate-wise minimum.

if f is pseudoconvex in (x_k, x_i) , $\forall i, k \in s^1 \cup \dots, \cup s^{T-1}$

We have

$$f(z^{j-1}) \leq f(z^{j-1} + (0, \dots, d_{sj}, \dots, 0)), \quad j = 2, \dots, T$$

- (a). f is pseudoconvex in (x_k, x_i) for every i, k in $\{1, \dots, N\}$
- (b). f is pseudoconvex in (x_k, x_i) for every i, k in $\{1, \dots, N-1\}$
 \Rightarrow if f is pseudoconvex in (x_k, x_i) , $\forall i, k \in s^1 \cup \dots, \cup s^{T-1}$

Claim for $j = 1, \dots, T-1$,

$$f(z^j) \leq f(z^j + (0, \dots, d_k, \dots, 0)), \quad \forall d_k, \forall k = s^1, \dots, s^j. \quad (5)$$

Proof by Induction

- $j = 1$, automatically satisfied by the minimization.
- Suppose (5) holds for $j = 1, \dots, \ell - 1$ for $\ell \in \{2, \dots, T-1\}$, we'll show (5) holds for ℓ .

$$\begin{aligned}
& f(z^{j-1}) \leq f(z^{j-1} + (0, \dots, d_{sj}, \dots, 0)) \quad \forall d_{sj}, j = 2, \dots, T \\
\Rightarrow & f(z^{\ell-1}) \leq f(z^{\ell-1} + (0, \dots, d_{s^\ell}, \dots, 0)) \quad \forall d_{s^\ell} \\
\Rightarrow & f'(z^{\ell-1}; (0, \dots, z_{s^\ell}^\ell - z_{s^\ell}^{\ell-1}, \dots, 0)) \geq 0 \quad (\text{pseudoconvexity})
\end{aligned}$$

Based on Induction assumption, we have

$$\begin{aligned}
& f'(z^{\ell-1}; (0, \dots, d_k, \dots, 0)) \geq 0, \forall d_k, k = s^1, \dots, s^{\ell-1} \\
\Rightarrow & \underbrace{f'(z^{\ell-1}; (0, \dots, d_k, \dots, 0) + (0, \dots, z_{s^\ell}^\ell - z_{s^\ell}^{\ell-1}, \dots, 0))}_{\text{as } f \text{ is regular}} \geq 0 \quad (6)
\end{aligned}$$

$$\Rightarrow f(z^{\ell-1}) \leq f(z^\ell + (0, \dots, d_k, \dots, 0)) \quad (f \text{ is pseudoconvex}) \quad (7)$$

$$\Rightarrow f(z^\ell) = f(z^{\ell-1}) \leq f(z^\ell + (0, \dots, d_k, \dots, 0)) \quad k = s^1, \dots, s^{\ell-1} \quad (8)$$

$$\text{As } f(z^j) \leq f(z^j + (0, \dots, d_{sj}, \dots, 0)), \quad \forall d_{sj}, j = 1, \dots, T \quad (9)$$

$$\Rightarrow f(z^\ell) \leq f(z^\ell + (0, \dots, d_k, \dots, 0)) \quad k = s^1, \dots, s^\ell \quad (10)$$

$$\Rightarrow \text{Claim holds for } \ell. \quad (11)$$

$$\begin{aligned}
& f(z^{j-1}) \leq f(z^{j-1} + (0, \dots, d_{sj}, \dots, 0)) \quad \forall d_{sj}, j = 2, \dots, T \\
\Rightarrow & f(z^{\ell-1}) \leq f(z^{\ell-1} + (0, \dots, d_{s^\ell}, \dots, 0)) \quad \forall d_{s^\ell} \\
\Rightarrow & f'(z^{\ell-1}; (0, \dots, z_{s^\ell}^\ell - z_{s^\ell}^{\ell-1}, \dots, 0)) \geq 0 \quad (\text{pseudoconvexity})
\end{aligned}$$

