

CS 2601 Linear and Convex Optimization

6. Gradient descent (part 2)

Bo Jiang

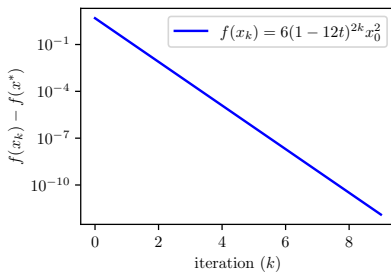
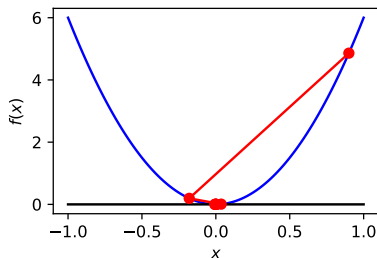
John Hopcroft Center for Computer Science
Shanghai Jiao Tong University

Fall 2022

Fast convergence

The following f is 12-smooth,

$$f(x) = 6x^2$$



For small enough step size t (e.g. 0.1),

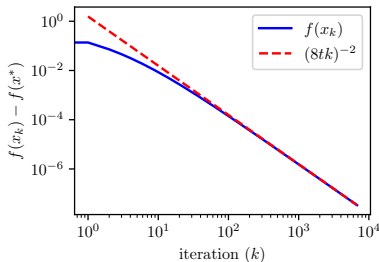
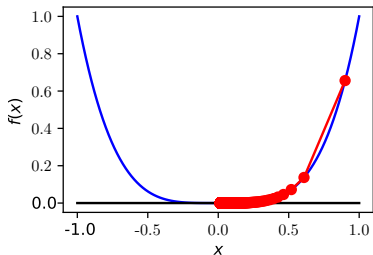
$$f(x_k) = 6x_0^2(1 - 12t)^{2k}$$

Need $O(\log \frac{1}{\epsilon})$ iterations to get within ϵ from optimal.

Slow convergence

The following f is also 12-smooth,

$$f(x) = \begin{cases} x^4, & \text{if } |x| \leq 1 \\ 4|x| - 3, & \text{if } |x| \geq 1 \end{cases}$$



For $x_0 \in (0, 1)$, small enough step size t (e.g. 0.1), and large k ,

$$x_k \sim \frac{1}{\sqrt{8tk}}, \quad f(x_k) \sim \frac{1}{(8tk)^2}$$

Need $O(1/\sqrt{\epsilon})$ iterations to get within ϵ from optimal value (i.e. 0).

Strong convexity

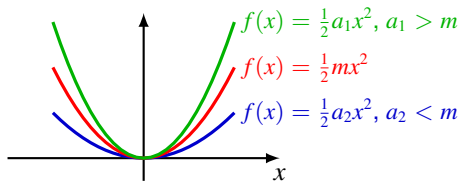
A function f is **strongly convex** with parameter $m > 0$, or simply **m -strongly convex**, if

$$\tilde{f}(\mathbf{x}) = f(\mathbf{x}) - \frac{m}{2}\|\mathbf{x}\|^2$$

is convex.

Note. $f(\mathbf{x}) = \frac{m}{2}\|\mathbf{x}\|^2 + \tilde{f}(\mathbf{x})$, i.e. f is $\frac{m}{2}\|\mathbf{x}\|^2$ plus an extra convex term.
Informally, “ m -strongly convex” means at least as “convex” as $\frac{m}{2}\|\mathbf{x}\|^2$.

