

# Probabilistic Generative Models — Convex analysis

Caio Corro

These notes contain the basic knowledge about convex analysis required for the Probabilistic Generative Models course. Some notions are not properly defined and can be ignored for the scope of this course (*e.g.* the relative interior of a set). The objective is to stay on simple notions to understand the methods studied in the course, without getting burdened with more complex notions that do not improve the clarity of the presentation. More details on convex analysis and optimization can be found in [Boyd et al., 2004], [Beck, 2017] and [Borwein and Lewis, 2006].

## 1 Convex sets and functions

Convex functions are interesting to study for (at least) two reasons:

1. any local minimizer of a convex function is also a global minimizer (if it exists!), which allows to derive “simple” optimality conditions that are both necessary and sufficient (Section 2);
2. we can derive parameterized bounds on convex functions that will be the basis of variational methods for inference.

Convex functions are defined as follows.

**Definition 1** (Convex set). Let  $U \subseteq \mathbb{R}^k$  be a set. The set  $U$  is convex if and only if:

$$\forall \mathbf{u}, \mathbf{v} \in U, \epsilon \in [0, 1] : \quad \epsilon \mathbf{u} + (1 - \epsilon) \mathbf{v} \in U.$$

**Definition 2** (Convex function). Let  $f : U \rightarrow \mathbb{R}$ ,  $U \subseteq \mathbb{R}^k$ , be a function. The function  $f$  is convex if and only if:

1.  $U$  is a convex set;
2. the function satisfies the following inequality:

$$\forall \mathbf{u}, \mathbf{v} \in U, \epsilon \in [0, 1] : \quad f(\epsilon \mathbf{u} + (1 - \epsilon) \mathbf{v}) \leq \epsilon f(\mathbf{u}) + (1 - \epsilon) f(\mathbf{v})$$

A function  $f$  is concave if and only if  $-f$  is convex.

If the domain of a function is an open set and it is (twice) differentiable everywhere, convexity can also be characterized in terms of its gradient (or Hessian).

**Proposition 1** (Convexity: first-order condition). Let  $U \subseteq \mathbb{R}^k$  be a convex open set and  $f : U \rightarrow \mathbb{R}$  be a differentiable function. Then,  $f$  is convex if and only if:

$$\forall \mathbf{u}, \mathbf{v} \in U : \quad f(\mathbf{v}) \geq f(\mathbf{u}) + \langle \nabla f(\mathbf{u}), \mathbf{v} - \mathbf{u} \rangle.$$

**Proposition 2** (Convexity: second-order condition). Let  $U \subseteq \mathbb{R}^k$  be a convex open set and  $f : U \rightarrow \mathbb{R}$  be a twice differentiable function. Then,  $f$  is convex if and only if the Hessian  $\nabla^2 f(\mathbf{u})$  is a positive semidefinite matrix, for all  $\mathbf{u} \in U$ .

Proof of Propositions 1 and 2 can be found in [Boyd et al., 2004, Section 3.1.3 and Exercise 3.8]. Note that the second order condition simplifies to non-negative second-order derivatives in the case of  $U \subseteq \mathbb{R}$ .

Two other important characteristics of functions are properness and closedness. A function  $f : U \rightarrow [-\infty, +\infty]$  is a proper function if and only if there is no  $\mathbf{u} \in U$  s.t.  $f(\mathbf{u}) = -\infty$  and there exists at least one  $\mathbf{u}' \in U$  s.t.  $f(\mathbf{u}') \neq +\infty$ . The closedness property is equivalent to lower semicontinuity. We won't study this property in this course, more information can be found in [Beck, 2017, Section 2.1].

Note that a convex function may not be differentiable everywhere. This is the case of the absolute value function  $f(u) = |u|$ , which is not differentiable at 0. In the case of convex functions, the concept of subgradient is a very useful generalization of the gradient,<sup>1</sup> that will be used to define optimality conditions, among others.

