# Chapter 7. Coordinate Descent (CD)

### 1. Cyclic CD

$$\min_{x \in \mathbb{R}^n} f(x) := f(x_1, x_2, x_3, \cdots, x_n)$$

For $k = 0, 1, \cdots$ (outer iterations)

For $i = 1, 2, \cdots, n$ (inner iterations)

$$x_i^{k+1} \leftarrow \operatorname{argmin}_{z} f(\underbrace{x_1^{k+1}, \cdots, x_{i-1}^{k+1}}_{\text{already done}}, \underset{\underset{\text{current}}{\underset{\text{coordinate}}{\uparrow}}}{z}, \underbrace{x_{i+1}^{k}, \cdots x_n^{k}}_{\text{haven't done}})$$

end

end

---

initial: $x^0 = [x_1^0, x_2^0, x_3^0, \cdots x_n^0]$

$\downarrow$

$[x_1^1, x_2^0, x_3^0, \cdots, x_n^0]$

$\downarrow$

$[x_1^1, x_2^1, x_3^0, \cdots, x_n^0]$

$\vdots$

$x^1 = [x_1^1, x_2^1, x_3^1, \cdots, x_n^1]$

$\downarrow$

$[x_1^2, x_2^1, x_3^1, \cdots, x_n^1]$

$\downarrow$

$($

- CD for constrained minimization

$$\min_x f(x) \quad st \quad x \in X_1 \otimes X_2 \times \cdots \times X_n$$

$$( x_1 \in X_1. \quad x_2 \in X_2 ), \cdots , x_n \in X_n )$$

For $k=0,1,\cdots$

For $i=1,2,\cdots,n$.

$$x^{k+1}_i \leftarrow \underset{\substack{\xi \in X_i}}{\arg\min} \ f(x^{k+1}_1, \cdots, x^{k+1}_{i-1}, \xi, x^k_{i+1}, \cdots, x^k_n)$$

end.

end

- Block Coordinate Descent $\Rightarrow$ each $X_i$ can be a $d_i$-dimensional vector )

- One variable update, equivalent to:

$$\delta \leftarrow \underset{\substack{\delta : x_i + \delta \in X_i}}{\arg\min} \ f(x + \delta e_i)$$

$$x_i \leftarrow x_i + \delta$$

current solution

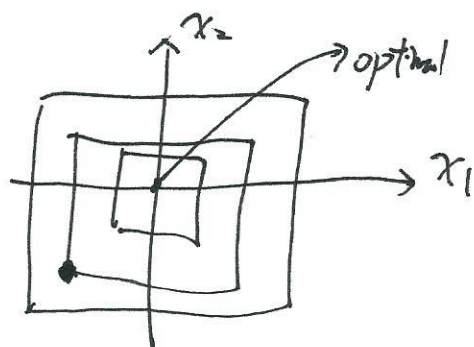$$\begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \leftarrow i\text{-th element.}$$

~~Does it~~

Does it converge to optimal solution?

Ex 1: $\min_x f(x) = \max(|x_1|, |x_2|)$
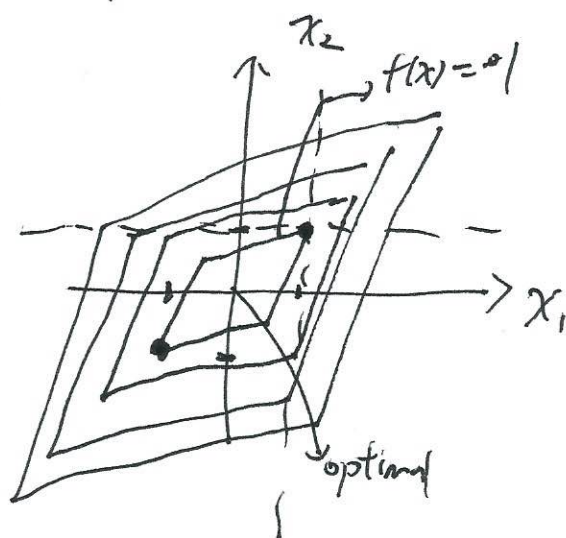


$x^0 = \begin{bmatrix} -1 \\ -1 \end{bmatrix}$ : $f(x^0) = 1$

update 1 variable

$x = \begin{bmatrix} \delta \\ -1 \end{bmatrix}$   $f(x) = \max(|\delta|, 1)$
$\geq 1$

Ex 2: $\min_x f(x) = \frac{1}{2}|x_1 + x_2| + |x_1 - x_2|$



$x^0 = \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

Cannot improve

by changing one coordinate.

