

First Order Descent Methods in Nonsmooth Convex Optimization

Juan PEYPOUQUET

Universidad Técnica Federico Santa María
Universidad de Chile

CIMPA School on Algorithmic Nonsmooth Optimization
Ibarra, September 19-21, 2017

NONSMOOTH CONVEX FUNCTIONS AND SUBDIFFERENTIAL CALCULUS

Extended real-valued functions

We consider $f : H \rightarrow \mathbb{R} \cup \{+\infty\}$ ¹, and write

$$\text{dom}(f) = \{x \in H : f(x) < +\infty\}.$$

The indicator function of a set $C \subset H$ is

$$\delta_C(x) = \begin{cases} 0 & \text{if } x \in C \\ +\infty & \text{if } x \notin C \end{cases}$$

Then,

$$\min\{f(x) : x \in C\} = \min\{f(x) + \delta_C(x) : x \in H\}.$$

¹We exclude here the case $f \equiv +\infty$.

Extended real-valued functions

We consider $f : H \rightarrow \mathbb{R} \cup \{+\infty\}$ ¹, and write

$$\text{dom}(f) = \{x \in H : f(x) < +\infty\}.$$

The **indicator function** of a set $C \subset H$ is

$$\delta_C(x) = \begin{cases} 0 & \text{if } x \in C \\ +\infty & \text{if } x \notin C \end{cases}$$

Then,

$$\min\{f(x) : x \in C\} = \min\{f(x) + \delta_C(x) : x \in H\}.$$

¹We exclude here the case $f \equiv +\infty$.

Extended real-valued functions

We consider $f : H \rightarrow \mathbb{R} \cup \{+\infty\}$ ¹, and write

$$\text{dom}(f) = \{x \in H : f(x) < +\infty\}.$$

The **indicator function** of a set $C \subset H$ is

$$\delta_C(x) = \begin{cases} 0 & \text{if } x \in C \\ +\infty & \text{if } x \notin C \end{cases}$$

Then,

$$\min\{f(x) : x \in C\} = \min\{f(x) + \delta_C(x) : x \in H\}.$$

¹We exclude here the case $f \equiv +\infty$.

Convex functions

A function $f : H \rightarrow \mathbb{R} \cup \{+\infty\}$ is **convex** if

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$$

for every $x, y \in H$ and every $\lambda \in (0, 1)$.

This is equivalent to the convexity of the epigraph of f :

$$\text{epi}(f) = \{(x, z) \in H \times \mathbb{R} : f(x) \leq z\},$$

and implies the convexity of $\text{dom}(f)$ and of every sublevel set

$$[f \leq \gamma] = \{x \in H : f(x) \leq \gamma\}, \quad \gamma \in \mathbb{R}.$$

Convex functions

A function $f : H \rightarrow \mathbb{R} \cup \{+\infty\}$ is **convex** if

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$$

for every $x, y \in H$ and every $\lambda \in (0, 1)$.

This is **equivalent** to the convexity of the **epigraph** of f :

$$\text{epi}(f) = \{(x, z) \in H \times \mathbb{R} : f(x) \leq z\},$$

and **implies** the convexity of $\text{dom}(f)$ and of every **sublevel set**

$$[f \leq \gamma] = \{x \in H : f(x) \leq \gamma\}, \quad \gamma \in \mathbb{R}.$$

Lower-semicontinuity

A function $f : H \rightarrow \mathbb{R} \cup \{+\infty\}$ is **lower-semicontinuous** at $x \in \text{dom}(f)$ if, for every $\varepsilon > 0$, there is $\delta > 0$ such that

$$\|x - y\| < \delta \implies f(y) > f(x) - \varepsilon.$$

Equivalently, if $x_k \rightarrow x \implies f(x) \leq \liminf_{k \rightarrow +\infty} f(x_k)$.

We say f is lower-semicontinuous if it is so at every $x \in \text{dom}(f)$.

Example

- Every continuous function is lower-semicontinuous.
- δ_C is lower-semicontinuous if, and only if, C is closed.

