# Lec13-new

Saturday, February 21, 2009     8:58 PM

HW4 is on the web

Bring convergence/projection handout.

# Optimization algorithms

- We have now discussed

    - convex problems
    - optimality conditions

- Sometimes we can explicitly solve for the optimal solution from the optimality condition

    e.g.          $\min \ \|Ax - b\|_2^2$

    $\Rightarrow \quad A^T A x = A^T b$
    $\qquad\quad x^* = (A^T A)^{-1} A^T b$

- At other times a closed-form solution may not be possible

- We may then use numerical algorithms.

---

## Numerical Algorithms

- A numerical algorithm starts from some initial estimate $x_0$, and iteratively generate new estimates by

    $$x_k = T(x_{k-1}) \qquad k = 1, 2, \cdots$$

- Hopefully, as $k \to +\infty$, $x_k \to x^*$ the optimal

Solution.

① When will such a sequence converge to the optimal solution?

Ⓐ often by showing that the quality of $x_k$ improves in each iteration

    ⓐ The distance between $x_k$ and $x^*$ improves in each iteration. e.g

$$\| x_k - x^* \| \leq \alpha \| x_{k-1} - x^* \| \quad , \quad \alpha < 1$$

$$\text{or} \quad \| x_k - x^* \| \leq \| x_{k-1} - x^* \| - \beta \quad , \quad \beta > 0$$

    ⓑ The function value improves in each iteration. e.g

$$f_0 (x_k) \leq f_0 (x_{k-1}) - \gamma . \quad \gamma > 0.$$

A trivial example

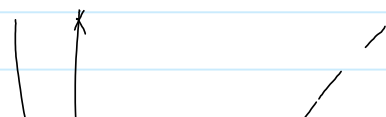- How to compute $\sqrt{2}$ using only $+, -, \times, /$

- $x = \sqrt{2} \iff x^2 - 1 = 1$
  $$\iff x = \frac{1}{x+1} + 1$$

- The iteration
  $$x_k = \frac{1}{x_{k-1}+1} + 1$$
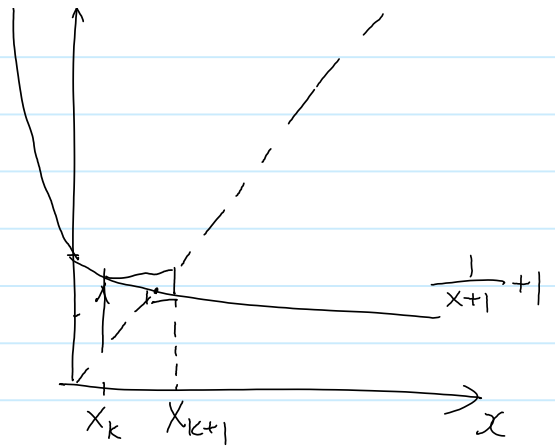
- For any two points

— For any two points $x, y$, let

$$Tx = \frac{1}{x+1} + 1$$
$$Ty = \frac{1}{y+1} + 1$$

$$Tx - Ty = \frac{1}{x+1} - \frac{1}{y+1}$$

$$= -\frac{x-y}{(x+1)(y+1)}$$

$$\|Tx - Ty\| \leq \frac{\|x-y\|}{4} \qquad \text{assuming that } x, y \geq 1. \qquad (\ast)$$

— Note that $\sqrt{2}$ is a fixed point of the mapping

$$T(\sqrt{2}) = \sqrt{2}$$

— Hence, the distance $\|x_k - \sqrt{2}\|$ is cut by $\frac{1}{4}$ in each iteration

$$\Rightarrow \quad x_k \to \sqrt{2}.$$

— The inequality $(\ast)$ describes a "contraction mapping":

