

# OPTIMIZATION FOR MACHINE LEARNING

September 19, 2024

Today: Intro to optimization (2/2), with exercises

Coming up: Sep. 23 : Gradient descent (Gabriel Peyré)  
Sep 26 : LAB on gradient descent (C. Royer)

## Back where we left off

"Find a  $C^1$  function  $f$  such that  $f$  is not convex

$$\nabla f(\bar{x}) = 0 \Leftrightarrow \bar{x} \in \underset{x}{\text{argmin}} f(x) "$$

NB: Such functions are called **invex** functions  
(or, more rarely, pseudo-convex functions)

$$\{\text{invex functions}\} \supset \{\text{convex functions}\}$$

In 1 dimension

$$x \mapsto x^2 + 2 \sin(x)^2$$

$$\varphi: x \mapsto -e^{-x^2}$$

$$\varphi'(x) = 2x e^{-x^2} \\ = 0 \text{ iff } x=0$$

In 2 dimensions:

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \mapsto x_1^2 - e^{-x_2^2}$$

## One application of invex functions

Data:  $\{(a_i, y_i)\}_{i=1..m}$      $a_i \in \mathbb{R}^d$      $y_i \in \{0, 1\}$

One optimization problem for binary classification:

$$\begin{aligned} \min_{x \in \mathbb{R}^d} f(x) &= \frac{1}{2m} \|\sigma(Ax) - y\|^2 \\ &= \frac{1}{2m} \sum_{i=1}^m (\sigma(a_i^T x) - y_i)^2 \end{aligned}$$

$\ell_2$  loss, model  $a \mapsto \sigma(a^T x) = \frac{1}{1 + e^{-a^T x}}$   
(sigmoid function)

Regularization:

1)  $\ell_2$  regularization

$$\min_{x \in \mathbb{R}^d} f(x) + \underbrace{\frac{\lambda}{2} \|x\|^2}_{\text{convex}}$$

For sufficiently large  $\lambda$ , we expect  $f + \frac{\lambda}{2} \|\cdot\|^2$  to behave (almost) like a convex function

2) Inverse regularization

$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{2m} \|\sigma(Ax) - y\|^2$$

$$\begin{aligned} \min_{\substack{x \in \mathbb{R}^d \\ p \in \mathbb{R}^m}} \hat{f}(x, p) &= \frac{1}{2m} \|\sigma(Ax) - y + p\|^2 \\ &= \frac{1}{2m} \|\sigma(Ax) - y\|^2 + \frac{\lambda^2}{2m} \|p\|^2 + \frac{\lambda}{m} p^T (\sigma(Ax) - y) \end{aligned}$$

The function  $\hat{f}: \mathbb{R}^{d+m} \rightarrow \mathbb{R}$  is invex  $\forall \lambda > 0$

Final note about invex functions

If  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is  $C^1$  and invex, then there exists a  $C^1$  function  $\eta: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^d$  such that

$$\forall (x, y) \in (\mathbb{R}^d)^2, \quad f(y) \geq f(x) + \nabla f(x)^T \eta(x, y)$$

$\Rightarrow \eta(x, y) = y - x$ , recover  $C^1$  convex functions

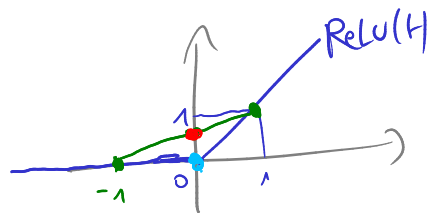
$$(f \text{ } C^1 \text{ convex} \Leftrightarrow f(y) \geq f(x) + \nabla f(x)^T (y-x) \quad \forall (x,y) \in \text{Dom})$$

## MORE ON CONVEX FUNCTIONS

↳ We defined convexity only for  $C^1$  functions using the gradient  
 $\Rightarrow$  The notion of a convex function does not involve derivatives

Def: A function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is convex if  
 $\forall (x,y) \in (\mathbb{R}^d)^2, \forall \alpha \in [0,1], f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y)$

Ex) ReLU:  $t \mapsto \max(t, 0)$



$$x = -1$$

$$y = 1 \quad \alpha = \frac{1}{2}$$

$$f(\alpha x + (1-\alpha)y) = f(0) = 0$$

$$\alpha f(x) + (1-\alpha)f(y) = \frac{1}{2} \times 0 + \frac{1}{2} \times 1 = \frac{1}{2}$$

↳ Checking convexity for functions with more than 2 variables is difficult!

