# Exercise List: Properties and examples of convexity and smoothness

Robert M. Gower and Francis Bach.

## February 7, 2019

Time to get familiarized with convexity, smoothness and a bit of strong convexity.

**Notation:** For every  $x, y \in \mathbb{R}^d$  let  $\langle x, y \rangle \stackrel{\text{def}}{=} x^\top y$  and let  $||x||_2 = \sqrt{\langle x, x \rangle}$ . Let  $\sigma_{\min}(A)$  and  $\sigma_{\max}(A)$  be the smallest and largest singular values of A defined by

$$\sigma_{\min}(A) \stackrel{\text{def}}{=} \min_{x \in \mathbb{R}^d} \frac{\|Ax\|_2}{\|x\|_2} \quad \text{and} \quad \sigma_{\max}(A) \stackrel{\text{def}}{=} \max_{x \in \mathbb{R}^d} \frac{\|Ax\|_2}{\|x\|_2}.$$
 (1)

Thus clearly

$$\frac{\|Ax\|_2^2}{\|x\|_2^2} \le \sigma_{\max}(A)^2, \quad \forall x \in \mathbb{R}^d.$$

$$\tag{2}$$

Let  $||A||_F^2 \stackrel{\text{def}}{=} \text{Tr}(A^{\top}A)$  denote the Frobenius norm of A. Finally, a result you will need, for every symmetric matrix G the L2 induced matrix norm can be equivalently defined by

$$||G||_{2} = \sigma_{\max}(G) = \sup_{x \in \mathbb{R}^{d}, x \neq 0} \frac{|\langle Gx, x \rangle|}{\|x\|_{2}^{2}} = \max_{x \in \mathbb{R}^{d}, x \neq 0} \frac{\|Gx\|_{2}}{\|x\|_{2}}.$$
(3)

# 1 Convexity

We say that a twice differentiable function  $f : \mathbb{R}^d \to \mathbb{R}$  is convex if

$$f(\lambda x + (1 - \lambda)y) \le \lambda f(x) + (1 - \lambda)f(y), \quad \forall x, y \in \mathbb{R}^d, \lambda \in [0, 1].$$
(4)

or equivalently

$$v^{\top} \nabla^2 f(x) v \ge 0, \quad \forall x, v \in \mathbb{R}^d.$$
 (5)

We say that f is  $\mu$ -strongly convex if

$$v^{\top} \nabla^2 f(x) v \ge \mu \|v\|_2^2, \quad \forall x, v \in \mathbb{R}^d.$$
(6)

**Ex. 1** — We say that  $\|\cdot\| \to \mathbb{R}_+$  is a norm over  $\mathbb{R}^d$  if it satisfies the following three properties

- 1. **Point separating:**  $||x|| = 0 \Leftrightarrow x = 0, \forall x \in \mathbb{R}^d$ .
- 2. **Subadditive:**  $||x + y|| \le ||x|| + ||y||, \forall x, y \in \mathbb{R}^d$
- 3. Homogeneous:  $||ax|| = |a|||x||, \forall x \in \mathbb{R}^d, a \in \mathbb{R}.$

#### Part I

Prove that  $x \mapsto ||x||$  is a convex function.

## Part II

For every convex function  $f: y \in \mathbb{R}^m \mapsto f(y)$ , prove that  $g: x \in \mathbb{R}^d \mapsto f(Ax - b)$  is a convex function, where  $A \in \mathbb{R}^{n \times d}$  and  $b \in \mathbb{R}^n$ .

## Part III

Let  $f_i : \mathbb{R}^d \to \mathbb{R}$  be convex for i = 1, ..., n. Prove that  $\sum_{i=1}^n f_i$  is convex.

## Part IV

For given scalars  $y_i \in \mathbb{R}$  and vectors  $a_i \in \mathbb{R}^d$  for i = 1, ..., m prove that the *logistic* regression function  $f(x) = \frac{1}{n} \sum_{i=1}^n \ln(1 + e^{-y_i(x, a_i)})$  is convex.

