

CD

Coordinate Descent Algorithms

Liu Jingchang

March 23, 2017

University of Science and Technology of China

Table of contents

1. Algorithms

2. Analysis

3. Conclusion

4. Q & A

Algorithms

Stochastic Gradient Descent

$$n \gg p \quad f(x) = \sum_{i=1}^n f_i(x) \quad (1)$$

Coordinate Descent

p may be big, such as in computational biology and health care problem.

Formulation

Unconstrained minimization

$$\min_x f(x) \quad (2)$$

Structured formulation

$$\min_x h(x) := f(x) + \lambda \Omega(x) \quad (3)$$

- f : smooth
- Ω : regularization function, may be nonsmooth, often convex, usually assumed to be separable or block-separable.

$$\Omega(x) = \sum_{i=1}^n \Omega_i(x_i) \quad (4)$$

- $\Omega(x) = \|x\|_1, \Omega_i(x_i) = \|x_i\|_1$
- Ω : box constraints, $\Omega_i(x_i) = I_{[l_i, u_i]}(x_i)$

Algorithm 1 Coordinate Descent for unconstrained case

- 1: Set $k \leftarrow 0$ and choose $x^0 \in \mathbb{R}^n$
 - 2: **repeat**
 - 3: Choose index $i_k \in \{1, 2, \dots, n\}$
 - 4: $x^{k+1} \leftarrow x^k - \alpha_k \nabla_{i_k} f(x^k) e_{i_k}$
 {equal to $x_{i_k}^{k+1} \leftarrow x_{i_k}^k - \alpha_k \nabla_{i_k} f(x^k)$ }
 - 5: $k \leftarrow k + 1$
 - 6: **until** termination test satisfied
-

Formulation

$$\min_x h(x) := f(x) + \lambda \Omega(x) \quad (5)$$

f : smooth, Ω : may be nonsmooth.

Iteration

$$x^{k+1} = \text{prox}_{\Omega}^{\eta} \left(x^k - \eta \nabla f(x^k) \right) \quad (6)$$

$$= \underset{y}{\operatorname{argmin}} \left\{ \Omega(y) + \frac{1}{2\eta} \left\| y - \left(x^k - \eta \nabla f(x^k) \right) \right\|_2^2 \right\} \quad (7)$$

Indicate function

$$x^{k+1} = \text{prox}_{\Omega}^{\eta} \left(x^k - \eta \nabla f(x^k) \right) \quad (8)$$

$$= \underset{y}{\operatorname{argmin}} \left\{ \Omega(y) + \frac{1}{2\eta} \left\| y - \left(x^k - \eta \nabla f(x^k) \right) \right\|_2^2 \right\} \quad (9)$$

$$= \underset{y \in \mathbb{C}}{\operatorname{argmin}} \left\{ \left\| y - \left(x^k - \eta \nabla f(x^k) \right) \right\|_2^2 \right\} \quad (10)$$

