

Modern Optimization Techniques

2. Unconstrained Optimization / 2.6. Coordinate Descent

Lars Schmidt-Thieme

Information Systems and Machine Learning Lab (ISMLL)
Institute of Computer Science
University of Hildesheim, Germany

original slides by Lucas Rego Drumond (ISMLL)

Syllabus

Tue. 18.10.	(0)	0. Overview
		1. Theory
Tue. 25.10.	(1)	1. Convex Sets and Functions
		2. Unconstrained Optimization
Tue. 1.11	(2)	2.1 Gradient Descent
Tue. 8.11.	(3)	2.2 Stochastic Gradient Descent
Tue. 15.11.	(4)	(ctd.)
Tue. 22.11.	(5)	2.3 Newton's Method
Tue. 29.11.	(6)	2.4 Quasi-Newton Methods
Tue. 6.12.	(7)	2.5 Subgradient Methods
Tue. 13.12.	(8)	2.6 Coordinate Descent
		3. Equality Constrained Optimization
Tue. 20.12.	(9)	3.1 Duality
	—	— Christmas Break —
Tue. 10.1.	(10)	3.2 Methods
		4. Inequality Constrained Optimization
Tue. 17.1.	(11)	4.1 Interior Point Methods
Tue. 24.1.	(11)	4.2 Cutting Plane Method
		5. Distributed Optimization
Tue. 31.1.	(12)	5.1 Alternating Direction Method of Multipliers

Outline

1. Coordinate Descent

2. Convergence

3. Examples

Outline

1. Coordinate Descent

2. Convergence

3. Examples

Coordinate Descent

- ▶ Gradient Descent and Stochastic Gradient Descent:
 - ▶ update **all** variables simultaneously.
 - ▶ use the gradient to do so.
(first order methods)
- ▶ **Coordinate Descent:**
 - ▶ update **one** variable at a time.
 - ▶ use an **analytic solver** to do so
(derivative-free method)
 - ▶ if not possible: use one-dimensional gradient steps
(first order method; often slow)

Coordinate Descent

We start with an initial guess $\mathbf{x}^{(0)}$

For $k = 1, 2, 3, \dots$:

Coordinate Descent

We start with an initial guess $\mathbf{x}^{(0)}$

For $k = 1, 2, 3, \dots$:

$$\mathbf{x}_1^{(k)} \leftarrow \arg \min_{\mathbf{x}_1} f(\mathbf{x}_1, \mathbf{x}_2^{(k-1)}, \mathbf{x}_3^{(k-1)}, \dots, \mathbf{x}_n^{(k-1)})$$

Coordinate Descent

We start with an initial guess $\mathbf{x}^{(0)}$

For $k = 1, 2, 3, \dots$:

$$\mathbf{x}_1^{(k)} \leftarrow \arg \min_{\mathbf{x}_1} f(\mathbf{x}_1, \mathbf{x}_2^{(k-1)}, \mathbf{x}_3^{(k-1)}, \dots, \mathbf{x}_n^{(k-1)})$$

$$\mathbf{x}_2^{(k)} \leftarrow \arg \min_{\mathbf{x}_2} f(\mathbf{x}_1^{(k)}, \mathbf{x}_2, \mathbf{x}_3^{(k-1)}, \dots, \mathbf{x}_n^{(k-1)})$$

Coordinate Descent

We start with an initial guess $\mathbf{x}^{(0)}$

For $k = 1, 2, 3, \dots$:

$$\mathbf{x}_1^{(k)} \leftarrow \arg \min_{\mathbf{x}_1} f(\mathbf{x}_1, \mathbf{x}_2^{(k-1)}, \mathbf{x}_3^{(k-1)}, \dots, \mathbf{x}_n^{(k-1)})$$

$$\mathbf{x}_2^{(k)} \leftarrow \arg \min_{\mathbf{x}_2} f(\mathbf{x}_1^{(k)}, \mathbf{x}_2, \mathbf{x}_3^{(k-1)}, \dots, \mathbf{x}_n^{(k-1)})$$

$$\mathbf{x}_3^{(k)} \leftarrow \arg \min_{\mathbf{x}_3} f(\mathbf{x}_1^{(k)}, \mathbf{x}_2^{(k)}, \mathbf{x}_3, \dots, \mathbf{x}_n^{(k-1)})$$

Coordinate Descent

We start with an initial guess $\mathbf{x}^{(0)}$

For $k = 1, 2, 3, \dots$:

$$\mathbf{x}_1^{(k)} \leftarrow \arg \min_{\mathbf{x}_1} f(\mathbf{x}_1, \mathbf{x}_2^{(k-1)}, \mathbf{x}_3^{(k-1)}, \dots, \mathbf{x}_n^{(k-1)})$$