**Definition 3** (Subgradient). Let  $U \subseteq \mathbb{R}^k$  be a set and  $f : U \rightarrow \mathbb{R}$  a convex function. Then, a vector  $\mathbf{g} \in \mathbb{R}^k$  is a subgradient of  $f$  at  $\mathbf{u} \in U$  if and only if:

$$\forall \mathbf{v} \in U : f(\mathbf{v}) \geq f(\mathbf{u}) + \langle \mathbf{g}, \mathbf{v} - \mathbf{u} \rangle.$$

The set of subgradients at a given point is called the subdifferential and is denoted as  $\partial f(\mathbf{u})$ .

**Proposition 3** (Subgradient properties). Let  $U$  be a convex set and  $f : U \rightarrow \mathbb{R}$  be a convex function. The subgradients of  $f$  satisfy the following properties:

1.  $\partial f(\mathbf{u}), \mathbf{u} \in U$ , is a convex set;
2. if  $f$  is a proper function and  $\mathbf{u} \in \text{int}(\text{dom } f)$ , then  $\partial f(\mathbf{u}) \neq \emptyset$ , *i.e.* there exists at least one subgradient at  $\mathbf{u}$  [Beck, 2017, Theorem 3.14];
3. if  $f$  is a proper function,  $\mathbf{u} \in \text{int}(\text{dom } f)$  and  $f$  differentiable at  $\mathbf{u}$ , then  $\partial f(\mathbf{u}) = \{\nabla f(\mathbf{u})\}$ .

In practice, it is often useful to work with extended-value extensions of a function. Let  $U \subseteq \mathbb{R}^k$  be a set and  $f : U \rightarrow \mathbb{R}$  be a function. The extended-value extensions of  $f$ , denoted  $\tilde{f} : \mathbb{R}^k \rightarrow \mathbb{R} \cup \{\infty\}$ , is defined as follows:

$$\tilde{f}(\mathbf{u}) = \begin{cases} f(\mathbf{u}) & \text{if } \mathbf{u} \in \text{dom } f, \\ \infty & \text{otherwise.} \end{cases}$$

The domain of a function is then defined as the set of input whose images are not equal to infinity (but potentially equal to minus infinity), *i.e.*  $\text{dom } \tilde{f} = \{\mathbf{u} \in \mathbb{R}^k \mid \tilde{f}(\mathbf{u}) \neq \infty\}$ . To simplify notation, it is usual to simply refer to the extended-value extension as  $f$ .

**Definition 4** (Indicator function). Let  $S \subseteq \mathbb{R}^k$  be a set. The indicator function  $\delta_S : \mathbb{R}^k \rightarrow \mathbb{R} \cup \{\infty\}$  of

<sup>1</sup>Note that for a nonconvex function, there may exist points where the function is differentiable but where no subgradient exists. Hence, the term generalization is an exaggeration in general.

the set  $S$  is defined as follows:

$$\delta_S(\mathbf{u}) = \begin{cases} 0 & \text{if } \mathbf{u} \in S, \\ \infty & \text{otherwise.} \end{cases}$$

Indicator functions can be used to introduce constraints in an optimization problem as a term in the objective. For example, the two following problems are equivalent:

$$\begin{aligned} \min_{\mathbf{u}} f(\mathbf{u}) &= \min_{\mathbf{u}} f(\mathbf{u}) + \delta_S(\mathbf{u}) \\ \text{s.t. } \mathbf{u} &\in S. \end{aligned}$$

## 2 Optimality conditions

In this section, we present optimality conditions for optimization (un)constrained optimization problems. They will be useful to compute closed-form expression of minimizers of these optimization problems.

**Proposition 4** (Fermat's optimality conditions). Let  $f : U \rightarrow \mathbb{R} \cup \{\infty\}$  be a proper convex function. Then,  $\hat{\mathbf{u}} \in \arg \min_{\mathbf{u} \in U} f(\mathbf{u})$  is a minimizer of  $f$  if and only if  $\mathbf{0} \in \partial f(\hat{\mathbf{u}})$ .