Example. $f(\mathbf{x}) = \frac{a}{2}\|\mathbf{x}\|^2$ is m -strongly convex iff $a \geq m$



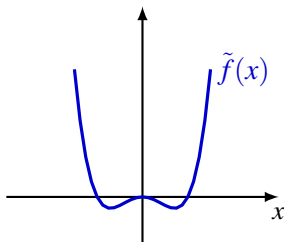
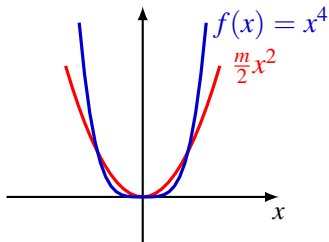
Strong convexity (cont'd) stronger than strict

Example. $f(\mathbf{x}) = \mathbf{a}^T \mathbf{x}$ is not m -strongly convex for any $m > 0$, as $\tilde{f}(\mathbf{x}) = \mathbf{a}^T \mathbf{x} - \frac{m}{2} \|\mathbf{x}\|^2$ is concave.

Example. $f(x) = x^4$ is not m -strongly convex for any $m > 0$, as $\tilde{f}(x) = x^4 - \frac{m}{2}x^2$ is not convex,

$$\tilde{f}''(x) = 12x^2 - m < 0$$

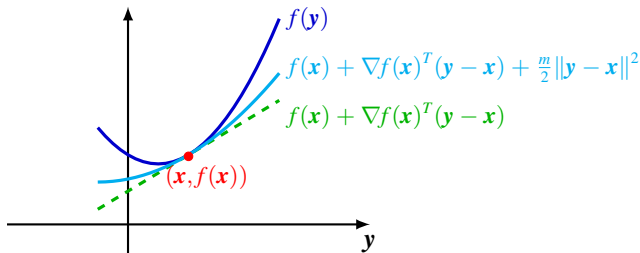
for $|x| < \sqrt{m/12}$.



First-order condition for strong convexity

A differentiable f is m -strongly convex iff

$$f(y) \geq f(x) + \nabla f(x)^T(y - x) + \frac{m}{2}\|x - y\|^2, \quad \forall x, y$$



- strong convexity \implies strict convexity \implies convexity
- m -strong convexity and L -smoothness together imply

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

Proof

1. By definition,

$$f \text{ is } m\text{-strongly convex} \iff \tilde{f}(\mathbf{x}) = f(\mathbf{x}) - \frac{m}{2}\|\mathbf{x}\|^2 \text{ is convex}$$

2. By first-order condition for convexity,

$$\iff \tilde{f}(\mathbf{y}) \geq \tilde{f}(\mathbf{x}) + \nabla \tilde{f}(\mathbf{x})^T(\mathbf{y} - \mathbf{x}), \quad \forall \mathbf{x}, \mathbf{y}$$

3. Noting $\nabla \tilde{f}(\mathbf{x}) = \nabla f(\mathbf{x}) - m\mathbf{x}$,

$$\iff f(\mathbf{y}) - \frac{m}{2}\|\mathbf{y}\|^2 \geq f(\mathbf{x}) - \frac{m}{2}\|\mathbf{x}\|^2 + (\nabla f(\mathbf{x}) - m\mathbf{x})^T(\mathbf{y} - \mathbf{x}), \quad \forall \mathbf{x}, \mathbf{y}$$

4. Rearranging and using $\mathbf{y}^T\mathbf{y} - \mathbf{x}^T\mathbf{x} - 2\mathbf{x}^T(\mathbf{y} - \mathbf{x}) = (\mathbf{y} - \mathbf{x})^T(\mathbf{y} - \mathbf{x})$,

$$\iff f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{m}{2}\|\mathbf{x} - \mathbf{y}\|^2, \quad \forall \mathbf{x}, \mathbf{y}$$

Second-order condition for strong convexity

A twice continuously differentiable f is m -strongly convex iff

$$\nabla^2 f(\mathbf{x}) \succeq m\mathbf{I}, \quad \forall \mathbf{x}$$

or equivalently, the smallest eigenvalue of $\nabla^2 f(\mathbf{x})$ satisfies

$$\lambda_{\min}(\nabla^2 f(\mathbf{x})) \geq m, \quad \forall \mathbf{x}$$

Proof. $\tilde{f}(\mathbf{x}) = f(\mathbf{x}) - \frac{m}{2}\|\mathbf{x}\|^2$ is convex iff $\nabla^2 \tilde{f}(\mathbf{x}) = \nabla^2 f(\mathbf{x}) - m\mathbf{I} \succeq \mathbf{O}$

Example. With $\mathbf{Q} = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}$, we obtain $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{Q} \mathbf{x} = \frac{1}{2}x_1^2 + x_2^2$ is 1-strongly convex.

