# OPTIMIZATION FOR MACHINE LEARNING
## Regularized, large-scale and distributed optimization

November 20, 2023

Today : Coordinate descent methods

Two lectures remaining:  Nov 23   Constrained optimization (V. Duval)

Dec 4   Distributed optimization (C. Royer)

Exam:  December 14  (open book)

Course project: Deadline  January 19, 2024

# COORDINATE DESCENT METHODS

## ① Basics

Context:    minimize $f(x)$ , $f: \mathbb{R}^d \to \mathbb{R}$  $C^1$
            $x \in \mathbb{R}^d$

$$\boxed{d \gg 1}$$

### Recall GD

$$x_{k+1} = x_k - \alpha_k \nabla f(x_k) \qquad \alpha_k > 0$$

$\underset{\in \mathbb{R}^d}{\uparrow \nearrow} \qquad \qquad \underset{\in \mathbb{R}^d}{\uparrow}$

At every iteration, GD:

- Computes a d-dimensional vector ($\nabla f(x_k)$)
- Updates a d-dimensional vector ($x_k$)

### Basic coordinate descent iteration

$$x_{k+1} = x_k - \alpha_k \nabla_{j_k} f(x_k) e_{j_k} \qquad \begin{array}{l} \alpha_k > 0 \\ j_k \in \{1, \dots d\} \end{array}$$

$$\nabla f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \vdots \\ \frac{\partial f(x)}{\partial x_d} \end{bmatrix}$$

$j_k$ th coordinate of $\nabla f(x_k)$

$$= \frac{\partial f}{\partial x_{j_k}}(x_k) \in \mathbb{R}$$

$$\begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \} j_k \quad \in \mathbb{R}^d$$

The iteration reduces to a scalar update

$$[x_{k+1}]_{j_k} = [x_k]_{j_k} - \alpha_k \nabla_{j_k} f(x_k) \qquad \qquad j_k \in \{1, \dots d\}$$

$\Rightarrow$ Only update 1 coordinate of the iterate at a time

$\Rightarrow$ $\nabla_{j_k} f(x_k)$ might depend on all coordinates of $x_k$ in general, but when the function is partially separable the cost of evaluating $\nabla_{j_k} f(x_k)$ can be much lower than the cost of evaluating $\nabla f(x_k)$

Def.. A function $f: \mathbb{R}^d \to \mathbb{R}$ is called separable if
$$\forall x \in \mathbb{R}^d, \quad f(x) = \sum_{j=1}^{d} f_j(x_j), \quad \text{where} \quad f_j : \mathbb{R} \to \mathbb{R}$$
and $f_j(x_j)$ only depends on $x_j$

$$\text{Ex)} \quad f(x) = \|x\|_1 = \sum_{j=1}^{d} |x_j|$$
$$f(x) = \frac{1}{2} \|x\|_2^2 = \sum_{j=1}^{d} \left( \frac{1}{2} x_j^2 \right)$$

. A function $f: \mathbb{R}^d \to \mathbb{R}$ is called <u>partially separable</u>
if $\quad f(x) = \sum_{g \in G} f_g(x_g)$, where $\quad f_g : \mathbb{R}^{|g|} \to \mathbb{R}$
$$g \subseteq \{1, \dots, d\}$$

$\underset{g \in G}{\bigcup} = \{1, \dots, d\}$ and every $f_g$ depends on

<span style="color:orange">optional but common $\longrightarrow$</span> a subset of the coordinates of $x$ $\left( x_g = \left[ x_j \right]_{j \in g} \right)$

$$\text{Ex)} \quad f(x) = \sum_{g \in G} \|x_g\|_2 \quad \text{is partially separable}$$

# Classical strategies for choosing $j_k$

a) **Cyclic coordinate descent**   Cycle through $\{1, -, d\}$ in that order

$$j_0 = 1, \; j_1 = 2, \; -, \; j_{d-1} = d, \; j_d = 1, \; -, \; j_{2d-1} = d, \; j_{2d} = 1 \ldots$$

