

Convex Optimization and Applications

12 - First Order Methods

Guillaume Sagnol



Nonsmooth Convex Optimization

Many convex problems that arise in machine learning / signal processing are

- Non-smooth
- Unconstrained (or constrained over a very simple set)
- When the dimension of the problem is very large, the Newton steps of interior point methods become too expensive in practice.
- → preference given to first-order algorithms, that quickly converge to a reasonably good solution.

Examples (1/2)

■ Lasso regression

$$\underset{\boldsymbol{\theta} \in \mathbb{R}^n}{\text{minimize}} \quad \|X\boldsymbol{\theta} - \mathbf{y}\|^2 + \lambda \|\boldsymbol{\theta}\|_1,$$

■ Soft-margin SVM

$$\underset{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}}{\text{minimize}} \quad \sum_{i=1}^m \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i - b)) + \lambda \|\mathbf{w}\|^2.$$

■ D-optimal design

$$\begin{aligned} \underset{\mathbf{w} \in \mathbb{R}^n}{\text{maximize}} \quad & \det^{\frac{1}{n}} \left(\sum_{i=1}^m w_i \mathbf{x}_i \mathbf{x}_i^T \right) \\ \text{s.t.} \quad & \mathbf{w} \geq \mathbf{0}, \quad \sum_i w_i = 1. \end{aligned}$$

Examples (2/2)

- Low-rank Matrix completion

$$\underset{Y \in \mathbb{R}^{m \times n}}{\text{minimize}} \quad \sum_{(i,j) \in \Omega} (X_{ij} - Y_{ij})^2 + \lambda \|Y\|_*,$$

The nuclear-norm $\|Y\|_* := \text{trace}(Y^T Y)^{\frac{1}{2}}$ serves as a convex approximation for **rank** Y .

- Total-Variation denoising

$$\underset{Y \in \mathbb{R}^{m \times n}}{\text{minimize}} \quad \|X - Y\|_F^2 + \lambda \text{TV}(Y),$$

The total-variation $\text{TV}(Y) := \sum_{i,j} \left\| \begin{bmatrix} y_{i+1,j} - y_{i,j} \\ y_{i,j+1} - y_{i,j} \end{bmatrix} \right\|_2$ penalizes the pixels with a high local variation (*noise*).

Outline

- 1 Gradient & Subgradient methods
- 2 Strong convexity and L-smoothness
- 3 The proximal operator
- 4 The proximal gradient method
- 5 The FISTA accelerated method
- 6 Optimality of accelerated gradient methods

Gradient descent

- A first order method is an algorithm to minimize a function F , that only uses first-order derivatives
- The typical algorithm is *the gradient descent* [Cauchy, mid-19th]

$$\mathbf{x}^{(k)} = \mathbf{x}^{(k-1)} - t_k \nabla F(\mathbf{x}^{(k-1)}),$$

where the stepsize t_k is selected with a line search procedure.

- Obviously, we cannot use this method for *Nonsmooth optimization*.
- Hence, the most natural idea is to use *subgradients* instead.

Subgradient

Definition (Subgradient).

The vector \mathbf{g} is a *subgradient* of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ at $\mathbf{x} \in \mathbf{dom} f$, if

$$f(\mathbf{z}) \geq f(\mathbf{x}) + \mathbf{g}^T(\mathbf{z} - \mathbf{x}), \quad \forall \mathbf{z} \in \mathbf{dom} f.$$

Geometrically, this means that the vector $[\mathbf{g}, -1]^T$ defines a supporting hyperplane to $\mathbf{epi} f$ at $(\mathbf{x}, f(\mathbf{x}))$.

The *subdifferential* of f at \mathbf{x} is the set of all subgradients:

$$\partial f(\mathbf{x}) := \{\mathbf{g} \in \mathbb{R}^n : \mathbf{g} \text{ is a subgradient of } f \text{ at } \mathbf{x}\}.$$

Properties

Proposition

- (i) $\partial f(\mathbf{x})$ is always a closed convex set.
- (ii) f convex $\implies \partial f(\mathbf{x}) \neq \emptyset, \forall \mathbf{x} \in \mathbf{int\,dom\,} f$.
- (iii) Let f be convex and $\mathbf{x} \in \mathbf{int\,dom\,} f$. Then,

$$f \text{ differentiable at } \mathbf{x} \iff \partial f(\mathbf{x}) = \{\nabla f(\mathbf{x})\}.$$

- For (i), we use that

$$\partial f(\mathbf{x}) = \bigcap_{z \in \mathbf{dom\,} f} \{\mathbf{g} : f(z) \geq f(\mathbf{x}) + \mathbf{g}^T(z - \mathbf{x})\}$$

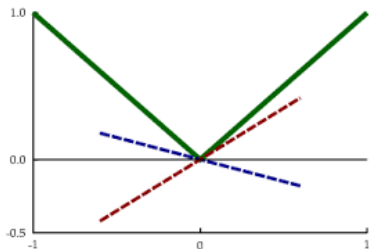
is an intersection of halfspaces.

- (ii) follows from the supporting hyperplane theorem.

Example

Let $f : \mathbb{R} \rightarrow \mathbb{R}, x \mapsto |x|$. Then, the subdifferential of f is given by

$$\partial f(x) = \begin{cases} \{-1\} & \text{if } x < 0; \\ [-1, 1] & \text{if } x = 0; \\ \{1\} & \text{if } x > 0. \end{cases}$$



Subgradient and optimality

Theorem

Let f be convex. Then, \mathbf{x}^* minimizes f over \mathbb{R}^n if and only if $\mathbf{0} \in \partial f(\mathbf{x}^*)$.

Proof:

$$\mathbf{0} \in \partial f(\mathbf{x}^*) \iff f(\mathbf{z}) \geq f(\mathbf{x}^*) + \underbrace{\mathbf{0}^T(\mathbf{z} - \mathbf{x}^*)}_{=0}, \forall \mathbf{z} \in \mathbf{dom} f.$$

Calculus rules for subdifferentials

Let f, f_1, \dots, f_m be **convex**.

■ Nonnegative scaling:

$$\partial(\alpha f)(\mathbf{x}) = \alpha \partial f(\mathbf{x}), \text{ for all } \alpha \geq 0.$$

Calculus rules for subdifferentials

Let f, f_1, \dots, f_m be **convex**.

- Nonnegative scaling:

$$\partial(\alpha f)(\mathbf{x}) = \alpha \partial f(\mathbf{x}), \text{ for all } \alpha \geq 0.$$

- Sum:

$$\partial(f_1 + \dots + f_m)(\mathbf{x}) = \partial f_1(\mathbf{x}) + \dots + \partial f_m(\mathbf{x}).$$

(Note: this is a Minkowski sum of convex sets).

Calculus rules for subdifferentials

Let f, f_1, \dots, f_m be **convex**.

- **Nonnegative scaling:**

$$\partial(\alpha f)(\mathbf{x}) = \alpha \partial f(\mathbf{x}), \text{ for all } \alpha \geq 0.$$

- **Sum:**

$$\partial(f_1 + \dots + f_m)(\mathbf{x}) = \partial f_1(\mathbf{x}) + \dots + \partial f_m(\mathbf{x}).$$

(Note: this is a Minkowski sum of convex sets).

- **Affine transformation:**

$$\partial(z \mapsto f(Az + b))(\mathbf{x}) = A^T \partial f(A\mathbf{x} + b).$$

Calculus rules for subdifferentials

Let f, f_1, \dots, f_m be **convex**.

- **Nonnegative scaling:**

$$\partial(\alpha f)(x) = \alpha \partial f(x), \quad \text{for all } \alpha \geq 0.$$

- **Sum:**

$$\partial(f_1 + \dots + f_m)(x) = \partial f_1(x) + \dots + \partial f_m(x).$$

(Note: this is a Minkowski sum of convex sets).

- **Affine transformation:**

$$\partial(z \mapsto f(Az + b))(x) = A^T \partial f(Ax + b).$$

- **Pointwise maximum:**

Let $g(x) = \max_{i=1, \dots, m} f_i(x)$. Then, $\partial g(x) = \mathbf{conv} \left(\bigcup_{j \in A(x)} \partial f_j(x) \right)$,

where $A(x)$ is the set of *active functions* at x , i.e.,

$$A(x) := \{j \in [m] : f_j(x) = g(x)\}.$$

(can be extended to pointwise supremums of infinitely many functions under additional technical conditions).