Based on Induction assumption, we have

$$f'(z^{\ell-1}; (0, \dots, d_k, \dots, 0)) \geq 0, \forall d_k, k = s^1, \dots, s^{\ell-1}$$

$$\Rightarrow \underbrace{f'(z^{\ell-1}; (0, \dots, d_k, \dots, 0) + (0, \dots, z_{s^\ell}^\ell - z_{s^\ell}^{\ell-1}, \dots, 0))}_{\text{as } f \text{ is regular}} \geq 0 \quad (6)$$

$$\Rightarrow f(z^{\ell-1}) \leq f(z^\ell + (0, \dots, d_k, \dots, 0)) \quad (f \text{ is pseudoconvex}) \quad (7)$$

$$\Rightarrow f(z^\ell) = f(z^{\ell-1}) \leq f(z^\ell + (0, \dots, d_k, \dots, 0)) \quad k = s^1, \dots, s^{\ell-1} \quad (8)$$

$$\text{As } f(z^j) \leq f(z^j + (0, \dots, d_{sj}, \dots, 0)), \quad \forall d_{sj}, j = 1, \dots, T \quad (9)$$

$$\Rightarrow f(z^\ell) \leq f(z^\ell + (0, \dots, d_k, \dots, 0)) \quad k = s^1, \dots, s^\ell \quad (10)$$

$$\Rightarrow \text{Claim holds for } \ell. \quad (11)$$

$$\begin{aligned}
& f(z^{j-1}) \leq f(z^{j-1} + (0, \dots, d_{sj}, \dots, 0)) \quad \forall d_{sj}, j = 2, \dots, T \\
\Rightarrow & f(z^{\ell-1}) \leq f(z^{\ell-1} + (0, \dots, d_{s^\ell}, \dots, 0)) \quad \forall d_{s^\ell} \\
\Rightarrow & f'(z^{\ell-1}; (0, \dots, z_{s^\ell}^\ell - z_{s^\ell}^{\ell-1}, \dots, 0)) \geq 0 \quad (\text{pseudoconvexity})
\end{aligned}$$

Based on Induction assumption, we have

$$\begin{aligned}
& f'(z^{\ell-1}; (0, \dots, d_k, \dots, 0)) \geq 0, \forall d_k, k = s^1, \dots, s^{\ell-1} \\
\Rightarrow & \underbrace{f'(z^{\ell-1}; (0, \dots, d_k, \dots, 0) + (0, \dots, z_{s^\ell}^\ell - z_{s^\ell}^{\ell-1}, \dots, 0))}_{\text{as } f \text{ is regular}} \geq 0 \quad (6)
\end{aligned}$$

$$\Rightarrow f(z^{\ell-1}) \leq f(z^\ell + (0, \dots, d_k, \dots, 0)) \quad (f \text{ is pseudoconvex}) \quad (7)$$

$$\Rightarrow f(z^\ell) = f(z^{\ell-1}) \leq f(z^\ell + (0, \dots, d_k, \dots, 0)) \quad k = s^1, \dots, s^{\ell-1} \quad (8)$$

$$\text{As } f(z^j) \leq f(z^j + (0, \dots, d_{sj}, \dots, 0)), \quad \forall d_{sj}, j = 1, \dots, T \quad (9)$$

$$\Rightarrow f(z^\ell) \leq f(z^\ell + (0, \dots, d_k, \dots, 0)) \quad k = s^1, \dots, s^\ell \quad (10)$$

$$\Rightarrow \text{Claim holds for } \ell. \quad (11)$$

$$\text{As } f(z^{j-1}) = f(z^j) \leq f(z^{j-1} + (0, \dots, d_{sj}, \dots, 0)) \quad \forall d_{sj}, j = 2, \dots, T,$$
$$f(z^{T-1}) \leq f(z^{T-1} + (0, \dots, d_k, \dots, 0)) \quad k = s^T$$

Combined with our induction proof, we have

$$f(z^{T-1}) \leq f(z^{T-1} + (0, \dots, d_k, \dots, 0)) \quad k = s^1, \dots, s^T$$

Recall that $z^{T-1} = z$, hence z is coordinate-wise minimum.

As f is regular, z is also a stationary point.

Unique Minimizer at Each Step \implies unique limiting point?