More generally, $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{Q} \mathbf{x}$ with $\mathbf{Q} \succ \mathbf{O}$ is $\lambda_{\min}(\mathbf{Q})$ -strongly convex, where $\lambda_{\min}(\mathbf{Q})$ is the smallest eigenvalue of \mathbf{Q} .

Bound on suboptimality gap

If f is m -strongly convex, then

$$f(\mathbf{x}) - f(\mathbf{x}^*) \leq \frac{1}{2m} \|\nabla f(\mathbf{x})\|^2$$

Note. When $\nabla f(\mathbf{x})$ is small, then $f(\mathbf{x})$ is close to optimal. (Do we have this property for general convex functions?) **No!**

Proof. By the first-order condition,

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{m}{2} \|\mathbf{y} - \mathbf{x}\|^2$$

Minimize over \mathbf{y} on both sides,

$$f(\mathbf{x}^*) = \min_{\mathbf{y}} f(\mathbf{y}) \geq \min_{\mathbf{y}} \left[f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{m}{2} \|\mathbf{y} - \mathbf{x}\|^2 \right]$$

RHS minimized at $\mathbf{y} = \mathbf{x} - \frac{1}{m} \nabla f(\mathbf{x})$, so

$$f(\mathbf{x}^*) \geq f(\mathbf{x}) - \frac{1}{2m} \|\nabla f(\mathbf{x})\|^2$$

Convergence: 1D example

$f(x) = \frac{1}{2}mx^2$ with $m > 0$ is both m -smooth and m -strongly convex..

Recall the gradient descent step is

$$x_{k+1} = x_k - tf'(x_k) = (1 - mt)x_k$$

and $x_k \rightarrow x^* = 0$ iff $t \in (0, \frac{2}{m})$.

If $t = \frac{1}{m}$, it gets to x^* in one step.

For $t \in (0, \frac{1}{m}) \cup (\frac{1}{m}, \frac{2}{m})$,

$$x_k = (1 - mt)^k x_0$$

so both $x_k \rightarrow x^*$ and $f(x_k) \rightarrow f(x^*)$ exponentially fast,

$$|x_k - x^*| = (1 - mt)^k \cdot |x_0 - x^*|$$

$$|f(x_k) - f(x^*)| = \frac{m(1 - mt)^{2k}}{2} |x_0 - x^*|^2$$

Convergence analysis

Theorem. If f is m -strongly convex and L -smooth, and \mathbf{x}^* is a minimum of f , then for step size $t \in (0, \frac{1}{L}]$, the sequence $\{\mathbf{x}_k\}$ produced by the gradient descent algorithm satisfies

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \frac{L(1 - mt)^k}{2} \|\mathbf{x}_0 - \mathbf{x}^*\|^2$$
$$\|\mathbf{x}_k - \mathbf{x}^*\|^2 \leq (1 - mt)^k \|\mathbf{x}_0 - \mathbf{x}^*\|^2$$

Notes.

- $0 \leq 1 - \frac{m}{L} \leq 1 - mt < 1$, so $\mathbf{x}_k \rightarrow \mathbf{x}^*$ and $f(\mathbf{x}_k) \rightarrow f(\mathbf{x}^*)$ exponentially fast
- The number of iterations to reach $f(\mathbf{x}_k) - f(\mathbf{x}^*) \leq \epsilon$ is $O(\log \frac{1}{\epsilon})$. For $\epsilon = 10^{-p}$, $k = O(p)$, linear in the number of significant digits!
- Since $\nabla f(\mathbf{x}^*) = 0$, the bounds on slide 5 yield

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

relating the bounds on $\|\mathbf{x}_k - \mathbf{x}^*\|^2$ and those on $f(\mathbf{x}_k) - f(\mathbf{x}^*)$

Proof

Similar to proof without strong convexity, with difference highlighted.

