

# The Gradient Method (GM)

The GM is one of the most common algorithms used in optimization.

Starting from an initialization  $x^0 \in \mathbb{R}^n$ , the GM iteratively generates the sequence

$$x^t = x^{t-1} - \gamma_t (\nabla f)(x^{t-1}), \quad t = 1, 2, 3, \dots$$

where  $\gamma_t > 0$  is the step size. The algorithm requires the user to select an appropriate step size.

Step-size Selection:

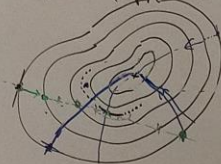
★ • Fixed step:  $\gamma_t \equiv \gamma > 0$ .  $\gamma$  small to ensure convergence  
large enough for GM to be fast

• Decreasing step:  $\gamma_t = \frac{C}{t}$  or  $\frac{C}{\sqrt{t}}$ ,  $C > 0$  constant

★ • Exact Line Search: Selecting  $\gamma_t$  via

$$\gamma_t = \arg \min_{\gamma \geq 0} \phi(\gamma) := f(x^{t-1} - \gamma (\nabla f)(x^{t-1})).$$

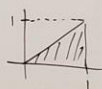
$$\dot{x}(t) = -(\nabla f)(x(t))$$



Recall: Theorem: Let  $f \in C_1^1(\mathbb{R}^n)$  be the cost that has finite minimum at  $x^* \in \mathbb{R}^n$ . Then for any step size  $\gamma \in (0, \frac{2}{L}]$ , the iterates generated by the gradient method satisfy  $\lim_{t \rightarrow \infty} \|(\nabla f)(x^t)\| = 0$ .

Lemma: For any  $f \in C_1^1(\mathbb{R}^n)$ , we have that  $f(y) \leq f(x) + (\nabla f)(x)^T (y-x) + \frac{L}{2} \|y-x\|^2$ ,  $\forall x, y \in \mathbb{R}^n$ .

Proof: Define  $\varphi(t) := f(x+tv)$  where  $v := y-x$ . Fundamental thm of calculus:  $\int_0^1 \varphi'(t) dt = \varphi(1) - \varphi(0)$   $\left| -\varphi'(0) \right.$



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

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

$$f(y) = f(x) + (\nabla f)(x)^T v + \int_0^1 ((\nabla f)(x+tv) - (\nabla f)(x))^T v dt$$

$$\leq f(x) + (\nabla f)(x)^T v + \int_0^1 \underbrace{\|(\nabla f)(x+tv) - (\nabla f)(x)\|}_{\leq L\|tv-x\|} \|v\| dt \quad \text{Cauchy Schwarz}$$

$$\leq f(x) + (\nabla f)(x)^T v + \int_0^1 L\|tv-x\| \|v\| dt = L\|tv\| = L\|v\|$$

$$= f(x) + (\nabla f)(x)^T v + \frac{L}{2} \|v\|^2 \quad \leftarrow \text{recall } v = y-x$$

In the proof of the theorem, we managed to show:

$$\frac{\gamma}{2} \sum_{i=1}^t \|(\nabla f)(x^{i-1})\|^2 \leq f(x^0) - f(x^1) + \dots + f(x^{t-1}) - f(x^t) = \dots = f(x^0) - f(x^t) \leftarrow \text{finite}$$

We concluded that  $\lim_{i \rightarrow \infty} \|(\nabla f)(x^i)\| = 0$ ,

which proves the theorem. We can, in addition, say a little more. Define

$$g_t = \min_{1 \leq i \leq t} \|(\nabla f)(x^{i-1})\|^2.$$

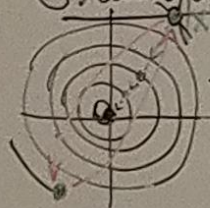
Then  $\frac{\gamma}{2} t g_t \leq \frac{\gamma}{2} \sum_{i=1}^t \|(\nabla f)(x^{i-1})\|^2 \leq f(x^0) - f(x^t)$

$$\Rightarrow g_t \leq \frac{2}{\gamma t} (f(x^0) - f(x^t)) \Rightarrow \min_{1 \leq i \leq t} \|(\nabla f)(x^{i-1})\|^2 \leq \sqrt{\frac{2}{\gamma t} (f(x^0) - f(x^t))} \leq \epsilon, \text{ e.g. } 10^{-6}$$

Thus, GM finds an  $\epsilon$ -approximate stationary point in  $O(\frac{1}{\epsilon^2})$  iterations.  $(\nabla f)(x) \leq \epsilon \dots \epsilon$  small

$$\Rightarrow \frac{2}{\gamma t} \text{const.} \leq \epsilon^2 \Rightarrow t \geq \frac{2 \text{const.}}{\gamma} \epsilon^2.$$

Examples: Fixed step  $\|(\nabla f)(x) - (\nabla f)(y)\| \leq L \|x - y\|$  ( $L=2$ )



$$f(x_1, x_2) = x_1^2 + x_2^2, \quad (\nabla f)(x) = \begin{pmatrix} 2x_1 \\ 2x_2 \end{pmatrix} = 2x, \quad x^0 = \begin{pmatrix} 3 \\ 4 \end{pmatrix}$$

$$x^1 = x^0 - \gamma (\nabla f)(x^0)$$

(i)  $\gamma = 1$ :

bad  
oscillating!

$$= \begin{pmatrix} 3 \\ 4 \end{pmatrix} - 1 \times \begin{pmatrix} 6 \\ 8 \end{pmatrix} = \begin{pmatrix} -3 \\ -4 \end{pmatrix}$$

(ii)  $\gamma = \frac{1}{2}$ :

works  
(maybe too well)

$$x^1 = \begin{pmatrix} 3 \\ 4 \end{pmatrix} - \frac{1}{2} \times \begin{pmatrix} 6 \\ 8 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

(iii)  $\gamma = \frac{1}{4}$ :

$$x^1 = \begin{pmatrix} 3 \\ 4 \end{pmatrix} - \frac{1}{4} \times \begin{pmatrix} 6 \\ 8 \end{pmatrix} = \begin{pmatrix} 3/2 \\ 2 \end{pmatrix}$$

$$x^2 = \begin{pmatrix} 3/2 \\ 2 \end{pmatrix} - \frac{1}{4} \times \begin{pmatrix} 3 \\ 4 \end{pmatrix} = \begin{pmatrix} 3/4 \\ 1 \end{pmatrix}$$

$$x^3 = \begin{pmatrix} 3/4 \\ 1 \end{pmatrix} - \frac{1}{4} \times \begin{pmatrix} 3/2 \\ 2 \end{pmatrix} = \begin{pmatrix} 3/8 \\ 1/2 \end{pmatrix}$$

Exact Line Search:

quadratic cost with  
ellipses as level sets