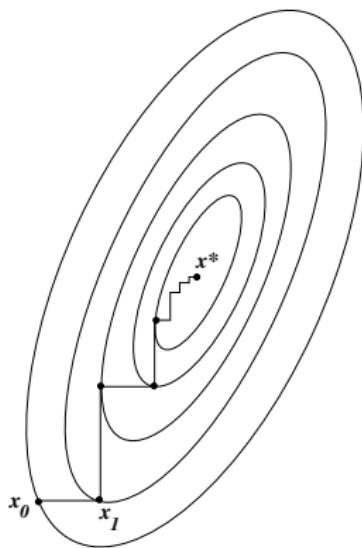
$$\mathbf{x}_2^{(k)} \leftarrow \arg \min_{\mathbf{x}_2} f(\mathbf{x}_1^{(k)}, \mathbf{x}_2, \mathbf{x}_3^{(k-1)}, \dots, \mathbf{x}_n^{(k-1)})$$

$$\mathbf{x}_3^{(k)} \leftarrow \arg \min_{\mathbf{x}_3} f(\mathbf{x}_1^{(k)}, \mathbf{x}_2^{(k)}, \mathbf{x}_3, \dots, \mathbf{x}_n^{(k-1)})$$

⋮

$$\mathbf{x}_n^{(k)} \leftarrow \arg \min_{\mathbf{x}_n} f(\mathbf{x}_1^{(k)}, \mathbf{x}_2^{(k)}, \mathbf{x}_3^{(k)}, \dots, \mathbf{x}_n)$$

Coordinate Descent Algorithm



[Nocedal and Wright, 2006, p.249]

Does Coordinate-wise Minimization Lead to Global Minima?

Question: If a point \mathbf{x} is minimal along each axis, is it a global minimum?

$$\begin{aligned} f(\mathbf{x} + t \mathbf{e}^{(n)}) &\geq f(\mathbf{x}) \quad \forall t \in \mathbb{R}, \forall n \in \{1, \dots, N\} \\ \xrightarrow{?} f(\mathbf{x}) &= \min_{\mathbf{y}} f(\mathbf{y}) \end{aligned}$$

where:

- $\mathbf{e}^{(n)} \in \mathbb{R}^N$ is the n -th unit vector with $\mathbf{e}_n^{(n)} := 1$ and $\mathbf{e}_m^{(n)} := 0$ for $m \neq n$.

Does Coordinate-wise Minimization Lead to Global Minima?

If f is **convex** and **differentiable**: yes.

Proof: $g^{(i)}(\mu) := f(\mathbf{x} + \mu \mathbf{e}^{(i)})$ are convex and differentiable.
 $\mu = 0$ is their minimum, thus

$$0 = \frac{\partial g^{(i)}}{\partial \mu}(\mu) = \frac{\partial f}{\partial \mathbf{x}_i}(\mathbf{x})$$

And then

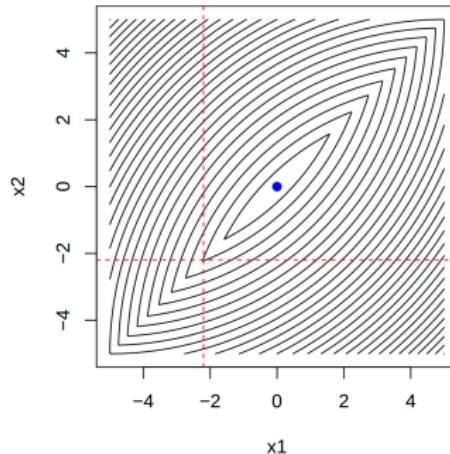
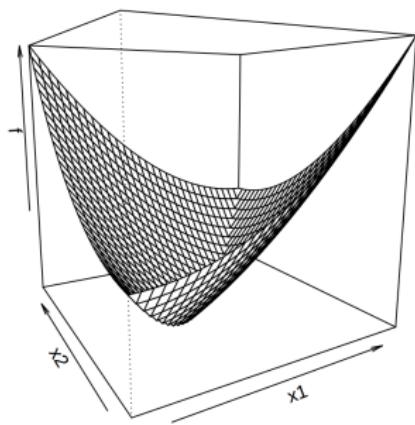
$$\nabla f(\mathbf{x}) = \left(\frac{\partial f(\mathbf{x})}{\partial \mathbf{x}_1}, \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}_2}, \dots, \frac{\partial f(\mathbf{x})}{\partial \mathbf{x}_n} \right)^T = \mathbf{0}$$

and thus x is the global optimum.

Does Coordinate-wise Minimization Lead to Global Minima?

If f is **convex**, but **not differentiable**: in general, no.

Counter example:



<https://www.cs.cmu.edu/~ggordon/10725-F12/slides/25-coord-desc.pdf>

Does Coordinate-wise Minimization Lead to Global Minima?

If f is a sum of a differentiable convex and a separable convex function:
yes.

$$f(\mathbf{x}) = g(\mathbf{x}) + \sum_{n=1}^N h_n(\mathbf{x}_n)$$

with

- i. g is differentiable and convex
- ii. all h_n are convex

Does Coordinate-wise Minimization Lead to Global Minima? / Proof

Proof: For any $\mathbf{y} \in \text{dom } f$:

$$\begin{aligned} f(\mathbf{y}) - f(\mathbf{x}) &\geq \sum_i \nabla g(\mathbf{x})^T (\mathbf{y}_i - \mathbf{x}_i) + \sum_{n=1}^N (h_n(\mathbf{y}_n) - h_n(\mathbf{x}_n)) \\ &= \sum_{n=1}^N \underbrace{\nabla g(\mathbf{x})(\mathbf{y}_n - \mathbf{x}_n)}_{\geq 0} + h_n(\mathbf{y}_n) - h_n(\mathbf{x}_n) \geq 0 \end{aligned}$$

where \geq is argued as follows: as x is minimal along axis n ,

- $\rightsquigarrow f(X_n; X_{-n})$ is convex and has minimum at $X_n = x_n$
- $\rightsquigarrow 0$ is a subgradient of $f(X_n; X_{-n})$ at x_n :
- $0 \in \partial f(X_n; X_{-n}) = \nabla g(x) + \partial h_n(x_n)$
- $\rightsquigarrow \exists s \in \partial h_n(x_n) : \nabla g(x) + s = 0$
- $$\begin{aligned} &\nabla g(x)(\mathbf{y}_n - \mathbf{x}_n) + h_n(\mathbf{y}_n) - h_n(\mathbf{x}_n) \\ &\geq \nabla g(x)(\mathbf{y}_n - \mathbf{x}_n) + s(\mathbf{y}_n - \mathbf{x}_n) = 0 \end{aligned}$$