  — The distance between $Tx$ & $Ty$ is less than that of $x$ & $y$

  $$\|Tx - Ty\| \leq c \|x-y\|, \quad 0 \leq c \leq 1$$

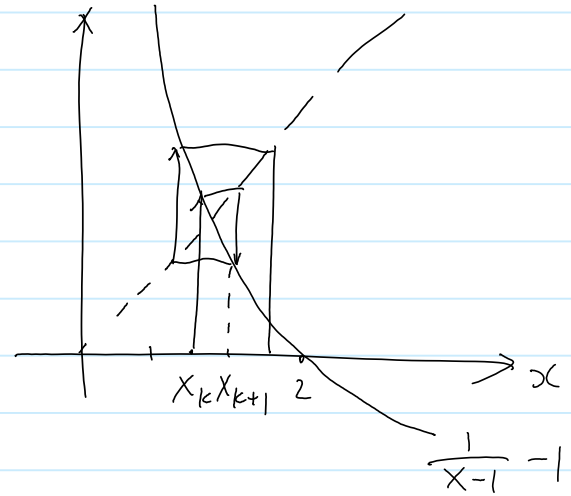  — When $c < 1$ (the contraction mapping is strict, we immediately have

  ① There is a unique fixed point $x^\ast$ with $Tx^\ast = x^\ast$

  ② $x_{k+1} = T(x_k)$ converges "geometrically" to $x^\ast$

② $X_{k+1} = T(X_k)$ converges "geometrically" to $x^*$.

However, not all iterative algorithms converge.

$$X = \sqrt{2} \iff X = \frac{1}{X-1} - 1$$



$$\frac{1}{X-1} - 1$$

⑩

- $\min f(x)$     $f$ is convex
- Optimality condition
  - If $f$ is differentiable
  $$\nabla f(x^*) = 0$$
  - If $f$ is not differentiable
  $$f'(x; x - x^*) \geq 0 \qquad \forall x$$

  i.e, the directional derivative is positive in all directions.

---

- Assume that $f(\cdot)$ is differentiable.
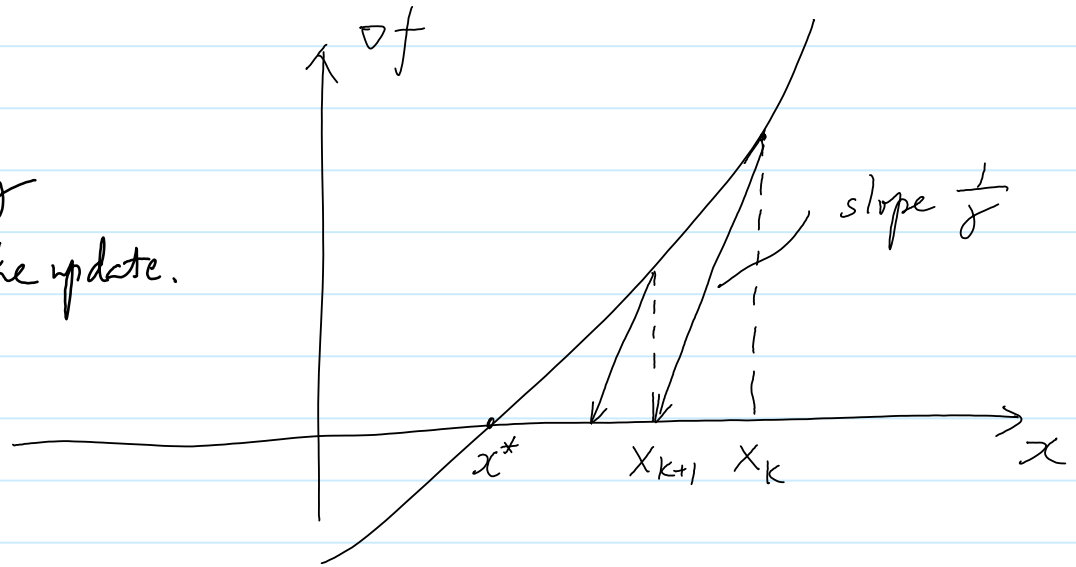- Consider the iteration of the type
$$x(t+1) = x(t) - \gamma \nabla f(x(t)) \overset{\triangle}{=} T(x(t))$$
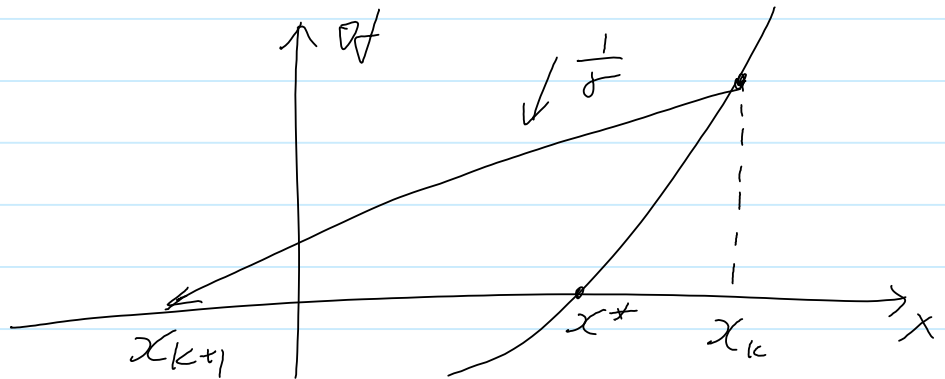$$\uparrow$$
gradient