The subgradient method

To minimize a non-smooth convex function F over \mathbb{R}^n , we can use the subgradient method:

$$\mathbf{x}^{(k)} = \mathbf{x}^{(k-1)} - t_k \mathbf{g}^{(k-1)}, \quad \text{for some } \mathbf{g}^{(k-1)} \in \partial F(\mathbf{x}^{(k-1)}).$$

A few properties of this algorithm:

- Not a *descent method* (we can have $F(\mathbf{x}^{(k)}) > F(\mathbf{x}^{(k-1)})$, even for arbitrarily small step sizes.)

The subgradient method

To minimize a non-smooth convex function F over \mathbb{R}^n , we can use the subgradient method:

$$\mathbf{x}^{(k)} = \mathbf{x}^{(k-1)} - t_k \mathbf{g}^{(k-1)}, \quad \text{for some } \mathbf{g}^{(k-1)} \in \partial F(\mathbf{x}^{(k-1)}).$$

A few properties of this algorithm:

- Not a *descent method* (we can have $F(\mathbf{x}^{(k)}) > F(\mathbf{x}^{(k-1)})$, even for arbitrarily small step sizes.)
- Method can fail to converge if we use exact or backtracking line search.

The subgradient method

To minimize a non-smooth convex function F over \mathbb{R}^n , we can use the subgradient method:

$$\mathbf{x}^{(k)} = \mathbf{x}^{(k-1)} - t_k \mathbf{g}^{(k-1)}, \quad \text{for some } \mathbf{g}^{(k-1)} \in \partial F(\mathbf{x}^{(k-1)}).$$

A few properties of this algorithm:

- Not a *descent method* (we can have $F(\mathbf{x}^{(k)}) > F(\mathbf{x}^{(k-1)})$, even for arbitrarily small step sizes.)
- Method can fail to converge if we use exact or backtracking line search.
- Convergence can be proved for some *offline rules*, e.g.
 - Constant step sizes ($t_k = t > 0, \forall k \in \mathbb{N}$).
 - Nonsummable diminishing ($t_k \rightarrow 0, \sum_{k \in \mathbb{N}} t_k = \infty$)

The subgradient method

To minimize a non-smooth convex function F over \mathbb{R}^n , we can use the subgradient method:

$$\mathbf{x}^{(k)} = \mathbf{x}^{(k-1)} - t_k \mathbf{g}^{(k-1)}, \quad \text{for some } \mathbf{g}^{(k-1)} \in \partial F(\mathbf{x}^{(k-1)}).$$

A few properties of this algorithm:

- Not a *descent method* (we can have $F(\mathbf{x}^{(k)}) > F(\mathbf{x}^{(k-1)})$, even for arbitrarily small step sizes.)
- Method can fail to converge if we use exact or backtracking line search.
- Convergence can be proved for some *offline rules*, e.g.
 - Constant step sizes ($t_k = t > 0, \forall k \in \mathbb{N}$).
 - Nonsummable diminishing ($t_k \rightarrow 0, \sum_{k \in \mathbb{N}} t_k = \infty$)
- Convergence typically *slow*: after k iteration, the best iterate seen so far satisfies $f(\mathbf{x}_{\text{best}}^{(k)}) \leq f(\mathbf{x}^*) + O\left(\frac{1}{\sqrt{k}}\right)$.

Outline

- 1 Gradient & Subgradient methods
- 2 Strong convexity and L-smoothness**
- 3 The proximal operator
- 4 The proximal gradient method
- 5 The FISTA accelerated method
- 6 Optimality of accelerated gradient methods

Strong convexity

Definition (ν -strong convexity).

f is ν -strongly convex for some $\nu > 0$ iff $\mathbf{x} \mapsto f(\mathbf{x}) - \frac{\nu}{2}\|\mathbf{x}\|^2$ is convex.

Remark: If f is twice diff., then f is ν -strongly convex iff

$$\nabla^2 f(\mathbf{x}) \succeq \nu I, \quad \forall \mathbf{x} \in \mathbf{dom} f.$$

Strong convexity

Definition (ν -strong convexity).

f is ν -strongly convex for some $\nu > 0$ iff $x \mapsto f(x) - \frac{\nu}{2}\|x\|^2$ is convex.

Remark: If f is twice diff., then f is ν -strongly convex iff

$$\nabla^2 f(x) \succeq \nu I, \quad \forall x \in \mathbf{dom} f.$$

Proposition

Let f be ν -strongly convex. Then, $\forall x_0 \in \mathbf{dom} \partial f, \forall g \in \partial f(x_0)$,

$$f(x) \geq f(x_0) + \langle g, x - x_0 \rangle + \frac{\nu}{2}\|x - x_0\|^2, \quad \forall x \in \mathbf{dom} f.$$

Remark: The converse statement is also true.

Strong convexity

Proposition

Let f be ν -strongly convex. Then, $\forall \mathbf{x}_0 \in \mathbf{dom} \partial f, \forall \mathbf{g} \in \partial f(\mathbf{x}_0)$,

$$f(\mathbf{x}) \geq f(\mathbf{x}_0) + \langle \mathbf{g}, \mathbf{x} - \mathbf{x}_0 \rangle + \frac{\nu}{2} \|\mathbf{x} - \mathbf{x}_0\|^2, \quad \forall \mathbf{x} \in \mathbf{dom} f.$$

Proof:

- Let $F(\mathbf{x}) := f(\mathbf{x}) - \frac{\nu}{2} \|\mathbf{x}\|^2$; this is a convex function.
- Rule for sum of subdifferentials of convex functions:

$$\mathbf{g} \in \partial f(\mathbf{x}_0) = \partial F(\mathbf{x}_0) + \frac{\nu}{2} \nabla(\mathbf{x} \mapsto \|\mathbf{x}\|^2) = \partial F(\mathbf{x}_0) + \nu \mathbf{x}_0$$

- So $\mathbf{g} - \nu \mathbf{x}_0$ is a subgradient of F at \mathbf{x}_0 : $\forall \mathbf{x}_0 \in \mathbf{dom} f$,

$$f(\mathbf{x}) - \nu/2 \|\mathbf{x}\|^2 \geq f(\mathbf{x}_0) - \nu/2 \|\mathbf{x}_0\|^2 + \langle \mathbf{g} - \nu \mathbf{x}_0, \mathbf{x} - \mathbf{x}_0 \rangle.$$

- Re-arranging yields the proposition.

Minimizer of strong convex function

Theorem

Let f be a closed, ν -strongly convex. Then f has a unique minimizer \mathbf{x}^* , and

$$f(\mathbf{x}) \geq f(\mathbf{x}^*) + \frac{\nu}{2} \|\mathbf{x} - \mathbf{x}^*\|^2, \quad \forall \mathbf{x} \in \mathbf{dom} f.$$

Proof:

See blackboard.

L -smoothness

Definition (L -smoothness).

A differentiable function f is called L -smooth for some $L > 0$ if its gradient is L -Lipschitz:

$$\|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\| \leq L\|\mathbf{x} - \mathbf{y}\|, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbf{dom} f,$$

Remark: If f is twice diff., then f is L -smooth iff

$$\nabla^2 f(\mathbf{x}) \preceq LI, \quad \forall \mathbf{x} \in \mathbf{dom} f.$$

L-smoothness

Proposition

For a differentiable function f , consider the following statements:

- (i) f is L -smooth (i.e., ∇f is L -Lipschitz)
- (ii) $f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2, \quad \forall \mathbf{x}, \mathbf{y} \in \mathbf{dom} f.$
- (iii) $\mathbf{x} \mapsto \frac{L}{2} \|\mathbf{x}\|^2 - f(\mathbf{x})$ is convex

It holds: (i) \implies (ii) \iff (iii).

If moreover f is convex, then (i) \iff (ii) \iff (iii).

We prove (i) \implies (ii) on the next slide.