*Proof.* For any  $\mathbf{g} \in \partial f(\hat{\mathbf{u}})$ , the subgradient inequality holds:

$$\forall \mathbf{v} \in U : f(\mathbf{v}) \geq f(\hat{\mathbf{u}}) + \langle \mathbf{g}, \mathbf{v} - \hat{\mathbf{u}} \rangle$$

In particular, we have  $\mathbf{0} \in \partial f(\hat{\mathbf{u}})$  by assumption, therefore setting  $\mathbf{g} = \mathbf{0}$  we obtain:

$$f(\mathbf{v}) \geq f(\hat{\mathbf{u}}),$$

hence  $\hat{\mathbf{u}}$  is a minimizer of  $f$ . □

In the constrained case, we restrict ourselves to problems of the following form:

$$\begin{aligned} \min_{\mathbf{u} \in \mathbb{R}^k} f(\mathbf{u}) \\ \text{s.t. } \mathbf{A}\mathbf{u} &= \mathbf{b} \\ \mathbf{C}\mathbf{u} &\leq \mathbf{d} \end{aligned}$$

where  $f$  is a convex function and  $\mathbf{A} \in \mathbb{R}^{m \times k}$ ,  $\mathbf{b} \in \mathbb{R}^m$ ,  $\mathbf{C} \in \mathbb{R}^{n \times k}$  and  $\mathbf{d} \in \mathbb{R}^n$  defines a set of  $m$  equality and  $n$  inequality constraints. We assume there exists at least one feasible solution to the problem. The Lagrangian  $L : \mathbb{R}^k \times \mathbb{R}^m \times \mathbb{R}_+^n \rightarrow \mathbb{R}$  associated with this problem is defined as:

$$L(\mathbf{u}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{u}) + \langle \boldsymbol{\lambda}, \mathbf{A}\mathbf{u} - \mathbf{b} \rangle + \langle \boldsymbol{\mu}, \mathbf{C}\mathbf{u} - \mathbf{d} \rangle$$

where  $\boldsymbol{\lambda} \in \mathbb{R}^m$  and  $\boldsymbol{\mu} \in \mathbb{R}_+^n$  are dual variables associated with inequalities and equalities, respectively. Importantly,  $\boldsymbol{\mu}$  must contain only non-negative values. By analogy, the vector  $\mathbf{u}$  is called the primal variable. Although outside the scope of this course, for any couple of dual variables, the minimizing the Lagrangian over the primal variable defines an lower bound on the original problem:

$$\begin{aligned} \forall \boldsymbol{\lambda} \in \mathbb{R}^m, \boldsymbol{\mu} \in \mathbb{R}_+^n : \min_{\mathbf{u} \in \mathbb{R}^k} L(\mathbf{u}, \boldsymbol{\lambda}, \boldsymbol{\mu}) &\leq \min_{\mathbf{u} \in \mathbb{R}^k} f(\mathbf{u}) \\ \text{s.t. } \mathbf{A}\mathbf{u} &= \mathbf{b} \\ \mathbf{C}\mathbf{u} &\leq \mathbf{d} \end{aligned}$$

**Proposition 5** (Karush–Kuhn–Tucker optimality conditions). Let  $f : \mathbb{R}^k \rightarrow \mathbb{R}$  be a convex function and assume the following optimization problem:

$$\begin{aligned} \min_{\mathbf{u} \in \mathbb{R}^k} \quad & f(\mathbf{u}) \\ \text{s.t.} \quad & \mathbf{A}\mathbf{u} = \mathbf{b} \\ & \mathbf{C}\mathbf{u} \leq \mathbf{d} \end{aligned}$$

where  $\mathbf{A} \in \mathbb{R}^{m \times k}$ ,  $\mathbf{b} \in \mathbb{R}^m$ ,  $\mathbf{C} \in \mathbb{R}^{n \times k}$  and  $\mathbf{d} \in \mathbb{R}^m$ . Then, optimal primal and dual variable  $\hat{\mathbf{u}} \in \mathbb{R}^k$ ,  $\hat{\boldsymbol{\lambda}} \in \mathbb{R}^m$  and  $\hat{\boldsymbol{\mu}} \in \mathbb{R}_{++}^k$  satisfy the following constraints:

$$\begin{aligned} \text{(stationarity)} \quad & \partial_{\mathbf{u}} L(\mathbf{u}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \ni \mathbf{0} \quad \iff \quad \partial f(\mathbf{u}) + \mathbf{A}^\top \boldsymbol{\lambda} + \mathbf{C}^\top \boldsymbol{\mu} \ni \mathbf{0} \\ \text{(primal feasibility)} \quad & \mathbf{A}\mathbf{u} = \mathbf{b} \\ & \mathbf{C}\mathbf{u} \leq \mathbf{d} \\ \text{(dual feasibility)} \quad & \boldsymbol{\mu} \geq \mathbf{0} \\ \text{(complementary slackness)} \quad & \langle \boldsymbol{\mu}, \mathbf{C}\mathbf{u} - \mathbf{d} \rangle = \mathbf{0} \end{aligned}$$

Note that the KKT conditions are in general only necessary conditions, but they become sufficient conditions in the case of convex problems as presented in Proposition 5

### 3 Fenchel conjugates

Fenchel conjugates are a very useful tools that appears in many area of machine learning (duality for convex learning problems like SVMs, loss functions, variational approximations, ...).

**Definition 5** (Fenchel conjugate). Let  $f : U \rightarrow \mathbb{R} \cup \{\infty\}$  be a function. The Fenchel conjugate of  $f$ , denoted  $f^*$ , is defined as follows:

$$f^*(\mathbf{t}) = \sup_{\mathbf{u} \in \text{dom } f} \langle \mathbf{u}, \mathbf{t} \rangle - f(\mathbf{u}).$$

**Proposition 6** (Fenchel-Young inequality). Let  $f : U \rightarrow \mathbb{R} \cup \{\infty\}$  be a function. Then, the following inequality holds:

$$\forall \mathbf{u} \in U, \mathbf{t} \in \text{dom } f^* : \quad f^*(\mathbf{t}) \geq \langle \mathbf{u}, \mathbf{t} \rangle - f(\mathbf{u}).$$

*Proof.* The inequality directly follows from the definition of the Fenchel conjugate. □

**Proposition 7** (Fenchel biconjugation). Let  $f : U \rightarrow \mathbb{R} \cup \{\infty\}$  be a closed convex function. Then:

$$f = f^{**},$$

that is conjugate of the conjugate of  $f$  is  $f$ .

## 4 Jensen's inequality

**Proposition 8** (Jensen's inequality (simplified)). Let  $f : U \rightarrow \mathbb{R} \cup \{\infty\}$  be a convex function,  $\mathbf{u}^{(1)} \dots \mathbf{u}^{(n)}$  be  $n$  datapoints and  $\boldsymbol{\mu} \in \Delta(n)$ . Then:

$$\sum_{i=1}^n \mu_i f(\mathbf{u}^{(i)}) \geq f\left(\sum_{i=1}^n \mu_i \mathbf{u}^{(i)}\right).$$

*Proof.* Let  $\mathbf{u} = \sum_{i=1}^n \mu_i \mathbf{u}^{(i)}$  and  $\mathbf{g} \in \partial f(\mathbf{u})$ . By convexity of  $f$ , we have:

$$\forall i \in \{1 \dots n\} : f(\mathbf{u}^{(i)}) \geq f(\mathbf{u}) + \langle \mathbf{g}, \mathbf{u}^{(i)} - \mathbf{u} \rangle.$$

As  $\boldsymbol{\mu}$  contains only positive values, by multiplying and then summing all inequalities we obtain:

$$\begin{aligned} \sum_i \mu_i f(\mathbf{u}^{(i)}) &\geq \sum_i \mu_i \left( f(\mathbf{u}) + \langle \mathbf{g}, \mathbf{u}^{(i)} - \mathbf{u} \rangle \right), \\ \sum_i \mu_i f(\mathbf{u}^{(i)}) &\geq \sum_i \mu_i f(\mathbf{u}) + \left\langle \mathbf{g}, \sum_i \mu_i (\mathbf{u}^{(i)} - \mathbf{u}) \right\rangle, \end{aligned}$$