- Note that $x^*$ is a fixed point of the iteration / mapping

— If $x(t) = x^*$, then $x(t+1) = x^*$.

The smaller $\gamma$
~~the smaller the update.~~

slope $\frac{1}{\gamma}$

$$x_k \to x^* \quad \text{if} \quad \gamma \text{ is small.}$$

$\frac{1}{\gamma}$

May not converge if $\gamma$ is too large.

$x_k$ $\quad$ $x$

May be more difficult when $\nabla f$ is not "smooth"

— when $f$ has sharp corners.

(20)

# Conditions related to convergence
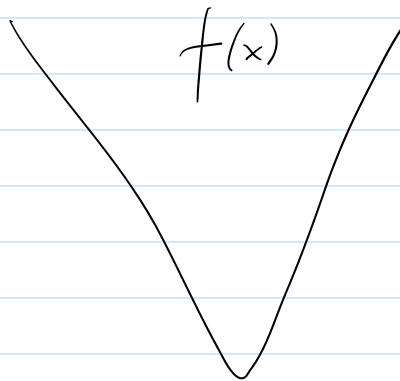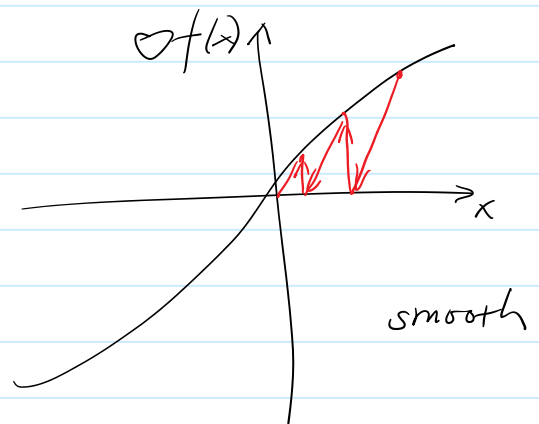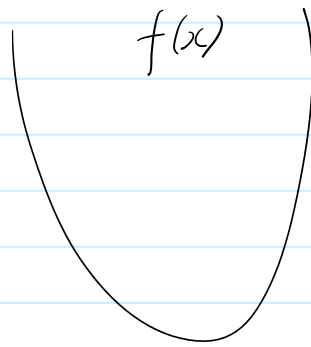
- How can we study the convergence of the gradient algorithms?

- We will usually encounter two types of conditions on the gradient $\nabla f(x)$

① Smoothness

- $\nabla f(x)$ does not abruptly change as $x$ change

$$\|\nabla f(x) - \nabla f(y)\|_2 \leq L \, \|x - y\|_2$$

$f(x)$

$\nabla f(x)$

smooth

$f(x)$

$\nabla f(x)$

non-smooth

— For non-smooth func, even when you are already in a neighborhood of $x^*$, it would be difficult for the gradient algorithm to converge exactly to $x^*$

— Instead, smoothness implies convergence, provided that the stepsize is sufficiently small.

(2) Strong convexity

- $\nabla f(x)$ increases quickly as $x$ lies away from $x^*$

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

$f(x)$

$\nabla f(x)$

strongly convex

$f(x)$

linear

$\nabla f(x)$

non-strongly-convex

- If the function is not strongly convex, then starting from $x(\circ)$ that is far away from $\bar{x}^*$,

the ⌣improvement of the gradient algorithm will be slow

— Instead, strong-convexity & smoothness combined implies geometric <u>descent</u>

$$\| x(t+1) - x^* \|_2 \leq \rho \, \| x(t) - x^* \|_2$$
$$\rho < 1$$

— This is usually referred to as "linear convergence" in the optimization literature

— In contrast to "quadratic convergence" by Newton's algorithm.

