

$f \in C'_L(\mathbb{R}^n) \Rightarrow$  Continuously differentiable functions with Lipschitz gradients

Lemma: Quadratic upper bound on  $f \in C'_L(\mathbb{R}^n)$

Let  $f \in C'_L(\mathbb{R}^n)$  and  $\text{dom } f$  is a convex set.

Then;  $\forall x, y \in \text{dom } f$

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

$\rightarrow$  quadratic function in  $y$



Proof:

Define  $g(t) = f(\underline{x + t(y-x)})$  for  $t \in [0, 1]$

Since  $\text{dom } f$  is convex  $(x + t(y-x))$  is in the  $\text{dom } f$ .

$$\xrightarrow{\quad} x + ty - tx = (1-t)x + ty$$

Note:  $g(1) = f(y)$

$g(0) = f(x)$

$$g'(t) = \nabla f(x + t(y-x))^T (y-x)$$

$$\boxed{g'(0) = \nabla f(x)^T (y-x)}$$

$$g'(0) = \nabla f(x)^T (y-x)$$

$$\begin{aligned} g'(t) - g'(0) &= \nabla f(x + t(y-x))^T (y-x) - \nabla f(x)^T (y-x) \\ &= \left[ \nabla f(x + t(y-x)) - \nabla f(x) \right]^T (y-x) \end{aligned}$$

Use Cauchy-Schwarz inequality ( $a^T b \leq \|a\|_2 \|b\|_2$ )

$$\leq \underbrace{\|\nabla f(x + t(y-x)) - \nabla f(x)\|_2}_{\text{Lipschitz continuous gradients}} \|y-x\|_2$$

Lipschitz continuous gradients

$$\leq L \|x + t(y-x) - x\|_2 \|y-x\|_2$$

$$g'(t) - g'(0) \leq tL \|y-x\|_2^2$$

$$\int_0^1 g'(t) dt = g(1) - g(0)$$

$$\Rightarrow g(1) = g(0) + \int_0^1 g'(t) dt$$

$$\leq g(0) + \int_0^1 (tL \|y-x\|_2^2 + g'(0)) dt$$

$$f(y) \leq f(x) + L \|y-x\|_2^2 \frac{t^2}{2} \Big|_0^1 + g'(0)$$

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



Descent Lemma (when  $\Delta x = \nabla f(x)$ )

## Descent Lemma (when $\Delta x = -\nabla f(x)$ )

Let  $f \in C'_L(\mathbb{R}^n)$  with  $\text{dom } f$  being convex. Let us consider the iterative method:

$$x^{(k+1)} = x^{(k)} + t^{(k)} \Delta x^{(k)}$$

and focus on the case  $\Delta x^{(k)} = -\nabla f(x^{(k)})$ . Then, as long as  $0 < t^{(k)} < \frac{2}{L}$ ,  $\Delta x^{(k)}$  is a descent direction, i.e.  $f(x^{(k+1)}) < f(x^{(k)})$  as long as  $x^{(k)}$  is not  $x^*$ .

Remarks: A range of step sizes guarantee descent. The larger  $L$  is, the smaller is the range.

Proof: Let  $y = x^{(k)} + t^{(k)} \Delta x^{(k)}$

and  $x = x^{(k)}$  in the Quadratic upper bound.

$$\begin{aligned} f(x^{(k)} + t^{(k)} \Delta x^{(k)}) &\leq f(x^{(k)}) + \nabla f(x^{(k)})^T (x^{(k)} + t^{(k)} \Delta x^{(k)} - x^{(k)}) \\ &\quad + \frac{L}{2} \|x^{(k)} + t^{(k)} \Delta x^{(k)} - x^{(k)}\|_2^2 \end{aligned}$$

$$f(x^{(k+1)}) \leq f(x^{(k)}) + t^{(k)} \nabla f(x^{(k)})^T \Delta x^{(k)} + t^{(k)2} \frac{L}{2} \|\Delta x^{(k)}\|_2^2$$

Now put  $\Delta x^{(k)} = -\nabla f(x^{(k)})$

$$f(x^{(k+1)}) \leq f(x^{(k)}) - t^{(k)} \|\nabla f(x^{(k)})\|_2^2 + t^{(k)2} \frac{L}{2} \|\nabla f(x^{(k)})\|_2^2$$

need  $< 0$

$$\leq f(x^{(k)}) - \underbrace{\left(t^{(k)} - t^{(k)^2} \frac{L}{2}\right)}_{\text{most reduction when this is the largest.}} \|\nabla f(x)\|_2^2$$