Lower-semicontinuity

A function $f : H \rightarrow \mathbb{R} \cup \{+\infty\}$ is **lower-semicontinuous** at $x \in \text{dom}(f)$ if, for every $\varepsilon > 0$, there is $\delta > 0$ such that

$$\|x - y\| < \delta \implies f(y) > f(x) - \varepsilon.$$

Equivalently, if $x_k \rightarrow x \implies f(x) \leq \liminf_{k \rightarrow +\infty} f(x_k)$.

We say f is **lower-semicontinuous** if it is so at every $x \in \text{dom}(f)$.

Example

- Every continuous function is lower-semicontinuous.
- δ_C is lower-semicontinuous if, and only if, C is closed.

Lower-semicontinuity

A function $f : H \rightarrow \mathbb{R} \cup \{+\infty\}$ is **lower-semicontinuous** at $x \in \text{dom}(f)$ if, for every $\varepsilon > 0$, there is $\delta > 0$ such that

$$\|x - y\| < \delta \implies f(y) > f(x) - \varepsilon.$$

Equivalently, if $x_k \rightarrow x \implies f(x) \leq \liminf_{k \rightarrow +\infty} f(x_k)$.

We say f is **lower-semicontinuous** if it is so at every $x \in \text{dom}(f)$.

Example

- *Every continuous function is lower-semicontinuous.*
- δ_C is lower-semicontinuous if, and only if, C is closed.

Two important properties

Proposition

Let $f : H \rightarrow \mathbb{R} \cup \{+\infty\}$. The following are equivalent:

- f is lower-semicontinuous.
- $\text{epi}(f)$ is closed.
- $[f \leq \gamma]$ is closed for every $\gamma \in \mathbb{R}$.

Proposition

Let $f : H \rightarrow \mathbb{R} \cup \{+\infty\}$ be convex and lower-semicontinuous. For every $x \in \text{dom}(f)$,

$$x_k \rightarrow x \implies f(x) \leq \liminf_{k \rightarrow +\infty} f(x_k).$$

Two important properties

Proposition

Let $f : H \rightarrow \mathbb{R} \cup \{+\infty\}$. The following are equivalent:

- f is lower-semicontinuous.
- $\text{epi}(f)$ is closed.
- $[f \leq \gamma]$ is closed for every $\gamma \in \mathbb{R}$.

Proposition

Let $f : H \rightarrow \mathbb{R} \cup \{+\infty\}$ be convex and lower-semicontinuous. For every $x \in \text{dom}(f)$,

$$x_k \rightarrow x \implies f(x) \leq \liminf_{k \rightarrow +\infty} f(x_k).$$

Existence of minimizers

A function $f : H \rightarrow \mathbb{R} \cup \{+\infty\}$ is **coercive** if

$$\lim_{\|x\| \rightarrow +\infty} f(x) = +\infty.$$

In other words, if every sublevel set $[f \leq \gamma]$ is bounded.

Theorem (Weierstrass-Hilbert-Tonelli)

If $f : H \rightarrow \mathbb{R} \cup \{+\infty\}$ is convex, lower-semicontinuous and coercive, then $S \neq \emptyset$.

Existence of minimizers

A function $f : H \rightarrow \mathbb{R} \cup \{+\infty\}$ is **coercive** if

$$\lim_{\|x\| \rightarrow +\infty} f(x) = +\infty.$$

In other words, if every sublevel set $[f \leq \gamma]$ is bounded.

Theorem (Weierstrass-Hilbert-Tonelli)

If $f : H \rightarrow \mathbb{R} \cup \{+\infty\}$ is convex, lower-semicontinuous and coercive, then $S \neq \emptyset$.

Existence of minimizers

A function $f : H \rightarrow \mathbb{R} \cup \{+\infty\}$ is **coercive** if

$$\lim_{\|x\| \rightarrow +\infty} f(x) = +\infty.$$

In other words, if every sublevel set $[f \leq \gamma]$ is bounded.