Skip

Two proof techniques:

① Show decrease of $f(x(t+1))$ : More intuitive but ∧less powerful results.

② Show decrease of $\| x(t) - x^* \|$

First Approach:

Lemma: Assume that function $f$ is continuously differentiable & there exists a constant $L$ such that

$$\| \nabla f(x) - \nabla f(y) \|_2 \leq L \|x - y\|_2$$

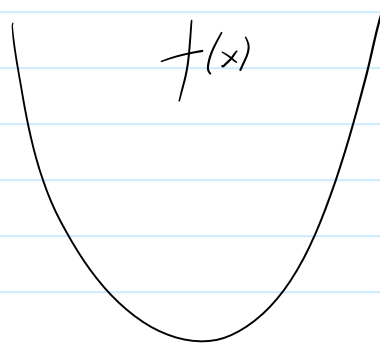<span style="color:red">Smoothness</span>

Then

$$f(y) \leq f(x) + \nabla f(x)^T (y-x) + \frac{L}{2} \|y - x\|_2^2$$

Note: The Lipshitz condition bounds how smooth $\nabla f(\cdot)$ is.

<span style="color:red">$g(t) = f(x + t(y-x))$</span>
<span style="color:red">$g'(t) = \nabla f(x + t(y-x)) \cdot (y-x)$</span>

Proof:

$$f(y) = f(x + (y-x)) \overset{\color{red}= g(\cdot)}{}$$

$$= \overset{\color{red}= g(\cdot)}{f(x)} + \int_0^1 \nabla f(x + t(y-x))^T \cdot (y-x) \, dt$$

$$= f(x) + \nabla f(x)^T \cdot (y-x)$$

$$\quad + \int_0^1 \left[ \nabla f(x + t(y-x)) - \nabla f(x) \right]^T (y-x) \, dt$$

$$\leq f(x) + \nabla f(x)^T (y-x)$$

$$\quad + \int_0^1 L t \|y-x\|^2 \, dt$$

$$= f(x) + \nabla f(x)^T (y-x) + \frac{L}{2} \|y-x\|^2.$$

Theorem: (Bertsekas & Tsitsiklis  p 203)

Assume that function $f$ is continuously differentiable & there exists a constant $L$ such that

$$\| \nabla f(x) - \nabla f(y) \|_2 \le L \| x - y \|_2 \qquad \forall x, y$$

Suppose $f$ is bounded from below by $f^*$.

Assume further that $0 < \gamma < \frac{2}{L}$. If the sequence of points $x(t)$ generated by

$$x(t+1) = x(t) - \gamma \nabla f(x(t))$$

has a limit point $x^*$, then

$$\nabla f(x^*) = 0$$

---

Start with the basic Taylor-series expansion:

$$f(x(t+1)) \approx f(x(t)) + \nabla f(x(t))^T (x(t+1) - x(t))$$
$$= f(x(t)) + \underbrace{\nabla f(x(t))^T \cdot [-\gamma \nabla f(x(t))]}_{\le 0}$$

— However, there may be other higher-order terms, which may create problems when $\nabla f(x(t))$ is already small

— The smoothness condition controls these higher-order terms.

Proof: Use the above Lemma. Let $x = x(t)$,
$$y = x(t) - \gamma \nabla f(x(t)) = x(t+1)$$

$$\Rightarrow f(x(t+1)) \le f(x(t)) - \gamma \| \nabla f(x(t)) \|^2$$
$$+ \frac{L}{2} \cdot \gamma^2 \| \nabla f(x(t)) \|^2 \qquad (*)$$

Since $\gamma < \frac{2}{L}$, we have $\gamma - \frac{L}{2} \gamma^2 > 0$.

Therefore, $f(x(t))$ is non-increasing. But it has a lower bound $f(x^*)$. Hence, the decrement must go to zero.