clearly,  $f(x^{(k+1)}) < f(x^{(k)})$

if and only if

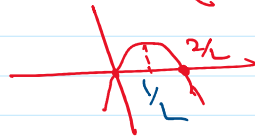
$$t^{(k)} - t^{(k)^2} \frac{L}{2} > 0$$

$$\Leftrightarrow t^{(k)^2} \frac{L}{2} < t^{(k)}$$

$$\Leftrightarrow t^{(k)} < \frac{2}{L}$$

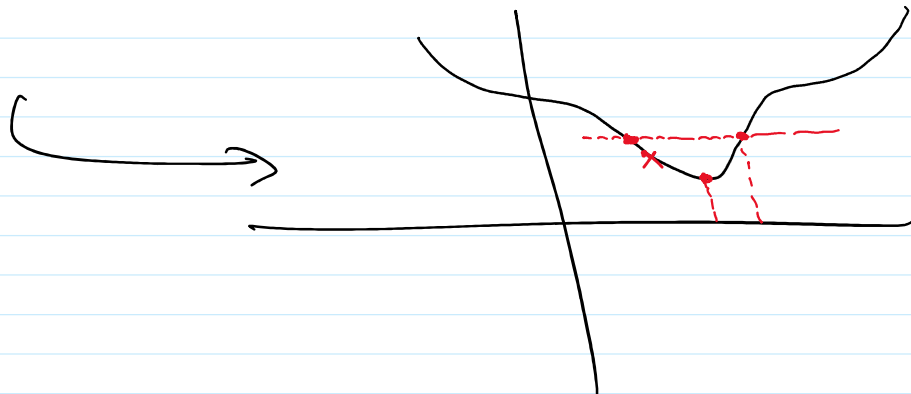


what is the  $t^{(k)}$  when 'L' is known?



The best  $t^{(k)}$ , when it comes to most reduction for gradient descent is  $t^{(k)} = \frac{1}{L}$ .

Ex:  $f(x)$



General Descent Method

Initialize:  $x^{(0)} \in \text{dom} f$   
 $k \leftarrow 0$

Repeat

1. Determine a descent direction  $\Delta x^{(k)}$   
(i.e.,  $-\nabla f(x^{(k)})^\top \Delta x^{(k)} > 0$ )

2. Line Search: Choose a step size  $t^{(k)} > 0$

3. Update the iterate:  $x^{(k+1)} \leftarrow x^{(k)} + t^{(k)} \Delta x^{(k)}$

Until stopping criterion is satisfied.

Special Cases:

① Gradient descent:  $\Delta x^{(k)} = -\nabla f(x^{(k)})$

② Newton's method:  $\Delta x^{(k)} = -\nabla^2 f(x^{(k)})^{-1} \nabla f(x^{(k)})$

Line Search: How to pick the step size in a descent method?

① Choose a fixed step size  $t^{(k)} = \eta \ \forall \ k \geq 0$ .

e.g., If  $f \in C'_L(\mathbb{R}^n)$  and  $L$  is known or can be computed efficiently then pick

$$\eta = \frac{1}{L} \text{ for gradient descent.}$$

In other cases, trial and error helps pick a step size.

Issues with this approach:

- Sometimes it is too costly to compute  $L$ .
- Sometimes functions are not  $C'_L$ .
- Even when one has descent in some iterations, does not mean the step size would work in all points in dom  $f$ .

Unless  $L$  is known, the larger the step size, the better 'perhaps'.

## ② Variable step size $t^{(k)}$

↳ The most common approach in the literature.

(a) Decaying step size policy

$$t^{(k)} \rightarrow 0 \text{ as } k \rightarrow \infty$$

Typical policy

$$(i) \quad t^{(k)} \rightarrow 0 \text{ as } k \rightarrow \infty$$

$$(ii) \quad \sum_{k=0}^{\infty} t^{(k)} = \infty$$

$$(iii) \quad \sum_{k=0}^{\infty} [t^{(k)}]^2 < \infty$$

e.g.:  $t^{(k)} = \frac{\text{const}}{k}$

↳ Often used in machine learning / stochastic

↳ Often used in machine learning / Stochastic optimization.

(b) Search for a nice step size that reduces the objective function using a subroutine in each iteration  $k$ .

↳ often used in practical deterministic optimization problems.

↳ we will study this in detail under 'Exact line search' and 'Inexact line search'.

### Exact line search

$$x^{(k+1)} = x^{(k)} + t^{(k)} \Delta x^{(k)}$$

When we have fixed  $\Delta x^{(k)}$  and are looking for  $t^{(k)}$ , we are effectively looking at a one-dimensional function:

$$\tilde{f}(t) = f(x^{(k)} + t \Delta x^{(k)})$$

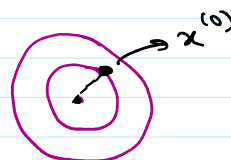
$$\text{Pick } t^{(k)} = \arg \min_{t \geq 0} \tilde{f}(t)$$

Example: Say  $f(x) = x_1^2 + x_2^2$

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

$$\text{Let } x^{(0)} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} = \underline{1}$$

$$\Rightarrow x^{(1)} = x^{(0)} - 2t^{(1)} x^{(0)}$$



$$\begin{aligned}
 &= (1 - 2t^{(1)}) \mathbb{1} \quad \leftarrow \quad = \left(1 - \frac{1}{2} \times 2\right) \mathbb{1} = 0 \cdot \mathbb{1} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \\
 \tilde{f}(t) &= f(x^{(1)}) = (1 - 2t)^2 \cdot 1 + (1 - 2t)^2 \cdot 1 \\
 &= 2(1 - 2t)^2 \\
 \boxed{t^{(1)}} &\stackrel{1}{=} \arg \min_{t \geq 0} \tilde{f}(t) = \arg \min_{t \geq 0} 2(1 - 2t)^2
 \end{aligned}$$

$$\tilde{f}'(t) = 4(1 - 2t) \cdot -2 = -8(1 - 2t)$$

$$\tilde{f}'(t) = 0 \Rightarrow -8(1 - 2t) = 0 \\ t = \frac{1}{2}$$

Exact line search, in which one solves a one-dim optimization problem in each iteration, works in cases where:

- ① The solution has an analytical form.
- ② It might be computationally feasible to numerically solve the 1-D problem.

But if the cost of exact line search is too much, we resort to 'inexact line search'.

↳ Backtracking (Armijo-Goldstein line search)