---

Thm: If f is ① continuous differentiable over X

② the minimum of

$$\min_{\mathcal{z}} f(x_1, \cdots, x_{i-1}, \mathcal{z}, x_{i+1}, \cdots, x_n)$$

is uniquely attained, $\forall i$, $\forall x$

Then every limit point of $\{x^{(k)}\}_{k=1}^{\infty}$ is a stationary point.

(convex ⇒ stationary point = global optimizer)

Thm: If $f$ is continuously differentiable and there are only two blocks.

$\Rightarrow$ Every limit points are stationary points.

---

Examples:

— Linear regression.

$$\min_{x \in \mathbb{R}^n} \left\{ \frac{1}{2} \| Ax - y \|^2 \right\} = f(x)$$

$$A = \begin{bmatrix} \overbrace{a_1 \mid a_2 \mid \cdots \mid a_n}^{n} \end{bmatrix} \}m$$

$$= \begin{bmatrix} \cdots \overline{a_1} \cdots \\ \vdots \\ \overline{a_m} \end{bmatrix}$$

$$\min_{\delta} f(x + \delta e_i)$$

Define $g(\delta) = f(x + \delta e_i)$

$g'(\delta) = \nabla_i f(x + \delta e_i)$ (Want to be 0)

$$= \left( A^T (A(x + \delta e_i) - y) \right)_i$$

$$= \left( A^T (Ax - y) + A^T A \delta e_i \right)_i$$

$$= \left( \begin{bmatrix} \cdots a_1 \cdots \\ a_2 \\ \vdots \\ a_n \end{bmatrix} (Ax - y) \right)_i + \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix} a_i$$

$$= a_i^T (Ax - y) + a_i^T a_i \cdot \delta = 0$$

$$\delta^* = a_i^T (y - Ax) / \| a_i \|^2$$

For $k=0, 1, \cdots$

    For $i=1, \cdots, n$

$$x_i \leftarrow x_i + \frac{\boxed{a_i^T(y-Ax)}}{\|a_i\|^2}$$

    end

end.

compute $Ax \Rightarrow O(mn)$



compute $y - Ax \Rightarrow O(m)$

$a_i^T(y-Ax) \Rightarrow O(m)$

$\|a_i\|^2 \Rightarrow O(m)$

$\Rightarrow$ Each inner iteration:

$$O(mn)$$

---

Initial $x^0$, $r = y - Ax^0$

For $k=0, 1, \cdots$

    For $i=1, \cdots, n$

$$\begin{cases} \delta = \dfrac{a_i^T(r)}{\|a_i\|^2} \\[2mm] x_i \leftarrow x_i + \delta \\[2mm] r \leftarrow r - \delta a_i \end{cases}$$

    end

end

$O(m) \qquad O(m)$

Each inner iteration

$$\Rightarrow O(m)$$

Each outer iteration

$$O(mn)$$

GD: $x \leftarrow x - \nabla f(x) = x - A^T \underbrace{(Ax - y)}_{y_m}$

$O(mn)$

// Same time complexity

$$x^k = \begin{bmatrix} x_1^k \\ x_2^k \\ \vdots \\ x_n^k \end{bmatrix} \xrightarrow{\ GD\ } \begin{bmatrix} x_1^{k+1} \\ x_2^{k+1} \\ \vdots \\ x_n^{k+1} \end{bmatrix} \qquad \begin{bmatrix} x_1^k \\ x_2^k \\ \vdots \\ x_n^k \end{bmatrix} \rightarrow \begin{bmatrix} x_1^{k+1} \\ x_2^{k+1} \end{bmatrix}$$

Advantages :

- Each iteration is cheap (& simple)

- No step size tuning ✷

Disadvantages:

- Convergence proof is hard.

- Slow in Matlab or R. ✷

- Slow when near optimal. (depeds on functn)

3. Other variations :

3.1 Different ways to select variables :

General Form of CD:

For $k = 0, 1, \cdots$

- Pick an index $i \in \{1, \cdots, n\}$

- ~~$x_i \leftarrow \arg$~~ Update $x_i$ (by solving the one-variable problem)
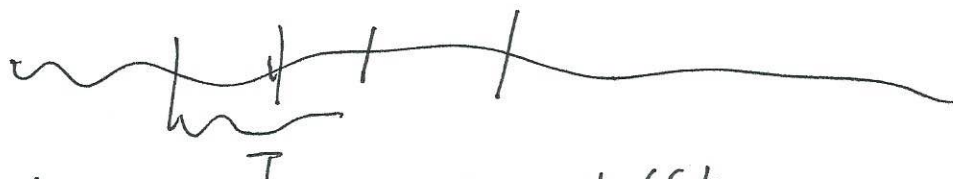
end