↳ After $d$ iterations, all coordinates of the iterate have been updated

b) **Randomized cyclic coordinate descent**

↳ Every $d$ iterations, choose a random permutation of $\{1, -, d\} \longrightarrow \{\sigma(1), -, \sigma(d)\}$

↳ Choose the indices for the next $d$ iterations as $\sigma(1), \ldots, \sigma(d)$

$\Rightarrow$ For <u>separable functions</u>, $d$ iterations of $\begin{vmatrix} \text{Cyclic CD} \\ \text{Randomized cyclic CD} \end{vmatrix}$

are equivalent to $1$ iteration of GD

c) **Randomized CD:**   $j_k$ chosen at random in $\{1, -, d\}$

↳ The randomized techniques have better theoretical guarantees than cyclic CD!

# Connection between stochastic gradient and coordinate descent

- First viewpoint

$$\nabla f(x_k) = \sum_{j=1}^{d} \nabla_j f(x_k) e_j = \frac{1}{d} \sum_{j=1}^{d} \left( d \, \nabla_j f(x_k) \right) e_j$$

$$(*) \qquad \boxed{\nabla f(x_k) = \frac{1}{d} \sum_{j=1}^{d} \nabla f_j(x_k)} \qquad \text{gradient is a finite sum}$$

where we define

$$\nabla f_j(x_k) \text{ to be } \left( d \, \nabla_j f(x_j) \right) e_j$$

**Recall:** For SG, we had $\quad \nabla f(x_k) = \frac{1}{m} \sum_{i=1}^{m} \nabla f_i(x_k)$

$\longrightarrow$ Can view randomized CD as a special case of SG applied to $(*)$

- Second viewpoint

Consider a finite-sum problem

$$(?) \qquad \underset{x \in \mathbb{R}^d}{\text{minimize}} \quad \frac{1}{m} \sum_{i=1}^{m} l_i\left( a_i^T x \right) + \lambda \Omega(x)$$

where 
$a_i \in \mathbb{R}^d$ data vector

$l_i : \mathbb{R} \to \mathbb{R}$ loss function that may depend on the $i$-th data point

$\Omega : \mathbb{R}^d \to \mathbb{R}$ regularization term, $\lambda > 0$

$l_i, \Omega$ convex

The Fenchel dual (see V. Duval's lectures) of (P) is

(D)   minimize   $\frac{1}{m} \sum_{j=1}^{m} \ell_j^*(-[y]_j) + \lambda \, \Omega^*\left(\frac{1}{\lambda m} A^\top y\right)$
$y \in \mathbb{R}^m$

$\underset{\text{separable}}{\uparrow}$

where   $A^\top = [a_1 \cdots a_m] \in \mathbb{R}^{d \times m}$   ($A \in \mathbb{R}^{m \times d}$)

and $\forall$ convex function $\phi : \mathbb{R}^d \to \mathbb{R}$, the convex conjugate

of $\phi$ is       $\phi^* : \mathbb{R}^d \to \mathbb{R}$
                    $y \mapsto \sup_{z \in \mathbb{R}^d} \left( z^\top y - \phi(z) \right)$

• We can apply SG to (P) using

$$\nabla\left[ \ell_{i_k}\left(a_{i_k}^\top x_k\right) \right] = \nabla \ell_{i_k}\left(a_{i_k}^\top x_k\right) a_{i_k}$$

$\Rightarrow$ Produces a sequence $\{x_k\}$          $i_k \in \{1, \dots, m\}$

• We can also apply CD to (D) $\Rightarrow$ Produces a sequence $\{y_k\}$

using

$$\nabla\left[ \ell_{j_k}^*\left(-[y]_j\right) \right] = \nabla_{j_k}\left( \frac{1}{m} \sum_{j=1}^{m} \ell_j^*\left(-[y]_j\right) \right)$$