$\Rightarrow$  We use a set of rules to certify convexity of most convex functions used in ML

Rules: I.  $f: \mathbb{R}^d \rightarrow \mathbb{R}^m$  is convex  
 $x \mapsto Wx + b$  where  $W \in \mathbb{R}^{m \times d}$  and  $b \in \mathbb{R}^m$

II. All norms in  $\mathbb{R}^d$  are convex  
 $(\sqrt{\sum_{i=1}^d x_i^2}) = \|\cdot\|$ ,  $\|x\|_1 = \sum_{i=1}^d |x_i|$ ,  $\|x\|_\infty (= \max_{1 \leq i \leq d} |x_i|)$

III. For any convex function  $f$  and any  $\alpha \geq 0$ ,  $\alpha f$  is convex

IV.  $f, g$  convex  $\Rightarrow f+g$  convex

V.  $f, g$  convex  $\Rightarrow \max(f, g)$  convex

VI.  $f: \mathbb{R}^m \rightarrow \mathbb{R}$  convex  $\Rightarrow x \mapsto f(Wx+b)$  convex  
 $W \in \mathbb{R}^{m \times d}$ ,  $b \in \mathbb{R}^m$

Proof of V:

$\forall (x, y) \in (\mathbb{R}^d)^2$ ,  $\forall \alpha \in [0, 1]$ ,

$\max(f(\alpha x + (1-\alpha)y), g(\alpha x + (1-\alpha)y))$

If the max is attained at  $f$

$\max(f(\alpha x + (1-\alpha)y), g(\alpha x + (1-\alpha)y))$

$= f(\alpha x + (1-\alpha)y) \stackrel{\text{convexity of } f}{\leq} \alpha f(x) + (1-\alpha)f(y)$

$\leq \alpha \max(f(x), g(x)) + (1-\alpha) \max(f(y), g(y))$

$\downarrow$   
convexity inequality

The case where the maximum is attained for  $g$  is handled the same way.

Exercise: Show that the objective function for SVM

$x \mapsto \frac{1}{n} \sum_{i=1}^n \max(1 - y_i a_i^T x, 0)$  is convex  
 (using the rules)

$$(a_i \in \mathbb{R}^d, y_i \in \{-1, 1\}, x \in \mathbb{R}^d)$$

Solution By I,  $x \mapsto 0$  is convex and  $x \mapsto 1 - y_i a_i^T x$  is convex  $\forall i$

By V,  $\forall i$ ,  $x \mapsto \max(1 - y_i a_i^T x, 0)$

By IV (applied  $n-1$  times),  $x \mapsto \sum_{i=1}^n \max(1 - y_i a_i^T x, 0)$  is convex

By III,  $x \mapsto \frac{1}{n} \sum_{i=1}^n \max(1 - y_i a_i^T x, 0)$  is convex

↳ Another characterization of convexity for  $C^2$  functions

Def:  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is  $C^2$  (twice continuously differentiable)

if it is  $C^1$  and  $\nabla f: \mathbb{R}^d \rightarrow \mathbb{R}^d$  is also  $C^1$

$f \in C^2 \Rightarrow \exists \nabla^2 f(x) \in \mathbb{R}^{d \times d}$  symmetric matrix

such that

2<sup>nd</sup>-order Taylor expansion of  $f$  around  $x$

$$f(y) \approx f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} (y-x)^T \nabla^2 f(x) (y-x)$$

when  $\|y-x\|$  is small

Note:  $\nabla f(x) = \begin{bmatrix} \frac{\partial f}{\partial x_1}(x) \\ \vdots \\ \frac{\partial f}{\partial x_d}(x) \end{bmatrix}$  and  $\nabla^2 f(x) = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2}(x) & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_d}(x) \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_d \partial x_1}(x) & \dots & \frac{\partial^2 f}{\partial x_d^2}(x) \end{bmatrix}$

