

# COMPLEXITY IN CONTINUOUS OPTIMIZATION

Course 1: January 30, 2025

Detailed proofs:

$$\|x\| = \sqrt{\sum_{i=1}^m x_i^2} \quad \forall i=1..m$$

Global minimum  $\Rightarrow$  local minimum  $\Rightarrow \|\nabla f(\bar{x})\| = 0$

---

Key Lemma

$\nexists x_k \in \mathbb{R}^m$  is such that  $\|\nabla f(x_k)\| \neq 0$ ,

then for any  $\alpha > 0$ ,

$$\begin{aligned} f(x_k - \alpha \nabla f(x_k)) &\leq f(x_k) + \nabla f(x_k)^T (x_k - \alpha \nabla f(x_k) - x_k) \\ &\quad + \frac{L}{2} \|x_k - \alpha \nabla f(x_k) - x_k\|^2 \\ &= f(x_k) - \alpha \|\nabla f(x_k)\|^2 + \frac{L}{2} \alpha^2 \|\nabla f(x_k)\|^2 \\ &= f(x_k) - \underbrace{\left(\alpha - \frac{L}{2} \alpha^2\right)}_{> 0 \text{ when } \alpha < \frac{2}{L}} \|\nabla f(x_k)\|^2 \end{aligned}$$

Hence

$$f(x_k - \alpha \nabla f(x_k)) < f(x_k) \text{ when } \alpha < \frac{2}{L}$$


---

*Iteration complexity result*

$$\text{Set } \alpha_k = \frac{1}{L} \quad \forall k$$

(why?)

$$\alpha - \frac{L}{2} \alpha^2$$

maximised for  $\alpha = \frac{1}{L}$

Then,  $\forall k$ ,

$$\begin{aligned} f(x_{k+1}) &\leq f(x_k) - \left(\frac{1}{L} - \frac{L}{2} \left(\frac{1}{L}\right)^2\right) \|\nabla f(x_k)\|^2 \\ &= f(x_k) - \frac{1}{2L} \|\nabla f(x_k)\|^2 \end{aligned}$$

Suppose that  $\forall k \geq 1$ ,  $\|\nabla f(x_k)\| > \varepsilon \quad \forall k=0 \dots K-1$

(false  $\Rightarrow \|\nabla f(x_0)\| \leq \varepsilon$ , so the method stops at the first iteration)

Then,  $\forall k=0 \dots K-1$ ,  $f(x_{k+1}) \leq f(x_k) - \frac{1}{2L} \|\nabla f(x_k)\|^2$

$$\text{—————} \quad f(x_{k+1}) < f(x_k) - \frac{1}{2L} \varepsilon^2$$

$$\Leftrightarrow \text{—————} \quad , \quad \frac{1}{2L} \varepsilon^2 < f(x_k) - f(x_{k+1})$$

Sum the inequalities:

$$\underbrace{\sum_{k=0}^{K-1} \frac{1}{2L} \varepsilon^2}_{\frac{K\varepsilon^2}{2L}} < \underbrace{\sum_{k=0}^{K-1} (f(x_k) - f(x_{k+1}))}_{\text{"telescopic sum"}} = f(x_0) - f(x_K)$$

$$\frac{K\varepsilon^2}{2L} < f(x_0) - f(x_K) \leq f(x_0) - \underbrace{f_{\text{low}}}_{\substack{\text{Lower bound} \\ \text{of } f}}$$

$$\Leftrightarrow K < 2L (f(x_0) - f_{\text{low}}) \varepsilon^{-2}$$

We have shown that if  $\|\nabla f(x_k)\| > \varepsilon \quad \forall k=0 \dots K-1$ ,  
then  $K < 2L (f(x_0) - f_{\text{low}}) \varepsilon^{-2}$

Thus, if we run the method for  $\lceil 2L (f(x_0) - f_{\text{low}}) \varepsilon^{-2} \rceil$  iterations,

Then necessarily there will be an iterate  $x_{\frac{1}{k}}$  such that  $\|\nabla f(x_{\frac{1}{k}})\| \leq \varepsilon$ .

Condition: The algorithm will terminate in at most

$\lceil 2L (f(x_0) - f_{\text{low}}) \varepsilon^{-2} \rceil$  iterations.