One-dimensional Subproblems

Let $f : X \rightarrow \mathbb{R}$, $X \subseteq \mathbb{R}^N$ be a function, $x \in \text{dom } f$ a point.

n -th one-dimensional subproblem at x :

$$g_n^{(x)} : T_n^{(x)} \rightarrow \mathbb{R}$$

$$g_n^{(x)}(t) := f(x + t e^{(n)})$$

$$T_n^{(x)} := \{t \in \mathbb{R} \mid x + t e^{(n)} \in \text{dom } f\}$$

n -th one-dimensional subproblem solver at x :

$$h_n(x) := \arg \min_{t \in T_n^{(x)}} g_n^{(x)}(t)$$

One-dimensional Subproblems / Example

Solving a linear system of equations / least squares / linear regression:

$$f(x) := x^T A x - b^T x, \quad A \in \mathbb{R}^{N \times N} \text{ sym. pos.def., } b \in \mathbb{R}^N$$

$$g_n^{(x)}(t) = f(x_n = t; x_{-n}), \quad T_n^{(x)} = \mathbb{R}$$

$$\begin{aligned} f(x_n; x_{-n}) &= A_{n,n}x_n^2 + (2(x_{-n}^T A_{-n,.})_n - b_n)x_n + x_{-n}^T A_{-n,-n}x_{-n} - b_{-n}^T x_{-n} \\ &= A_{n,n}x_n^2 + (2A_{n,-n}x_{-n} - b_n)x_n + x_{-n}^T A_{-n,-n}x_{-n} - b_{-n}^T x_{-n} \end{aligned}$$

analytic minimum:

$$f'(x_n; x_{-n}) = 2A_{n,n}x_n + 2A_{n,-n}x_{-n} - b_n \stackrel{!}{=} 0$$

$$\rightsquigarrow x_n = \frac{b_n - 2A_{n,-n}x_{-n}}{2A_{n,n}}$$

$$\text{i.e., } h_n(x) := \frac{b_n - 2A_{n,-n}x_{-n}}{2A_{n,n}}$$

Coordinate Descent / Random Coordinate

```

1 min-cd( $f, h, x^{(0)}, K, \epsilon$ ):
2   for  $k := 1, \dots, K$ :
3     draw  $n^{(k)} \sim \text{unif}(\{1, \dots, N\})$ 
4      $x_m^{(k)} := x_m^{(k-1)}$  for  $m \in \{1, \dots, N\}, m \neq n^{(k)}$ 
5      $x_{n^{(k)}}^{(k)} := h_{n^{(k)}}(x^{(k-1)})$ 
6     if  $k \geq N$  and  $\|x^{(k)} - x^{(k-N)}\| \leq \epsilon$ 
7       return  $x^{(k)}$ 
8   return "not converged"

```

where

- ▶ h solvers for one-dimensional subproblems $g_n^{(x)}$.
- ▶ ϵ convergence threshold step in an epoch

Coordinate Descent / Cyclic Epochs

```

1 min-cd( $f, h, x^{(0)}, K, \epsilon$ ):
2   for  $k := 1, \dots, K$ :
3      $x_1^{(k)} := h_1(x_1^{(k-1)}, \dots, x_N^{(k-1)})$ 
4      $x_2^{(k)} := h_2(x_1^{(k)}, x_2^{(k-1)}, \dots, x_N^{(k-1)})$ 
5      $x_3^{(k)} := h_3(x_1^{(k)}, x_2^{(k)}, x_3^{(k-1)}, \dots, x_N^{(k-1)})$ 
6     :
7      $\vdots$ 
8      $x_N^{(k)} := h_N(x_1^{(k)}, \dots, x_{N-1}^{(k)}, x_N^{(k-1)})$ 
9     if  $\|x^{(k)} - x^{(k-1)}\| \leq \epsilon$ 
10    return  $x^{(k)}$ 
11  return "not converged"

```

where

- ▶ h solvers for one-dimensional subproblems $g_n^{(x)}$.
- ▶ ϵ convergence threshold step in an epoch

Coordinate Descent / Cyclic Epochs

```

1 min-cd( $f, h, x^{(0)}, K, \epsilon$ ):
2    $x^{(0,N)} := x^{(0)}$ 
3   for  $k := 1, \dots, K$ :
4      $x^{(k,0)} := x^{(k-1,N)}$ 
5     for  $n := 1, \dots, N$ :
6        $x_m^{(k,n)} := x_m^{(k,n-1)}$  for  $m \in \{1, \dots, N\}, m \neq n$ 
7        $x_n^{(k,n)} := h_n(x^{(k,n-1)})$ 
8     if  $\|x^{(k,N)} - x^{(k-1,N)}\| \leq \epsilon$ 
9       return  $x^{(k,N)}$ 
10  return "not converged"