## Part V

Let  $A \in \mathbb{R}^{n \times d}$  have full column rank. Prove that  $f(x) = \frac{1}{2} ||Ax - b||_2^2$  is  $\sigma_{\min}^2(A)$ -strongly convex.

## Part VI

Now suppose that the function f(x) is  $\mu$ -strongly convex, that is, it satisfies

$$f(y) \ge f(x) + \langle \nabla f(x), y - x \rangle + \frac{\mu}{2} \|y - x\|_2^2, \quad \forall x, y \in \mathbb{R}^d.$$

$$\tag{7}$$

Prove that f(x) satisfies the Polyak-Lojasiewicz condition, that is

$$\|\nabla f(x)\|_2^2 \ge 2\mu(f(x) - f(x^*)), \quad \forall x.$$
 (8)

**Answer (Ex. I)** — Let  $x, y \in \mathbb{R}^d$  and  $\lambda \in [0, 1]$ . It follows that

$$\begin{aligned} \|\lambda x + (1-\lambda)y\| & \stackrel{\text{item } 2}{\leq} & \|\lambda x\| + \|(1-\lambda)y\| \\ & \text{item } 3 \\ & \stackrel{1}{\leq} & \lambda\|x\| + (1-\lambda)\|y\|. \end{aligned}$$

**Answer (Ex. II)** — Let  $x, y \in \mathbb{R}^d$  and  $\lambda \in [0, 1]$ . It follows that

$$g(\lambda x + (1 - \lambda)y) = f(A(\lambda x + (1 - \lambda))y - b)$$
  
=  $f(\lambda(Ax - b) + (1 - \lambda)(Ay - b))$  (9)  
$$f \text{ is conv.}$$
  
=  $\lambda f(Ax - b) + (1 - \lambda)f(Ay - b).$ 

Answer (Ex. III) — Immediate through either definition.

**Answer (Ex. IV)** — From exercise V we need only prove that  $f(x) = \ln(1 + e^{-y\langle x, w \rangle})$ is convex for a given  $y \in \mathbb{R}$  and  $w \in \mathbb{R}^d$ . From exercise II we need only prove that  $\phi(\alpha) = \ln(1 + e^{\alpha})$  is convex, since  $x \mapsto -y \langle x, w \rangle$  is a linear function. The convexity of  $f(\alpha)$  now follows by differentiating once

$$\phi'(\alpha) = \frac{e^{\alpha}}{1 + e^{\alpha}},$$

then differentiating again

$$\phi''(\alpha) = \frac{e^{\alpha}}{1+e^{\alpha}} - \frac{e^{2\alpha}}{(1+e^{\alpha})^2} = \frac{e^{\alpha}}{(1+e^{\alpha})^2} \ge 0, \quad \forall \alpha.$$
(10)

We can now call upon the definition (5), but since  $\alpha \in \mathbb{R}$  is a scalar, the above already proves that  $\phi(\alpha)$  is convex.

Answer (Ex. V) — Differentiating twice we have that

$$\nabla^2 f(x) = A^\top A.$$

Consequently

$$v^{\top} \nabla^2 f(x) v = v^{\top} A^{\top} A v = \|Av\|_2^2 \ge \sigma_{\min}(A)^2 \|v\|_2^2$$

Answer (Ex. VI) — Multiplying (7) by minus and substituting  $y = x^*$  we have that

$$\begin{split} f(x) - f(x^*) &\leq \langle \nabla f(x), x - x^* \rangle - \frac{\mu}{2} \| x^* - x \|_2^2 \\ &= -\frac{1}{2} \| \sqrt{\mu} (x - x^*) - \frac{1}{\sqrt{\mu}} \nabla f(x) \|_2^2 + \frac{1}{2\mu} \| \nabla f(x) \|_2^2 \\ &\leq \frac{1}{2\mu} \| \nabla f(x) \|_2^2. \end{split}$$

# 2 Smoothness

We say that a function  $f : \mathbb{R}^d \to \mathbb{R}$  is *L*-smooth if

$$\|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|$$
(11)

or equivalently if f is twice differentiable then

$$v^{\top} \nabla^2 f(x) v \le L \|v\|_2^2, \quad \forall x, v \in \mathbb{R}^d.$$
(12)

## Ex. 2 — Part I

Prove that  $x \mapsto \frac{1}{2} ||x||^2$  is 1-smooth.

