

# Lec14-new

Monday, February 28, 2011 5:01 PM

# Geometric convergence

Monday, February 6, 2023 11:34 AM

See Bubeck p278

- Sébastien Bubeck, "Convex Optimization: Algorithms and Complexity, in Foundations and Trends in Machine Learning, Vol. 8, No. 3-4 (2015)

- With smoothness, we can ensure convergence, but the speed of convergence may be slow

$$\|x(t+1) - x^*\|^2$$

$$\leq \|x(t) - x^*\|^2 - \left(\frac{2\gamma}{L} - \gamma^2\right) \|\nabla f(x(t))\|^2$$

- Adding strong convexity enables us to prove geometric descent.

$$\|\nabla f(x) - \nabla f(y)\|_2 \geq \alpha \|x - y\|_2$$

- Intuitively, this ensures that when  $x(t)$  is far away from  $x^*$ , the improvement of  $\|x(t+1) - x^*\|_2^2$  is directly related to  $\|x(t) - x^*\|_2^2$ .

$$\|x(t+1) - x^*\|^2 \leq \left(1 - \left(\frac{2\gamma}{L} - \gamma^2\right)\alpha\right) \|x(t) - x^*\|^2 \quad (*)$$

- The result below is stronger (i.e., faster descent).

---

Skip

- Lemma: Let  $f$  be a  $L$ -smooth &  $\alpha$ -strongly convex function on  $\mathbb{R}^n$ . Then for all  $x, y \in \mathbb{R}^n$ ,

$$(\nabla f(x) - \nabla f(y))^T (x - y) \geq \frac{\alpha L}{\alpha + L} \|x - y\|^2 + \frac{1}{\alpha + L} \|\nabla f(x) - \nabla f(y)\|^2$$

Proof: Since  $f$  is  $\alpha$ -strongly-convex, we can show that

$\varphi(x) = f(x) - \frac{\alpha}{2} \|x\|^2$   
 is still convex. Further, we can show that  $\varphi(x)$  is  $(L-\alpha)$ -smooth. Thus, using the earlier lemma for smooth functions, we have

$$\begin{aligned} [\nabla\varphi(x) - \nabla\varphi(y)]^T(x-y) &\geq \frac{1}{L-\alpha} \|\nabla\varphi(x) - \nabla\varphi(y)\|^2 \\ \Leftrightarrow [\nabla f(x) - \nabla f(y)]^T(x-y) - \alpha(x-y)^T(x-y) \\ &\geq \frac{1}{L-\alpha} \|(\nabla f(x) - \nabla f(y)) - \alpha(x-y)\|^2 \\ &= \frac{1}{L-\alpha} \|\nabla f(x) - \nabla f(y)\|^2 + \frac{\alpha^2}{L-\alpha} \|x-y\|^2 \\ &\quad - \frac{2\alpha}{L-\alpha} (\nabla f(x) - \nabla f(y))^T(x-y) \end{aligned}$$

$$\begin{aligned} \Leftrightarrow \frac{L+\alpha}{L-\alpha} [\nabla f(x) - \nabla f(y)]^T(x-y) \\ &\geq \frac{1}{L-\alpha} \|\nabla f(x) - \nabla f(y)\|^2 \\ &\quad + \frac{L\alpha}{L-\alpha} \|x-y\|^2 \end{aligned}$$

The result of the lemma then follows. #

Theorem: Let  $f$  be a  $L$ -smooth &  $\alpha$ -strongly-convex function on  $\mathbb{R}^n$ . Then, by setting a stepsize

$$\gamma < \frac{2}{\alpha+L},$$

gradient descent satisfies

$$\|x(t) - x^*\|^2 \leq \left(1 - \frac{2\gamma\alpha L}{\alpha+L}\right)^t \|x(0) - x^*\|^2$$

Proof: Starting with the norm approach again.