1. By the basic gradient step $\mathbf{x}_{k+1} = \mathbf{x}_k - t\nabla f(\mathbf{x}_k)$,

$$\begin{aligned}\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 &= \|\mathbf{x}_k - t\nabla f(\mathbf{x}_k) - \mathbf{x}^*\|^2 \\ &= \|\mathbf{x}_k - \mathbf{x}^*\|^2 + t^2\|\nabla f(\mathbf{x}_k)\|^2 + 2t\nabla f(\mathbf{x}_k)^T(\mathbf{x}^* - \mathbf{x}_k)\end{aligned}$$

2. By L -smoothness, the second term is upper bounded by

$$t^2\|\nabla f(\mathbf{x}_k)\|^2 \leq 2t[f(\mathbf{x}_k) - f(\mathbf{x}_{k+1})]$$

3. By m -strong convexity,

$$\nabla f(\mathbf{x}_k)^T(\mathbf{x}^* - \mathbf{x}_k) \leq f(\mathbf{x}^*) - f(\mathbf{x}_k) - \frac{m}{2}\|\mathbf{x}_k - \mathbf{x}^*\|^2$$

4. Plugging 2 and 3 into 1,

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 \leq (1 - mt)\|\mathbf{x}_k - \mathbf{x}^*\|^2 + 2t[f(\mathbf{x}^*) - f(\mathbf{x}_{k+1})]$$

5. Since $f(\mathbf{x}^*) \leq f(\mathbf{x}_{k+1})$,

$$\|\mathbf{x}_{k+1} - \mathbf{x}^*\|^2 \leq (1 - mt)\|\mathbf{x}_k - \mathbf{x}^*\|^2$$

Convergence: 2D quadratic function

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x}, \quad \mathbf{Q} = \begin{pmatrix} m & 0 \\ 0 & L \end{pmatrix}$$

where $L > m > 0$. f is L -smooth and m -strongly convex. $\mathbf{x}^* = \mathbf{0}$.

The gradient descent step is

$$\mathbf{x}_{k+1} = \mathbf{x}_k - t \nabla f(\mathbf{x}_k) = (\mathbf{I} - t \mathbf{Q}) \mathbf{x}_k$$

so

$$\mathbf{x}_k = (\mathbf{I} - t \mathbf{Q})^k \mathbf{x}_0 = \begin{bmatrix} (1 - mt)^k x_{01} \\ (1 - Lt)^k x_{02} \end{bmatrix}$$

and

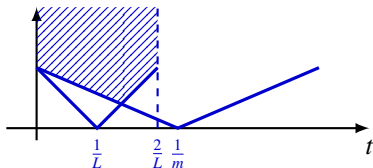
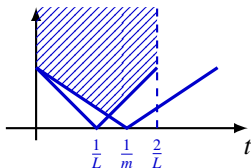
$$f(\mathbf{x}_k) = \frac{m}{2} (1 - mt)^{2k} x_{01}^2 + \frac{L}{2} (1 - Lt)^{2k} x_{02}^2$$

To ensure convergence, $t < \frac{2}{L}$. The convergence rate is determined by the slower of $(1 - Lt)^{2k}$ and $(1 - mt)^{2k}$.

