Examples:  $f(x) = \frac{1}{2}\|x\|_2^2$ is $1$-strongly convex

$$(\nabla^2 f(x) = I \geq 1 \cdot I)$$

$f(x) = \frac{1}{2}\|Ax-b\|_2^2$ is $\lambda_{min}(A^T A)$-strongly convex

(HW 2)

Lemma  If $\{f_i\}$ is $\mu_i$-strongly convex, then

$$\sum_i f_i \text{ is } \sum_i \mu_i \text{-strongly convex.}$$

Proof.  ~~We know~~ (Lets assume each $f_i$ is $C^1$).

Then for any $x, y$, we have

$$f_i(y) \geq f_i(x) + \nabla f_i(x)^T (y-x) + \frac{\mu_i}{2}\|y-x\|_2^2.$$

Summing over $i$

$$\sum_i f_i(y) \geq \sum_i f_i(x) + \left(\sum_i \nabla f_i(x)\right)^T (y-x) + \frac{\sum_i \mu_i}{2}\|y-x\|_2^2$$

$$\underset{\nabla(\sum_i f_i)(x)}{}$$

This equivalent to $\sum f_i$ being $\sum \mu_i$-strongly convex.   □

Then we know the following are strongly convex.

1. Training Support Vector Machines: $\pm 1$ as a label.

Given observations $(x_i, y_i)$, we want $w^T x \approx \text{sign}(y)$

(feature vector)

$$\min_{w} \sum_{i} \left( \max\{0, 1 - y_i\, w^T x_i\} \right) + \frac{\lambda}{2} \|w\|_2^2$$

$\lambda$-strongly convex

2. Sparse Regression (LASSO)

$$\min_{x} \|Ax - b\|_2^2 + \lambda \|x\|_1, \quad = \sum |x_i|$$

$\lambda_{\min}(A^T A)$ - strongly convex

"0

---

__Theorem__ Let $f$ be $\mu$-strongly convex with $L$-Lipschitz gradient and has a minimizer $x^*$. Then GD with $\alpha_K = \frac{2}{\mu + L}$ has

$$\|x_{K} - x^*\|^2 \leq \left(1 - \frac{4\mu^2}{(\mu+L)^2}\right)^K \|x_0 - x^*\|^2.$$

**Proof.** We want to a contraction at each step.

$$\|x_{k+1} - x^*\|^2 = \|x_k - \frac{2}{\mu+L}\nabla f(x_k) - x^*\|^2$$

$$= \|x_k - x^*\|^2 + \frac{4}{(\mu+L)^2}\|\nabla f(x_k)\|^2 - \frac{4}{\mu+L}\nabla f(x_k)^T(x_k - x^*)$$

Recall $\quad (\nabla f(x_k) - \nabla f(x^*))^T(x_k - x^*) \geq \frac{1}{L}\|\nabla f(x_k) - \nabla f(x^*)\|^2$

$\frac{L}{\mu+L}\nearrow$

(by smoothness equiv condition)

$\frac{\mu}{\mu+L}\longrightarrow (\nabla f(x_k) - \nabla f(x^*))^T(x_k - x^*) \geq \mu\|x_k - x^*\|^2$

$$\implies (\nabla f(x_k) - \nabla f(x^*))^T(x_k - x^*)$$

$$\geq \frac{L}{\mu+L}\cdot\frac{1}{L}\|\nabla f(x_k) - \nabla f(x^*)\|^2$$

$$+ \frac{\mu}{\mu+L}\mu\|x_k - x^*\|^2$$

$$= \frac{1}{\mu+L}\|\nabla f(x_k) - \nabla f(x^*)\|^2$$

$$+ \frac{\mu^2}{\mu+L}\|x_k - x^*\|^2.$$

$$\leq \|x_k - x^*\|^2 + \frac{4}{(\mu+L)^2}\|\nabla f(x_k) - \nabla f(x^*)\|^2$$

$$- \frac{4}{\mu+L}\left(\frac{1}{\mu+L}\|\nabla f(x_k) - \nabla f(x^*)\|^2 + \frac{\mu^2}{(\mu+L)^2}\|x_k - x^*\|^2\right)$$

$$= \|x_k - x^*\|^2 - \frac{4\mu^2}{(\mu+L)^2}\|x_k - x^*\|^2. \qquad \square$$

# 6+7   Complexity Lowerbounds and Acceleration

We have $\overset{\text{shown}}{\wedge}$ GD has

$$f(x_k) - f(x^*) \leq \frac{2L \|x_0 - x^*\|^2}{k}.$$

Lets imagine a wider class of algorithms following gradient directions

<u>Assumption 1</u>   The given method produce points $x_k$ satisfying

$$x_k \in x_0 + \text{Lin}\{\nabla f(x_0), \ldots, \nabla f(x_{k-1})\}.$$

$$\left(\text{For example,} \quad x_k = x_0 - \sum_{i=0}^{k-1} \alpha_{\blacksquare i} \nabla f(x_i)\right)$$

is the gradient descent sequence.