$j_k \in \{1, \dots, m\}$

With the same sequence of random indices   ( $\{i_k\}$ for SG
                                              $\{j_k\}$ for CD
                                              $i_k = j_k$ ),
                                              $\forall k$

then the two methods are equivalent

and       $x_k$ is equivalent to $\frac{1}{\lambda m} A^\top y_k$
              $\underset{\text{equal up to}}{\underbrace{\qquad}}$
              a constant factor

$\rightarrow$ Useful in overparameterized settings where $d \gg n$

Remark: Randomized coordinate descent is sometimes called stochastic dual descent because of this connection

# Block coordinate descent ($\sim$ Batch SG)

$$x_{k+1} = x_k - \alpha_k \sum_{j \in B_k} \nabla_j f(x_k) e_j \qquad \alpha_k > 0$$

where $B_k \subseteq \{1, -, d\}$ is a block of coordinates

⚠ Difference with SG: $B_k$ never contains duplicates of indices

$\rightarrow$ Randomized block CD: draw indices without replacement

$$|B_k| = 1 \implies \text{Randomized CD}$$
$$|B_k| = d \implies \text{GD}$$

# Proximal coordinate descent

$\hookrightarrow$ Applies to $\underset{x \in \mathbb{R}^d}{\text{minimize}} \ f(x) + \lambda \Omega(x)$

$\hookrightarrow$ Particularly interesting when $\Omega$ is separable

$$\Omega(x) = \sum_{j=1}^{d} \Omega_i(x_i)$$

$\Rightarrow$ The iteration of a proximal CD then becomes

The objective function of the subproblem is separable in $x$

$$x_{k+1} \in \underset{x \in \mathbb{R}^d}{\arg\min} \left\{ f(x_k) + \left[ \nabla_{j_k} f(x_k) e_{j_k} \right]^T (x - x_k) \right.$$
$$\left. + \frac{1}{2\alpha_k} \| x - x_k \|^2 + \lambda r_{j_k}\left( [x]_{j_k} \right) \right\}$$

replace $\nabla f(x_k)$ by a coordinate vector $\nabla_{j_k} f(x_k) e_{j_k}$

$\uparrow$ regularization w.r.t. $j_k$th coordinate

$\Rightarrow$ this can be rewritten as a 1-dimensional problem

$$[x_{k+1}]_{j_k} \in \underset{c \in \mathbb{R}}{\arg\min} \left\{ f(x_k) + \nabla_{j_k} f(x_k) \left( c - [x_k]_{j_k} \right) \right.$$
$$\left. + \frac{1}{2\alpha_k} \left( c - [x_k]_{j_k} \right)^2 + \lambda r_{j_k}(c) \right\}$$

(The other coordinates of $x_{k+1}$ are identical to those of $x_k$)

$\Rightarrow$ cheap updates, easy to extend to block CD

$\Rightarrow$ CD methods are used in sparse optimization because many sparsity-inducing regularizers are (partially) separable

2) Analysis of Coordinate descent

Focus on the basic variant

$$x_{k+1} = x_k - \alpha_k \nabla_{j_k} f(x_k) e_{j_k}$$

$$j_k \in \{1, -, d\}$$

Theorem ( Powell, 1973) :   Cyclic CD doesn't work.!

$\quad\quad\quad\quad$ ↳ Powell gave a counterexample

$\quad\quad\quad\quad d = 3$

$$f\left(x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}\right) = -\left(x_1 x_2 + x_1 x_3 + x_2 x_3\right)$$

$$+ \sum_{i=1}^{3} \max\left(|x_i| - 1, 0\right)$$

$$f \subset^1 \quad, \quad \underset{x \in \mathbb{R}^3}{\operatorname{argmin}} f(x) = \left\{ \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}, \begin{bmatrix} -1 \\ -1 \\ -1 \end{bmatrix} \right\}$$