- (c). f has at most one minimum in x_k for $k = 2, \dots, N - 1$, and if the cycle rule is used. Then every cluster point z of $\{x^r\}_{r \equiv (N-1) \pmod N}$, is a coordinatewise minimum point of f . If f is regular at z , then it's also a stationary point.

Proof

Define a function as $d_{sj} \rightarrow f(z^j + (0, \dots, d_{sj}, \dots, 0))$

$$f(z^{j-1}) = f(z^j) \leq f(z^{j-1} + (0, \dots, d_{sj}, \dots, 0)) \quad \forall d_{sj}, j = 2, \dots, T \quad (12)$$

attains its minimum at both 0 and $z_{sj}^{j-1} - z_{sj}^j$.

$$\implies z_{sj}^{j-1} - z_{sj}^j = 0 \quad (\text{uniqueness of minimization function})$$

$$\implies z^{j-1} = z^j \implies z^1 = z^2 = \dots, z^{T-1} = z$$

Plus, $f(z^{j-1}) = f(z^j) \leq f(z^{j-1} + (0, \dots, d_{sj}, \dots, 0)) \quad \forall d_{sj}, j = 2, \dots, T$

Hence, z is the coordinate-wise minimizer.

Recap the Theorem

Assuming f continuous, without using the Special Structure

Theorem 4.1 Assume the level set $X^0 = \{x : f(x) \leq f(x^0)\}$ is compact and that f is continuous on X^0 . Then, the sequence generated by BCD is defined and bounded. Moreover,

- (a) If $f(x_1, \dots, x_N)$ is pseudoconvex in (x_k, x_i) for every $i, k \in \{1, \dots, N\}$, and if f is regular at every $x \in X^0$, then every cluster point of $\{x^r\}$ is a stationary point of f .
- (b) If $f(x_1, \dots, x_N)$ is pseudoconvex in (x_k, x_i) for every $i, k \in \{1, \dots, N-1\}$, if f is regular at every $x \in X^0$, and if the cyclic rule is used, then every cluster point of $\{x^r\}_{r \equiv (N-1) \pmod N}$ is a stationary point of f .
- (c) If $f(x_1, \dots, x_N)$ has at most one minimum in x_k for $k = 2, \dots, N-1$, and if the cyclic rule is used, then every cluster point z of $\{x^r\}_{r \equiv (N-1) \pmod N}$ is a coordinatewise minimum point of f . In addition, if f is regular at z , then z is a stationary point of f .

Summary & Comments

- if f is pseudoconvex, then f is pseudoconvex in (x_k, x_i) for all k, i
- if f is quasiconvex and hemivariate in x_k , then f has at most one minimum in x_k . Some papers refer it as **strict quasiconvex**.
- If f is continuous, and only **2-blocks** are involved. Then it does not require unique minimizer to converge to a stationary point. (This result is used in the convergence proof of alternating least-square proof in NMF)
- The previous proof does not take advantage of the special structure and assume f to be continuous on a bounded level set.
- Next we show that considering the special structure **without requiring f to be smooth**.

Summary & Comments

- if f is pseudoconvex, then f is pseudoconvex in (x_k, x_i) for all k, i
- if f is quasiconvex and hemivariate in x_k , then f has at most one minimum in x_k . Some papers refer it as **strict quasiconvex**.
- If f is continuous, and only **2-blocks** are involved. Then it does not require unique minimizer to converge to a stationary point. (This result is used in the convergence proof of alternating least-square proof in NMF)
- The previous proof does not take advantage of the special structure and assume f to be continuous on a bounded level set.
- Next we show that considering the special structure **without requiring f to be smooth**.

Sleepy? Shall we continue?



Assumptions

(B1) f_0 is continuous on $\text{dom } f_0$

(B2) for each $k \in \{1, \dots, N\}$, and $(x_j)_{j \neq k}$, the function $x_k \rightarrow f(x_1, \dots, x_N)$ is quasiconvex and hemivariate.

(B3) f_0, f_1, \dots, f_N is lower semi-continuous.