$\nabla^2 f(x)$ : Hessian matrix

$$\frac{\partial^2 f}{\partial x_i \partial x_j}(x) = \frac{\partial^2 f}{\partial x_j \partial x_i}(x)$$

Theorem: Let  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  be  $C^2$ .

$$\left[ f \text{ is convex} \right] \Leftrightarrow \left[ \forall x \in \mathbb{R}^d, \nabla^2 f(x) \succeq 0 \right]$$

↑  
"positive semidefinite"

$\succeq$ :  $A \in \mathbb{R}^{d \times d}$  is positive semidefinite (or PSD)

if  $A = A^T$  (A symmetric) and  $\forall x \in \mathbb{R}^d, \underbrace{x^T A x}_{1 \times d \quad d \times d} \geq 0$

$$M \in \mathbb{R}^{d \times m}$$

$$M = [M_{ij}]_{\substack{1 \leq i \leq d \\ 1 \leq j \leq m}}$$

$$M^T \in \mathbb{R}^{m \times d}$$

$$[M^T]_{ji} = M_{ij}$$

$$x = \begin{bmatrix} \phantom{x} \\ \phantom{x} \\ \phantom{x} \end{bmatrix}$$

$$x^T = \begin{bmatrix} \phantom{x} & \phantom{x} & \phantom{x} \end{bmatrix}$$

(Ex)  $A \in \mathbb{R}^{m \times d}$ ,  $y \in \mathbb{R}^m$ ,  $x \in \mathbb{R}^d$

•  $f(x) = \frac{1}{2m} \|Ax - y\|^2$  linear least squares regression

$$\nabla f(x) = \frac{1}{m} A^T (Ax - y) \in \mathbb{R}^d$$

$$\nabla^2 f(x) = \frac{1}{m} A^T A$$

$$\forall v \in \mathbb{R}^d, v^T \nabla^2 f(x) v = v^T \left( \frac{A^T A}{m} \right) v = \frac{1}{m} \overbrace{v^T A^T A v}^{(Av)^T} = \frac{1}{m} \|Av\|^2 \geq 0$$

•  $f(x) = \frac{1}{m} \sum_{i=1}^m \log(1 + \exp(-y_i a_i^T x))$

logistic loss

$a_i \in \mathbb{R}^d, y_i \in \{0, 1\}$

$$\Rightarrow \forall i=1..m, f_i(x) = \log(1 + \exp(-y_i a_i^T x))$$

$$\frac{-y_i \exp(-y_i a_i^T x)}{1 + \exp(-y_i a_i^T x)} = \nabla f_i(x) = \frac{\overbrace{-y_i}^{\in \mathbb{R}}}{1 + \exp(y_i a_i^T x)} \underbrace{a_i}_{\in \mathbb{R}^d}$$

$$\nabla^2 f_i(x) = \frac{y_i^2 \exp(y_i a_i^T x)}{(1 + \exp(y_i a_i^T x))^2} \underbrace{a_i a_i^T}_{\substack{d \times d \\ \underbrace{d \times 1} \quad \underbrace{1 \times d}}}$$

NB:  $(a_i a_i^T)^T = (a_i^T)^T a_i^T = a_i a_i^T$

$$\forall v \in \mathbb{R}^d, v^T \nabla^2 f_i(x) v = \frac{y_i^2 \exp(y_i a_i^T x)}{(1 + \exp(y_i a_i^T x))^2} \underbrace{v^T a_i}_{1 \times d} \underbrace{a_i^T v}_{d \times 1}$$

$$= \frac{y_i^2 \exp(y_i a_i^T x)}{(1 + \exp(y_i a_i^T x))^2} \underbrace{a_i^T v}_{1 \times 1} \underbrace{v^T a_i}_{1 \times 1}$$

$$= \frac{y_i^2 \exp(y_i a_i^T x)}{(1 + \exp(y_i a_i^T x))^2} (a_i^T v)^2 \geq 0$$