L -smoothness

- Let f be L -smooth, $\mathbf{x}, \mathbf{y} \in \mathbf{dom} f$.
- From the fundamental theorem of calculus,

$$\begin{aligned} f(\mathbf{y}) &= f(\mathbf{x}) + \int_{t=0}^1 \langle \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})), \mathbf{y} - \mathbf{x} \rangle dt. \\ &= f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \int_{t=0}^1 \langle \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle dt. \end{aligned}$$

L-smoothness

- Let f be L -smooth, $\mathbf{x}, \mathbf{y} \in \text{dom } f$.
- From the fundamental theorem of calculus,

$$\begin{aligned} f(\mathbf{y}) &= f(\mathbf{x}) + \int_{t=0}^1 \langle \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})), \mathbf{y} - \mathbf{x} \rangle dt. \\ &= f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \int_{t=0}^1 \langle \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle dt. \end{aligned}$$

- Hence,

$$\begin{aligned} |f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle| &= \left| \int_{t=0}^1 \langle \nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle dt \right|. \\ &\stackrel{\text{[Cauchy-Schwartz]}}{\leq} \int_{t=0}^1 \|\nabla f(\mathbf{x} + t(\mathbf{y} - \mathbf{x})) - \nabla f(\mathbf{x})\| \cdot \|\mathbf{y} - \mathbf{x}\| dt \\ &\stackrel{\text{[L-smoothness]}}{\leq} \int_{t=0}^1 Lt \|\mathbf{y} - \mathbf{x}\| \cdot \|\mathbf{y} - \mathbf{x}\| dt = \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2. \end{aligned}$$

Outline

- 1 Gradient & Subgradient methods
- 2 Strong convexity and L-smoothness
- 3 The proximal operator**
- 4 The proximal gradient method
- 5 The FISTA accelerated method
- 6 Optimality of accelerated gradient methods

Prox operator

Definition (Prox operator).

Let $g : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ be a closed convex function. We define the proximal mapping of g by

$$\mathbf{prox}_g(\mathbf{x}) := \operatorname{argmin}_{\mathbf{u} \in \mathbb{R}^n} g(\mathbf{u}) + \frac{1}{2} \|\mathbf{x} - \mathbf{u}\|^2.$$

The proximal operator generalizes the notion of *projection*:

- If C is convex set, define the convex indicator function

$$I_C(\mathbf{x}) = \begin{cases} 0 & \text{if } \mathbf{x} \in C \\ \infty & \text{otherwise.} \end{cases}$$

- Then,

$$P_C(\mathbf{x}) = \operatorname{argmin}_{\mathbf{u} \in C} \|\mathbf{u} - \mathbf{x}\|^2 = \operatorname{argmin}_{\mathbf{u}} \|\mathbf{u} - \mathbf{x}\|^2 + I_C(\mathbf{u}) = \mathbf{prox}_{I_C}(\mathbf{x}).$$

Properties of prox

Theorem

Let $g : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ be a closed convex function. Then,

- (i) $\mathbf{prox}_g(x) \in \mathbb{R}^n$ is well defined over $\mathbf{dom} g$.
(i.e., for all $x \in \mathbf{dom} g$, there is a unique minimizer).
- (ii) $u = \mathbf{prox}_g(x) \iff x - u \in \partial g(u)$
- (iii) x^* is a minimizer of $g \iff x^* = \mathbf{prox}_g(x^*)$.

Properties of prox

Theorem

Let $g : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ be a closed convex function. Then,

- (i) $\mathbf{prox}_g(x) \in \mathbb{R}^n$ is well defined over $\mathbf{dom} g$.
(i.e., for all $x \in \mathbf{dom} g$, there is a unique minimizer).
- (ii) $u = \mathbf{prox}_g(x) \iff x - u \in \partial g(u)$
- (iii) x^* is a minimizer of $g \iff x^* = \mathbf{prox}_g(x^*)$.

Proof:

(i) Let $h(u) := g(u) + \frac{1}{2}\|x - u\|^2$. Then,

$$h(u) - \frac{1}{2}\|u\|^2 = g(u) + \frac{1}{2}\|x - u\|^2 - \frac{1}{2}\|u\|^2 = g(u) + \frac{1}{2}\|x\|^2 - x^T u$$

is convex. So g is strongly convex with parameter $\nu = 1$ and has a single minimizer.

Properties of prox

Theorem

Let $g : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ be a closed convex function. Then,

- (i) $\text{prox}_g(x) \in \mathbb{R}^n$ is well defined over $\text{dom } g$.
(i.e., for all $x \in \text{dom } g$, there is a unique minimizer).
- (ii) $u = \text{prox}_g(x) \iff x - u \in \partial g(u)$
- (iii) x^* is a minimizer of $g \iff x^* = \text{prox}_g(x^*)$.

Proof:

(ii) Let $h(u) := g(u) + \frac{1}{2}\|u - x\|^2$. The subdifferential of h is $\partial h(u) = \partial g(u) + u - x$. So u minimizes h iff

$$0 \in \partial g(u) + u - x \iff x - u \in \partial g(u)$$

Properties of prox

Theorem

Let $g : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ be a closed convex function. Then,

- (i) $\mathbf{prox}_g(x) \in \mathbb{R}^n$ is well defined over $\mathbf{dom} g$.
(i.e., for all $x \in \mathbf{dom} g$, there is a unique minimizer).
- (ii) $u = \mathbf{prox}_g(x) \iff x - u \in \partial g(u)$
- (iii) x^* is a minimizer of $g \iff x^* = \mathbf{prox}_g(x^*)$.

Proof:

(ii) Let $h(u) := g(u) + \frac{1}{2}\|u - x\|^2$. The subdifferential of h is $\partial h(u) = \partial g(u) + u - x$. So u minimizes h iff

$$0 \in \partial g(u) + u - x \iff x - u \in \partial g(u)$$

(iii) $x = \mathbf{prox}_g(x) \iff x - x = 0 \in \partial g(x) \iff x \in \mathbf{argmin} g$.

Computing the prox

- In general, computing $\mathbf{prox}_g(\mathbf{x})$ can be as hard as minimizing g ...
- Good news: for many functions, the prox operator can be computed efficiently, i.e. in $O(n)$ or $O(n \log(n))$, by using a closed-form formula, or by reducing to a one-dimensional problem.
- A catalog of known prox. operators can be found at <http://proximity-operator.net>
- Simple rule for separable sums:

$$\text{if } f(\mathbf{x}_1, \dots, \mathbf{x}_n) = \sum_i f_i(\mathbf{x}_i), \text{ then } \mathbf{prox}_f(\mathbf{x}) = \begin{bmatrix} \mathbf{prox}_{f_1}(\mathbf{x}_1) \\ \vdots \\ \mathbf{prox}_{f_n}(\mathbf{x}_n) \end{bmatrix}$$

Example of proximal operators

We usually need to compute the proximal operator of a scaling $t \cdot g$ of a convex function g , for some $t > 0$.

$g(x)$	$\text{prox}_{tg}(x)$
$\ x\ _2$	$\left(1 - \frac{t}{\max(\ x\ _2, t)}\right) x$
$x^T Q x + p^T x$	$(tQ + I)^{-1}(x - tp)$
$\ x\ _1$	$[x - t\mathbf{1}]_+ \odot \text{sign}(x)$
$\sum_{i=1}^n [x_i]_+$	$\left[x - \frac{t}{2}\mathbf{1} - \frac{t}{2}\mathbf{1}\right]_+ \odot \text{sign}(x)$
$\sum_{i=1}^n x_i \log(x_i)$	$t W(t^{-1} e^{\frac{x}{t}-1})$
$\max_{i=1, \dots, n} x_i$	$\{\min(x_i, s)\}_{i=1, \dots, n}$ where s solves $\sum_i [x_i - s]_+ = t$

Example: Prox of ℓ_1 -norm

- $f(\mathbf{x}) = t\|\mathbf{x}\|_1 = \sum_i t|x_i|$ is a separable sum, so we can focus on the 1-dimensional function $g : \mathbb{R} \rightarrow \mathbb{R}, x \rightarrow t|x|$ and apply the prox *elementwise*.

Example: Prox of ℓ_1 -norm

- $f(x) = t\|x\|_1 = \sum_i t|x_i|$ is a separable sum, so we can

focus on the 1-dimensional function $g : \mathbb{R} \rightarrow \mathbb{R}, x \rightarrow t|x|$ and apply the prox *elementwise*.

- $u^* = \mathbf{prox}_{x \mapsto t|x|}(x) \iff x - u^* \in \partial g(u^*) \iff$
 $(x - u^* = -t \wedge u^* < 0) \text{ or } (x - u^* \in [-t, t] \wedge u^* = 0) \text{ or } (x - u^* = t \wedge u^* > 0).$

Example: Prox of ℓ_1 -norm

- $f(x) = t\|x\|_1 = \sum_i t|x_i|$ is a separable sum, so we can

focus on the 1-dimensional function $g : \mathbb{R} \rightarrow \mathbb{R}, x \rightarrow t|x|$ and apply the prox *elementwise*.

- $u^* = \mathbf{prox}_{x \mapsto t|x|}(x) \iff x - u^* \in \partial g(u^*) \iff$
 $(x - u^* = -t \wedge u^* < 0)$ or $(x - u^* \in [-t, t] \wedge u^* = 0)$ or $(x - u^* = t \wedge u^* > 0)$.

