

# Introductory Course on Non-smooth Optimisation

## Lecture 02 - Proximal gradient descent

---

Jingwei Liang

Department of Applied Mathematics and Theoretical Physics

- 1 Subgradient descent
- 2 Proximal gradient descent
- 3 Proximal mapping
- 4 Inertial proximal gradient
- 5 Fast iterative shrinkage-thresholding algorithm (FISTA)
- 6 Restarting FISTA
- 7 Numerical experiments

## Unconstrained non-smooth optimisation

Consider minimising

$$\min_{x \in \mathbb{R}^n} R(x),$$

where  $R : \mathbb{R}^n \rightarrow ]-\infty, +\infty]$  is proper convex and **lower semi-continuous**.

- $\Gamma_0$ : the class of proper convex and lower semi-continuous functions on  $\mathbb{R}^n$ .
- The set of minimisers, *i.e.*

$$\text{Argmin}(R) = \{x \in \mathbb{R}^n : R(x) = \min_{x \in \mathbb{R}^n} R(x)\}$$

is non-empty.

- $R(x)$  is non-differentiable...

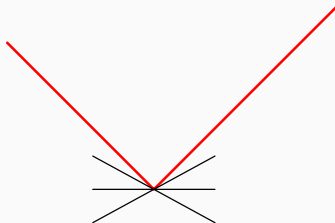
## Subdifferential

Let  $R \in \Gamma_0$ , the subdifferential of  $R$  at  $x \in \text{dom}(R)$  is defined by

$$\partial R : \mathbb{R}^n \rightrightarrows \mathbb{R}^n, x \mapsto \{g \in \mathbb{R}^n \mid R(y) \geq R(x) + \langle g, y - x \rangle, \forall y \in \mathbb{R}^n\}.$$

**Example** absolute value function

$$\partial|x| = \begin{cases} +1 : x > 0 \\ [-1, 1] : x = 0 \\ -1 : x < 0 \end{cases}$$



## Subdifferential

Let  $R \in \Gamma_0$ , the subdifferential of  $R$  at  $x \in \text{dom}(R)$  is defined by

$$\partial R : \mathbb{R}^n \rightrightarrows \mathbb{R}^n, x \mapsto \{g \in \mathbb{R}^n \mid R(y) \geq R(x) + \langle g, y - x \rangle, \forall y \in \mathbb{R}^n\}.$$

## Convexity

Let  $R \in \Gamma_0$  and  $x \in \text{dom}(R)$ , then

- $\partial R(x) = \{g \in \mathbb{R}^n : R(y) \geq R(x) + \langle g, y - x \rangle\}.$
- $\partial R(x)$  is closed and convex.

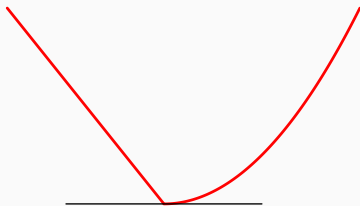
## Monotonicity

Let  $R \in \Gamma_0$ , then  $\forall x, y \in \text{dom}(R)$ ,

$$\langle u - v, x - y \rangle \geq 0, \forall u \in \partial R(x), v \in \partial R(y).$$

Given  $x^* \in \mathbb{R}^n$ , it minimises  $R(x)$  if and only if

$$0 \in \partial R(x^*).$$



$$R(y) \geq R(x^*) + \langle g, y - x^* \rangle \text{ holds for all } y \in \text{dom}(R) \iff 0 \in \partial R(x^*).$$

## Subgradient descent

**initial** :  $x_0 \in \text{dom}(R)$ ;

**repeat** :

1. Choose step-size  $\gamma_k > 0$  and a subgradient  $g_k \in \partial R(x_k)$
2. Update  $x_{k+1} = x_k - \gamma_k g_k$

**until** : stopping criterion is satisfied.

## Subgradient descent

**initial** :  $x_0 \in \text{dom}(R)$ ;

**repeat** :

1. Choose step-size  $\gamma_k > 0$  and a subgradient  $g_k \in \partial R(x_k)$
2. Update  $x_{k+1} = x_k - \gamma_k g_k$

**until** : stopping criterion is satisfied.

Step-size rule:

- Fixed step-size:  $\gamma_k$  is constant.
- Fixed length:  $\gamma_k \|g_k\| = \|x_{k+1} - x_k\|$  is a constant.
- Diminishing step-size:  $\gamma_k \rightarrow 0$ ,  $\sum_i \gamma_i = +\infty$ .



Assumptions:

- $R$  has minimiser  $x^*$  and finite optimal value  $R(x^*)$ .
- $R$  is convex,  $\text{dom}(R) = \mathbb{R}^n$ .
- $R$  is Lipschitz continuous with constant  $L$ :

$$|R(x) - R(y)| \leq L\|x - y\|, \quad \forall x, y \in \text{dom}(R).$$

**NB:** the Lipschitz continuity implies  $\|g\| \leq L$  for all  $x \in \text{dom}(R)$ .

Subgradient descent is **NOT** a descent method.