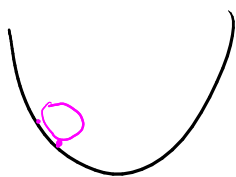
↳ Convexity:

→ can be checked for  $C^2$ ,  $C^1$ , and even not  $C^1$  functions

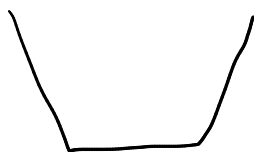
(even true for functions  $\mathbb{R}^d \rightarrow \mathbb{R} \cup \{\infty\}$ )

→ for convex functions,  $\|\nabla f(\bar{x})\| = 0$

$$\Leftrightarrow \bar{x} \in \underset{x \in \mathbb{R}^d}{\text{argmin}} f(x)$$



Convex



Convex



Convex



Def: (Strongly convex function)

•  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is  $\mu$ -strongly convex ( $\mu > 0$ ) if  
 $\forall (x, y) \in (\mathbb{R}^d)^2, f(\alpha x + (1-\alpha)y) \leq \alpha f(x) + (1-\alpha)f(y) - \frac{\mu \alpha(1-\alpha)}{2} \|y-x\|^2$

•  $f: \mathbb{R}^d \rightarrow \mathbb{R} \ C^1$  is  $\mu$ -strongly convex ( $\mu > 0$ ) if  
 $\forall (x, y) \in (\mathbb{R}^d)^2, f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{\mu}{2} \|y-x\|^2$

•  $f: \mathbb{R}^d \rightarrow \mathbb{R} \ C^2$  is  $\mu$ -strongly convex ( $\mu > 0$ ) if

$$\nabla^2 f(x) \succeq \mu I \quad (\Leftrightarrow) \quad \nabla^2 f(x) - \mu I \succeq 0$$

$$I = \begin{bmatrix} 1 & & \\ & \ddots & \\ & & 1 \end{bmatrix} \in \mathbb{R}^{d \times d}$$

• Strongly convex functions are convex

• If  $f$  is  $\mu$ -strongly convex, then  $g: x \mapsto f(x) - \frac{\mu}{2} \|x\|^2$  is convex

$$\nabla^2 g(x) = \nabla^2 f(x) - \mu I \succeq 0$$

$$\varphi: x \mapsto \frac{\mu}{2} \|x\|^2$$

$$\nabla \varphi(x) = \mu x$$

$$\nabla^2 \varphi(x) = \mu I$$

• Examples of strongly convex functions

$$x \mapsto \frac{\mu}{2} \|x\|^2$$

• If  $f$  is convex,  $f + \frac{\mu}{2} \|x\|^2$  is  $\mu$ -strongly convex

  $f \ C^1, \mu\text{-strongly convex}$

$$y \mapsto f(y) + \nabla f(x)^T (y-x) + \frac{\mu}{2} \|y-x\|^2$$

## Properties of strongly convex functions for optimization

• If  $f$  is strongly convex, it has at most one global minimum, i.e.  $|\arg\min_x f(x)| \leq 1$ .

• If  $f$  is strongly convex and continuous on  $\mathbb{R}^d$ , then it has a unique minimum,

• If  $f$  is strongly convex and  $C^1$  on  $\mathbb{R}^d$ , then it has a unique minimum and that minimum is the only solution to  $\nabla f(x) = 0$

NB: Strict convexity  $\neq \mu$ -strong convexity

$$f \in C^2 \text{ strictly convex} \Leftrightarrow \nabla^2 f(x) \succ 0 \quad \forall x \in \mathbb{R}^d$$
$$\Leftrightarrow v^T \nabla^2 f(x) v > 0 \quad \forall x \in \mathbb{R}^d, \forall v \in \mathbb{R}^d \setminus \{0\}$$

Ex)  $x \mapsto e^{-x}$



strictly convex but not strongly convex

$$\nabla^2 f(x) = -e^{-x} > 0 \quad \forall x \in \mathbb{R}$$

↳ There are many inequalities that hold for convex and strongly convex inequalities (cheat sheet by Fabian Pedregosa)