- We can solve this system by analyzing the sign of u^* :

$$\mathbf{prox}_{x \mapsto t|x|}(x) = u^* := \begin{cases} x + t & \text{if } x < -t \\ 0 & \text{if } x \in [-t, t] \\ x - t & \text{if } x > t \end{cases} = [|x| - t]_+ \text{sign}(x),$$

Example: Prox of ℓ_1 -norm

- $f(x) = t\|x\|_1 = \sum_i t|x_i|$ is a separable sum, so we can

focus on the 1-dimensional function $g : \mathbb{R} \rightarrow \mathbb{R}, x \rightarrow t|x|$ and apply the prox *elementwise*.

- $u^* = \mathbf{prox}_{x \mapsto t|x|}(x) \iff x - u^* \in \partial g(u^*) \iff$
 $(x - u^* = -t \wedge u^* < 0)$ or $(x - u^* \in [-t, t] \wedge u^* = 0)$ or $(x - u^* = t \wedge u^* > 0)$.

- We can solve this system by analyzing the sign of u^* :

$$\mathbf{prox}_{x \mapsto t|x|}(x) = u^* := \begin{cases} x + t & \text{if } x < -t \\ 0 & \text{if } x \in [-t, t] \\ x - t & \text{if } x > t \end{cases} = [|x| - t]_+ \text{sign}(x),$$

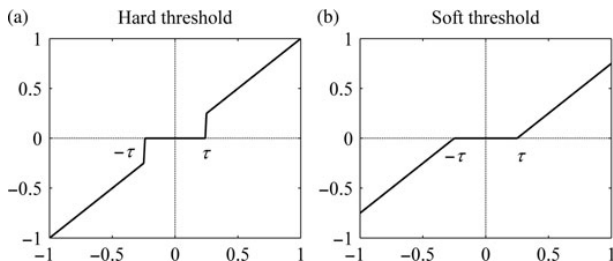
- Finally, apply the formula *componentwise*:

$$\mathbf{prox}_{t\mathbf{f}}(x) = \mathcal{T}_t(x) := [|x| - t\mathbf{1}]_+ \odot \text{sign}(x).$$

Prox of ℓ_1 -norm: Soft thresholding

The proximal operator $\mathcal{T}_t(x)$ of $x \mapsto t\|x\|_1$ is called the *soft thresholding operator* (a level t).

$\mathcal{T}_\tau(x)$ acts on each coordinate as a thresholding operator that zeroes values $|x| < \tau$, but the function is shifted to make it continuous:



Outline

- 1 Gradient & Subgradient methods
- 2 Strong convexity and L-smoothness
- 3 The proximal operator
- 4 The proximal gradient method**
- 5 The FISTA accelerated method
- 6 Optimality of accelerated gradient methods

Composite model

Definition

From now on we consider a *Composite convex model*

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad F(x) := f(x) + g(x),$$

where

- $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ is convex and L -smooth.
- $g : \mathbb{R}^n \rightarrow \mathbb{R}$ is closed, convex, possibly nonsmooth, but it has a *cheap proximal operator*.

Composite model

Definition

From now on we consider a *Composite convex model*

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad F(x) := f(x) + g(x),$$

where

- $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ is convex and L -smooth.
- $g : \mathbb{R}^n \rightarrow \mathbb{R}$ is closed, convex, possibly nonsmooth, but it has a *cheap proximal operator*.

Special cases

- $g = 0$: Smooth convex, unconstrained optimization
- $g = I_C$: Minimization of f over the convex set C (for a “simple” convex set C such that the projection over C ($= \text{prox}_{I_C}$) can be computed easily).

Basic idea

- $F(x) := f(x) + g(x)$, with f L -smooth, g proximable.

Basic idea

- $F(x) := f(x) + g(x)$, with f L -smooth, g proximable.
- f is L -smooth, so

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

Basic idea

- $F(x) := f(x) + g(x)$, with f L -smooth, g proximable.
- f is L -smooth, so

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

- At iteration k , we use this to obtain an overestimator of F around the current iterate $x^{(k)}$: Given $0 < t_k \leq \frac{1}{L}$,

$$F(y) \leq \hat{F}(y) := f(x^{(k)}) + \langle \nabla f(x^{(k)}), y - x^{(k)} \rangle + \frac{1}{2t_k} \|y - x^{(k)}\|^2 + g(y)$$

Basic idea

- $F(x) := f(x) + g(x)$, with f L -smooth, g proximable.
- f is L -smooth, so

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

- At iteration k , we use this to obtain an overestimator of F around the current iterate $x^{(k)}$: Given $0 < t_k \leq \frac{1}{L}$,

$$F(y) \leq \hat{F}(y) := f(x^{(k)}) + \langle \nabla f(x^{(k)}), y - x^{(k)} \rangle + \frac{1}{2t_k} \|y - x^{(k)}\|^2 + g(y)$$

- The next iterate is determined by computing

$$x^{(k+1)} := \underset{y}{\operatorname{argmin}} \hat{F}(y).$$

This reduces to evaluating the $\operatorname{prox}_{t_k g}$ operator !

Proximal Gradient iteration

$$\begin{aligned} \mathbf{x}^{(k+1)} &:= \operatorname{argmin}_y \hat{F}(y) \\ &= \operatorname{argmin}_y g(y) + \mathbf{y}^T \nabla f(\mathbf{x}^{(k)}) + \frac{1}{2t_k} \|\mathbf{y} - \mathbf{x}^{(k)}\|^2 \\ &= \operatorname{argmin}_y g(y) + \mathbf{y}^T \nabla f(\mathbf{x}^{(k)}) + \frac{1}{2t_k} (\|\mathbf{y}\|^2 - 2\mathbf{y}^T \mathbf{x}^{(k)}) \\ &= \operatorname{argmin}_y t_k g(y) + \frac{1}{2} \|\mathbf{y}\|^2 - \mathbf{y}^T (\mathbf{x}^{(k)} - t_k \nabla f(\mathbf{x}^{(k)})) \\ &= \operatorname{argmin}_y t_k g(y) + \frac{1}{2} \|\mathbf{y} - (\mathbf{x}^{(k)} - t_k \nabla f(\mathbf{x}^{(k)}))\|^2 \\ &= \mathbf{prox}_{t_k g}(\mathbf{x}^{(k)} - t_k \nabla f(\mathbf{x}^{(k)})) \end{aligned}$$

Proximal Gradient iteration

$$\begin{aligned} \mathbf{x}^{(k+1)} &:= \operatorname{argmin}_y \hat{F}(y) \\ &= \operatorname{argmin}_y g(y) + \mathbf{y}^T \nabla f(\mathbf{x}^{(k)}) + \frac{1}{2t_k} \|\mathbf{y} - \mathbf{x}^{(k)}\|^2 \\ &= \operatorname{argmin}_y g(y) + \mathbf{y}^T \nabla f(\mathbf{x}^{(k)}) + \frac{1}{2t_k} (\|\mathbf{y}\|^2 - 2\mathbf{y}^T \mathbf{x}^{(k)}) \\ &= \operatorname{argmin}_y t_k g(y) + \frac{1}{2} \|\mathbf{y}\|^2 - \mathbf{y}^T (\mathbf{x}^{(k)} - t_k \nabla f(\mathbf{x}^{(k)})) \\ &= \operatorname{argmin}_y t_k g(y) + \frac{1}{2} \|\mathbf{y} - (\mathbf{x}^{(k)} - t_k \nabla f(\mathbf{x}^{(k)}))\|^2 \\ &= \mathbf{prox}_{t_k g}(\mathbf{x}^{(k)} - t_k \nabla f(\mathbf{x}^{(k)})) \end{aligned}$$

Definition Proximal Gradient iteration

For some step size $t > 0$, update $\mathbf{x}^+ \leftarrow \mathbf{prox}_{tg}(\mathbf{x} - t\nabla f(\mathbf{x}))$.

Analysis of the proximal gradient method

Theorem (Prox-grad inequality).