More precisely, summing $(*)$ over $t$, we have

$$f^* \leq f(x(t+1)) \leq f(x(0)) - \sum_{k=0}^{T} \left(\gamma - \frac{L}{2}\gamma^2\right) \|\nabla f(x(k))\|^2$$

$$\therefore \quad \sum_{k=0}^{t} \|\nabla f(x(k))\|^2 < f(x(0)) - f^*, \quad \forall t$$

$$\Rightarrow \quad \nabla f(x(t)) \rightarrow 0 \quad \text{as} \quad t \rightarrow +\infty$$

Since $x^*$ is a limit point of $x(t)$, and $\nabla f$ is continuous, we must have

$$\nabla f(x^*) = 0$$

(40)

— Does $x(t) \rightarrow x^*$ ?

Lemma: If $f$ is convex, then

$$(\nabla f(x) - \nabla f(y))^T (x - y) \geq 0 \qquad\qquad (*)$$

Note: This holds even if $\nabla f(\cdot)$ is a sub-gradient. A mapping $\nabla f(\cdot)$ that satisfy $(*)$ is called a monotone mapping.

- For 1-dim, $(f'(x) - f'(y))(x - y) \geq 0$
  $\Longleftrightarrow$ $f'(x) \geq f'(y)$ when $x > y$ (monotonicity)
  $\Longleftrightarrow$ $f''(x) \geq 0$.

Proof: From first-order condition of convexity;

$$f(y) \geq f(x) + \nabla f(x)(y - x)$$

$$f(x) \geq f(y) + \nabla f(y)(x - y)$$

Summing them,

$$\Rightarrow \quad [\nabla f(y) - \nabla f(x)]^T (y - x) \geq 0 \qquad\qquad \#$$
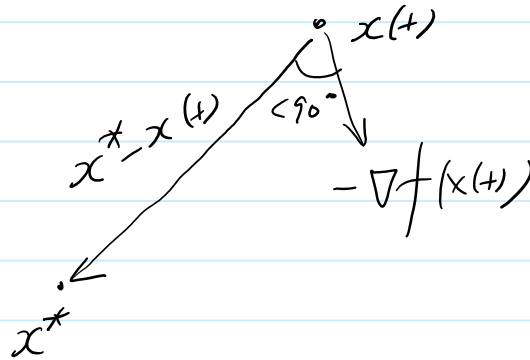
---

- Later on, we will apply this to

$$y = x(t)$$
$$x = x^*$$

$$\Rightarrow \quad \nabla f(x(t))^T \cdot \underline{(x(t) - x^*)} \geq 0$$

$$\Rightarrow \quad \nabla f(x(t))^T \cdot (x(t) - x^*) \geq 0$$
$$\left[ - \nabla f(x(t)) \right]^T (x^* - x(t)) \geq 0$$



If the derivative is Lipschitz-continuous, then a stronger version can be shown.

Lemma: Let $f$ be a convex & differentiable function such that

$$\| \nabla f(x) - \nabla f(y) \|_2 \leq L \| x - y \|_2 \quad \forall x, y \in R^n$$

Then

$$(\nabla f(x) - \nabla f(y))^T (x - y) \geq \frac{1}{L} \| \nabla f(x) - \nabla f(y) \|_2^2$$

$$\forall x, y \in R^n \quad (**)$$

Note: If a mapping $\nabla f(\cdot)$ that satisfies $(**)$, the inverse mapping is called strongly monotone.
$(\nabla f(x) \rightarrow x)$

We will prove this Lemma later.

Theorem: Assume that

$$\left(\nabla f(x) - \nabla f(y)\right)^T (x - y) \geq \frac{1}{L} \|\nabla f(x) - \nabla f(y)\|_2^2$$

$$\forall x, y \in \mathbb{R}^n$$

Further, assume that $0 < \gamma < \frac{2}{L}$ and there exists at least one point $x^*$ with $\nabla f(x^*) = 0$.

Then the sequence of $x(t)$ generated by

$$x(t+1) = x(t) - \gamma \nabla f(x(t))$$

converges, and the limit $x_0$ satisfies $\nabla f(x_0) = 0$.