Meanwhile, f_0 satisfy the one of the following assumption:

(C1) $\text{dom } f_0$ is open and f_0 tends to ∞ at every boundary point of $\text{dom } f_0$

(C2) $\text{dom } f_0 = Y_1 \times \dots \times Y_N$ for some $Y_k \subseteq R^{n_k}, k = 1, \dots, N$

- C2 allows a finite value at boundary point.
- We'll show that Assumption B1-B3, together with either C1 or C2, ensure that every cluster point of the iterates generated by the BCD methods is a coordinate minimum point of f .

Proposition 5.1

Suppose that f, f_0, \dots, f_N satisfy B1-B3 and f_0 satisfy C1 or C2. Then, either $\{f(x^r)\} \downarrow -\infty$ or else every cluster point $z = (z_1, \dots, z_N)$ is a coordinatewise minimum point of f .

Proof Strategy

As $f(x^0) < \infty$, and $f(x^{r+1}) \leq f(x^r)$

$\Rightarrow \{f(x^r)\} \downarrow -\infty$

or $\{f(x^r)\}$ converges to some limit and $\{f(x^{r+1}) - f(x^r)\} \rightarrow 0$

Let z be any cluster point of $\{x^r\}$

$\Rightarrow f(z) \leq \lim_{r \rightarrow \infty} f(x^r) \leq \infty$ (as f is lower semi-continuous)

- First, we show that for any convergent sequence $\{x^r\} \rightarrow z$, we have $\{x^{r+1}\} \rightarrow z$;
- We'll prove this by contradiction.
- Then, we prove z is a coordinate-wise minimum.

Claim of convergence for x^r

Claim: for any convergent subsequence $\{x^r\}_{r \in R} \rightarrow z$, we have

$$\{x^{r+1}\} \rightarrow z$$

Sketch of the Proof

- Proof by contradiction
- If $\{x^{r+1}\}$ converges to a different value z' , then all the values between z and z' have

$$f(\lambda z + (1 - \lambda)z') = f(z) = f(z')$$

contradicting to the uniqueness of each minimization of coordinate block.

Claim of convergence for x^r

Claim: for any convergent subsequence $\{x^r\}_{r \in R} \rightarrow z$, we have

$$\{x^{r+1}\} \rightarrow z$$

Prove by Contradiction

Suppose the above is not true, then there exists an infinite subsequence $R' \subseteq R$ and a scalar $\epsilon > 0$ such that

$$\|x^{r+1} - x^r\| \geq \epsilon, \quad \text{for all } r \in R'$$

So we can assume that there is some nonzero vector d such that

$$\{(x^{r+1} - x^r) / \|x^{r+1} - x^r\|\}_{r \in R'} \rightarrow d \quad (\text{not quite sure why?})$$

and the same coordinate block, say x_s is chosen at the $r + 1$ -th iteration. So

$$\{f_0(x^r) + f_s(x_s^r)\}_{r \in R'} \rightarrow \theta$$

Fix any $\lambda \in [0, \epsilon]$, Let $\hat{z} = z + \lambda d$, and for each $r \in R'$, let

$$\hat{x}^r = x^r + \lambda(x^{r+1} - x^r) / \|x^{r+1} - x^r\| \quad (13)$$

$$\Rightarrow \{\hat{x}^r\}_{r \in R'} \rightarrow \hat{z} \quad (14)$$

\hat{x}^r lies in the segment of x^{r+1} and x^r , thus

$$f(\hat{x}^r) \leq f(x^r) \quad \forall r \in R' \quad (\text{f is quasiconvex}) \quad (15)$$

$$\Rightarrow f_0(\hat{x}^r) + f_s(\hat{x}_s^r) \leq f_0(x^r) + f_s(x_s^r) \rightarrow \theta \quad (16)$$

$$\Rightarrow \lim_{r \rightarrow \infty, r \in R'} \sup \{f_0(\hat{x}^r) + f_s(\hat{x}_s^r)\} \leq \theta \quad (17)$$

$$\text{As } \{f(x^{r+1}) - f(x^r)\} \rightarrow 0 \quad (18)$$

$$\Rightarrow \{f_0(x^{r+1}) + f_s(x_s^{r+1}) - f_0(x^r) - f_s(x_s^r)\}_{r \in R'} \rightarrow 0 \quad (19)$$