**Fixed step-size**  $\gamma_k \equiv \gamma$

$$R_{k,best} - R(x^*) \leq \frac{\|x_0 - x^*\|^2}{2k\gamma} + \frac{\gamma L^2}{2}.$$

- Does not guarantee the convergence of  $R_{k,best}$ .
- For large  $k$ ,  $R_{k,best}$  is approximately  $\frac{\gamma L^2}{2}$  suboptimal.

Subgradient descent is **NOT** a descent method.

**Fixed step-length**  $\gamma_k = c/\|g_k\|$

$$R_{k,best} - R(x^*) \leq \frac{\|x_0 - x^*\|^2}{2kc} + \frac{cL}{2}.$$

- Does not guarantee the convergence of  $R_{k,best}$ .
- For large  $k$ ,  $R_{k,best}$  is approximately  $\frac{cL}{2}$  suboptimal.

Subgradient descent is **NOT** a descent method.

**Diminishing step-size**  $\gamma_k \rightarrow 0$ ,  $\sum_i \gamma_i = +\infty$ :

$$R_{k,best} - R(x^*) \leq \frac{\|x_0 - x^*\|^2 + L^2 \sum_{i=1}^k \gamma_i^2}{\sum_{i=1}^k \gamma_i}.$$

- If  $\sum_{i=1}^k \gamma_i^2 / \sum_{i=1}^k \gamma_i \rightarrow 0$ , then  $R_{k,best} \rightarrow R(x^*)$ .
- Choice of  $\gamma_k$ :  $\gamma_k = c/k^q$ ,  $q \in ]1/2, 1[$ .

**For fixed number of iterations** if  $c_i = \gamma_i \|g_i\|$  and  $\|x_0 - x^*\| \leq D$ ,

$$R_{k,best} - R(x^*) \leq \frac{D^2 + L^2 \sum_{i=1}^k c_i^2}{2 \sum_{i=1}^k \gamma_i / L}.$$

- For given  $k$ , rhs is minimised by  $c_i = c = D/\sqrt{k}$ .
- Hence the rate

$$R_{k,best} - R(x^*) \leq \frac{LD}{\sqrt{k}}.$$

- Iteration complexity: reach  $R_{k,best} - R(x^*) < \epsilon$  in  $O(1/\epsilon^2)$  steps.

**When  $R(x^*)$  is available** step-size

$$\gamma_k = \frac{R(x_k) - R(x^*)}{\|g_k\|^2}.$$

Convergence rate

$$R_{k,best} - R(x^*) \leq \frac{LD}{\sqrt{k}}.$$

**NB:**  $O(1/\sqrt{k})$  is the best rate can be obtained by subgradient method.

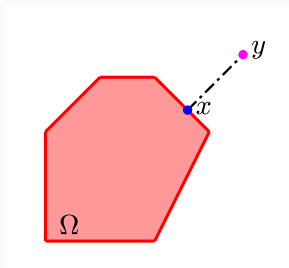
- Handles non-smooth problem
- Simple iterative scheme
- Slow convergence rate
- No clear stopping criterion

**NB:** need a better approach to handle non-smoothness...

- 1 Subgradient descent
- 2 Proximal gradient descent**
- 3 Proximal mapping
- 4 Inertial proximal gradient
- 5 Fast iterative shrinkage-thresholding algorithm (FISTA)
- 6 Restarting FISTA
- 7 Numerical experiments

**Indicator function** : let  $\Omega \subseteq \mathbb{R}^n$

$$\iota_{\Omega}(x) = \begin{cases} 0 & : x \in \Omega, \\ +\infty & : x \notin \Omega. \end{cases}$$



Projection of  $y$  onto  $\Omega$ :

$$\min_{x \in \Omega} \|x - y\|.$$

## Projection

Projection mapping onto a set is defined by

$$\mathcal{P}_{\Omega}(y) \stackrel{\text{def}}{=} \operatorname{argmin}_{x \in \Omega} \|x - y\|.$$



## Constrained smooth optimisation

Let  $F \in C_L^1$  and  $\Omega \subseteq \mathbb{R}^n$  be a closed and convex set

$$\min_{x \in \Omega} F(x).$$

## Projected gradient descent

**initial** :  $x_0 \in \Omega$ ;

**repeat** :

1. Choose step-size  $\gamma_k \in ]0, 2/L[$
2. Gradient descent  $x_{k+1/2} = x_k - \gamma_k \nabla F(x_k)$
3. Projection  $x_{k+1} = \mathcal{P}_\Omega(x_{k+1/2})$