Lipschitz  
constant

Distance (in function value) between  $x_0$   
and the lower bound

(Feature of nonlinearity)

→ Bound is  $O(\varepsilon^{-2})$  ( $\leq C \varepsilon^{-2}$  where  $C > 0$  that does not depend on  $\varepsilon$ )

→  $(\varepsilon \rightarrow \frac{\varepsilon}{10}) \Rightarrow (O(\varepsilon^{-2}) \rightarrow O(100\varepsilon^{-2}))$

→ The result says that  $\min_{0 \leq k \leq K-1} \|\nabla f(x_k)\| \leq \varepsilon$

after at most  $K = O(\varepsilon^{-2})$  iterations

Equivalently, after  $K \geq 1$  iterations, we have

$$\min_{0 \leq k \leq K-1} \|\nabla f(x_k)\| \leq O\left(\frac{1}{\sqrt{K}}\right)$$

"Convergence rate"

NB: 
$$f(x) = \frac{1}{2} \sum_{i=1}^N f_i(x)^2 = \frac{1}{2} \|r(x)\|^2$$

$$\|\nabla f(x)\| \leq \varepsilon \Rightarrow \begin{cases} \|\nabla f(x)\| < \varepsilon \|r(x)\| \\ \text{or } \|r(x)\| < \varepsilon \end{cases}$$

where  $r(x) = \begin{bmatrix} f_1(x) \\ \vdots \\ f_N(x) \end{bmatrix} \in \mathbb{R}^N$

### Line search

→ Want  $\alpha > 0$  such that  $f(x_k - \alpha \nabla f(x_k)) < f(x_k) - c\alpha \|\nabla f(x_k)\|^2$

→ Try  $\alpha = 1$ , then  $\alpha = 0$  if  $\alpha = 1$  fails, ...

(1)  $\alpha = 1$  satisfies the condition (✓)

(2) If  $\alpha = 1$  does not satisfy the condition, then

$$(1) \quad f(x_k) - c\alpha \|\nabla f(x_k)\|^2 \leq f(x_k - \alpha \nabla f(x_k)) \quad \text{for some } \alpha \in \{1, 0, 0^2, \dots\}$$

Since  $f \in C^2$ ,

$$(2) \quad f(x_k - \alpha \nabla f(x_k)) \leq f(x_k) - \alpha \|\nabla f(x_k)\|^2 + \frac{L\alpha^2}{2} \|\nabla f(x_k)\|^2$$

Putting (1) and (2) together,

$$\cancel{f(x_k)} - c\alpha \|\nabla f(x_k)\|^2 \leq \cancel{f(x_k)} - \alpha \|\nabla f(x_k)\|^2 + \frac{L\alpha^2}{2} \|\nabla f(x_k)\|^2$$

$$(1-c) \alpha \|\nabla f(x_k)\|^2 \leq \frac{L}{2} \alpha^2 \|\nabla f(x_k)\|^2 \quad \left( \|\nabla f(x_k)\| \neq 0 \right)$$

$$\boxed{\alpha \geq \frac{2(1-c)}{L}}$$

Thus, if  $\alpha$  doesn't satisfy the condition, then  $\alpha \geq \frac{2(1-c)}{L}$

$$\alpha = \theta^j \text{ for } j \in \mathbb{N} \Rightarrow \theta^j \geq \frac{2(1-c)}{L}$$

$$\Leftrightarrow j \leq \log_{\theta} \left( \frac{2(1-c)}{L} \right)$$

As a result, the line search must terminate after  $j_h \leq 1 + \left\lfloor \log_{\theta} \left( \frac{2(1-c)}{L} \right) \right\rfloor$  tries for  $\alpha$  and

$$\alpha_k \text{ is chosen as } \theta^{j_h} \text{ with } \theta^{j_h} \geq \frac{2\theta(1-c)}{L}$$

*Iteration complexity*

At every iteration,

$$\begin{aligned} f(x_{k+1}) &< f(x_k) - c \alpha_k \|\nabla f(x_k)\|^2 \\ &< f(x_k) - \frac{2\theta c(1-c)}{L} \varepsilon^2 \end{aligned}$$

---

In practice: If  $\alpha_{k-1} = \theta^{j_{k-1}}$ , then at iteration  $k$  start line search with  $\frac{\theta^{j_{k-1}}}{\theta} = \theta^{j_{k-1}-1}$

$\hookrightarrow$  Other conditions:
 

- $\alpha_k = \underset{\alpha > 0}{\operatorname{argmin}} f(x_k - \alpha \nabla f(x_k))$   
"Exact line search"
- Wolfe conditions  
 $f(x_k - \alpha \nabla f(x_k)) < f(x_k) - c_1 \alpha \|\nabla f(x_k)\|^2$   
 $-\nabla f(x_k - \alpha \nabla f(x_k))^T \nabla f(x_k) \leq -c_2 \|\nabla f(x_k)\|^2$

Give  $O(\varepsilon^{-2})$  iteration complexity

---

Bad function for gradient descent (simplified version,  
 Curtis Gold Toink 2022  
 "Evaluation complexity of nonconvex optimization methods")

