

## TD - Gradient descent : solution

Let  $f : \mathbb{R}^n \rightarrow \mathbb{R}$  be a differentiable function whose gradient is  $L$ -Lipschitz continuous i.e.  $\|\nabla f(y) - \nabla f(x)\| \leq L\|y - x\|$  for all  $x, y$ .

1. Prove that for all  $x, y$ ,  $\langle \nabla f(y) - \nabla f(x), y - x \rangle \leq L\|y - x\|^2$ .

► We are using Cauchy-Schwartz inequality :

$$\begin{aligned} \langle \nabla f(y) - \nabla f(x), y - x \rangle &\leq |\langle \nabla f(y) - \nabla f(x), y - x \rangle| \leq \|\nabla f(y) - \nabla f(x)\| \cdot \|y - x\| \\ &\leq L \|y - x\|^2 \end{aligned}$$

2. Set  $\varphi(t) = f(x + t(y - x))$  for all  $t \in [0, 1]$ . Prove that

$$f(y) - f(x) - \langle \nabla f(x), y - x \rangle = \varphi(1) - \varphi(0) - \varphi'(0).$$

► It is clear that  $\varphi(0) = f(x)$  and  $\varphi(1) = f(y)$ . Note that  $\varphi(t) = f(g(t))$  where  $g(t) = x + t(y - x)$ . By the theorem of derivation of composite functions,

$$\varphi'(t) = \langle \nabla f(g(t)), g'(t) \rangle = \langle \nabla f(x + t(y - x)), y - x \rangle \quad (1)$$

So  $\varphi'(0) = \langle \nabla f(x), y - x \rangle$ . Combining the three equalities yields

$$\varphi(1) - \varphi(0) - \varphi'(0) = f(y) - f(x) - \langle \nabla f(x), y - x \rangle.$$

3. Deduce that

$$f(y) - f(x) - \langle \nabla f(x), y - x \rangle = \int_0^1 \langle \nabla f(x + t(y - x)) - \nabla f(x), y - x \rangle dt$$

► As  $\varphi$  is a primitive of  $\varphi'$ ,

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

Hence, using (1)

$$\begin{aligned} f(y) - f(x) - \langle \nabla f(x), y - x \rangle &= \varphi(1) - \varphi(0) - \varphi'(0) = \int_0^1 \varphi'(t) dt - \varphi'(0) \\ &= \int_0^1 \langle \nabla f(x + t(y - x)), y - x \rangle dt - \int_0^1 \langle \nabla f(x), y - x \rangle dt \\ &= \int_0^1 \langle \nabla f(x + t(y - x)) - \nabla f(x), y - x \rangle dt \end{aligned}$$

4. Using the first question, conclude that

$$f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{L}{2} \|y - x\|^2.$$

► We know that for all  $x, z$ ,  $\langle \nabla f(z) - \nabla f(x), z - x \rangle \leq L\|z - x\|^2$ . We use this inequality with  $z = x + t(y - x)$  :  $\langle \nabla f(x + t(y - x)) - \nabla f(x), t(y - x) \rangle \leq L\|t(y - x)\|^2$ . Dividing by  $t$  and integrating between 0 and 1, we get

$$\int_0^1 \langle \nabla f(x + t(y - x)) - \nabla f(x), y - x \rangle dt \leq \int_0^1 tL \|y - x\|^2 dt = \left[ \frac{t^2}{2} \right]_0^1 L \|y - x\|^2 = \frac{L}{2} \|y - x\|^2$$

We conclude using Question 3.

Consider the gradient algorithm i.e., the sequence  $(x_k)$  defined by  $x_{k+1} = x_k - \gamma \nabla f(x_k)$  where  $\gamma > 0$  is a constant step size.