Ex) If  $f$  convex  $C^1$ , then

$$\forall (x, y) \in (\mathbb{R}^d)^2, (\nabla f(x) - \nabla f(y))^T (x - y) \geq 0$$

“co-coercivity”

Proof:

$$f(y) \geq f(x) + \nabla f(x)^T (y-x)$$

and

$$f(x) \geq f(y) + \nabla f(y)^T (x-y)$$

$$= f(y) - \nabla f(y)^T (y-x)$$

Sum the two inequalities

$$f(x) + f(y) \geq f(y) + f(x) + (\nabla f(x) - \nabla f(y))^T (y-x)$$

$$0 \geq (\nabla f(x) - \nabla f(y))^T (y-x)$$

$$0 \leq (\nabla f(x) - \nabla f(y))^T (x-y)$$

Exercise: Let  $f$  be a  $C^1$ ,  $\mu$ -strongly convex function in  $\mathbb{R}^d$   
and let  $x^* = \operatorname{argmin}_x f(x)$  ( $\{x^*\} = \operatorname{argmin}_x f(x)$ )

Goal:  $\forall x \in \mathbb{R}^d$ ,  $\|\nabla f(x)\|^2 \geq 2\mu (f(x) - f(x^*))$

Sometimes denoted  
by  $\min_{x \in \mathbb{R}^d} f(x)$

1) Start with

$$f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{\mu}{2} \|y-x\|^2 \quad \text{with fixed } x \in \mathbb{R}^d \text{ and } \forall y \in \mathbb{R}^d$$

2) Use  $\min_y f(y) = f(x^*)$  and that  $g: y \mapsto f(x) + \nabla f(x)^T (y-x) + \frac{\mu}{2} \|y-x\|^2$  is strongly convex

$$\nabla g(y) = \nabla f(x) + \mu(y-x)$$

$$f(y) \geq g(y) \quad \forall y \in \mathbb{R}^d \Rightarrow \min_{y \in \mathbb{R}^d} f(y) \geq \min_{y \in \mathbb{R}^d} g(y) \quad (1)$$

$$\underbrace{\hspace{10em}}_{= f(x^*)}$$

$g$  is a strongly convex function (sum of  $y \mapsto f(x) + \nabla f(x)^T(y-x)$ , that is convex, and  $y \mapsto \frac{\mu}{2} \|y-x\|^2$ , that is  $\mu$ -strongly convex)

hence it has a unique global minimum, which is the solution

$$\begin{aligned} \text{of } \nabla g(y) = 0 &\Leftrightarrow \nabla f(x) + \mu(y-x) = 0 \quad \mu > 0 \\ &\Leftrightarrow y = x - \frac{1}{\mu} \nabla f(x) \end{aligned}$$

Let  $y^* := x - \frac{1}{\mu} \nabla f(x)$ . Then (1) becomes

$$f(x^*) = \min_{y \in \mathbb{R}^d} f(y) \geq \min_{y \in \mathbb{R}^d} g(y) = g(y^*) \quad (2)$$

$$\begin{aligned} g(y^*) &= f(x) + \nabla f(x)^T(y^* - x) + \frac{\mu}{2} \|y^* - x\|^2 \\ &= f(x) + \nabla f(x)^T\left(-\frac{1}{\mu} \nabla f(x)\right) + \frac{\mu}{2} \left\|-\frac{1}{\mu} \nabla f(x)\right\|^2 \\ &= f(x) - \frac{1}{\mu} \|\nabla f(x)\|^2 + \frac{\mu}{2} \times \frac{1}{\mu^2} \|\nabla f(x)\|^2 \\ &= f(x) - \frac{1}{2\mu} \|\nabla f(x)\|^2 \end{aligned}$$

$$a^T a = \|a\|^2$$

Plugging this expression into (2), we get

$$f(x^*) \geq f(x) - \frac{1}{2\mu} \|\nabla f(x)\|^2$$

$$\|\nabla f(x)\|^2 \geq 2\mu \underbrace{(f(x) - f(x^*))}_{\geq 0}$$

Takeaway: Strongly convex functions are lower bounded by a family of strongly convex quadratic functions ( $y \mapsto f(x) + \nabla f(x)^T(y-x) + \frac{\mu}{2} \|y-x\|^2$ ), and these are very nice functions to optimize