If $\alpha_k$ is chosen through exact minimization

$\left(\text{i.e.} \quad \alpha_k = \underset{\alpha > 0}{\operatorname{argmin}} f\left(x_k - \alpha \frac{\partial f}{\partial_k}(x_k) e_{j_k}\right)\right)$

and start from $x_0 = \begin{bmatrix} -1 \\ -1 \\ 1 \end{bmatrix}$, then the method cycles between 6 points, none of which is a minimum



$x_0 \rightarrow x_1 \quad i=2$
$x_1 \rightarrow x_2 \quad i=3$
$x_2 \rightarrow x_3 \quad i=1$
$x_3 \rightarrow x_4 \quad i=2$
$x_4 \rightarrow x_5 \quad i=3$
$x_5 \rightarrow x_6 \quad i=1$

↳ Very pathological example: choosing another $x_0$ (other than the red vertices) leads to convergence
↳ But explains why CD methods were not investigated

much until late 2000s where they became of
interest in ML

$\longrightarrow$ ML point of view: • Apply CD on problems that
do not fall into the bad cases identified by
Powell (and others)

• Focus on randomized CD (because of
its ties to SG)

## Some theoretical results

We consider $\underset{x \in \mathbb{R}^d}{\text{minimize}} \; f(x)$, where $f$ is $C_L^{1,1}$

$\forall (x,y) \in (\mathbb{R}^d)^2, \; \|\nabla f(x) - \nabla f(y)\| \leq L \|x-y\|$

$\nabla f$ L-Lipschitz

$f$ is $\mu$-strongly convex

$\nabla f$ is coordinate-wise Lipschitz
continuous, i.e.

$\forall j=1..d, \; \forall (x,y) \in (\mathbb{R}^d)^2$

$$|\nabla_j f(x) - \nabla_j f(y)| \leq L_j \|x-y\|$$

$$L_j > 0$$

If $L_{max} = \underset{1 \leq j \leq d}{\max} L_j$, we have $\boxed{1 \leq \dfrac{L}{L_{max}} \leq d}$

$\varphi : \mathbb{R}^d \to \mathbb{R}^m$

$L$-Lip.

$L \approx \underset{\{(x_i, y_i)\}}{\max} \dfrac{\|\varphi(x_i) - \varphi(y_i)\|}{\|x_i - y_i\|}$

**Th)** Consider $K$ iterations of randomized CD with

- $j_k$ drawn uniformly at random in $\{1, \ldots, d\}$ $\forall k$
- $\alpha_k = \dfrac{1}{L_{j_k}}$ $\forall k$

Then

$$\underbrace{\mathbb{E}_{j_0, \ldots, j_{K-1}}}\left[ f(x_k) \right] - \min_{x \in \mathbb{R}^d} f(x) \leq \left( 1 - \frac{\mu}{d\, L_{max}} \right)^K \left( f(x_0) - \min_{x \in \mathbb{R}^d} f(x) \right)$$

$\longrightarrow 0$ as $K \to \infty$

$\uparrow$ Convergence rate in expected value (randomized CD!)

$1 - \dfrac{\mu}{d L_{max}} \in [0, 1)$

$\uparrow$

For GD, would get

$$1 - \frac{\mu}{L} \leq 1 - \frac{\mu}{d L_{max}}$$

$\Rightarrow$ Better rate in the worst case for GD

$d L_{max}$ : "price" for using a single coordinate at every iterate

$\ominus$ Worse than GD, expected value

$\oplus$ Better rate for large $d$ in terms of updates of coordinates of $x_k$, still get convergence to the optimum (unlike SG)

   <u>NB</u>: CD decreases the objective at every iteration (unlike SG)

**Th)** Consider $K$ iterations of cyclic CD under the same assumptions, then

$$f(x_k) - \min_{x \in \mathbb{R}^d} f(x) \leq \left( 1 - \frac{\mu}{2 L_{max}\left(1 + d \frac{L^2}{L_{max}}\right)} \right)^{K/d}$$

$$\times \left( f(x_0) - \min_{x \in \mathbb{R}^d} f(x) \right)$$

⊖ Rate of CV is worse than GD

$$\left(1 - \frac{\mu}{L}\right)^K \qquad vs \qquad \left(1 - \frac{\mu}{\ldots}\right)^{K/d}$$