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

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

$$= \|x(t) - x^*\|^2 + 2\left(x(t+1) - x(t)\right)^T \left(x(t) - x^*\right)$$
$$+ \|x(t+1) - x(t)\|^2$$

Note that

$$\left(x(t+1) - x(t)\right)^T \left(x(t) - x^*\right)$$

$$= -\gamma \left(\nabla f(x(t)) - \nabla f(x^*)\right) \left(x(t) - x^*\right) \overset{\substack{\| \\ 0}}{}$$

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

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

$$\leq \|x(t) - x^*\|^2 - \frac{2\sigma}{L} \|\nabla f(x(t))\|^2 + \sigma^2 \|\nabla f(x(t))\|^2$$

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

If $\sigma < \frac{2}{L}$, then $\alpha \overset{\Delta}{=} \frac{2\sigma}{L} - \sigma^2 > 0$

Summing over $t$

$$\|x(t+1) - x^*\|^2 \leq \|x(0) - x^*\|^2 - \alpha \sum_{s=0}^{t} \|\nabla f(x(s))\|^2$$

$$\overset{\vee}{\underset{0}{}}$$

$$\sum_{s=0}^{t} \|\nabla f(x(s))\|^2 < +\infty$$

$$\Rightarrow$$

$$\Rightarrow \quad \nabla f(x(t)) \to 0$$

(Q) why does $x(t)$ converge?

(A) (key technique) Assume that a subsequence converges, i.e,

$$x(t_h) \to x_0, \quad \text{as} \quad h \to +\infty$$

where

$$t_h \to +\infty \quad \text{as} \quad h \to +\infty.$$

We must then have $\nabla f(x_0) = 0$ due to the continuity of $\nabla f(\cdot)$.

Now, replace $x^*$ by $x_0$,

$$\lim_{t \to +\infty} \|x(t) - x_0\|^2 \leq \lim_{h \to +\infty} \|x(t_h) - x_0\|^2 = 0$$

$$\Rightarrow \quad x(t) \rightarrow x_0 \quad \text{as} \quad t \rightarrow +\infty.$$

$$\boxed{10}$$

- Assume that $\nabla f$ is Lipschitz

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x-y\| \qquad \forall x, y$$

We want to show that

$$\left(\nabla f(x) - \nabla f(y)\right)^T (x-y) \geq \frac{1}{L} \|\nabla f(x) - \nabla f(y)\|^2.$$

- Note that by the growth Lemma

(see the earlier section "First approach: Decrease of f")

$$f(x) \leq f(y) + \nabla f(y)(x-y) + \frac{L}{2}\|x-y\|^2 \qquad (*)$$

for any convex function and $\nabla f$ is Lipschitz.

We will show that

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

- Fix $y$.

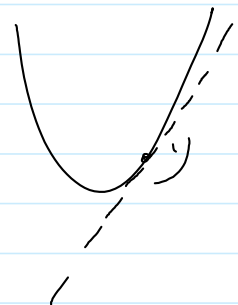Let $\hat{f}(x) = f(x) - f(y) - \nabla f(y)^T (x-y)$

We have

$$\hat{f}(y) = 0$$

$$\nabla \hat{f}(x) = \nabla f(x) - \nabla f(y)$$

$$\nabla \hat{f}(y) = 0$$

Hence, $y$ is the minimum point of $\hat{f}(\cdot)$.

- Let $z = x - \frac{1}{L}\nabla \hat{f}(x)$. Using $(*)$ on $\hat{f}(\cdot)$
have

$$0 \le \hat{f}(z) \le \hat{f}(x) + \nabla \hat{f}(x) \cdot \left( -\frac{1}{L} \nabla \hat{f}(x) \right)$$

$$+ \frac{L}{2} \cdot \frac{1}{L^2} \cdot \| \nabla \hat{f}(x) \|^2$$

— Using the definition of $\hat{f}(x)$, we have

$$f(x) - f(y) - \nabla f(y)^T (x-y) - \frac{1}{2L} \| \nabla f(x) - \nabla f(y) \|^2 \ge 0$$

Interchange the role of $x$ and $y$, we can similarly show that

$$f(y) - f(x) - \nabla f(x)^T (y-x) - \frac{1}{2L} \| \nabla f(x) - \nabla f(y) \|^2 \ge 0$$

Combine these two inequalities, we have

$$(\nabla f(y) - \nabla f(x))^T (y-x) \ge \frac{1}{L} \| \nabla f(y) - \nabla f(x) \|^2 .$$

$$\#$$