```

where

- ▶ h solves for one-dimensional subproblems $g_n^{(x)}$.
- ▶ ϵ convergence threshold step in an epoch

Coordinate Descent / 1-dim Gradient Steps

```
1 min-cd( $f, \nabla f, x^{(0)}, K, \epsilon$ ):
2   for  $k := 1, \dots, K$ :
3     draw  $n^{(k)} \sim \text{unif}(\{1, \dots, N\})$ 
4      $x^{(k)} := x^{(k-1)} - \mu^{(k)} (\nabla f(x^{(k-1)}))_{n^{(k)}} e^{(n^{(k)})}$ 
5     if  $k \geq N$  and  $\|x^{(k)} - x^{(k-N)}\| \leq \epsilon$ 
6       return  $x^{(k)}$ 
7   return "not converged"
```

where

- ▶ ∇f gradients of objective function.
- ▶ $\mu \in (\mathbb{R}^+)^*$ step length schedule (or controller)
- ▶ ϵ convergence threshold step in an epoch

Coordinate Descent - Considerations

- ▶ The order in which we cycle through the coordinates is arbitrary.
 - ▶ e.g., cyclic
 - ▶ better should be randomized.
- ▶ We may update blocks of coordinates at a time instead of only one
(block coordinate descent)
- ▶ No need to adjust a step-size!
 - ▶ if we have exact solvers h for the 1-dim. subproblems.
- ▶ Does not work in general with non-differentiable functions

Outline

1. Coordinate Descent

2. Convergence

3. Examples

Coordinate Lipschitz Constant

- standard Lipschitz constant L :

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\| \quad \forall x, y \in \text{dom } f$$

or equivalently

$$\|\nabla f(x + d) - \nabla f(x)\| \leq L\|d\| \quad \forall x \in \text{dom } f, d \in \mathbb{R}^N : x + d \in \text{dom } f$$

- component Lipschitz constant L_n : ($n \in \{1, \dots, N\}$)

$$\|\nabla f(x + t e^{(n)}) - \nabla f(x)\| \leq |t| L_n \quad \forall x \in \text{dom } f, t \in \mathbb{R} : x + t e^{(n)} \in \text{dom } f$$

- coordinate Lipschitz constant L_{\max} :

$$L_{\max} := \max_{n \in \{1, \dots, N\}} L_n$$

Note: In the following $\|\cdot\|$ denotes the L2-norm.

Lipschitz Continuous Functions / Bounded Derivative

Lemma

A differentiable function $f : X \rightarrow \mathbb{R}$, $X \subseteq \mathbb{R}$ is Lipschitz-continuous iff its derivative is bounded:

$$\begin{aligned} |f(x) - f(y)| \leq L|x - y| &\iff |f'(x)| \leq L \\ \forall x, y \in \text{dom } f & \quad \forall x \in \text{dom } f \end{aligned}$$

Proof: “ \Rightarrow ”:

$$\begin{aligned} \left| \frac{f(x) - f(x + t)}{|t|} \right| &= \frac{|f(x) - f(x + t)|}{|t|} \leq \frac{|t|L}{|t|} \leq L \quad |t| \rightarrow 0 \\ |f'(x)| &\leq L \end{aligned}$$

“ \Leftarrow ”:

$$\begin{aligned} f(x) &= f(y) + f'(\xi)(x - y) \quad \text{for a } \xi \in [x, y] \\ &\leq f(y) + L(x - y) \end{aligned}$$

$$|f(x) - f(y)| \leq L|x - y|$$

Theorem (convergence of coordinate descent)

If

- i. $f : X \rightarrow \mathbb{R}$, $X \subseteq \mathbb{R}^N$ is convex and differentiable,
- ii. ∇f is uniformly Lipschitz-continuous:

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|, \quad L \in \mathbb{R}_0^+$$

- iii. the sublevel set of $x^{(0)}$ is bounded:

$$\max_{x^*: f(x^*) = \min_x f(x)} \max_{x \in S_f(x^{(0)})} \|x - x^*\| \leq R, \quad R \in \mathbb{R}_0^+$$

- iv. constant steplength $\mu^{(k)} := 1/L_{\max}$ is used,

then coordinate descent converges and

$$\mathbf{E}(f(x^{(k)})) - f(x^*) \leq \frac{2NL_{\max}R^2}{k}$$

Convergence / Proof (1/3)

$$f(x^{\text{next}}) = f(x - \mu(\nabla f(x))_n e^{(n)})$$

$$\stackrel{\text{Taylor}}{=} f(x) - \nabla f(x)^T \mu(\nabla f(x))_n e^{(n)}$$

$$+ \frac{1}{2} (\mu(\nabla f(x))_n e^{(n)})^T \underbrace{\nabla^2 f(x - \xi \mu(\nabla f(x))_n e^{(n)})}_{\stackrel{\text{ii. bound.deriv.}}{\leq}} \mu(\nabla f(x))_n e^{(n)}$$

$$\leq f(x) - (\nabla f(x))_n^2 \mu + \frac{1}{2} \mu^2 (\nabla f(x))_n^2 L_n$$

$$= f(x) - (\nabla f(x))_n^2 \mu \left(1 - \frac{1}{2} \mu L_n\right)$$

$$\leq f(x) - (\nabla f(x))_n^2 \mu \left(1 - \frac{1}{2} \mu L_{\max}\right)$$

$$= f(x) - \frac{(\nabla f(x))_n^2}{2L_{\max}}$$

iv.