Let  $x^*$  be one optimal solution, i.e.,  $\nabla f(x^*) = 0$

$$\begin{aligned} & \|x(t+1) - x^*\|^2 \\ &= \|x(t) - x^*\|^2 + 2(x(t+1) - x(t))^T (x(t) - x^*) \\ &\quad + \|x(t+1) - x(t)\|^2 \end{aligned}$$

Note that

$$\begin{aligned} & (x(t+1) - x(t))^T (x(t) - x^*) \\ &= -\delta (\nabla f(x(t)) - \underbrace{\nabla f(x^*)}_0) (x(t) - x^*) \\ &\leq -\frac{\gamma \alpha L}{\alpha + L} \|x(t) - x^*\|^2 - \frac{\delta}{\alpha + L} \|\nabla f(x) - \nabla f(x^*)\|^2 \end{aligned}$$

Hence,

$$\begin{aligned} & \|x(t+1) - x^*\|^2 \\ &\leq \left(1 - \frac{2\delta \alpha L}{\alpha + L}\right) \|x(t) - x^*\|^2 + \underbrace{\left(\delta^2 - \frac{2\delta}{\alpha + L}\right)}_{(\leq 0 \text{ if } \delta < \frac{2}{\alpha + L})} \|\nabla f(x) - \nabla f(x^*)\|^2 \\ &\leq \left(1 - \frac{2\delta \alpha L}{\alpha + L}\right) \|x(t) - x^*\|^2 \end{aligned}$$

The result of the Theorem then follows. #

---

Note:

- For a faster convergence speed, would want

$$1 - \frac{2\delta \alpha L}{\alpha + L}$$

to be smaller, or  $\delta$  to be larger.

$$\Rightarrow \text{ set } \delta = \frac{2}{\alpha + L}$$

- Compare with (\*) at  $\delta = \frac{2}{\alpha + L}$

$$\begin{aligned} & 1 - \delta \left( \frac{2}{L} - \delta \right) \alpha \\ &= 1 - \frac{2\delta\alpha^2}{L(\alpha + L)} \end{aligned}$$

- When  $\frac{\alpha}{L}$  is small, this is much closer to 1.

---

- Using  $f(x) - f(x^*) \leq \frac{L}{2} \|x - x^*\|^2$

we have

$$f(x(t)) - f(x^*) \leq \frac{L}{2} \left( 1 - \frac{2\delta\alpha L}{\alpha + L} \right)^t \|x(0) - x^*\|^2$$

---

- This is usually referred to as "linear convergence"

- in contrast to "quadratic convergence" for second-order (e.g. Newton's) algorithms.

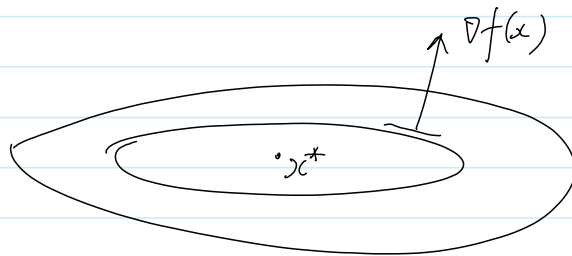
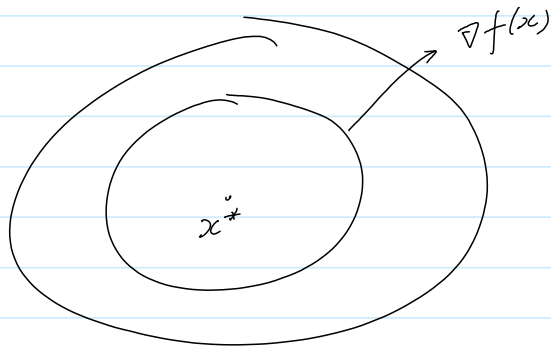
## Scaled gradient descent algorithm

Saturday, February 21, 2009 5:50 PM

- Standard gradient algorithm

$$x^{(t+1)} = x^{(t)} - \sigma \nabla f(x)$$

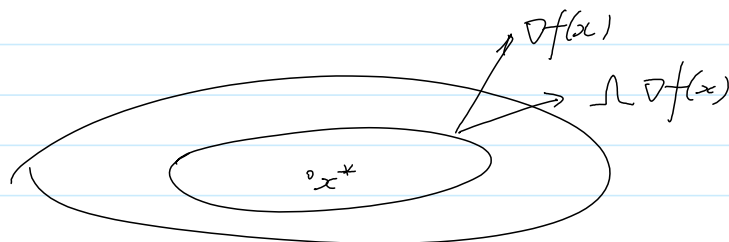
If  $\nabla f(x)$  is Lipschitz &  $\sigma < \frac{2}{L}$ , then the algorithm converges.