# Lipschitz continuous function and optimization

Def.  $\varphi: \mathbb{R}^d \rightarrow \mathbb{R}^m$  is a  $L$ -Lipschitz continuous function if  
 $\forall (x, y) \in (\mathbb{R}^d)^2, \quad \|\varphi(x) - \varphi(y)\| \leq L \|x - y\| \quad (L > 0)$

Ex)  $\varphi: x \mapsto Wx + b$  is  $\|W\|$ -Lipschitz (continuous)  
where  $\|W\| = \max_{x \neq 0} \frac{\|Wx\|}{\|x\|}$

NB:  $L$  is not unique, but typically the smallest possible value is chosen.

Def. A  $C^{1,1}$  function  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is called  $L$ -smooth with  $L > 0$

if  $\forall (x, y) \in (\mathbb{R}^d)^2, \quad \|\nabla f(x) - \nabla f(y)\| \leq L \|x - y\|$ ,  
i.e.  $\nabla f: \mathbb{R}^d \rightarrow \mathbb{R}^d$  is  $L$ -Lipschitz

Ex) Any quadratic function  $x \mapsto c + g^T x + \frac{1}{2} x^T H x$   
 $c \in \mathbb{R}, g \in \mathbb{R}^d, H \in \mathbb{R}^{d \times d}$   
is  $C^{1,1}_L$  with  $L = \frac{\|H + H^T\|}{2}$

• In particular,  $x \mapsto \frac{1}{2m} \|Ax - y\|^2 = \frac{1}{2m} y^T y - \frac{1}{m} y^T A x + \frac{1}{2m} x^T A^T A x$

is  $C^{1,1}_L$  with  $L = \frac{\|A^T A\|}{m}$

•  $x \mapsto \frac{1}{m} \sum_{i=1}^m \log(1 + e^{-y_i a_i^T x})$  is  $C^{1,1}_L$ ,  $L = \frac{\|A^T A\|}{4m}$

Theorem: Let  $f: \mathbb{R}^d \rightarrow \mathbb{R}$  is  $C_{L}^{1,1}$  ( $\Delta$  not necessarily convex)  $L > 0$

Then,  $\forall (x, y) \in (\mathbb{R}^d)^2$ ,

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

Additional properties

$\rightarrow$  If  $f \in C_{L}^{2,1}$  ( $= C^2 \cap C_{L}^{1,1}$ ),  $\|\nabla^2 f(x)\| \leq L$   
 $\nabla^2 f(x) \preceq LI$

$\rightarrow$  If  $f$  is  $C_{L}^{1,1}$  and  $\mu$ -strongly convex, then  $\forall (x, y) \in (\mathbb{R}^d)^2$ ,

$$\underbrace{f(x) + \nabla f(x)^T (y-x) + \frac{\mu}{2} \|y-x\|^2}_{\mu\text{-strong convexity}} \leq f(y) \leq \underbrace{f(x) + \nabla f(x)^T (y-x) + \frac{L}{2} \|y-x\|^2}_{L\text{-smoothness}}$$

Hence  $\mu \leq L$   
and  $f$  can be "sandwiched" between two (simple) quadratics (aka quadratic functions)

Summary:

- Convex functions are good (for optimization)
  - $\rightarrow$  Nice optimality conditions for global optimality
  - $\rightarrow$  Even nicer results for strongly convex functions
  - $\rightarrow$  Useful inequalities

Time for strongly convex quadratic functions  $\leftarrow$

• Best case:  $C_{L}^{1,1}$  +  $\mu$ -strongly convex function

Last-minute inequality:

$f$   $\mu$ -strongly convex and  $C^1$ , then

$$\forall x \in \mathbb{R}^d, \quad f(x) - f(x^0) \geq \frac{\mu}{2} \|x - x^0\|^2$$