$$\Rightarrow \{f_0(x^{r+1}) + f_s(x_s^{r+1})\} \rightarrow \theta \quad (20)$$

$$\text{Define } \delta = f_0(\hat{z}) + f_s(\hat{z}_s) - \theta \quad (21)$$

$$\text{Then } \delta \leq 0, \text{ actually } \delta = 0 \quad (22)$$

$$\text{As } \underbrace{\{(x_1^r, \dots, x_{s-1}^r, \hat{z}_s, x_{s+1}^r, \dots, x_N^r)\}}_{\hat{x}^r \text{ and } x^r \text{ only differ in } s\text{-th block}} \rightarrow \hat{z} \quad (23)$$

$$\lim_{r \rightarrow \infty, r \in R'} \sup \{f_0(\hat{x}^r) + f_s(\hat{x}_s^r)\} \leq \theta \quad (24)$$

if $\delta \neq 0$, then for r sufficiently large

$$f_0(x_1^r, \dots, x_{s-1}^r, \hat{z}_s, x_{s+1}^r, \dots, x_N^r) \leq f_0(x^{r+1}) + f_s(x_s^{r+1}) + \delta/2 \quad (25)$$

$$f(x_1^r, \dots, x_{s-1}^r, \hat{z}_s, x_{s+1}^r, \dots, x_N^r) \leq f(x^{r+1}) + \delta/2 \quad (26)$$

A contradiction to the fact that x^{r+1} is obtained from x^r by minimizing f with respect to the s -th coordinate block. Hence

$$\delta = 0 \text{ so } f_0(\hat{z}) + f_s(\hat{z}_s) = \theta \quad (27)$$

$$f_0(z + \lambda d) + f_s(z_s + \lambda d_s) = \theta, \forall \lambda \in [0, \epsilon] \quad (28)$$

A contradiction to B2 that f is hemivariate in each block. Therefore,

$$\{x^{r+1}\}_{r \in R} \rightarrow z$$

$$\{x^{r+j}\}_{r \in R} \rightarrow z, \quad \forall j = 0, 1, \dots, T \quad (29)$$

all converge to the same value, but the sequence could be different?

With (29) and Assumption C1 or C2,

$$f_0(z) + f_k(z_k) \leq f_0(z_1, \dots, z_{k-1}, x_k, z_{k+1}, \dots, z_N) + f_k(x_k)$$

$$f_0(x^{r+j}) + f_k(x_k^{r+j}) \leq f_0(x_1^{r+1}, \dots, x_{k-1}^{r+j}, x_k^{r+j}, x_{k+1}^{r+j}, \dots, x_N^{r+j}) + f_k(x_k) \forall x_k$$

Based on the continuity of f_0 and lower-semi-continuous property of f_k , we can push the above inequality to the limit and obtain the solution.

Theorem 5.1

Suppose that f, f_0, \dots, f_N satisfy Assumptions B1-B3 and that f_0 satisfies Assumption C1 or C2. Also, assume that $\{x : f(x) \leq f(x^0)\}$ is bounded. Then the sequence $\{x^r\}$ generated by the BCD method using the essentially cyclic rule is **defined, bounded, and every cluster point is a coordinate-wise minimum point of f .**

(B1) f_0 is continuous on $\text{dom } f_0$

(B2) for each $k \in \{1, \dots, N\}$, and $(x_j)_{j \neq k}$, the function $x_k \rightarrow f(x_1, \dots, x_N)$ is quasiconvex and hemivariate.

(B3) f_0, f_1, \dots, f_N is lower semi-continuous.

(C1) $\text{dom } f_0$ is open and f_0 tends to ∞ at every boundary point of $\text{dom } f_0$

(C2) $\text{dom } f_0 = Y_1 \times \dots \times Y_N$ for some $Y_k \subseteq R^{n_k}, k = 1, \dots, N$

- Does BCD always converge on a compact subset?
- If BCD converges, are all the sequence converging to the same value?
- What if those assumption are not satisfied, could we make any conclusion?