Let $\mathbf{x} \in \text{int dom } f$ denote the current iterate, and \mathbf{x}^+ be the next iterate, obtained after a step of size $t > 0$, i.e.,

$$\mathbf{x}^+ = \mathbf{prox}_{tg}(\mathbf{x} - t\nabla f(\mathbf{x})).$$

If

$$f(\mathbf{x}^+) \leq f(\mathbf{x}) + \nabla f(\mathbf{x})^T(\mathbf{x}^+ - \mathbf{x}) + \frac{1}{2t}\|\mathbf{x}^+ - \mathbf{x}\|^2,$$

(so in particular, if $t \leq \frac{1}{L}$), then for all $\xi \in \mathbb{R}^n$ it holds:

$$F(\xi) - F(\mathbf{x}^+) \geq \frac{1}{2t}(\|\xi - \mathbf{x}^+\|^2 - \|\xi - \mathbf{x}\|^2).$$

Proof of the prox-grad inequality (1/2)

- $\mathbf{x}^+ = \mathbf{prox}_{tg}(\mathbf{x} - t\nabla f(\mathbf{x})) \iff \mathbf{x} - t\nabla f(\mathbf{x}) - \mathbf{x}^+ \in \partial(tg)(\mathbf{x}^+)$
- So, by definition of a subgradient, for all $\xi \in \mathbb{R}^n$,

$$tg(\xi) \geq tg(\mathbf{x}^+) + \langle \mathbf{x} - t\nabla f(\mathbf{x}) - \mathbf{x}^+, \xi - \mathbf{x}^+ \rangle$$
$$\iff g(\xi) - g(\mathbf{x}^+) \geq \frac{1}{t} \langle \mathbf{x} - \mathbf{x}^+, \xi - \mathbf{x}^+ \rangle - \nabla f(\mathbf{x})^T (\xi - \mathbf{x}^+).$$

Proof of the prox-grad inequality (1/2)

- $\mathbf{x}^+ = \mathbf{prox}_{tg}(\mathbf{x} - t\nabla f(\mathbf{x})) \iff \mathbf{x} - t\nabla f(\mathbf{x}) - \mathbf{x}^+ \in \partial(tg)(\mathbf{x}^+)$
- So, by definition of a subgradient, for all $\xi \in \mathbb{R}^n$,

$$tg(\xi) \geq tg(\mathbf{x}^+) + \langle \mathbf{x} - t\nabla f(\mathbf{x}) - \mathbf{x}^+, \xi - \mathbf{x}^+ \rangle$$
$$\iff g(\xi) - g(\mathbf{x}^+) \geq \frac{1}{t} \langle \mathbf{x} - \mathbf{x}^+, \xi - \mathbf{x}^+ \rangle - \nabla f(\mathbf{x})^T (\xi - \mathbf{x}^+).$$

- And our assumption on the stepsize t can be rewritten as:

$$f(\xi) - f(\mathbf{x}^+) \geq f(\xi) - f(\mathbf{x}) - \nabla f(\mathbf{x})^T (\mathbf{x}^+ - \mathbf{x}) - \frac{1}{2t} \|\mathbf{x}^+ - \mathbf{x}\|^2$$

Proof of the prox-grad inequality (1/2)

- $x^+ = \mathbf{prox}_{tg}(x - t\nabla f(x)) \iff x - t\nabla f(x) - x^+ \in \partial(tg)(x^+)$
- So, by definition of a subgradient, for all $\xi \in \mathbb{R}^n$,

$$tg(\xi) \geq tg(x^+) + \langle x - t\nabla f(x) - x^+, \xi - x^+ \rangle$$
$$\iff g(\xi) - g(x^+) \geq \frac{1}{t} \langle x - x^+, \xi - x^+ \rangle - \nabla f(x)^T (\xi - x^+).$$

- And our assumption on the stepsize t can be rewritten as:

$$f(\xi) - f(x^+) \geq f(\xi) - f(x) - \nabla f(x)^T (x^+ - x) - \frac{1}{2t} \|x^+ - x\|^2$$

- We sum the above two inequalities:

$$F(\xi) - F(x^+) \geq \underbrace{f(\xi) - f(x) - \nabla f(x)^T (\xi - x)}_{\epsilon_f(x, \xi) \geq 0} + \frac{1}{t} \langle x - x^+, \xi - x^+ \rangle - \frac{1}{2t} \|x^+ - x\|^2.$$

Proof of the prox-grad inequality (2/2)

$$F(\xi) - F(\mathbf{x}^+) \geq \underbrace{f(\xi) - f(\mathbf{x}) - \nabla f(\mathbf{x})^T (\xi - \mathbf{x})}_{\epsilon_f(\mathbf{x}, \xi) \geq 0} + \frac{1}{t} \langle \mathbf{x} - \mathbf{x}^+, \xi - \mathbf{x}^+ \rangle - \frac{1}{2t} \|\mathbf{x}^+ - \mathbf{x}\|^2.$$

So,

$$F(\xi) - F(\mathbf{x}^+) \geq \frac{1}{2t} (2 \langle \mathbf{x} - \mathbf{x}^+, \xi - \mathbf{x}^+ \rangle - \|\mathbf{x}^+ - \mathbf{x}\|^2).$$

Proof of the prox-grad inequality (2/2)

$$F(\xi) - F(\mathbf{x}^+) \geq \underbrace{f(\xi) - f(\mathbf{x}) - \nabla f(\mathbf{x})^T (\xi - \mathbf{x})}_{\epsilon_f(\mathbf{x}, \xi) \geq 0} + \frac{1}{t} \langle \mathbf{x} - \mathbf{x}^+, \xi - \mathbf{x}^+ \rangle - \frac{1}{2t} \|\mathbf{x}^+ - \mathbf{x}\|^2.$$

So,

$$F(\xi) - F(\mathbf{x}^+) \geq \frac{1}{2t} (2 \langle \mathbf{x} - \mathbf{x}^+, \xi - \mathbf{x}^+ \rangle - \|\mathbf{x}^+ - \mathbf{x}\|^2).$$

Finally, we use the identity

$$\begin{aligned} \|\xi - \mathbf{x}\|^2 &= \|(\xi - \mathbf{x}^+) - (\mathbf{x} - \mathbf{x}^+)\|^2 \\ &= \|\xi - \mathbf{x}^+\|^2 + \|\mathbf{x} - \mathbf{x}^+\|^2 - 2 \langle \xi - \mathbf{x}^+, \mathbf{x} - \mathbf{x}^+ \rangle, \end{aligned}$$

Proof of the prox-grad inequality (2/2)

$$F(\xi) - F(\mathbf{x}^+) \geq \underbrace{f(\xi) - f(\mathbf{x}) - \nabla f(\mathbf{x})^T (\xi - \mathbf{x})}_{\epsilon_f(\mathbf{x}, \xi) \geq 0} + \frac{1}{t} \langle \mathbf{x} - \mathbf{x}^+, \xi - \mathbf{x}^+ \rangle - \frac{1}{2t} \|\mathbf{x}^+ - \mathbf{x}\|^2.$$

So,

$$F(\xi) - F(\mathbf{x}^+) \geq \frac{1}{2t} (2 \langle \mathbf{x} - \mathbf{x}^+, \xi - \mathbf{x}^+ \rangle - \|\mathbf{x}^+ - \mathbf{x}\|^2).$$

Finally, we use the identity

$$\begin{aligned} \|\xi - \mathbf{x}\|^2 &= \|(\xi - \mathbf{x}^+) - (\mathbf{x} - \mathbf{x}^+)\|^2 \\ &= \|\xi - \mathbf{x}^+\|^2 + \|\mathbf{x} - \mathbf{x}^+\|^2 - 2 \langle \xi - \mathbf{x}^+, \mathbf{x} - \mathbf{x}^+ \rangle, \end{aligned}$$

and we obtain the result:

$$F(\xi) - F(\mathbf{x}^+) \geq \frac{1}{2t} (\|\xi - \mathbf{x}^+\|^2 - \|\xi - \mathbf{x}\|^2).$$

Sufficient decrease

Corollary

If the step size t is “well chosen” (i.e., it satisfies the condition of the previous theorem), then

$$F(\mathbf{x}) - F(\mathbf{x}^+) \geq \frac{1}{2t} \|\mathbf{x} - \mathbf{x}^+\|^2.$$

In particular, the proximal gradient method is a *descent method*.

Convergence Analysis

- We assume that $\mathbf{x}^{(0)} \in \text{int dom } f$ and constant step sizes $t_k = \frac{1}{L}$ are used:

$$\mathbf{x}^{(k+1)} := \text{prox}_{\frac{1}{L}g}(\mathbf{x}^{(k)} - \frac{1}{L}\nabla f(\mathbf{x}^{(k)})).$$

- If L is unknown, one can use backtracking line search to find a step size t_k that satisfies the condition of the previous theorem – then, similar analysis.