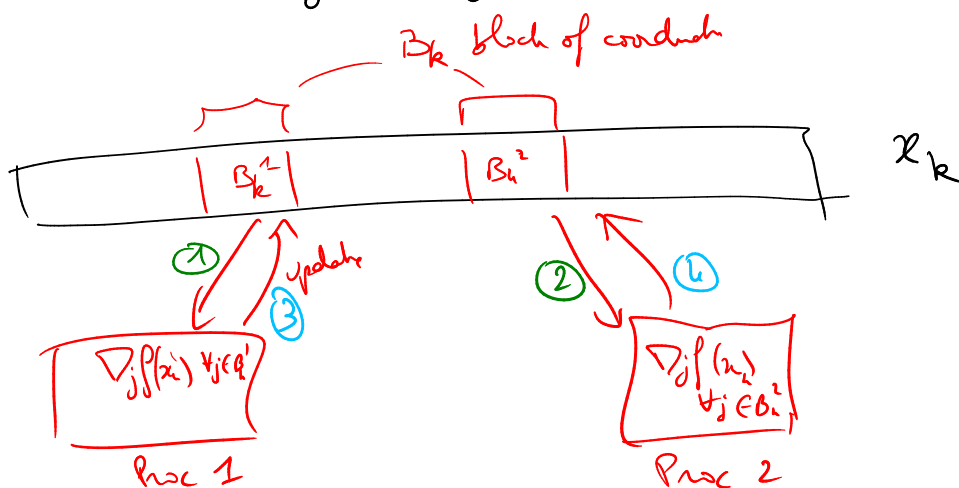⟹ Even in terms of coordinate updates, the rate will not be better than that of gradient descent

↳ These results can be improved on specific problems which is usually how CD methods are analyzed in ML

3) CD for parallel/distributed optimization

Setup: $x_k \in \mathbb{R}^d$ is too large to be updated via GD
⟹ Stored in memory

Iteration $k$ of CD
⟶ Distribute steps along (blocks of) coordinates to different processors
⟶ Every processor updates their coordinates
⟶ Synchronization after every update



$B_k$ block of coordinates

$x_k$

$B_k^1$     $B_k^2$

①   update ③     ②   ④

$\nabla_j f(x_k) \ \forall j \in B_k^1$     $\nabla_j f(x_k) \ \forall j \in B_k^2$

Proc 1     Proc 2

① ② } run in parallel

③ ④ } run in parallel

→ Synchronization: wait until ③ and ④ are completed before sending the coordinates again

⟹ Because of synchronization, this process does not bring a lot of benefit

⟹ Still efficient when the problem is separable (synchronizat° is not an issue)

# Asynchronous CD : Do not wait for the other processors!

- Initialization: $x_0 \in \mathbb{R}^d$ stored, shared iteration count $k=0$.

- Repeat loop (for all processors)

  choose $j_k \in \{1, -, d\}$

  $$[x_{k+1}]_{j_k} \leftarrow [x_k]_{j_k} - \alpha_k \nabla_{j_k} f(\hat{x}_k)$$

  $\hat{x}_k$: value of $x_k$ when $\nabla_{j_k} f(x_k)$ was computed

  $\alpha_k > 0$

  Asynchronous update

  $k \leftarrow k+1$

  Synchronization only for $k$

↳ Surprisingly, this can work!

⟹ For convex $f$, if
  - every coordinate is updated infinitely often    of $x_k$
  - every coordinate of $\hat{x}_k$ is updated infinitely often

  then $x_k \xrightarrow[k \to \infty]{} x^* \in \underset{x \in \mathbb{R}^d}{\operatorname{argmin}} f(x)$

↳ Follow-up of asynchronous CD: Asynchronous SG!

⟹ Hogwild! : Asynchronous SG with theory and good practical success

⟹ NeurIPS 2011, won the test-of-time award in 2020