Note that as  $\boldsymbol{\mu} \in \Delta(n)$ , we have  $\sum_i \mu_i f(\mathbf{u}) = f(\mathbf{u}) \sum_i \mu_i = f(\mathbf{u})$ . A similar argument cancels the second term of the right-hand side of the inequality. Replacing  $\boldsymbol{\mu}$  by its definition, we obtain:

$$\sum_{i=1}^n \mu_i f(\mathbf{u}^{(i)}) \geq f\left(\sum_{i=1}^n \mu_i \mathbf{u}^{(i)}\right).$$

□

From the proposition above, we can see that it really looks like an expectation. This is indeed the case: the Jensen's inequality holds can be defined over expectation, including continuous random variables and transformed random variables. However, we don't prove this in these notes.

**Proposition 9** (Jensen's inequality). Let  $f : X \rightarrow \mathbb{R} \cup \{\infty\}$  be a convex function,  $\mathcal{X}$  a random variables taking values in  $X$  and  $p$  a probability distribution over  $\mathcal{X}$ . Then:

$$\mathbb{E}_{p(\mathcal{X})} [f(\mathcal{X})] \geq f\left(\mathbb{E}_{p(\mathcal{X})} [\mathcal{X}]\right).$$

Under mild conditions, the inequality also holds if  $\mathcal{X}$  is transformed via a (possibly nonconvex) function  $g$ :

$$\mathbb{E}_{p(\mathcal{X})} [f(g(\mathcal{X}))] \geq f\left(\mathbb{E}_{p(\mathcal{X})} [g(\mathcal{X})]\right).$$

## 5 Exercises

In the following,  $\sigma$  refers to the sigmoid function and  $\Delta(n)$  to the simplex of dimension  $n$ .

1. Let  $U$  be a convex set and  $n$  a strictly positive integer. Prove that:

$$\forall \mathbf{u}^{(1)} \dots \mathbf{u}^{(n)}, \boldsymbol{\mu} \in \Delta(n) : \sum_{i=1}^n \mu_i \mathbf{u}^{(i)} \in U.$$

2. Prove that the following functions are convex:
  - (a)  $f : \mathbb{R}^k \rightarrow \mathbb{R}$  defined as  $f(\mathbf{u}) = \langle \mathbf{a}, \mathbf{u} \rangle + b$  where  $\mathbf{a} \in \mathbb{R}^k$  and  $b \in \mathbb{R}$
  - (b)  $f : \mathbb{R} \rightarrow \mathbb{R}$  defined as  $f(u) = \exp(u)$
  - (c)  $f : \mathbb{R}_{++} \rightarrow \mathbb{R}$  defined as  $f(u) = -\log(u)$
  - (d)  $f : \mathbb{R}^k \rightarrow \mathbb{R}$  defined as  $f(\mathbf{u}) = \frac{1}{2} \|\mathbf{u}\|_2^2$
  - (e)  $f : \mathbb{R}^k \rightarrow \mathbb{R}$  defined as  $f(\mathbf{u}) = \log \sum_i \exp(u_i)$
3. Compute (explicitly!) the conjugates and the biconjugates of the above functions.
4. Prove that the function  $-\log \sigma(u)$  is convex.
5. Compute the Fenchel conjugate of function  $-\log \sigma(u)$ .
6. Use results of two previous questions to derive a parameterized upper bound to the sigmoid function.

## References

- [Beck, 2017] Beck, A. (2017). *First-order methods in optimization*. SIAM.
- [Borwein and Lewis, 2006] Borwein, J. and Lewis, A. (2006). *Convex Analysis*. Springer.
- [Boyd et al., 2004] Boyd, S., Boyd, S. P., and Vandenberghe, L. (2004). *Convex optimization*. Cambridge university press.