Theorem

For any optimal solution \mathbf{x}^* of the composite convex optimization problem (**minimize** $f + g$),

$$F(\mathbf{x}^{(k)}) - F(\mathbf{x}^*) \leq \frac{L}{2k} \|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2, \quad \forall k \geq 1.$$

Convergence Analysis: Proof

- Prox-grad inequality at $\xi = \mathbf{x}^*$:

$$F(\mathbf{x}^*) - F(\mathbf{x}^{(i+1)}) \geq \frac{L}{2} (\|\mathbf{x}^* - \mathbf{x}^{(i+1)}\|^2 - \|\mathbf{x}^* - \mathbf{x}^{(i)}\|^2).$$

Convergence Analysis: Proof

- Prox-grad inequality at $\xi = \mathbf{x}^*$:

$$F(\mathbf{x}^*) - F(\mathbf{x}^{(i+1)}) \geq \frac{L}{2} (\|\mathbf{x}^* - \mathbf{x}^{(i+1)}\|^2 - \|\mathbf{x}^* - \mathbf{x}^{(i)}\|^2).$$

- Summing over $i = 0, \dots, k-1$,

$$k F(\mathbf{x}^*) - \sum_{i=0}^{k-1} F(\mathbf{x}^{(i+1)}) \geq \frac{L}{2} (\|\mathbf{x}^* - \mathbf{x}^{(k)}\|^2 - \|\mathbf{x}^* - \mathbf{x}^{(0)}\|^2)$$

$$\implies \sum_{i=1}^k F(\mathbf{x}^{(i)}) - k F(\mathbf{x}^*) \leq \frac{L}{2} \|\mathbf{x}^* - \mathbf{x}^{(0)}\|^2.$$

Convergence Analysis: Proof

- Prox-grad inequality at $\xi = \mathbf{x}^*$:

$$F(\mathbf{x}^*) - F(\mathbf{x}^{(i+1)}) \geq \frac{L}{2} (\|\mathbf{x}^* - \mathbf{x}^{(i+1)}\|^2 - \|\mathbf{x}^* - \mathbf{x}^{(i)}\|^2).$$

- Summing over $i = 0, \dots, k-1$,

$$\begin{aligned} k F(\mathbf{x}^*) - \sum_{i=0}^{k-1} F(\mathbf{x}^{(i+1)}) &\geq \frac{L}{2} (\|\mathbf{x}^* - \mathbf{x}^{(k)}\|^2 - \|\mathbf{x}^* - \mathbf{x}^{(0)}\|^2) \\ \implies \sum_{i=1}^k F(\mathbf{x}^{(i)}) - k F(\mathbf{x}^*) &\leq \frac{L}{2} \|\mathbf{x}^* - \mathbf{x}^{(0)}\|^2. \end{aligned}$$

- Since the algorithm is a descent method, $\sum_{i=1}^k F(\mathbf{x}^{(i)}) \geq k F(\mathbf{x}^{(k)})$.

Convergence Analysis: Proof

- Prox-grad inequality at $\xi = \mathbf{x}^*$:

$$F(\mathbf{x}^*) - F(\mathbf{x}^{(i+1)}) \geq \frac{L}{2} (\|\mathbf{x}^* - \mathbf{x}^{(i+1)}\|^2 - \|\mathbf{x}^* - \mathbf{x}^{(i)}\|^2).$$

- Summing over $i = 0, \dots, k-1$,

$$\begin{aligned} k F(\mathbf{x}^*) - \sum_{i=0}^{k-1} F(\mathbf{x}^{(i+1)}) &\geq \frac{L}{2} (\|\mathbf{x}^* - \mathbf{x}^{(k)}\|^2 - \|\mathbf{x}^* - \mathbf{x}^{(0)}\|^2) \\ \implies \sum_{i=1}^k F(\mathbf{x}^{(i)}) - k F(\mathbf{x}^*) &\leq \frac{L}{2} \|\mathbf{x}^* - \mathbf{x}^{(0)}\|^2. \end{aligned}$$

- Since the algorithm is a descent method, $\sum_{i=1}^k F(\mathbf{x}^{(i)}) \geq k F(\mathbf{x}^{(k)})$.

- Hence,

$$k(F(\mathbf{x}^{(k)}) - F(\mathbf{x}^*)) \leq \frac{L}{2} \|\mathbf{x}^* - \mathbf{x}^{(0)}\|^2.$$

Convergence Analysis

Theorem

For any optimal solution x^* of the composite convex optimization problem (**minimize** $f + g$),

$$F(x^{(k)}) - F(x^*) \leq \frac{L}{2k} \|x^{(0)} - x^*\|^2, \quad \forall k \geq 1.$$

Remark

- It can also be shown that the sequence $(x^{(k)})_{k \in \mathbb{N}}$ converges to an optimal solution.
- This can NOT be considered as a polytime algorithm if ϵ is part of the input: $O(1/\epsilon)$ iterations required to find an ϵ -suboptimal solution, which is exponential w.r.t. input size $\langle \epsilon \rangle := \lceil \log \epsilon \rceil$.

Fast convergence for *strongly convex* functions

Now, consider a composite model (f, g) in which f is ν -strongly convex. Then, a *linear convergence rate* can be achieved (this time, we have a polytime algorithm)

Theorem

If f is ν -strongly convex, then the proximal gradient method with constant step sizes ($t_k = \frac{1}{L}$) generates a sequence of points satisfying

$$(i) \quad \|\mathbf{x}^{(k)} - \mathbf{x}^*\|^2 \leq \left(1 - \frac{\nu}{L}\right)^k \|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2;$$

$$(ii) \quad F(\mathbf{x}^{(k)}) - F(\mathbf{x}^*) \leq \frac{L}{2} \left(1 - \frac{\nu}{L}\right)^k \|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2,$$

where \mathbf{x}^* denotes the *unique* optimal solution.

$\implies O\left(\frac{L}{\nu} \log(LR^2/\epsilon)\right)$ iterations to get ϵ -suboptimal solution.

Linear convergence: Proof sketch

- Recall proof of prox-grad inequality:

$$F(\xi) - F(\mathbf{x}^+) \geq \underbrace{\epsilon_f(\mathbf{x}, \xi)}_{\geq 0} + \frac{1}{2t} (\|\xi - \mathbf{x}^+\|^2 - \|\xi - \mathbf{x}\|^2).$$

Linear convergence: Proof sketch

- Recall proof of prox-grad inequality:

$$F(\xi) - F(\mathbf{x}^+) \geq \underbrace{\epsilon_f(\mathbf{x}, \xi)}_{\geq 0} + \frac{1}{2t} (\|\xi - \mathbf{x}^+\|^2 - \|\xi - \mathbf{x}\|^2).$$

- With f strongly convex, stronger bound: $\epsilon_f(\mathbf{x}, \xi) \geq \frac{\nu}{2} \|\xi - \mathbf{x}\|^2$.

Linear convergence: Proof sketch

- Recall proof of prox-grad inequality:

$$F(\xi) - F(\mathbf{x}^+) \geq \underbrace{\epsilon_f(\mathbf{x}, \xi)}_{\geq 0} + \frac{1}{2t} (\|\xi - \mathbf{x}^+\|^2 - \|\xi - \mathbf{x}\|^2).$$

- With f strongly convex, stronger bound: $\epsilon_f(\mathbf{x}, \xi) \geq \frac{\nu}{2} \|\xi - \mathbf{x}\|^2$.
- At $\xi = \xi^*$, with step size $t = 1/L$,

$$\begin{aligned} F(\mathbf{x}^*) - F(\mathbf{x}^+) &\geq \frac{L}{2} (\|\mathbf{x}^* - \mathbf{x}^+\|^2 - \|\mathbf{x}^* - \mathbf{x}\|^2) + \frac{\nu}{2} \|\mathbf{x}^* - \mathbf{x}\|^2 \\ &= \frac{L}{2} \|\mathbf{x}^* - \mathbf{x}^+\|^2 - \frac{L - \nu}{2} \|\mathbf{x}^* - \mathbf{x}\|^2. \end{aligned}$$

Linear convergence: Proof sketch

- Recall proof of prox-grad inequality:

$$F(\xi) - F(\mathbf{x}^+) \geq \underbrace{\epsilon_f(\mathbf{x}, \xi)}_{\geq 0} + \frac{1}{2t} (\|\xi - \mathbf{x}^+\|^2 - \|\xi - \mathbf{x}\|^2).$$

- With f strongly convex, stronger bound: $\epsilon_f(\mathbf{x}, \xi) \geq \frac{\nu}{2} \|\xi - \mathbf{x}\|^2$.
- At $\xi = \xi^*$, with step size $t = 1/L$,

$$\begin{aligned} F(\mathbf{x}^*) - F(\mathbf{x}^+) &\geq \frac{L}{2} (\|\mathbf{x}^* - \mathbf{x}^+\|^2 - \|\mathbf{x}^* - \mathbf{x}\|^2) + \frac{\nu}{2} \|\mathbf{x}^* - \mathbf{x}\|^2 \\ &= \frac{L}{2} \|\mathbf{x}^* - \mathbf{x}^+\|^2 - \frac{L - \nu}{2} \|\mathbf{x}^* - \mathbf{x}\|^2. \end{aligned}$$

- Then, we use $F(\mathbf{x}^*) - F(\mathbf{x}^+) \leq 0$:

$$\frac{L}{2} \|\mathbf{x}^* - \mathbf{x}^+\|^2 \leq \frac{L - \nu}{2} \|\mathbf{x}^* - \mathbf{x}\|^2 \iff \|\mathbf{x}^* - \mathbf{x}^+\|^2 \leq \left(1 - \frac{\nu}{L}\right) \|\mathbf{x}^* - \mathbf{x}\|^2.$$

- The rest of the proof follows by easy induction.