Convergence / Proof (2/3)

$$\begin{aligned}
 f(x^{\text{next}}) &\leq f(x) - \frac{(\nabla f(x))_n^2}{2L_{\max}} \quad |\mathbb{E}_n(\dots) \\
 \mathbb{E}_n(f(x^{\text{next}})) &\leq \mathbb{E}_n(f(x) - \frac{(\nabla f(x))_n^2}{2L_{\max}}) \\
 &= f(x) - \frac{1}{N} \sum_{n=1}^N \frac{(\nabla f(x))_n^2}{2L_{\max}} = f(x) - \frac{1}{2NL_{\max}} \|\nabla f(x)\|^2 \quad (1)
 \end{aligned}$$

$\mathbb{E}_{n^{(0:K)}} := \mathbb{E}_{n^{(0)}, n^{(1)}, \dots, n^K}$ expectation over **all** $n^{(k)}$

$$\phi^{(k)} := \mathbb{E}_{n^{(0:K)}}(f(x^{(k)})) - f(x^*)$$

$$f(x) - f(x^*) \stackrel{i.}{\leq} \nabla f(x)^T (x - x^*) \leq \|\nabla f(x)\| \|x - x^*\| \stackrel{iii.}{\leq} R \|\nabla f(x)\|$$

$$\mathbb{E}_{n^{(0:K)}}(\|\nabla f(x^{(k)})\|) \geq \frac{1}{R} \phi^{(k)} \quad (2)$$

Convergence / Proof (3/3)

$$\begin{aligned}
 & \mathbb{E}_{n^{(k)}}(f(x^{(k+1)})) \stackrel{(1)}{\leq} f(x^{(k)}) - \frac{1}{2NL_{\max}} \|\nabla f(x^{(k)})\|^2 \quad | \mathbb{E}_{n^{(0:k)}} \\
 & \phi^{(k+1)} \leq \phi^{(k)} - \frac{1}{2NL_{\max}} \mathbb{E}_{n^{(0:k)}}(\|\nabla f(x^{(k)})\|^2) \\
 & \stackrel{\text{Jensen's Ineq.}}{\leq} \phi^{(k)} - \frac{1}{2NL_{\max}} (\mathbb{E}_{n^{(0:k)}}(\|\nabla f(x^{(k)})\|))^2 \\
 & \stackrel{(2)}{\leq} \phi^{(k)} - \frac{1}{2NL_{\max}R^2} (\phi^{(k)})^2 \tag{3}
 \end{aligned}$$

$$\begin{aligned}
 \frac{1}{\phi^{(k+1)}} - \frac{1}{\phi^{(k)}} &= \frac{\phi^{(k)} - \phi^{(k+1)}}{\phi^{(k)}\phi^{(k+1)}} \geq \frac{\phi^{(k)} - \phi^{(k+1)}}{(\phi^{(k)})^2} \stackrel{(3)}{\geq} \frac{1}{2NL_{\max}R^2} \\
 \frac{1}{\phi^{(k+1)}} &\stackrel{\text{rec}}{\geq} \frac{1}{\phi^{(0)}} + \frac{k+1}{2NL_{\max}R^2} \geq \frac{k+1}{2NL_{\max}R^2}
 \end{aligned}$$

Outline

1. Coordinate Descent

2. Convergence

3. Examples

Coordinate Descent - Example

For $\mathbf{x} \in \mathbb{R}^2$

$$\min_{\mathbf{x}} (a_1x_1 - a_2x_2 + a_3)^2$$

Algorithm:

- ▶ Initialize $\mathbf{x}_1^{(0)}, \mathbf{x}_2^{(0)}$
- ▶ Repeat until convergence:
 - ▶ $\mathbf{x}_1^{(k)} \leftarrow \arg \min_{\mathbf{x}_1} (a_1\mathbf{x}_1 - a_2\mathbf{x}_2^{(k-1)} + a_3)^2$
 - ▶ $\mathbf{x}_2^{(k)} \leftarrow \arg \min_{\mathbf{x}_2} (a_1\mathbf{x}_1^{(k)} - a_2\mathbf{x}_2 + a_3)^2$

Coordinate Descent - Example

$$\arg \min_{x_1} (a_1 x_1 - a_2 x_2 + a_3)^2$$

Find a closed form solution:

$$0 = \frac{d}{dx_1} (a_1 x_1 - a_2 x_2 + a_3)^2$$

Coordinate Descent - Example

$$\arg \min_{x_1} (a_1 x_1 - a_2 x_2 + a_3)^2$$

Find a closed form solution:

$$0 \stackrel{!}{=} \frac{d}{dx_1} (a_1 x_1 - a_2 x_2 + a_3)^2 = 2(a_1 x_1 - a_2 x_2 + a_3) a_1$$