Theorem (Weierstrass-Hilbert-Tonelli)

If $f : H \rightarrow \mathbb{R} \cup \{+\infty\}$ is convex, lower-semicontinuous and coercive, then $S \neq \emptyset$.

Subdifferential calculus

Recall that if $f : H \rightarrow \mathbb{R}$ is convex and differentiable, then

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle \quad \forall x, y \in H.$$

Keeping this in mind, we define the subdifferential of $f : H \rightarrow \mathbb{R} \cup \{+\infty\}$ at $x \in H$ by

$$\partial f(x) = \{x^* \in H : f(y) \geq f(x) + \langle x^*, y - x \rangle \quad \forall y \in H\}.$$

We write $\text{dom}(\partial f) = \{x \in H : \partial f(x) \neq \emptyset\}$.

Subdifferential calculus

Recall that if $f : H \rightarrow \mathbb{R}$ is convex and differentiable, then

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle \quad \forall x, y \in H.$$

Keeping this in mind, we define the **subdifferential of $f : H \rightarrow \mathbb{R} \cup \{+\infty\}$ at $x \in H$** by

$$\partial f(x) = \{x^* \in H : f(y) \geq f(x) + \langle x^*, y - x \rangle \quad \forall y \in H\}.$$

We write $\text{dom}(\partial f) = \{x \in H : \partial f(x) \neq \emptyset\}$.

Subdifferential calculus

Recall that if $f : H \rightarrow \mathbb{R}$ is convex and differentiable, then

$$f(y) \geq f(x) + \langle \nabla f(x), y - x \rangle \quad \forall x, y \in H.$$

Keeping this in mind, we define the **subdifferential of $f : H \rightarrow \mathbb{R} \cup \{+\infty\}$ at $x \in H$** by

$$\partial f(x) = \{x^* \in H : f(y) \geq f(x) + \langle x^*, y - x \rangle \quad \forall y \in H\}.$$

We write $\text{dom}(\partial f) = \{x \in H : \partial f(x) \neq \emptyset\}$.

Some properties

Proposition

If f is differentiable at x , then $\partial f(x) = \{\nabla f(x)\}$.

Proposition

If $x^ \in \partial f(x)$ and $y^* \in \partial f(y)$, then $\langle x^* - y^*, x - y \rangle \geq 0$.*

Proposition (Fermat's Rule)

$$x \in S \iff 0 \in \partial f(x).$$

Some properties

Proposition

If f is differentiable at x , then $\partial f(x) = \{\nabla f(x)\}$.

Proposition

If $x^ \in \partial f(x)$ and $y^* \in \partial f(y)$, then $\langle x^* - y^*, x - y \rangle \geq 0$.*

Proposition (Fermat's Rule)

$$x \in S \iff 0 \in \partial f(x).$$

Some properties

Proposition

If f is differentiable at x , then $\partial f(x) = \{\nabla f(x)\}$.

Proposition

If $x^* \in \partial f(x)$ and $y^* \in \partial f(y)$, then $\langle x^* - y^*, x - y \rangle \geq 0$.

Proposition (Fermat's Rule)

$$x \in \mathcal{S} \iff 0 \in \partial f(x).$$

A calculus rule

Proposition (Moreau-Rockafellar)

For $f, g : H \rightarrow \mathbb{R} \cup \{+\infty\}$, we have

$$\partial f(x) + \partial g(x) \subseteq \partial(f + g)(x) \quad \forall x \in H.$$

If g is continuous *at some* $\bar{x} \in \text{dom}(f)$, then

$$\partial(f + g)(x) = \partial f(x) + \partial g(x) \quad \forall x \in H.$$

Example: The version of Fermat's Rule we saw yesterday.