$$= \text{Proj}_{\mathbb{C}} \left(x^k - \eta \nabla f(x^k) \right) \quad (11)$$

Other

- l_1 norm: soft-thresholding
- l_0 norm: hard-thresholding

Figure

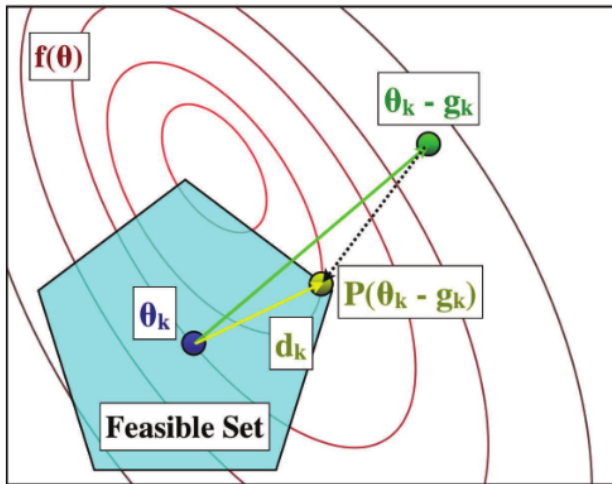


Figure 1: Illustration of projected gradient descent

Algorithm

Algorithm 2 Coordinate Descent for regularization case

- 1: Set $k \leftarrow 0$ and choose $x^0 \in \mathbb{R}^n$
 - 2: **repeat**
 - 3: Choose index $i_k \in \{1, 2, \dots, n\}$
 - 4: $z_{i_k}^{k+1} \leftarrow \text{prox}_{\Omega_{i_k}}^\eta \left(x_{i_k}^k - \eta \nabla_{i_k} f(x^k) \right)$
 - 5: $x^{k+1} \leftarrow x^k + \left(z_{i_k}^k - x_{i_k}^k \right) e_{i_k}$
 {4 and 5 is equal to $x_{i_k}^{k+1} \leftarrow \text{prox}_{\Omega_{i_k}}^\eta \left(x_{i_k}^k - \eta \nabla_{i_k} f(x^k) \right)$ }
 - 6: $k \leftarrow k + 1$
 - 7: **until** termination test satisfied
-

Property

$$\begin{aligned} \text{prox}_{\Omega}^\eta \left(x^k - \eta \nabla f(x^k) \right) = \\ \left(\text{prox}_{\Omega_1}^\eta \left(x_1^k - \eta \nabla_1 f(x^k) \right), \text{prox}_{\Omega_2}^\eta \left(x_2^k - \eta \nabla_2 f(x^k) \right), \dots, \text{prox}_{\Omega_n}^\eta \left(x_n^k - \eta \nabla_n f(x^k) \right) \right) \end{aligned} \quad (12)$$

Analysis

Gradient Descent

Full Gradient Descent

$$\min_x f(x) \quad x^{k+1} \leftarrow x^k - \eta_k \nabla f(x^k)$$

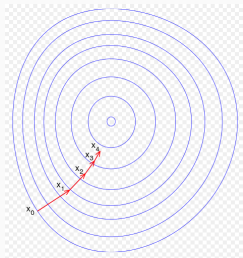
Stochastic Gradient Descent

$$\min_x \sum_{i=1}^n f_i(x) \quad x^{k+1} \leftarrow x^k - \eta_k \nabla f_{i_k}(x^k)$$

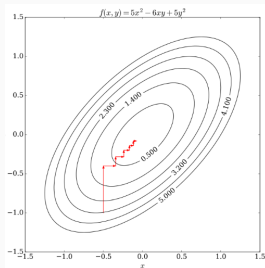
Stochastic Coordinate Descent

$$\min_x f(x) \quad x^{k+1} \leftarrow x^k - \eta_k \nabla_{i_k} f(x^k)$$

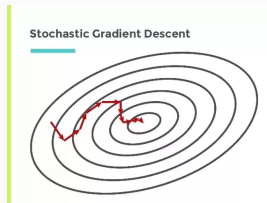
Gradient Descent



SCD



SGD



Assumption

Strong convexity

$$\exists \delta > 0, \text{ s.t. } f(y) \geq f(x) + \nabla f(x)^T (y - x) + \frac{\delta}{2} \|y - x\|_2^2, \quad \text{for all } x, y \quad (13)$$

Component Lipschitz Continuous

$$\|\nabla_i f(x + te_i) - \nabla_i f(x)\| \leq L_i \|t\| \quad (14)$$

Standard Lipschitz Continuous

$$\|\nabla f(x + d) - \nabla f(x)\| \leq L \|d\| \quad (15)$$

Assumption

f : Convex, Lipschitz continuous gradient, give x^0 , $\exists R_0$, s.t.

$$\max_{x^* \in \mathcal{S}} \max_x \{ \|x - x^*\| : f(x) \leq f(x^0) \} \leq R_0 \quad (16)$$