Coordinate Descent - Example

$$\arg \min_{x_1} (a_1 x_1 - a_2 x_2 + a_3)^2$$

Find a closed form solution:

$$0 \stackrel{!}{=} \frac{d}{dx_1} (a_1 x_1 - a_2 x_2 + a_3)^2 = 2(a_1 x_1 - a_2 x_2 + a_3) a_1$$

$$x_1 = \frac{a_2 x_2 - a_3}{a_1}$$

Coordinate Descent - Example

$$\arg \min_{x_2} (a_1 x_1 - a_2 x_2 + a_3)^2$$

Find a closed form solution:

$$0 \stackrel{!}{=} \frac{d}{dx_2} (a_1 x_1 - a_2 x_2 + a_3)^2$$

Coordinate Descent - Example

$$\arg \min_{x_2} (a_1 x_1 - a_2 x_2 + a_3)^2$$

Find a closed form solution:

$$0 \stackrel{!}{=} \frac{d}{dx_2} (a_1 x_1 - a_2 x_2 + a_3)^2 = -2(a_1 x_1 - a_2 x_2 + a_3) a_2$$

Coordinate Descent - Example

$$\arg \min_{x_2} (a_1 x_1 - a_2 x_2 + a_3)^2$$

Find a closed form solution:

$$0 \stackrel{!}{=} \frac{d}{dx_2} (a_1 x_1 - a_2 x_2 + a_3)^2 = -2(a_1 x_1 - a_2 x_2 + a_3) a_2$$
$$x_2 = \frac{a_1 x_1 + a_3}{a_2}$$

Coordinate Descent - Example

For $\mathbf{x} \in \mathbb{R}^2$

$$\min_{\mathbf{x}} (a_1x_1 - a_2x_2 + a_3)^2$$

Algorithm:

- ▶ Initialize $\mathbf{x}_1^{(0)}, \mathbf{x}_2^{(0)}$
- ▶ Repeat until convergence:

$$\mathbf{x}_1^{(k)} \leftarrow \frac{a_2\mathbf{x}_2^{(k-1)} - a_3}{a_1}$$

$$\mathbf{x}_2^{(k)} \leftarrow \frac{a_1\mathbf{x}_1^{(k)} + a_3}{a_2}$$

Coordinate Descent - Example

For $\mathbf{x} \in \mathbb{R}^2$, $a_1 = 0.1$, $a_2 = 2$, $a_3 = 1$

$$\min_{\mathbf{x}} (0.1x_1 - 2x_2 + 1)^2$$
$$\mathbf{x}_1^{(k)} \leftarrow \frac{2\mathbf{x}_2^{(k-1)} - 1}{0.1} \quad \mathbf{x}_2^{(k)} \leftarrow \frac{0.1\mathbf{x}_1^{(k)} + 1}{2}$$

Coordinate Descent - Example

For $\mathbf{x} \in \mathbb{R}^2$, $a_1 = 0.1$, $a_2 = 2$, $a_3 = 1$

$$\min_{\mathbf{x}} (0.1x_1 - 2x_2 + 1)^2$$

$$\mathbf{x}_1^{(k)} \leftarrow \frac{2\mathbf{x}_2^{(k-1)} - 1}{0.1} \quad \mathbf{x}_2^{(k)} \leftarrow \frac{0.1\mathbf{x}_1^{(k)} + 1}{2}$$

Start with $\mathbf{x}_1^{(0)} = 1$, $\mathbf{x}_2^{(0)} = 2$

- ▶ $\mathbf{x}_1^{(1)} \leftarrow \frac{2 \cdot 2 - 1}{0.1} = 30$
- ▶ $\mathbf{x}_2^{(1)} \leftarrow \frac{0.1 \cdot 30 + 1}{2} = 2$

Coordinate Descent - Example

For $\mathbf{x} \in \mathbb{R}^2$, $a_1 = 0.1$, $a_2 = 2$, $a_3 = 1$

$$\min_{\mathbf{x}} (0.1x_1 - 2x_2 + 1)^2$$

$$\mathbf{x}_1^{(k)} \leftarrow \frac{2\mathbf{x}_2^{(k-1)} - 1}{0.1} \quad \mathbf{x}_2^{(k)} \leftarrow \frac{0.1\mathbf{x}_1^{(k)} + 1}{2}$$

Start with $\mathbf{x}_1^{(0)} = 1$, $\mathbf{x}_2^{(0)} = 2$

- ▶ $\mathbf{x}_1^{(1)} \leftarrow \frac{2 \cdot 2 - 1}{0.1} = 30$
- ▶ $\mathbf{x}_2^{(1)} \leftarrow \frac{0.1 \cdot 30 + 1}{2} = 2$

$$\mathbf{x}_1^{(1)} = 30, \mathbf{x}_2^{(1)} = 2$$

- ▶ $\mathbf{x}_1^{(2)} \leftarrow \frac{2 \cdot 2 - 1}{0.1} = 30$
- ▶ $\mathbf{x}_2^{(2)} \leftarrow \frac{0.1 \cdot 30 + 1}{2} = 2$

Coordinate Descent for Linear Regression

For the problem

$$\text{minimize} \quad \sum_{i=1}^m (y_i - \mathbf{x}^T \mathbf{a}_i)^2 = \sum_{i=1}^m \left(y_i - \sum_{j=1}^n x_j a_{ij} \right)^2$$