## Part II

Let  $f : \mathbb{R}^d \to \mathbb{R}$  be twice differentiable and *L*-smooth. Show that

$$\sigma_{\max}(\nabla^2 f(x)) = \|\nabla^2 f(x)\|_2 \le L.$$

## Part III

For every twice differentiable *L*-smooth function  $f : y \in \mathbb{R}^n \mapsto f(y)$ , prove that  $g : x \in \mathbb{R}^d \mapsto f(Ax - b)$  is a smooth function, where  $A \in \mathbb{R}^{n \times d}$  and  $b \in \mathbb{R}^n$ . Find the smoothness constant of g.

## Part IV

Let  $f_i : \mathbb{R}^d \to \mathbb{R}$  be a twice differentiable and  $L_i$ -smooth for i = 1, ..., n. Prove that  $\frac{1}{n} \sum_{i=1}^{n} f_i$  is  $\sum_{i=1}^{n} \frac{L_i}{n}$ -smooth.

#### Part V

For given scalars  $y_i \in \mathbb{R}$  and vectors  $a_i \in \mathbb{R}^d$  for i = 1, ..., n prove that the *logistic regression* function  $f(x) = \frac{1}{n} \sum_{i=1}^n \ln(1 + e^{-y_i \langle x, a_i \rangle})$  is smooth. Find the smoothness constant!

## Part VI

Let  $A \in \mathbb{R}^{n \times d}$  be any matrix. Prove that  $||Ax - b||_2^2$  is  $\sigma_{\max}^2(A)$ -smooth.

#### Part VII

Let M > 0 be a positive constant. Let  $f(x) = \frac{1}{n} \sum_{i=1}^{n} \phi_i(a_i^{\top} x)$  where  $\phi_i : \mathbb{R} \to \mathbb{R}$  is a scalar function such that  $\phi''_i(t) \leq M$  for all  $t \in \mathbb{R}$ . Prove that f(x) is  $\frac{M}{n} \sigma_{\max}^2(A)$ -smooth.

With this result, can you find a better estimate of the smoothness constant of the logistic regression loss? *Hint 1: ...* 

Part VIII

Co-coercivity. Let f be L-smooth, show that

$$\langle \nabla f(y) - \nabla f(x), y - x \rangle \ge \frac{1}{L} \| \nabla f(x) - \nabla f(y) \|_2^2$$

*Hint:* Start by showing that  $f(y) - f(x) \le \langle \nabla f(y), y - x \rangle - \frac{1}{2L} \| \nabla f(x) - \nabla f(y) \|_2^2$ .

**Answer (Ex. I)** — Clearly  $\nabla^2 \frac{1}{2} ||x||^2 = I$  and thus follows from definition (11).

Answer (Ex. II) — Using that the induced norm for symmetric matrices is given by

$$\|\nabla^2 f(x)\|_2 = \sup_{v \neq 0} \frac{|v^\top \nabla^2 f(x)v|}{\|v\|_2^2} \stackrel{(12)}{\leq} \sup_{v \neq 0} \frac{L\|v\|_2^2}{\|v\|_2^2} = L.$$