Convergence rate of unconstrained case

Th

Under above assumption. Suppose that $\alpha_k = \frac{1}{L_{\max}}$, then

$$E\left(f\left(x^k\right)\right) - f^* \leq \frac{2nL_{\max}R_0^2}{k} \quad (17)$$

When δ strongly-convex,

$$E\left(f\left(x^k\right)\right) - f^* \leq \left(1 - \frac{\delta}{nL_{\max}}\right)^k (f(x^0) - f^*) \quad (18)$$

$$\begin{aligned}
 f(x^{k+1}) &= f(x^k - \alpha_k \nabla_{i_k} f(x^k) e_{i_k}) \\
 &\leq f(x^k) - \alpha_k [\nabla_{i_k} f(x^k)]^2 + \frac{1}{2} \alpha_k^2 L_{i_k} [\nabla_{i_k} f(x^k)]^2 \\
 &\leq f(x^k) - \alpha_k \left(1 - \frac{L_{\max}}{2} \alpha_k\right) [\nabla_{i_k} f(x^k)]^2 \\
 &= f(x^k) - \frac{1}{2L_{\max}} [\nabla_{i_k} f(x^k)]^2
 \end{aligned} \tag{19}$$

Taking the expectation of both sides over the random index i_k ,

$$\begin{aligned}
 E_{i_k} f(x^{k+1}) &\leq f(x^k) - \frac{1}{2L_{\max}} \frac{1}{n} \sum_{i=1}^m [\nabla_i f(x^k)]^2 \\
 &\leq f(x^k) - \frac{1}{2nL_{\max}} [\nabla f(x^k)]^2
 \end{aligned} \tag{20}$$

Take expectation with respect to i_1, i_2, \dots, i_{k-1} ,

$$E \left(f \left(x^{k+1} \right) \right) \leq E \left(f \left(x^k \right) \right) - \frac{1}{2nL_{\max}} E \left[\left\| \nabla f \left(x^k \right) \right\|^2 \right] \quad (21)$$

By the convexity of f

$$f \left(x^k \right) - f^* \leq \nabla f \left(x^k \right)^T \left(x^k - x^* \right) \leq \left\| \nabla f \left(x^k \right) \right\| \left\| x^k - x^* \right\| \leq R_0 \left\| \nabla f \left(x^k \right) \right\| \quad (22)$$

Combine 21 and 22, we can get 17.

Proof for strong convex

Strong convex

$$\exists \delta > 0, \text{ s.t. } f(y) \geq f(x) + \nabla f(x)^T (y - x) + \frac{\delta}{2} \|y - x\|_2^2, \quad \text{for all } x, y$$

Taking the minimum of both sides with respect to y , and setting $x = x^k$

$$f^* \geq f(x^k) - \frac{1}{2\alpha} \|\nabla f(x^k)\|^2 \quad (23)$$

Combine with 18, we obtain

$$E_{i_k} f(x^{k+1}) - f^* \leq f(x^k) - \frac{\delta}{nL_{\max}} f(x^k) \quad (24)$$

Take expectation with respect to i_1, i_2, \dots, i_{k-1} ,

$$E(f(x^{k+1})) - f^* \leq E(f(x^k)) - f^* - \frac{\delta}{nL_{\max}} (E(f(x^k)) - f^*) \quad (25)$$

Separable Regularized Case

Assumption

- f : Lipschitz continuous gradient, Strongly convex with $\delta > 0$
- $\Omega_i, i = 1, 2, \dots, n$: convex

Th

Under above assumption, $\alpha_k = \frac{1}{L_{\max}}$, Then

$$E\left(h\left(x^k\right)\right) - h^* \leq \left(1 - \frac{\delta}{nL_{\max}}\right)^k \left(h\left(x^0\right) - h^*\right) \quad (26)$$

Conclusion

Conclusion

- Coordinate Descent is an approximation to Full Gradient Descent.
- Proximal methods is designed for regularized optimization.
- CD can get $O\left(\frac{1}{n}\right)$ convergence; And linear convergence when strongly convex.

Q & A