- Scaled gradient algorithm

$$x^{(t+1)} = x^{(t)} - \sigma \Omega \nabla f(x)$$

where  $\Omega$  is a positive-definite matrix



Some choice of  $\Omega$  can make the algorithm converge faster.

- For example,  $\Omega = \begin{bmatrix} \lambda_1 & & & \\ & \lambda_2 & & \\ & & \ddots & \\ & & & \lambda_m \end{bmatrix}$

- each processor chooses different step-size.
- Regardless the use of  $\Omega$ , the scaled gradient algorithm will converge under a similar Lipschitz condition.

$$\gamma < \frac{2}{L \lambda_{\max}}$$

- where  $\lambda_{\max}$  is the largest eigenvalue of  $\Omega$ , i.e.  
 $x^T \Omega x \leq \lambda_{\max} \|x\|^2$  for all  $x$

Sketch of proof:

- choose a different norm.

$$\begin{aligned} & (x(t+1) - x^*)^T \Omega^{-1} (x(t+1) - x^*) \\ &= (x(t) - x^*)^T \Omega^{-1} (x(t) - x^*) \\ & \quad + 2 (x(t+1) - x(t))^T \Omega^{-1} (x(t) - x^*) \\ & \quad + (x(t+1) - x(t))^T \Omega^{-1} (x(t+1) - x(t)) \end{aligned}$$

Note that

$$\begin{aligned} & (x(t+1) - x(t))^T \Omega^{-1} (x(t) - x^*) \\ &= -\gamma \nabla f(x(t))^T (x(t) - x^*) \end{aligned}$$

$$\begin{aligned} &= -\sigma [\nabla f(x^{(+)}) - \nabla f(x^*)] [x^{(+)} - x^*] \\ &\leq -\frac{\sigma}{L} \|\nabla f(x^{(+)})\|^2 \end{aligned}$$

Work out the rest in hw.

(75)



# Constrained optimization

Saturday, February 21, 2009 6:06 PM

- Consider the problem

$$\begin{array}{l} \min f(x) \\ \text{sub to } x \in \mathbb{X} \end{array}$$

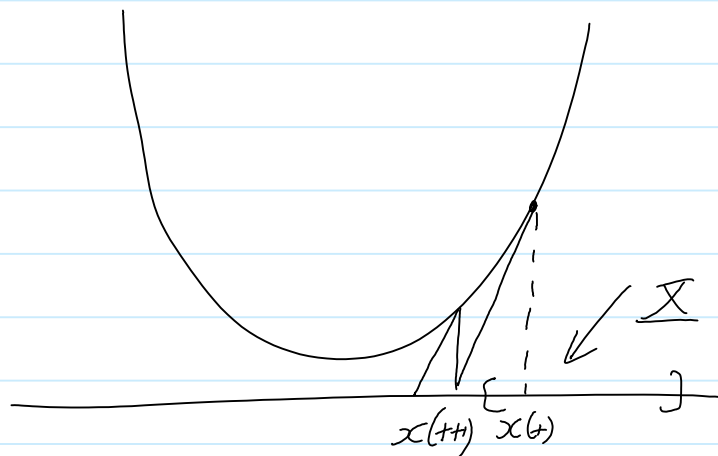
Optimality condition

$$f'(x; y-x) \geq 0 \quad \text{for all } y \in \mathbb{X}$$

If  $f$  is differentiable, then

$$(\nabla f(x))^T (y-x) \geq 0 \quad \text{for all } y \in \mathbb{X}.$$

- However, the normal gradient projection does not work any more because the new  $x^{(t+1)}$  can go outside  $\mathbb{X}$ .