Outline

- 1 Gradient & Subgradient methods
- 2 Strong convexity and L-smoothness
- 3 The proximal operator
- 4 The proximal gradient method
- 5 The FISTA accelerated method**
- 6 Optimality of accelerated gradient methods

History

- For smooth optimization, *accelerated* gradient methods were proposed by Nesterov [80's]
- Convergence in $\epsilon = O(1/k^2)$ instead of $O(1/k)$.
- The rate is *optimal* in some sense over the class of first-order methods
- Idea: Update $x^{(k+1)}$ by taking a gradient step at point $y^{(k)}$, where $y^{(k)}$ is a well-chosen linear combination of previous 2 iterates, $x^{(k)}$ and $x^{(k-1)}$.
- Generalized to nonsmooth composite models by Beck & Teboulle [2009]. Method called FISTA for *fast iterative shrinkage-thresholding algorithm*, which describes the proximal gradient steps when $g(x) = \|x\|_1$.

FISTA

FISTA (here, with constant step sizes $t_k = \frac{1}{L}, \forall k$)

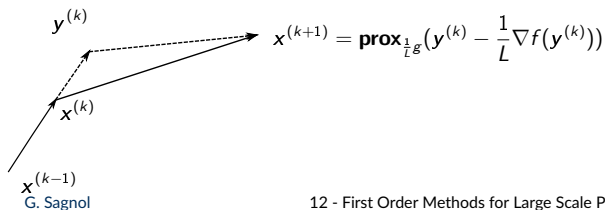
Initialization: $\mathbf{y}^{(0)} = \mathbf{x}^{(0)} \in \mathbf{int\,dom}\,f, \tau_0 = 1.$

For $k = 0, 1, 2, \dots,$

1 $\mathbf{x}^{(k+1)} = \mathbf{prox}_{\frac{1}{L}g}(\mathbf{y}^{(k)} - \frac{1}{L}\nabla f(\mathbf{y}^{(k)}))$

2 $\tau_{k+1} = \frac{1 + \sqrt{1 + 4\tau_k^2}}{2}$

3 $\mathbf{y}^{(k+1)} = \mathbf{x}^{(k+1)} + \left(\frac{\tau_k - 1}{\tau_{k+1}}\right) (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})$



FISTA

FISTA (here, with constant step sizes $t_k = \frac{1}{L}, \forall k$)

Initialization: $\mathbf{y}^{(0)} = \mathbf{x}^{(0)} \in \text{int dom } f, \tau_0 = 1.$

For $k = 0, 1, 2, \dots,$

$$\mathbf{1} \quad \mathbf{x}^{(k+1)} = \text{prox}_{\frac{1}{L}g}(\mathbf{y}^{(k)} - \frac{1}{L}\nabla f(\mathbf{y}^{(k)}))$$

$$\mathbf{2} \quad \tau_{k+1} = \frac{1 + \sqrt{1 + 4\tau_k^2}}{2}$$

$$\mathbf{3} \quad \mathbf{y}^{(k+1)} = \mathbf{x}^{(k+1)} + \left(\frac{\tau_k - 1}{\tau_{k+1}} \right) (\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})$$

■ In fact, τ_{k+1} solves the equation $\tau_{k+1}^2 - \tau_{k+1} = \tau_k^2.$

■ Simple induction:

$$\tau_k \geq \frac{k+2}{2} \geq 1, \forall k \in \mathbb{N}.$$

Theorem

Consider the sequence of iterates $\mathbf{x}^{(k)}$ generated by FISTA (with constant step sizes $t_k = \frac{1}{L}, \forall k$). Then, for any optimal solution \mathbf{x}^* of the composite model (**minimize** $f + g$), it holds

$$F(\mathbf{x}^{(k)}) - F(\mathbf{x}^*) \leq \frac{2L \|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2}{(k+1)^2}.$$

Proof: See blackboard.

Example: Lasso regression (1/4)

```
In [1]: #import packages
import numpy as np
import picos
%matplotlib inline
import matplotlib
import matplotlib.pyplot as plt
import time
```

The goal of this Notebook is to implement FISTA to solve the Lasso-regression problem

$$\min_x \|Ax - y\|^2 + \lambda \|x\|_1.$$

```
In [2]: #generate data
#In this example, y = Ax_0 + noise, where x_0 is sparse
m,n = 5000,1000
x0 = np.random.randn(n)
x0[:int(3*n/4)]=0
A = np.random.rand(m,n)
y = A.dot(x0) + 0.01 * np.random.rand(m)
lbda = 0.05
```

Example: Lasso regression (2/4)

Define Proximal (soft-thresholding) operator of $g = \|\cdot\|_1$

```
In [3]: def prox_threshold(x,t):  
        return np.maximum(0,np.abs(x)-t) * np.sign(x)
```

Compute Lipschitz constant

```
In [4]: L = 2*lbda*max(np.linalg.svd(A)[1])**2
```

Proximal Gradient Method

```
In [ ]: Niter = 10000  
x = np.zeros(n)  
prox_grad = []  
for k in range(Niter):  
    #compute value of objective function  
    r = A.dot(x)-y  
    prox_grad.append(lbda*np.linalg.norm(r)**2 + np.linalg.norm(x,1))  
    #compute gradient of f  
    grad = 2 * lbda * A.T.dot(r)  
    #proximal step  
    x = prox_threshold(x - 1./L * grad,1./L)
```

Example: Lasso regression (3/4)

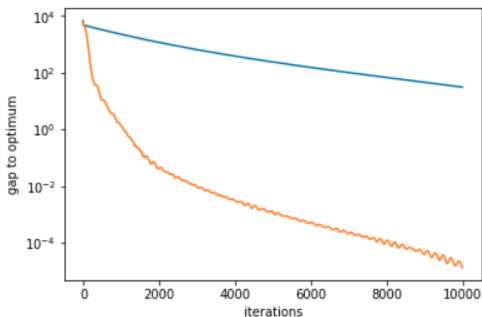
FISTA

```
In [ ]: x_current = np.zeros(n)
y_current = np.zeros(n)
tau_current = 1.
fista = []
for k in range(Niter):
    #compute value of objective function
    fista.append(lbda*np.linalg.norm(A.dot(x_current)-y)**2 + np.linalg.norm(x_current,1))
    #compute gradient of f at y
    grad = 2 * lbda * A.T.dot(A.dot(y_current)-y)
    x_new = prox_threshold(y_current - 1./L * grad,1./L)
    #update tau and y
    tau_new = (1+(1+4*tau_current**2)**0.5)/2.
    y_current = x_new + (tau_current-1.)/tau_new * (x_new-x_current)
    #update current values
    x_current = x_new
    tau_current = tau_new
```

Example: Lasso regression (4/4)

```
In [18]: opt = min(min(prox_grad),min(fista))
plt.semilogy(np.array(prox_grad)-opt)
plt.semilogy(np.array(fista[:Niter])-opt)
plt.xlabel("iterations")
plt.ylabel("gap to optimum")
print "time(proximal gradient)=",t_proxgrad
print "time(FISTA)=",t_fista
```

```
time(proximal gradient)= 44.049719
time(FISTA)= 59.7696934545
```



Example: Lasso regression

- On this example, FISTA \gg standard proximal gradient
- But FISTA *is not* a descent method
- The upper bounds for the gap $\delta_k \leq \frac{LR^2}{2k}$ and

$\delta_k \leq \frac{2LR^2}{(k+1)^2}$ are very pessimistic:

After $k = 10^4$ iterations, assuming the exact value of $R = \|\mathbf{x}^{(0)} - \mathbf{x}^*\|$ is known,

Algorithm	δ_k	upper bound
Proximal gradient	30.77	1421.41
FISTA	$1.29 \cdot 10^{-5}$	0.5684

- In practice, we can use much better duality bound on δ_k as stopping criterion.