① Cyclic order: (in practice: very bad)

$$i = 1, 2, \cdots n, 1, 2, \cdots n$$

② Almost cyclic order:

each coordinate picked at least once
within every $T_{>0}$ iterations



Useful example: random shuffling

For $k = 0, 1, \cdots$
    Generate a random permutation $\pi$

    For $i = 1, \cdots, n$
        Update $x_{\pi(i)}$
    end
end

✓③ Random Sampling (Stochatic Coordinte Descent)

For $k = 0, 1, \cdots$
    ⓪ random sample $i \in \{1, \cdots, n\}$
    Update $x_i$
end

④ Greedy Coordinate Descent (Gauss-Southwell)

For $k = 0, 1, \cdots$
    $i \leftarrow \arg\max_i |\nabla_i f(x)| \quad \leftarrow$
    Update $x_i$
end

**Example:** Quadratic minimization with constraints.

$$\min_{x \in \mathbb{R}^n} \frac{1}{2} x^T Q x + b^T x \quad \text{st} \quad \overset{(x \geq 0)}{x \geq 0} \quad (x_i \geq 0 \ \forall i)$$

Stochastic CD:

$$\boxed{\delta \geq -x_i}$$

$$\underset{\boxed{\delta}}{\arg\min} \ f(x + \delta e_i) \quad \overset{\left[\begin{smallmatrix} 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{smallmatrix}\right]}{\phantom{.}} \quad \text{st} \quad x_i + \delta \geq 0$$

$$= \frac{1}{2}(x + \delta e_i)^T Q (x + \delta e_i) + b^T(x + \delta e_i)$$

$$= \frac{1}{2} x^T Q x + x^T Q \delta e_i + \frac{1}{2}\delta^2 e_i^T Q e_i + b^T x + \delta b^T e_i$$

$$= \frac{1}{2} Q_{ii} \delta^2 + (x^T q_i)\delta + b_i \delta + \text{const}$$

$$= \frac{1}{2} Q_{ii} \delta^2 + (b_i + x^T q_i)\delta + \text{const}.$$
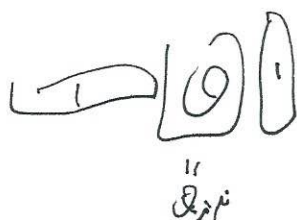
$$\ddot{\delta}_{,i}$$
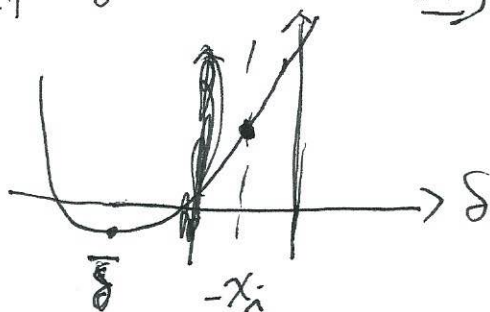
If no constraint:

$$\delta = -(b_i + x^T q_i)/Q_{ii} = \bar{\delta}$$

If $\bar{\delta} \geq -x_i \implies \delta^* = \bar{\delta}$



If $\bar{\delta} < -x_i \implies \delta^* = -x_i$

SCD for solving this

    For $k=0, 1, \ldots$

        randomly choose $i \in \{1, \ldots n\}$

        Compute $\bar{\delta} = -(b_i + x^T q_i) / Q_{ii}$

        If $\bar{\delta} \geq -x_i$

            $x_i \leftarrow x_i + \bar{\delta}$

        If $\bar{\delta} < -x_i$

            $x_i \leftarrow x_i - x_i = 0$

   end

---

$\min_{x} x^2 \quad st \quad x > 0 \quad (x \geq 0)$

Convergence for stochastic CD. $(x^k$

Thm : If $f$ is twice differentiable and $\left(\text{actually } L = \max_{\Lambda} \nabla_{ii}^2 f(x) \atop \forall x\right)$

$\checkmark \quad mI \preceq \nabla^2 f(x) \preceq LI$ , $m, L > 0$ , $\forall x$

then $\{x^k\}_{k=1}^{\infty}$ generated by SCD :

$$E[f(x^{k+1})] - f(x^*) \leq \left(1 - \frac{m}{nL}\right)\left(E[f(x^k)] - f(x^*)\right)$$

$\uparrow$
Iter k+1

$\uparrow$
iteration k

Pf: Strategy : $\nearrow^{\phi^k} \quad \nearrow^{\phi^{k+1}} \quad \frac{1}{2Ln}$

① $E[f(x^k) - f(x^{k+1})] \geq \frac{1}{2Ln} E[\|\nabla f(x^k)\|^2]$

② $f(x^k) - f(x^*) \leq \frac{1}{2m}\|\nabla f(x^k)\|^2$.