We can compute the update rule for a specific x_k :

$$\frac{\partial f(\mathbf{x})}{\partial x_k} \stackrel{!}{=} 0$$

$$\frac{\partial f(\mathbf{x})}{\partial x_k} = 2 \cdot \sum_{i=1}^m a_{ik} \cdot \left(y_i - \sum_{j=1}^n x_j a_{ij} \right)$$

Coordinate Descent for Linear Regression

$$2 \cdot \sum_{i=1}^m a_{ik} \cdot \left(y_i - \sum_{j=1}^n x_j a_{ij} \right) = 0$$

$$\sum_{i=1}^m a_{ik} \cdot y_i - \sum_{i=1}^m a_{ik} \cdot \sum_{j=1}^n x_j a_{ij} = 0$$

$$\sum_{i=1}^m a_{ik} \cdot y_i - \sum_{i=1}^m a_{ik} \cdot \left(x_k a_{ik} + \sum_{j=1, j \neq k}^n x_j a_{ij} \right) = 0$$

$$\sum_{i=1}^m a_{ik} \cdot y_i - \sum_{i=1}^m a_{ik} \cdot x_k a_{ik} - \sum_{i=1}^m a_{ik} \cdot \sum_{j=1, j \neq k}^n x_j a_{ij} = 0$$

$$\sum_{i=1}^m a_{ik} \cdot y_i - x_k \sum_{i=1}^m a_{ik}^2 - \sum_{i=1}^m a_{ik} \cdot \sum_{j=1, j \neq k}^n x_j a_{ij} = 0$$

Coordinate Descent for Linear Regression

$$\sum_{i=1}^m a_{ik} \cdot y_i - x_k \sum_{i=1}^m a_{ik}^2 - \sum_{i=1}^m a_{ik} \cdot \sum_{j=1, j \neq k}^n x_j a_{ij} = 0$$

$$x_k \cdot \sum_{i=1}^m a_{ik}^2 = \sum_{i=1}^m a_{ik} \cdot y_i - \sum_{i=1}^m a_{ik} \cdot \sum_{j=1, j \neq k}^n x_j a_{ij}$$
$$x_k = \frac{\sum_{i=1}^m a_{ik} \cdot \left(y_i - \sum_{j=1, j \neq k}^n x_j a_{ij} \right)}{\sum_{i=1}^m a_{ik}^2}$$

Linear Regression Coordinate Descent - Simple Algorithm

```
1: procedure LINEAR REGRESSION-CD
  input:  $f$ 
2:   Get initial point  $\mathbf{x}^{(0)}$ 
3:   repeat
4:     for  $k \in 1, \dots, n$  do
5:        $x_k \leftarrow \frac{\sum_{i=1}^m a_{ik} \cdot (y_i - \sum_{j=1, j \neq k}^n x_j a_{ij})}{\sum_{i=1}^m a_{ik}^2}$ 
6:     end for
7:   until convergence
8:   return  $\mathbf{x}, f(\mathbf{x})$ 
9: end procedure
```

Coordinate Descent for Linear Regression

For each parameter we have the following update rule:

$$x_k = \frac{\sum_{i=1}^m a_{ik} \cdot \left(y_i - \sum_{j=1, j \neq k}^n x_j a_{ij} \right)}{\sum_{i=1}^m a_{ik}^2}$$

One Coordinate descent epoch has a cost of $O(m \cdot n^2)$!

Can we do it faster?

CD for Linear Regression - Smart Update

For each parameter we have the following update rule:

$$x_k = \frac{\sum_{i=1}^m a_{ik} \cdot \left(y_i - \sum_{j=1, j \neq k}^n x_j a_{ij} \right)}{\sum_{i=1}^m a_{ik}^2}$$

We can rewrite:

$$\begin{aligned} \sum_{i=1}^m \left(y_i - \sum_{j=1, j \neq k}^n x_j a_{ij} \right) &= \sum_{i=1}^m \left(y_i - \sum_{j=1}^n x_j a_{ij} + x_k a_{ik} \right) \\ &= \sum_{i=1}^m \left(y_i - \sum_{j=1}^n x_j a_{ij} \right) + \sum_{i=1}^m x_k a_{ik} \\ &= \sum_{i=1}^m \left(y_i - \sum_{j=1}^n x_j a_{ij} \right) + x_k \sum_{i=1}^m a_{ik} \end{aligned}$$

CD for Linear Regression - Smart Update

From which we have:

$$\begin{aligned}
 x_k &= \frac{\sum_{i=1}^m a_{ik} \cdot \left(y_i - \sum_{j=1, j \neq k}^n x_j a_{ij} \right)}{\sum_{i=1}^m a_{ik}^2} \\
 &= \frac{\sum_{i=1}^m a_{ik} \cdot \left(y_i - \sum_{j=1}^n x_j a_{ij} + x_k^{old} a_{ik} \right)}{\sum_{i=1}^m a_{ik}^2} \\
 &= \frac{\sum_{i=1}^m a_{ik} \cdot \left(y_i - \sum_{j=1}^n x_j a_{ij} \right)}{\sum_{i=1}^m a_{ik}^2} + \frac{\sum_{i=1}^m a_{ik} \cdot (x_k^{old} a_{ik})}{\sum_{i=1}^m a_{ik}^2} \\
 &= \frac{\sum_{i=1}^m a_{ik} \cdot \left(y_i - \sum_{j=1}^n x_j a_{ij} \right)}{\sum_{i=1}^m a_{ik}^2} + \frac{x_k^{old} \cdot \sum_{i=1}^m a_{ik}^2}{\sum_{i=1}^m a_{ik}^2} \\
 &= \frac{\sum_{i=1}^m a_{ik} \cdot \left(y_i - \sum_{j=1}^n x_j a_{ij} \right)}{\sum_{i=1}^m a_{ik}^2} + x_k^{old}
 \end{aligned}$$