5. Show that

$$f(x_{k+1}) \leq f(x_k) - \gamma(1 - \frac{\gamma L}{2}) \|\nabla f(x_k)\|^2. \quad (2)$$

► We use the result of Question 4 :

$$f(x_{k+1}) \leq f(x_k) + \langle \nabla f(x_k), x_{k+1} - x_k \rangle + \frac{L}{2} \|x_{k+1} - x_k\|^2$$

But  $x_{k+1} = x_k - \gamma \nabla f(x_k)$  so

$$f(x_{k+1}) \leq f(x_k) - \gamma \langle \nabla f(x_k), \nabla f(x_k) \rangle + \frac{L}{2} \|\gamma \nabla f(x_k)\|^2 = f(x_k) - \gamma(1 - \frac{\gamma L}{2}) \|\nabla f(x_k)\|^2$$

6. Provide a condition on  $\gamma$  which ensures that  $f(x_{k+1}) \leq f(x_k)$ .

► If  $\gamma(1 - \frac{\gamma L}{2}) \|\nabla f(x_k)\|^2 \geq 0$ , then  $f(x_{k+1}) \leq f(x_k) - \gamma(1 - \frac{\gamma L}{2}) \|\nabla f(x_k)\|^2 \leq f(x_k)$ .  
Now  $\gamma(1 - \frac{\gamma L}{2}) \|\nabla f(x_k)\|^2 \geq 0 \Leftrightarrow \gamma(1 - \frac{\gamma L}{2}) \geq 0 \Leftrightarrow \gamma \leq \frac{2}{L}$  since  $\gamma > 0$ .

7. Based on Eq. (2), what value of  $\gamma$  would you suggest to choose ?

► There are two possible answers for this question.

- First answer : we may want to choose  $\gamma$  the largest possible that ensures a strict decrease of the objective value. Hence, we will take  $\gamma = \frac{2}{L} - \epsilon$  where  $\epsilon > 0$  is small. Taking long stepsizes helps obtaining an algorithm that goes faster to the minimum.
- Second answer : we may want to choose  $\gamma$  such that the bound we found is the smallest possible. Hence, we will take  $\gamma = \frac{1}{L}$ . Minimizing the bound helps guaranteeing a better decrease in the objective function.

For the rest of the exercise, we will take  $\gamma = \frac{1}{L}$  because the expressions are easier to manipulate.

From now on, we set  $\gamma$  equal to the value suggested above. Consider an arbitrary  $y \in \mathbb{R}^n$ .

8. Prove that

$$\langle \nabla f(x_k), y - x_k \rangle + \frac{L}{2} \|x_k - y\|^2 - \frac{L}{2} \|x_{k+1} - y\|^2 = -\frac{1}{2L} \|\nabla f(x_k)\|^2.$$

$$\begin{aligned} \blacktriangleright \quad \frac{L}{2} \|x_{k+1} - y\|^2 &= \frac{L}{2} \|x_k - \gamma \nabla f(x_k) - y\|^2 \\ &= \frac{L}{2} \|x_k - y\|^2 - L\gamma \langle x_k - y, \nabla f(x_k) \rangle + \frac{L}{2} \gamma^2 \|\nabla f(x_k)\|^2 \end{aligned}$$

Rearranging and simplifying  $\gamma L = 1$ , we get the expected result.

9. Deduce that

$$f(x_{k+1}) \leq f(x_k) + \langle \nabla f(x_k), y - x_k \rangle + \frac{L}{2} \|x_k - y\|^2 - \frac{L}{2} \|x_{k+1} - y\|^2. \quad (3)$$

$$\begin{aligned} \blacktriangleright \quad f(x_{k+1}) &\leq f(x_k) - \gamma(1 - \frac{\gamma L}{2}) \|\nabla f(x_k)\|^2 = f(x_k) - \frac{1}{2L} \|\nabla f(x_k)\|^2 \\ &= f(x_k) + \langle \nabla f(x_k), y - x_k \rangle + \frac{L}{2} \|x_k - y\|^2 - \frac{L}{2} \|x_{k+1} - y\|^2 \end{aligned}$$

We assume from now on that  $f$  is convex and admits (at least) one minimizer  $x^*$ .