Convergence: 2D quadratic function (cont'd)

To maximize convergence rate, solve

$$\begin{aligned} \min_t \quad & \max\{|1 - Lt|, |1 - mt|\} \\ \text{s. t.} \quad & 0 < t < 2/L \end{aligned}$$



Maximum rate achieved by $1 - mt = Lt - 1 \implies t = \frac{2}{m+L}$, in which case

$$\mathbf{x}_k = \left(\frac{L-m}{L+m}\right)^k \begin{bmatrix} x_{01} \\ (-1)^k x_{02} \end{bmatrix} \implies \|\mathbf{x}_k - \mathbf{x}^*\|_2 = \left(\frac{L-m}{L+m}\right)^k \|\mathbf{x}_0 - \mathbf{x}^*\|_2$$

$$f(\mathbf{x}_k) - f(\mathbf{x}^*) = \left(\frac{L-m}{L+m}\right)^{2k} [f(\mathbf{x}_0) - f(\mathbf{x}^*)]$$

Depends on $\kappa(\mathbf{Q}) = \frac{\lambda_{\max}(\mathbf{Q})}{\lambda_{\min}(\mathbf{Q})} = \frac{L}{m}$, the condition number of \mathbf{Q}

Condition number

For a matrix $\mathbf{Q} \in \mathbb{R}^{n \times n}$ s.t. $\mathbf{Q} \succ \mathbf{0}$, its **condition number**¹ is defined as

$$\kappa(\mathbf{Q}) = \frac{\lambda_{\max}(\mathbf{Q})}{\lambda_{\min}(\mathbf{Q})}$$

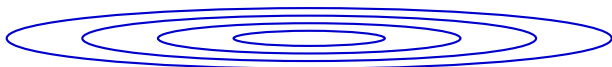
It characterizes **how stretched the level curves** of $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^T \mathbf{Q} \mathbf{x}$ are.

Example. $\mathbf{Q} = \text{diag}\{\gamma, 1\}$, $f(x_1, x_2) = \frac{\gamma}{2}x_1^2 + \frac{1}{2}x_2^2$



$$\mathbf{Q} = \text{diag}\{1, 1\}$$

$$\kappa(\mathbf{Q}) = 1$$



$$\mathbf{Q} = \text{diag}\{0.01, 1\}$$

$$\kappa(\mathbf{Q}) = 100$$

Nondiagonal case reduces to diagonal case in eigenbasis of \mathbf{Q} .

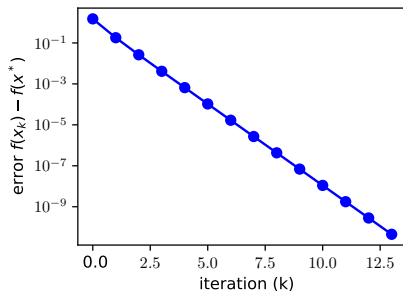
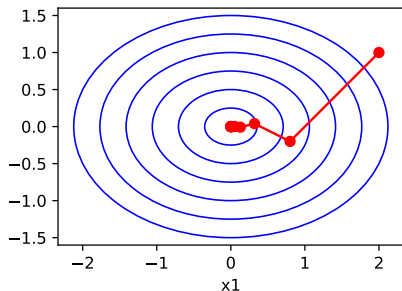
For nonquadratic case, $\kappa(\nabla^2 f(\mathbf{x}))$ plays a similar role.

¹For a general nonsingular matrix, **the condition number is the ratio between its largest and smallest singular values**, $\kappa(\mathbf{A}) = \sigma_{\max}(\mathbf{A})/\sigma_{\min}(\mathbf{A})$.

Well-conditioned Problem

The problem $\min_x \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x}$ is well-conditioned if $\kappa(\mathbf{Q})$ is small.

Example. $\mathbf{Q} = \text{diag}\{0.5, 1\}$, $f(x_1, x_2) = \frac{1}{4}x_1^2 + \frac{1}{2}x_2^2$, $\kappa(\mathbf{Q}) = 2$