Fixed stepsize:  $\alpha_k = \alpha > 0$

Goal: Design a function so that gradient descent

$n = 1$

$$x_{k+1} = x_k - \alpha_k g_k, \quad g_k = \nabla f(x_k)$$

takes exactly  $\varepsilon^{-2}$  iterations (assuming  $\frac{1}{\varepsilon}$  integer)

Define: •  $k_\varepsilon = \varepsilon^{-2}$

$$\bullet \forall k \in [0, k_\varepsilon], \quad g_k = -\varepsilon \left( 2 - \frac{k}{k_\varepsilon} \right)$$

$$\bullet g_0 = -2\varepsilon \leq g_1 < g_2 < \dots < g_{k_\varepsilon} = -\varepsilon$$

The first iteration for which  $|g_k| \leq \varepsilon$  is  $k_\varepsilon$

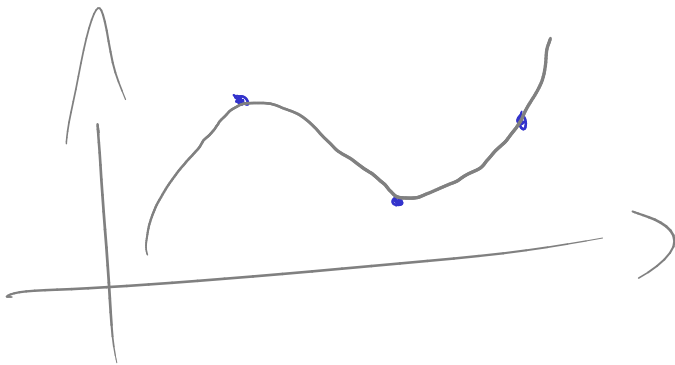
$$\bullet x_0 = 0, \quad \text{then } x_{k+1} = x_k - \alpha g_k \quad \forall k \in [0, k_\varepsilon - 1]$$

•  $f_0 = 4\alpha > 0$ , then  $f_{k+1} = f_k - \alpha / |g_k|^2$   
 $\forall k \in \{0, k_\varepsilon - 1\}$

$\Rightarrow$  Guarantees that line search would accept  $\alpha$  as a stepsize

Next step: Define a function  $f \in C^1$  such that  
 $f(x_k) = f_k \forall k \in \{0, k_\varepsilon\}$  and  $\nabla f(x_k) = g_k$   
 $\forall k \in \{0, k_\varepsilon\}$

$\Rightarrow$  Hermite interpolation guarantees that such a function exists



Rosenbrock function

$\forall x \in \mathbb{R}^2, f(x) = 100(x_2 - x_1^2)^2 + (1 - x_1)^2$   
 unique minimum  $\begin{bmatrix} 1 \\ 1 \end{bmatrix}$

Newton: 18 iterations,  $\|\nabla f(x_k)\| \approx 10^{-7}$

Gradient descent: 100 iterations,  $\|\nabla f(x_k)\| \approx 0.8$



Newton: Uses the second-order derivative of the function

$$x_{k+1} = x_k - \alpha_k \left[ \nabla^2 f(x_k) \right]^{-1} \nabla f(x_k)$$

$\uparrow$   
 matrix of derivatives  $\mathbb{R}^{n \times n}$

$\alpha_k = 1$  works most of the time  
 Typical: line search starting from 1

---

### Bad example for Newton

→ In  $\mathbb{R}^2$

→ Newton iteration well defined

$$x_0 = \begin{bmatrix} 0 \\ 0 \end{bmatrix}, x_{k+1} = x_k - \underbrace{H_k^{-1}}_{\begin{bmatrix} 1/f_k \\ \varepsilon/f_k \end{bmatrix}} g_k$$

$$f_k = 2 - k\varepsilon^2$$

$$g_k = \begin{pmatrix} -\varepsilon^2 f_k \\ -\varepsilon f_k \end{pmatrix}$$

$$H_k = \begin{bmatrix} \varepsilon^2 f_k^2 & 0 \\ 0 & f_k^2 \end{bmatrix}$$

$$\begin{aligned} \|g_k\| &= \sqrt{\varepsilon^4 f_k^2 + \varepsilon^2 f_k^2} = \sqrt{(\varepsilon^4 + \varepsilon^2) (2 - k\varepsilon^2)^2} \\ &= \varepsilon \underbrace{\sqrt{(\varepsilon^2 + 1)}}_{> 1} \underbrace{(2 - k\varepsilon^2)}_{\geq 1 \text{ when } k \leq \varepsilon^{-2}} > \varepsilon \text{ for } k \leq \varepsilon^{-2} \end{aligned}$$

The function

$$\sum_{k=0}^{k_\varepsilon - 1} \sigma(\|x - x_k\|) \left( f_k + g_k^T (x - x_k) + \frac{1}{2} (x - x_k)^T H_k (x - x_k) \right)$$

or Hermite interpolation