MINIMIZATION OF NONSMOOTH FUNCTIONS

The proximal point algorithm

Let $f : H \rightarrow \mathbb{R} \cup \{+\infty\}$ be convex and lower-semicontinuous. Given $z_k \in H$ y $\lambda > 0$, we define z_{k+1} as the unique solution of

$$\operatorname{argmin} \left\{ f(z) + \frac{1}{2\lambda} \|z - z_k\|^2 \right\}.$$

Remark

- *It is well defined.*
- *There is no need for f to be differentiable (not even continuous).*

The proximal point algorithm

Let $f : H \rightarrow \mathbb{R} \cup \{+\infty\}$ be convex and lower-semicontinuous. Given $z_k \in H$ y $\lambda > 0$, we define z_{k+1} as the unique solution of

$$\operatorname{argmin} \left\{ f(z) + \frac{1}{2\lambda} \|z - z_k\|^2 \right\}.$$

Remark

- *It is well defined.*
- *There is no need for f to be differentiable (not even continuous).*

Relationship with the previous systems

Steepest descent dynamics

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

Gradient method:

$$\frac{x_{k+1} - x_k}{\lambda} = -\nabla f(x_k) \quad \Leftrightarrow \quad x_{k+1} = x_k - \lambda \nabla f(x_k).$$

Proximal point algorithm:

$$\frac{z_{k+1} - z_k}{\lambda} \in -\partial f(z_{k+1}) \quad \Leftrightarrow \quad z_{k+1} + \lambda \partial f(z_{k+1}) \ni z_k.$$

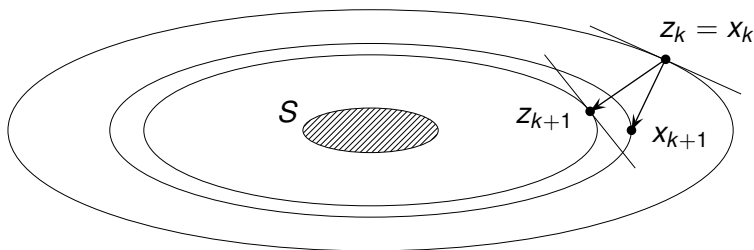
Geometric comparison

Gradient

$$x_{k+1} = x_k - \lambda \nabla f(x_k)$$

Proximal

$$z_{k+1} + \lambda \nabla f(z_{k+1}) = z_k$$



Convergence

- The sequence $(f(z_k))$ is nonincreasing.
- If f is bounded from below, then $\sum_{k \geq 0} \|z_{k+1} - z_k\|^2 < +\infty$.
- $\lim_{k \rightarrow +\infty} f(z_k) = \inf(f)$. In particular, if $z_{n_k} \rightharpoonup \bar{z}$, then $\bar{z} \in S$.
- If $S \neq \emptyset$, then $f(z_k) - \min(f) \leq \frac{\text{dist}(z_0, S)^2}{2\lambda k}$.
- For each $z^* \in S$, the sequence $(\|z_k - z^*\|)$ is nonincreasing.
- Let $S \neq \emptyset$. As $k \rightarrow +\infty$, z_k converges weakly to a point in S .

THE “SMOOTH + NONSMOOTH” STRUCTURE

We have two basic methods

Gradient method (for a smooth g):

$$x_{k+1} = \text{Grad}_{\lambda g}(x_k) = x_k - \lambda \nabla g(x_k)$$

Proximal point algorithm (for any h):

$$z_{k+1} = \text{Prox}_{\lambda h}(z_k) = \operatorname{argmin} \left\{ h(z) + \frac{1}{2\lambda} \|z - z_k\|^2 \right\}$$

The good and the bad

Gradient method

- + Easy to implement (explicit formula).
- For convergence, f is required to be smooth. The step size λ must be properly chosen.

Proximal point algorithm

- + Convergence does not rely on the differentiability of f or the step size.
- The implementation may be difficult or costly (the formula is implicit).

The good and the bad

Gradient method

- + Easy to implement (explicit formula).
- For convergence, f is required to be smooth. The step size λ must be properly chosen.

Proximal point algorithm

- + Convergence does not rely on the differentiability of f or the step size.
- The implementation may be difficult or costly (the formula is implicit).