Outline

- 1 Gradient & Subgradient methods
- 2 Strong convexity and L-smoothness
- 3 The proximal operator
- 4 The proximal gradient method
- 5 The FISTA accelerated method
- 6 Optimality of accelerated gradient methods**

$\Omega(1/k^2)$ lower bound

The following result basically states that $O(1/k^2)$ is the best convergence rate we can hope for in the class of first-order methods:

Theorem

There exists a function $f : \mathbb{R}^{2k+1} \rightarrow \mathbb{R}$ which is twice differentiable and L -smooth, such that for any sequence $(\mathbf{x}^{(i)})_{i \in \mathbb{N}}$ satisfying

$$\mathbf{x}^{(i+1)} \in \mathbf{x}^{(0)} + \text{span}(\nabla f(\mathbf{x}^{(0)}), \dots, \nabla f(\mathbf{x}^{(i)})), \quad \forall i \in \mathbb{N},$$

it holds

$$f(\mathbf{x}^{(k)}) - f(\mathbf{x}^*) \geq \frac{3L \|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2}{32(k+1)^2}.$$

$\Omega(1/k^2)$ lower bound: Proof sketch

$$f_k(\mathbf{x}) := \frac{L}{4} \left(\frac{1}{2} \mathbf{x}^T A \mathbf{x} - \mathbf{e}_1^T \mathbf{x} \right), \text{ where } A = \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{pmatrix} \in \mathbb{S}^k.$$

Assume (w.l.o.g.) that $\mathbf{x}^{(0)} = \mathbf{0}$. We can show that

$\Omega(1/k^2)$ lower bound: Proof sketch

$$f_k(\mathbf{x}) := \frac{L}{4} \left(\frac{1}{2} \mathbf{x}^T A \mathbf{x} - \mathbf{e}_1^T \mathbf{x} \right), \text{ where } A = \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{pmatrix} \in \mathbb{S}^k.$$

Assume (w.l.o.g.) that $\mathbf{x}^{(0)} = \mathbf{0}$. We can show that

- $A \succeq 0$ and $\lambda_{\max}(A) \leq 4$, hence f_k is convex and L -smooth, $\forall k$.

$\Omega(1/k^2)$ lower bound: Proof sketch

$$f_k(\mathbf{x}) := \frac{L}{4} \left(\frac{1}{2} \mathbf{x}^T A \mathbf{x} - \mathbf{e}_1^T \mathbf{x} \right), \text{ where } A = \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{pmatrix} \in \mathbb{S}^k.$$

Assume (w.l.o.g.) that $\mathbf{x}^{(0)} = \mathbf{0}$. We can show that

- $A \succeq 0$ and $\lambda_{\max}(A) \leq 4$, hence f_k is convex and L -smooth, $\forall k$.
- f_k is minimized over \mathbb{R}^k at $\mathbf{x}^* = A^{-1} \mathbf{e}_1$, and

$$f_k(\mathbf{x}^*) = -\frac{L}{8} \mathbf{e}_1^T A^{-1} \mathbf{e}_1 = -\frac{L}{8} \left(1 - \frac{1}{k+1} \right).$$

$\Omega(1/k^2)$ lower bound: Proof sketch

$$f_k(\mathbf{x}) := \frac{L}{4} \left(\frac{1}{2} \mathbf{x}^T A \mathbf{x} - \mathbf{e}_1^T \mathbf{x} \right), \text{ where } A = \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{pmatrix} \in \mathbb{S}^k.$$

Assume (w.l.o.g.) that $\mathbf{x}^{(0)} = \mathbf{0}$. We can show that

- $A \succeq 0$ and $\lambda_{\max}(A) \leq 4$, hence f_k is convex and L -smooth, $\forall k$.
- f_k is minimized over \mathbb{R}^k at $\mathbf{x}^* = A^{-1} \mathbf{e}_1$, and

$$f_k(\mathbf{x}^*) = -\frac{L}{8} \mathbf{e}_1^T A^{-1} \mathbf{e}_1 = -\frac{L}{8} \left(1 - \frac{1}{k+1} \right).$$

- Let $f = f_{2k+1}$. A simple induction shows that for all $i < 2k + 1$,
 $\text{span}(\nabla f(\mathbf{x}^{(0)}), \dots, \nabla f(\mathbf{x}^{(i)})) \subseteq \text{span}(\mathbf{e}_1, \dots, \mathbf{e}_{i+1})$.

$\Omega(1/k^2)$ lower bound: Proof sketch

$$f_k(\mathbf{x}) := \frac{L}{4} \left(\frac{1}{2} \mathbf{x}^T A \mathbf{x} - \mathbf{e}_1^T \mathbf{x} \right), \text{ where } A = \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{pmatrix} \in \mathbb{S}^k.$$

Assume (w.l.o.g.) that $\mathbf{x}^{(0)} = \mathbf{0}$. We can show that

- $A \succeq 0$ and $\lambda_{\max}(A) \leq 4$, hence f_k is convex and L -smooth, $\forall k$.
- f_k is minimized over \mathbb{R}^k at $\mathbf{x}^* = A^{-1} \mathbf{e}_1$, and

$$f_k(\mathbf{x}^*) = -\frac{L}{8} \mathbf{e}_1^T A^{-1} \mathbf{e}_1 = -\frac{L}{8} \left(1 - \frac{1}{k+1} \right).$$

- Let $f = f_{2k+1}$. A simple induction shows that for all $i < 2k + 1$,
 $\text{span}(\nabla f(\mathbf{x}^{(0)}), \dots, \nabla f(\mathbf{x}^{(i)})) \subseteq \text{span}(\mathbf{e}_1, \dots, \mathbf{e}_{i+1})$.

- So, $f(\mathbf{x}^{(k)}) \geq \inf_{\mathbf{z}} f_k(\mathbf{z}) = -\frac{L}{8} \left(1 - \frac{1}{k+1} \right)$ and $f(\mathbf{x}^*) = -\frac{L}{8} \left(1 - \frac{1}{2k+2} \right)$

$\Omega(1/k^2)$ lower bound: Proof sketch

$$f_k(\mathbf{x}) := \frac{L}{4} \left(\frac{1}{2} \mathbf{x}^T A \mathbf{x} - \mathbf{e}_1^T \mathbf{x} \right), \text{ where } A = \begin{pmatrix} 2 & -1 & & & \\ -1 & 2 & -1 & & \\ & \ddots & \ddots & \ddots & \\ & & -1 & 2 & -1 \\ & & & -1 & 2 \end{pmatrix} \in \mathbb{S}^k.$$

Assume (w.l.o.g.) that $\mathbf{x}^{(0)} = \mathbf{0}$. We can show that

- $A \succeq 0$ and $\lambda_{\max}(A) \leq 4$, hence f_k is convex and L -smooth, $\forall k$.
- f_k is minimized over \mathbb{R}^k at $\mathbf{x}^* = A^{-1} \mathbf{e}_1$, and

$$f_k(\mathbf{x}^*) = -\frac{L}{8} \mathbf{e}_1^T A^{-1} \mathbf{e}_1 = -\frac{L}{8} \left(1 - \frac{1}{k+1} \right).$$

- Let $f = f_{2k+1}$. A simple induction shows that for all $i < 2k + 1$,
 $\text{span}(\nabla f(\mathbf{x}^{(0)}), \dots, \nabla f(\mathbf{x}^{(i)})) \subseteq \text{span}(\mathbf{e}_1, \dots, \mathbf{e}_{i+1})$.

- So, $f(\mathbf{x}^{(k)}) \geq \inf_{\mathbf{z}} f_k(\mathbf{z}) = -\frac{L}{8} \left(1 - \frac{1}{k+1} \right)$ and $f(\mathbf{x}^*) = -\frac{L}{8} \left(1 - \frac{1}{2k+2} \right)$
- Finally, we can bound $\|\mathbf{x}^*\|^2 \leq \frac{2}{3}(k+1)$. Putting all together yields the desired bound.