**Answer (Ex. III)** — Differentiating g(x) once gives

$$\nabla g(x) = A^{\top} \nabla f(Ax - b).$$

First we prove the claim using the definition (11). Indeed note that

$$\begin{aligned} \|\nabla g(x) - \nabla g(y)\|_{2} &= \|A^{\top} (\nabla f(Ax - b) - \nabla f(Ay - b))\|_{2} \\ &\leq \|A^{\top}\|_{2} \|\nabla f(Ax - b) - \nabla f(Ay - b)\|_{2} \\ &\text{smooth. of } f \\ &\leq L \|A^{\top}\|_{2} \|Ax - b - (Ay - b)\|_{2} \\ &\leq L \|A^{\top}\|_{2} \|A\|_{2} \|x - y\|_{2}. \end{aligned}$$

This the smoothness parameter is given by  $L||A||_2^2$  where we used that  $||A^{\top}||_2 = ||A||_2$ . This completes the proof.

We can also prove the claim using (12). Differentiating again we have that

$$\nabla^2 g(x) = A^\top \nabla^2 f(Ax - b)A.$$

Consequently

$$\|\nabla^2 g(x)\|_2^2 \le \|A\|_2^2 \|\nabla^2 f(Ax - b)\|_2^2 \le L \|A\|_2^2.$$

We could further tighten this by considering the smoothness constant of f restricted to the set  $\{x \mid Ax - b\}$  which might be smaller then  $\mathbb{R}^d$ .

Answer (Ex. IV) — Clearly

$$\nabla^2(\frac{1}{n}\sum_{i=1}^n f_i(x)) = \frac{1}{n}\sum_{i=1}^n \nabla^2 f_i(x) \preceq \frac{1}{n}\sum_{i=1}^n L_i I.$$

You can also prove this using the definition (11) and applying repeatedly the subadditivity of the norm.

**Answer (Ex. V)** — First note that from (10) we can see that the function  $\phi(\alpha) = \ln(1 + e^{\alpha})$  is at least 1–smooth. Consequently from exercise II the function  $f_i(x) = \ln(1 + e^{-y_i \langle x, a_i \rangle})$  is  $y_i^2 ||a_i||_2^2$ –smooth. Finally from exercise III the logistic regression function is  $\sum_{i=1}^n \frac{y_i^2 ||a_i||_2^2}{n}$ –smooth.

But this is not the tightest smoothness constant, as we will see in the next two exercises!

Answer (Ex. VI) — Differentiating twice we have that

$$\nabla^2 f(x) = A^\top A.$$

Consequently

$$v^{\top} \nabla^2 f(x) v = v^{\top} A^{\top} A v \le \|Av\|_2^2 \le \sigma_{\max}(A)^2 \|v\|_2^2$$

**Answer (Ex. VII)** — By analysing directly the Hessian of  $f(x) = \frac{1}{n} \sum_{i=1}^{n} f_i(x)$  we see that

$$\nabla^2 f(x) = A^{\top} \Phi(x) A,$$

where  $\Phi(x) = \operatorname{diag}(\phi_1''(a_1^{\top}x), \dots, \phi_n''(a_n^{\top}x))$ , Consequently

$$\|\nabla^2 f(x)\|_2 = \frac{1}{n} \|A^{\top} \Phi(x)A\|_2 \le \frac{1}{n} \|A\|_2^2 \|\Phi(x)\|_2 \le M \|A\|_2^2 \stackrel{(1)}{=} \frac{M}{n} \sigma_{\max}(A)^2.$$

For the logistic function, note that  $\phi''(a_i^{\top}x) = \frac{e^{\alpha}}{(1+e^{\alpha})^2}$ , where  $\alpha = -y_i \langle a_i, x \rangle$ . Furthermore

$$\phi''(\alpha) = \frac{e^{\alpha}}{(1+e^{\alpha})^2} \le \frac{1}{4}, \quad \forall \alpha.$$
(13)

Consequently a better estimate of the smoothness constant is given by

$$L \le \frac{\sigma_{\max}(A)^2}{4n}.$$

This is a much tighter smoothness constant and the one that is used in practice.