Three solutions

## Three solutions

① Penalty function method:

choose  $g(x)$  such that

$$g(x) = 0 \quad \text{if } x \in \mathcal{X}$$

$$g(x) \geq 0 \quad \text{if } x \notin \mathcal{X}$$

- we can now minimize  $f(x) + \beta g(x)$

let the solution be  $x^*(\beta)$

- As  $\beta \uparrow +\infty$ , the penalty becomes larger & larger, the solution  $x^*(\beta)$  will approach  $x^*$  (the original constrained problem).

- We will discuss the engineering implication of this approach when we discuss TCP as an example.

---

② Interior point method (Barrier Method).

- Choose  $g(x)$  such that  $g(x) \rightarrow +\infty$  as  $x$  approaches the boundary of  $\mathcal{X}$  from inside.

Example:  $\bar{X} = \{x \geq 0\}$ ,  $g(x) = -\log x$

$\bar{X} = \{x \leq a\}$ ,  $g(x) = -\log(a-x)$

- We then minimize  $f(x) + \beta g(x)$ .
- Due to the barrier  $g(x)$ , the optimal solution  $x^*(\beta)$  must be in the interior of  $\bar{X}$ .
- As  $\beta \downarrow 0$ ,  $x^*(\beta) \rightarrow x^*$  (the original constrained problem).

In both the penalty-function method or the interior-point method, the problem is converted to a unconstrained problem. Hence, we can use gradient algorithm to solve  $x^*(\beta)$ .

- However, it may be difficult to ensure the Lipschitz condition of the gradient.

---

(3) Projection method.

(10)

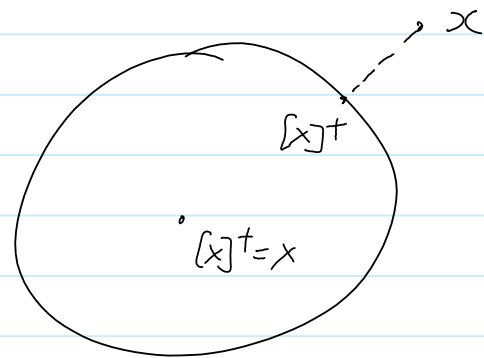
# Projection

Saturday, February 21, 2009 6:19 PM

## Projection $[x]^+$ : (for $L_2$ -norm)

The projection of  $x$  onto the convex and closed set  $\mathcal{X}$  is the point  $z$  in  $\mathcal{X}$  that is closest to  $x$ , i.e.,

$$[x]^+ = \underset{z \in \mathcal{X}}{\operatorname{argmin}} \|z - x\|_2$$



Example:

$$\textcircled{1} \quad \mathcal{X} = \bigotimes_{i=1}^m [a_i, b_i]$$

$$[x]^+ = \{ [x_i]^+ \}$$

$$\text{where } [x_i]^+ = \begin{cases} a_i & x_i \leq a_i \\ b_i & x_i \geq b_i \\ x_i & \text{o/w} \end{cases}$$

② Projection to a Polyhedra

Boyd p 390

$$\begin{aligned} \min \quad & \|x - x_0\|_2^2 \\ \text{sub to} \quad & Ax \leq b \end{aligned}$$

- A quadratic program.

Solution:

- On a hyperplane  $a^T x = b$

$$P_C(x_0) = x_0 + (b - a^T x_0) \cdot a / \|a\|_2^2$$

- On a half-space  $a^T x \leq b$

$$P_C(x_0) = \begin{cases} x_0 & a^T x_0 \leq b \\ x_0 + (b - a^T x_0) \cdot a / \|a\|_2^2 & \text{if } a^T x_0 > b. \end{cases}$$

② Can we derive these by optimality conditions?