<u>Theorem</u>   For any $1 \leq k \leq \frac{1}{2}(d-1)$ and $L \geq 0$, there exists a convex function $f: \mathbb{R}^d \to \mathbb{R}$ that L-smooth such any algorithm satisfying Assumption 1 has

$$f(x_k) - f(x^*) \geq \frac{3L \|x_0 - x^*\|^2}{32(k+1)^2}$$

$$\|x_k - x^*\|^2 \geq \frac{1}{32} \|x_0 - x^*\|^2.$$

where $x^*$ minimizes $f$.

It turns out we can give a faster method

Nesterov's Accelerated Gradient Method (1983)

Let $y_0 = x_0$. Then iterate

$$y_{k+1} = x_k - \frac{1}{L} \nabla f(x_k)$$

$$x_{k+1} = y_{k+1} + \left( \frac{\lambda_k - 1}{\lambda_{k+1}} \right) (y_{k+1} - y_k)$$
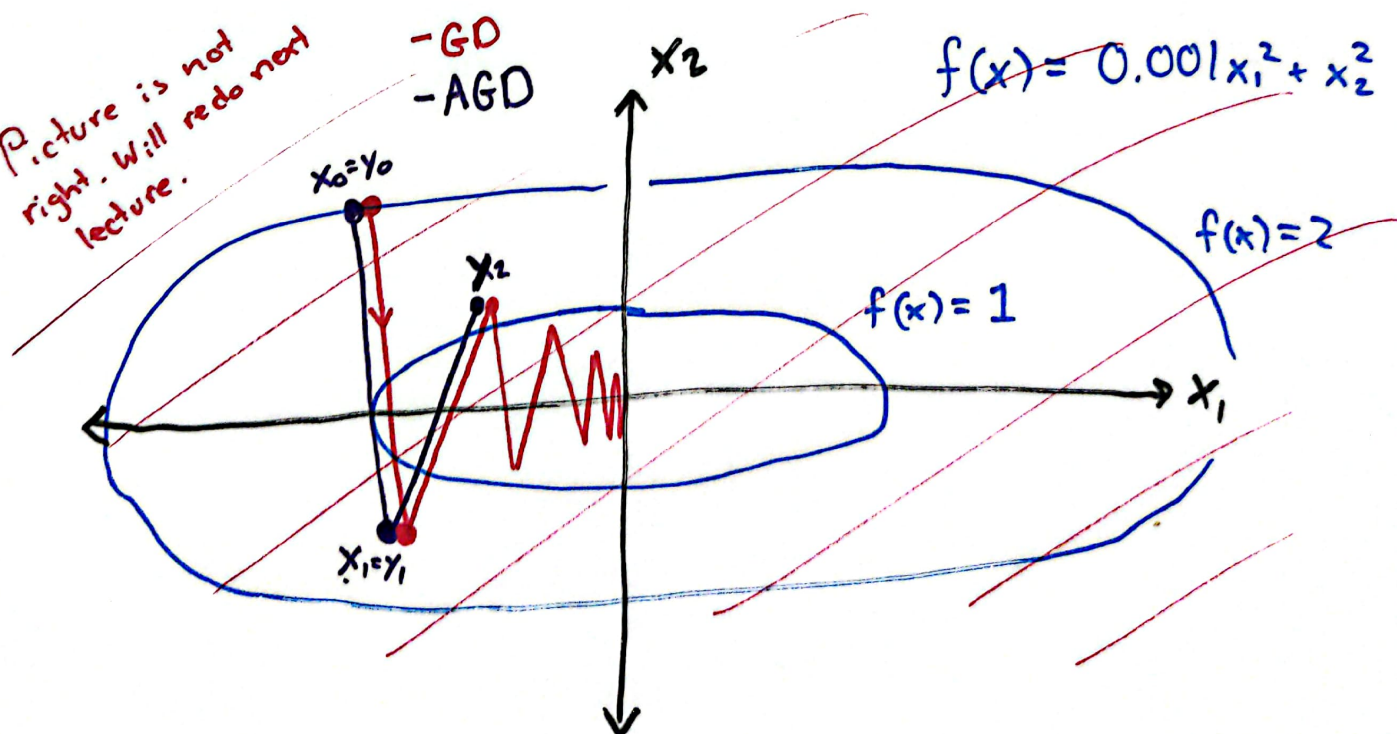
$$\approx \frac{k}{k+3}$$

where $\lambda_0 = 0$

$$\lambda_{k+1} = \frac{1 + \sqrt{1 + 4\lambda_k^2}}{2}.$$

Note $x_k \in x_0 + \text{Lin} \{ \nabla f(x_0), \ldots \nabla f(x_{k-1}) \}$

Picture is not right. Will redo next lecture.

−GD
−AGD

$x_2$

$f(x) = 0.001 x_1^2 + x_2^2$

$x_0 = y_0$

$x_2$

$f(x) = 2$

$f(x) = 1$

$x_1 = y_1$

$x_1$

**Theorem**  Let $f$ be convex with $L$-Lipschitz grad.
Then for any minimizer $x^*$,

$$f(y_k) - f(x^*) \leq \frac{2L \|x_0 - x^*\|^2}{k^2}.$$