Mixed structure

We want to minimize $f = g + h$, where g is smooth but h is not.

Example

The constrained problem

$$\min\{g(x) : x \in C\}$$

is equivalent to

$$\min\{g(x) + \delta_C(x) : x \in C\}.$$

Mixed structure

We want to minimize $f = g + h$, where g is smooth but h is not.

Example

The constrained problem

$$\min\{g(x) : x \in C\}$$

is equivalent to

$$\min\{g(x) + \delta_C(x) : x \in C\}.$$

Mixed structure

We want to minimize $f = g + h$, where g is smooth while h is not.

Example

Let $\mu > 0$, $b \in \mathbb{R}^M$ and A a matrix of size $M \times N$. Consider the functions $g, h : \mathbb{R}^N \rightarrow \mathbb{R}$ defined by

$$g(x) = \mu \sum_{n=1}^N |x_n| \quad \text{and} \quad h(x) = \frac{1}{2} \|Ax - b\|^2$$

The function $f = g + h$ appears in image processing, data analysis, among other contexts.

Strategy

Use the gradient method for the smooth part, and the proximal point algorithm for the nonsmooth one:

The proximal-gradient method:

$$x_{k+1} = T_\lambda(x_k) = \text{Prox}_{\lambda h}(\text{Grad}_{\lambda g}(x_k))$$

Case $h = 0$: gradient

Case $h = \delta_C$: projected gradient $x_{k+1} = \text{Proj}_C(\text{Grad}_{\lambda g}(x_k))$

Case $g = 0$: proximal

Strategy

Use the gradient method for the smooth part, and the proximal point algorithm for the nonsmooth one:

The proximal-gradient method:

$$x_{k+1} = T_\lambda(x_k) = \text{Prox}_{\lambda h}(\text{Grad}_{\lambda g}(x_k))$$

Case $h = 0$: gradient

Case $h = \delta_C$: projected gradient $x_{k+1} = \text{Proj}_C(\text{Grad}_{\lambda g}(x_k))$

Case $g = 0$: proximal

A fundamental property

Proposition

If ∇g is Lipschitz-continuous with constant L and $\lambda L \leq 1$, then

$$f(T_\lambda(x)) + \frac{\|y - T_\lambda(x)\|^2}{2\lambda} \leq f(y) + \frac{\|y - x\|^2}{2\lambda}.$$

Convergence

- The sequence $(f(x_k))$ is nonincreasing.
- If f is bounded from below, then $\sum_{k \geq 0} \|x_{k+1} - x_k\|^2 < +\infty$.
- $\lim_{k \rightarrow +\infty} f(x_k) = \inf(f)$. In particular, if $x_{n_k} \rightarrow \bar{x}$, then $\bar{x} \in S$.
- If $S \neq \emptyset$, then $f(x_k) - \min(f) \leq \frac{\text{dist}(x_0, S)^2}{2\lambda k}$.
- For each $x^* \in S$, the sequence $(\|x_k - x^*\|)$ is nonincreasing.
- Let $S \neq \emptyset$. As $k \rightarrow +\infty$, x_k converges weakly to a point in S .

EXERCISES

Proximal-gradient method and $\ell^1 + \ell^2$ minimization

- For $g : \mathbb{R}^N \rightarrow \mathbb{R}$ defined by $g(x) = \frac{1}{2}\|Ax - b\|^2$, compute

$$\text{Grad}_{\lambda g}(y) = y - \lambda \nabla g(y).$$

- For $h : \mathbb{R}^N \rightarrow \mathbb{R}$ defined by $h(x) = \mu\|x\|_1$, compute²

$$\text{Prox}_{\lambda h}(y) = \operatorname{argmin} \left\{ h(x) + \frac{1}{2\lambda}\|x - y\|^2 \right\}.$$

- Write down the corresponding proximal-gradient iteration $x_{k+1} = T_{\lambda}(x_k)$ explicitly.

²Study the case $N = 1$ first.