10. Show that

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

$\blacktriangleright$  As  $f$  is convex,  $f(x^*) \geq f(x_k) + \langle \nabla f(x_k), x^* - x_k \rangle$ . Combining this inequality with (3) applied with  $y = x^*$ , we get :

$$\begin{aligned} f(x_{k+1}) &\leq f(x_k) + \langle \nabla f(x_k), x^* - x_k \rangle + \frac{L}{2} \|x_k - x^*\|^2 - \frac{L}{2} \|x_{k+1} - x^*\|^2 \\ &\leq f(x^*) + \frac{L}{2} \|x_k - x^*\|^2 - \frac{L}{2} \|x_{k+1} - x^*\|^2 \end{aligned}$$

11. Deduce that for all  $k \geq 1$ ,

$$\sum_{i=1}^k f(x_i) \leq kf(x^*) + \frac{L}{2} \|x_0 - x^*\|^2.$$

$\blacktriangleright$  We sum the inequality in Question 10 for  $i$  between 0 and  $k-1$  :

$$\sum_{i=0}^{k-1} f(x_{i+1}) \leq kf(x^*) + \sum_{i=0}^{k-1} \left( \frac{L}{2} \|x_i - x^*\|^2 - \frac{L}{2} \|x_{i+1} - x^*\|^2 \right)$$

We recognise a telescoping sum. We also make a change of variable  $i \leftarrow i+1$  in the first sum.

$$\sum_{i=1}^k f(x_i) \leq kf(x^*) + \frac{L}{2} \|x_0 - x^*\|^2 - \frac{L}{2} \|x_k - x^*\|^2 \leq kf(x^*) + \frac{L}{2} \|x_0 - x^*\|^2$$

12. Show that

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

► Recall that  $f(x_k) \leq f(x_{k-1})$  for all  $k$ . So

$$f(x_k) - f(x^*) \leq \sum_{i=1}^k \frac{1}{k} (f(x_i) - f(x^*)) \leq \frac{L \|x_0 - x^*\|^2}{2k}$$

We assume from now on that  $f$  is  $\mu$ -strongly convex. Recall that a function  $f$  is said  $\mu$ -strongly convex if  $f - \frac{\mu}{2} \|\cdot\|^2$  is convex.

13. Prove that for any  $x, y$ ,

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2} \|y - x\|^2.$$

► The function  $g = f - \frac{\mu}{2} \|\cdot\|^2$  is convex so for all  $x, y$ ,

$$\begin{aligned} g(y) &\geq g(x) + \langle \nabla g(x), y - x \rangle \\ f(y) - \frac{\mu}{2} \|y\|^2 &\geq f(x) - \frac{\mu}{2} \|x\|^2 + \langle \nabla f(x) - \mu x, y - x \rangle \end{aligned}$$

We get the result since  $-\|x\|^2 - 2\langle x, y - x \rangle + \|y\|^2 = \|x\|^2 - 2\langle x, y \rangle + \|y\|^2 = \|y - x\|^2$ .

14. Using Eq. (3), prove that

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

► As  $f$  is  $\mu$ -strongly convex,  $f(x^*) \geq f(x_k) + \langle \nabla f(x_k), x^* - x_k \rangle + \frac{\mu}{2} \|x^* - x_k\|^2$ . Combining this inequality with (3) applied with  $y = x^*$ , we get :

$$\begin{aligned} f(x_{k+1}) &\leq f(x_k) + \langle \nabla f(x_k), x^* - x_k \rangle + \frac{L}{2} \|x_k - x^*\|^2 - \frac{L}{2} \|x_{k+1} - x^*\|^2 \\ &\leq f(x^*) + \frac{L - \mu}{2} \|x_k - x^*\|^2 - \frac{L}{2} \|x_{k+1} - x^*\|^2 \end{aligned}$$