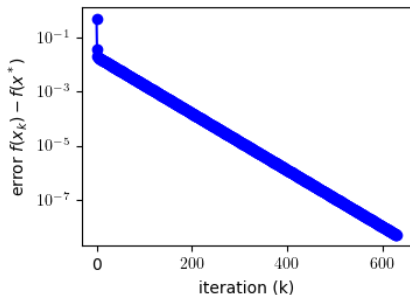
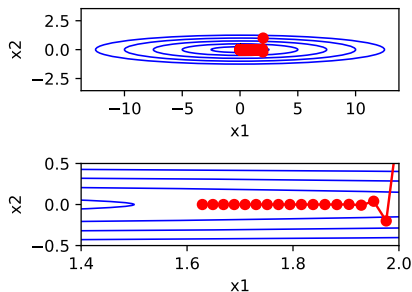
Fast convergence: for $\mathbf{x}_0 = (2, 1)^T$, $t = 1.2$, and large k ,

$$f(\mathbf{x}_k) \sim \frac{m}{2}(1 - mt)^{2k}x_{01}^2 = (0.4)^{2k}$$

Ill-conditioned problem

The problem $\min_x \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x}$ is **ill-conditioned** if $\kappa(\mathbf{Q})$ is large.

Example. $\mathbf{Q} = \text{diag}\{0.01, 1\}$, $f(x_1, x_2) = \frac{1}{200}x_1^2 + \frac{1}{2}x_2^2$, $\kappa(\mathbf{Q}) = 100$



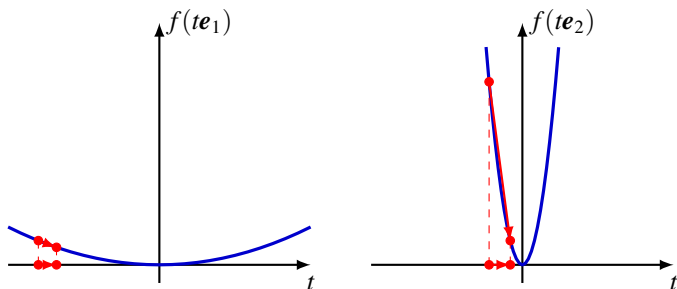
Slow convergence (relatively): for $\mathbf{x}_0 = (2, 1)^T$, $t = 1.2$, and large k ,

$$f(\mathbf{x}_k) \sim \frac{m}{2}(1 - mt)^{2k}x_{01}^2 = \frac{1}{50}(0.988)^{2k}$$

III-conditioned problem (cont'd)

$$f(x_1, x_2) = \frac{1}{2} \mathbf{x}^T \mathbf{Q} \mathbf{x} = \frac{1}{200} x_1^2 + \frac{1}{2} x_2^2, \quad \mathbf{Q} = \text{diag}\{0.01, 1\}, \quad \kappa(\mathbf{Q}) = 100$$

- 1-smooth \implies To guarantee convergence, step size² $t < 2$
- This limit is imposed by movement along e_2 direction
- Too pessimistic along other directions, e.g. along e_1 , can use $t < 200$



²We proved convergence for $t \in (0, 1/L]$. The proofs can be modified slightly to show convergence for $t \in (0, 2/L]$.

Ill-condition problem (cont'd)

The negative gradient direction is far away from the “ideal” direction for ill-conditioned problem.

For $\mathbf{Q} = \text{diag}\{\gamma, 1\}$, $f(x_1, x_2) = \frac{1}{2}\mathbf{x}^T \mathbf{Q} \mathbf{x} = \frac{\gamma}{2}x_1^2 + \frac{1}{2}x_2^2$,

negative gradient direction

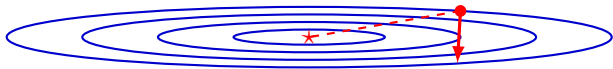
$$-\nabla f(\mathbf{x}) = -\mathbf{Q}\mathbf{x} = (-\gamma x_1, -x_2)^T$$

“ideal” direction

$$\mathbf{d} = -\mathbf{x}$$



$$\mathbf{Q} = \text{diag}\{1, 1\}$$
$$\kappa(\mathbf{Q}) = 1$$



$$\mathbf{Q} = \text{diag}\{0.01, 1\}$$
$$\kappa(\mathbf{Q}) = 100$$