CD for Linear Regression - Smart Update

Now we have:

$$x_k = \frac{\sum_{i=1}^m a_{ik} \cdot \left(y_i - \sum_{j=1}^n x_j a_{ij} \right)}{\sum_{i=1}^m a_{ik}^2} + x_k^{old}$$

So we can define our residual vector $\mathbf{r} \in \mathbb{R}^m$ such that

$$r_i = y_i - \sum_{j=1}^n x_j a_{ij}$$

CD for Linear Regression - Smart Update

After each update of x_k , we can maintain r_i instead of recomputing it given the old value \mathbf{x}_k^{old} :

$$\begin{aligned} r_i^{new} &= y_i - \left(\sum_{j=1}^n x_j a_{ij} - \mathbf{x}_k^{old} a_{ik} + x_k a_{ik} \right) \\ &= r_i^{old} + (\mathbf{x}_k^{old} - x_k) a_{ik} \end{aligned}$$

CD for Linear Regression - Smart Update

Now our algorithm looks like:

1. Initialize \mathbf{x}
2. Compute $r_i = y_i - \sum_{j=1}^n x_j a_{ij}$
3. While Not Converged
 - 3.1 For each $k = 1, \dots, n$
 - 3.1.1 $\mathbf{x}_k^{old} \leftarrow x_k$
 - 3.1.2 $x_k \leftarrow \frac{\sum_{i=1}^m a_{ik} \cdot (r_i)}{\sum_{i=1}^m a_{ik}^2} + \mathbf{x}_k^{old}$
 - 3.1.3 For all i $r_i \leftarrow r_i + (\mathbf{x}_k^{old} - x_k) a_{ik}$

This algorithm is now $O(m \cdot n)$!

Linear Regression Coordinate Descent Algorithm

```
1: procedure LINEAR REGRESSION-CD
  input:  $f$ 
2:   Get initial point  $\mathbf{x}^{(0)}$ 
3:    $\mathbf{r} \leftarrow \mathbf{y} - A\mathbf{x}^{(0)}$ 
4:   repeat
5:     for  $k \in 1, \dots, n$  do
6:        $\mathbf{x}_k^{old} \leftarrow x_k$ 
7:        $x_k \leftarrow \frac{\sum_{i=1}^m a_{ik} \cdot r_i}{\sum_{i=1}^m a_{ik}^2} + \mathbf{x}_k^{old}$ 
8:       for  $i \in 1, \dots, m$  do
9:          $r_i \leftarrow r_i + (\mathbf{x}_k^{old} - x_k) a_{ik}$ 
10:      end for
11:    end for
12:    until convergence
13:    return  $\mathbf{x}, f(\mathbf{x})$ 
14: end procedure
```

Real World Dataset: Body Fat prediction

We want to estimate the percentage of body fat based on various attributes:

- ▶ Age (years)
- ▶ Weight (lbs)
- ▶ Height (inches)
- ▶ Neck circumference (cm)
- ▶ Chest circumference (cm)
- ▶ Abdomen 2 circumference (cm)
- ▶ Hip circumference (cm)
- ▶ Thigh circumference (cm)
- ▶ Knee circumference (cm)
- ▶ ...

<http://lib.stat.cmu.edu/datasets/bodyfat>



Real World Dataset: Body Fat prediction

The data is represented it as:

$$A_{m,n} = \begin{pmatrix} 1 & a_{1,1} & a_{1,2} & \dots & a_{1,n} \\ 1 & a_{2,1} & a_{2,2} & \dots & a_{2,n} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & a_{m,1} & a_{m,2} & \dots & a_{m,n} \end{pmatrix} \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix}$$

with $m = 252$, $n = 14$

We can model the percentage of body fat y is a linear combination of the body measurements with parameters \mathbf{x} :

$$\hat{y}_i = \mathbf{x}^T \mathbf{a}_i = x_0 1 + x_1 a_{i,1} + x_2 a_{i,2} + \dots + x_n a_{i,n}$$

Coordinate Descent - Body fat dataset

Year Prediction Data Set

- ▶ Least Squares Problem
- ▶ Prediction of the release year of a song from audio features
- ▶ 90 features
- ▶ Experiments done on a subset of 1000 instances of the data

Coordinate Descent - Year Prediction

Further Readings

- ▶ The coordinate descent method is not covered by Boyd and Vandenberghe [2004]
- ▶ Coordinate descent:
 - ▶ very briefly [Nocedal and Wright, 2006, ch. 9.3]
 - ▶ A brief, but dense survey: Wright [2015]
- ▶ Convergence proof: Wright [2015]

References I

- Stephen Boyd and Lieven Vandenberghe. *Convex Optimization*. Cambridge Univ Press, 2004.
- Jorge Nocedal and Stephen J. Wright. *Numerical Optimization*. Springer, 2006.
- Stephen J. Wright. Coordinate descent algorithms. *Mathematical Programming*, 151(1):3–34, 2015.