15. Define  $\Delta_{k+1} = f(x_{k+1}) - f(x^*) + \frac{L}{2} \|x_{k+1} - x^*\|^2$ . Show that

$$\Delta_{k+1} \leq \left(1 - \frac{\mu}{L}\right) \Delta_k.$$

► Using Question 14 and the fact that  $f(x_k) \geq f(x^*)$ , we get

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

16. Conclude that

$$\begin{aligned} f(x_k) - f(x^*) &\leq \left(1 - \frac{\mu}{L}\right)^k \Delta_0 \\ \|x_k - x^*\|^2 &\leq \left(1 - \frac{\mu}{L}\right)^k \frac{2\Delta_0}{L}. \end{aligned}$$

► Iterating the relation  $\Delta_k \leq \left(1 - \frac{\mu}{L}\right)\Delta_{k-1}$ , we get that  $\Delta_k \leq \left(1 - \frac{\mu}{L}\right)^k \Delta_0$ . As  $\Delta_k$  is the sum of two nonnegative quantities, this means that  $f(x_k) - f(x^*) \leq \left(1 - \frac{\mu}{L}\right)^k \Delta_0$  and  $\frac{2}{L} \|x_k - x^*\|^2 \leq \left(1 - \frac{\mu}{L}\right)^k \Delta_0$ .

17. The ratio  $Q = L/\mu$  is called the condition number of  $f$ . Discuss the influence of  $Q$  on the convergence rate.

► The smaller the condition number the faster the convergence. If the problem is ill-conditioned, there will be a slow convergence.

*Application :* From now on, we define  $f(x) = \frac{1}{2}x^T Hx + c^T x$  where  $H$  is positive semidefinite  $n \times n$  matrix. We denote by  $\lambda_{\max}$  and  $\lambda_{\min}$  the largest and smallest eigenvalues of  $H$  respectively.

18. What is the Hessian matrix of  $f$ ? Deduce that  $f$  is convex.

► The Hessian matrix at  $x$  of  $f$  is  $H$  for all  $x$ . As it is positive semi-definite,  $f$  is convex.

19. Justify briefly that  $\nabla f$  is  $\lambda_{\max}$ -Lipschitz continuous.

►  $\|\nabla f(x) - \nabla f(y)\| = \|Hx - Hy\| \leq \|H\| \|x - y\| = \lambda_{\max} \|x - y\|$ .

20. Prove that  $f$  is  $\lambda_{\min}$ -strongly convex.

► The Hessian matrix of  $f - \frac{\lambda_{\min}}{2} \|\cdot\|^2$  is equal to  $H - \lambda_{\min}I$ , which is positive semi-definite by definition of  $\lambda_{\min}$ . Note that if  $\lambda_{\min} = 0$ , then  $f$  is not strongly convex.

21. Write the condition number  $Q$  of  $f$ . What kind of matrix  $H$  yields the smallest condition number?

►  $Q = \lambda_{\max}/\lambda_{\min}$ . The smallest condition number is for  $H = \alpha I$ ,  $\alpha > 0$ . In this case  $Q = 1$ .

22. Characterize the set of minimizers of  $f$ .

► As  $f$  is convex and differentiable,  $x$  is a minimizer of  $f$  if and only if  $\nabla f(x) = 0$ . This means  $Hx = -c$  and so we have three cases.

- If  $H$  is invertible, then  $x = H^{-1}c$  is the unique solution.
- If  $H$  is not invertible and  $-c \in \text{Range}(H)$ , then the set of solution is  $H^{-1}(\{-c\})$ , the pre-image of  $\{-c\}$  by  $H$ . It is a linear subspace of  $\mathbb{R}^m$ .
- If  $H$  is not invertible and  $-c \notin \text{Range}(H)$ , then there is no solution.