③ In general

$$X = \{x \mid f_i(x) \leq 0, h_i(x) = 0\}$$

$$[x]^+ = \arg \min_z \|z - x\|_2$$

$$\begin{aligned} \text{sub to} \quad & f_i(z) \leq 0 \\ & h_i(z) = 0 \end{aligned}$$

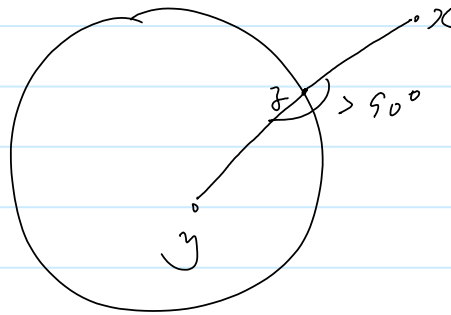
- Not a trivial operation! (Reason to use dual algorithm.)

Projection Theorem (Bertsekas & Tsitsiklis p211)

(a) For every  $x \in \mathbb{R}^n$ , there exists a unique  $z \in X$  that minimizes  $\|z - x\|_2$  over all  $z \in X$ , and will be denoted as  $[x]^+$

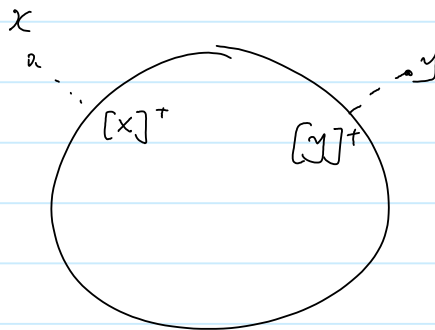
(b) Given some  $x \in \mathbb{R}^n$ , a vector  $z \in X$  is equal to  $[x]^+$  if & only if

$$(y - z)'(x - z) \leq 0 \text{ for all } y \in X$$



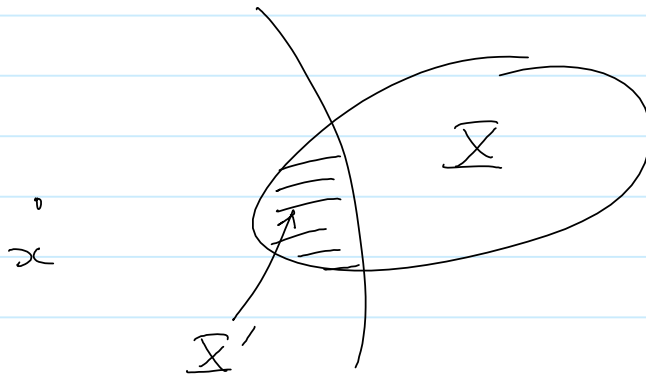
(c) The mapping  $f(x) = [x]^+$  is continuous and non-expansive, that is

$$\|[x]^+ - [y]^+\|_2 \leq \|x - y\|_2 \text{ for all } x, y \in \mathbb{R}^n$$



Proof:

(a) Can intersect  $X$  with a closed & bounded set



-  $\min \|z - x\|^2$  over  $\underline{X}'$  (closed, bounded)

$\Rightarrow$  exists a minimum point  $[x]^+$

- It is unique because  $\|z - x\|^2$  is strictly convex.

(b) Consider  $\min f(z) = \|z - x\|^2$

sub to  $z \in \underline{X}$

-  $\nabla f(z) = 2(z - x)$

- By optimality condition

$z$  optimal

$\Leftrightarrow 2(z - x)^T (y - z) \geq 0$  for all  $y \in \underline{X}$

$\Leftrightarrow (y - z)^T (x - z) \leq 0$

(c) From part (b)

$([y]^+ - [x]^+)^T (x - [x]^+) \leq 0$

Similarly

$$([x]^+ - [y]^+)^T (y - [y]^+) \leq 0$$

$$\Rightarrow ([y]^+ - [x]^+)^T ([y]^+ - [x]^+ - (y - x)) \leq 0$$

$$\| [y]^+ - [x]^+ \|^2 \leq ([y]^+ - [x]^+) (y - x)$$

$$\leq \| [y]^+ - [x]^+ \| \| y - x \|.$$

The result then follows.

(25)