① $g(\delta) = f(x + \delta e_i)$ , $\left(\begin{array}{c} \delta^* = \operatorname*{argmin}_{\delta} g(\delta) \\ x_i^{k+1} = x_i^k + \delta^* \end{array}\right)$

$g'(\delta) = \nabla_i f(x + \delta e_i)$

$g''(\delta) = \nabla_{ii}^2 f(x + \delta e_i) \leq L$

$\Rightarrow \quad m \leq g''(\delta) \leq L$

Taylor expansion.

$$g(\delta) \leq g(0) + g'(0) \cdot \delta + \frac{L}{2}\delta^2 \qquad \left(\begin{array}{c} g'(0) + L\delta = 0 \\ \delta = -\frac{g'(0)}{L} \end{array}\right)$$

take $\tilde{\delta} = -\frac{g'(0)}{L}$

$$g(\delta^*) \leq g(\tilde{\delta}) \leq g(0) - \frac{g'(0)^2}{L} + \frac{1}{2}\left(\frac{g'(0)^2}{L}\right) = g(0) - \frac{g'(0)^2}{2L}$$

$$\rightsquigarrow f(x^{k+1}) \leq f(x^k) - \frac{(\nabla_i f(x^k))^2}{2L}$$

$$E_i\left[f(x^{k}) - f(x^{(k+1)})\right] \geq E_i\left\{\frac{\nabla_i f(x^k)^2}{2L}\right\}$$

$$= \frac{1}{2L} \cdot \left(\frac{1}{n}\sum_i \nabla_i f(x^k)^2\right)$$

$$= \frac{1}{2Lh} \cdot \|\nabla_i f(x^k)\|^2$$

② Taylor Expansion

$$f(y) \geq \underbrace{f(x) + \nabla f(x)^\top(y-x) + \frac{m}{2}\|y-x\|^2}_{h(y)} \quad \forall y, x$$

$$\tilde{y} = \underset{y}{\arg\min}\, h(y) \quad \Rightarrow \quad \nabla f(x) + m(\tilde{y} - x) = 0$$

$$\tilde{y} = x - \frac{1}{m}\nabla f(x)$$

$$h(\tilde{y}) = f(x) - \frac{\|\nabla f(x)\|^2}{m} + \frac{\|\nabla f(x)\|^2}{2m}$$

$$= f(x) - \frac{\|\nabla f(x)\|^2}{2m}$$

$$f(x^*) \geq h(x^*) \geq h(\tilde{y}) = f(x) - \frac{\|\nabla f(x)\|^2}{2m}$$

$$\Rightarrow f(x^k) - f(x^*) \leq \frac{\|\nabla f(x^k)\|^2}{2m}$$

Combining ①, ②.

$$\phi^k = E[f(x^k)], \quad \phi^* = f(x^*),$$

$$\phi^{k+1} - \phi^* = \phi^{k+1} - \phi^k + \phi^k - \phi^*$$

$$\leq -\frac{1}{2Lh}E[\|\nabla f(x^k)\|^2] + \phi^k - \phi^*$$

$$\leq -\frac{1}{2Lh} \cdot 2m \cdot \left(\underbrace{E[f(x^k)]}_{\phi^k} - \underbrace{f(x^*)}_{\phi^*}\right) + \phi^k - \phi^*$$

$$= \left(1 - \frac{m}{Lh}\right)(\phi^k - \phi^*)$$

$$GD: \quad f(x^{k+1}) - f(x^*) \leq \left(1 - \frac{m}{L}\right)\left(f(x^k) - f(x^*)\right)$$

$$SCD: \quad E[f(x^{k+1})] - f(x^*) \leq \left(1 - \frac{m}{nL}\right)\left(E[f(x^k)] - f(x^*)\right)$$

$n$ iterations of SCD $\approx$ 1 iteration for GD

$$\left(1 - \frac{m}{nL}\right)^n \overset{\approx}{\underset{\geq}{\cancel{\phantom{xxx}}}} \left(1 - \frac{m}{L}\right)$$

if $a < 1$

$(1-a)^n$

$\geq 1 - na$

$$\geq 1 - n\cdot\frac{m}{nL} = 1 - \frac{m}{L} \quad // \qquad \Rightarrow CD \approx GD$$

$(1-a)^2 = \qquad (1-a)^2 = 1 - 2a + a^2$
$0.5$

$$E_{i^{k+1}}\left[f(x^k) - f(x^{k+1})\right] \geq 0 \cdot (f(x^k) - f^*)$$

$$E_{i^{k+1}} \, \bar{E}_{i^k} \, \bar{E}_{i^{k-1}} \cdots \bar{E}_{i^0}\left(f(x) - f^*\right)$$

$$\bar{E}(f(x^k)) \; //$$