**until** : stopping criterion is satisfied.

As  $\iota_\Omega \in \Gamma_0$ , the constrained optimisation problem is equivalent to

$$\min_{x \in \mathbb{R}^n} \iota_\Omega(x) + F(x).$$

## Composite optimisation

Consider the following optimisation problem

$$\min_{x \in \mathbb{R}^n} \{ \Phi(x) \stackrel{\text{def}}{=} R(x) + F(x) \}.$$

## Assumptions

- $F \in C_L^1$
- $R \in \Gamma_0$
- $\text{Argmin}(\Phi) \neq \emptyset$

**Examples** regularised LSE, image processing...

Projection

$$\begin{aligned}\mathcal{P}_{\Omega}(y) &\stackrel{\text{def}}{=} \operatorname{argmin}_{x \in \Omega} \|x - y\| \\ &= \operatorname{argmin}_{x \in \mathbb{R}^n} \iota_{\Omega}(x) + \frac{1}{2} \|x - y\|^2.\end{aligned}$$

Proximal mapping

$$\operatorname{prox}_R(y) \stackrel{\text{def}}{=} \operatorname{argmin}_{x \in \mathbb{R}^n} R(x) + \frac{1}{2} \|x - y\|^2.$$

## Projected gradient descent

**initial** :  $x_0 \in \Omega$ ;

**repeat** :

1. Choose step-size  $\gamma_k \in ]0, 2/L[$
2. Gradient descent  $x_{k+1/2} = x_k - \gamma_k \nabla F(x_k)$
3. Projection  $x_{k+1} = \operatorname{prox}_{\gamma_k R}(x_{k+1/2})$

**until** : stopping criterion is satisfied.

- A.K.A Forward-Backward splitting
  - Forward step: gradient descent of  $F$ .
  - Backward step: proximal mapping of  $R$ .

- Iteration in one line

$$x_{k+1} = \text{prox}_{\gamma_k R}(x_k - \gamma_k \nabla F(x_k)).$$

- Definition of  $\text{prox}_{\gamma R}$ ,

$$\begin{aligned} x_{k+1} &= \operatorname{argmin}_x \left\{ \gamma_k R(x) + \frac{1}{2} \|x - (x_k - \gamma_k \nabla F(x_k))\|^2 \right\} \\ &= \operatorname{argmin}_x \left\{ \gamma_k R(x) + \gamma_k \langle \nabla F(x_k), x - x_k \rangle + \frac{1}{2} \|x - x_k\|^2 \right\} \\ &= \operatorname{argmin}_x \left\{ R(x) + \boxed{F(x_k) + \langle \nabla F(x_k), x - x_k \rangle + \frac{1}{2\gamma_k} \|x - x_k\|^2} \right\}. \end{aligned}$$

**NB:**  $x_{k+1}$  minimises  $R(x)$  plus the majorisation of  $F(x)$  at  $x_k$  if  $\gamma_k \leq \frac{1}{L}$ .

- **Gradient descent**  $R = 0$

$$x_{k+1} = x_k - \gamma_k \nabla F(x_k).$$

- **Proximal point algorithm**  $F = 0$

$$x_{k+1} = \text{prox}_{\gamma_k R}(x_k).$$

- **Projected gradient descent**  $R = \iota_\Omega$

$$x_{k+1} = \mathcal{P}_\Omega(x_k - \gamma_k \nabla F(x_k)).$$

- **ISTA: iterative shrinkage-thresholding algorithm**  $R = \|x\|_1$

$$x_{k+1} = \mathcal{T}_\gamma(x_k - \gamma_k \nabla F(x_k)),$$

where

$$(\mathcal{T}_\gamma(y))_i = \begin{cases} \text{sign}(y_i) \cdot (|y_i| - \gamma) & : |y_i| > \gamma, \\ 0 & : |y_i| \leq \gamma. \end{cases}$$

Define

$$E_\gamma(x, y) \stackrel{\text{def}}{=} R(x) + F(y) + \langle \nabla F(y), x - y \rangle + \frac{1}{2\gamma} \|x - y\|^2$$

and  $y_+ \stackrel{\text{def}}{=} \operatorname{argmin}_x E_\gamma(x, y)$ .

### Lemma

Let  $y \in \mathbb{R}^n$  and  $\gamma \in ]0, 2/L[$  such that

$$\Phi(y_+) \leq E_\gamma(y_+, y),$$

then for any  $x \in \mathbb{R}^n$ ,

$$\Phi(x) - \Phi(y_+) \geq \frac{1}{2\gamma} (\|x - y_+\|^2 - \|x - y\|^2).$$

### Lemma

Given  $y \in \mathbb{R}^n$  and  $\gamma \in ]0, 1/L]$ , then for any  $x \in \mathbb{R}^n$ ,

$$\Phi(y_+) + \frac{1}{2\gamma} \|y_+ - x\|^2 \leq \Phi(y) + \frac{1}{2\gamma} \|y - x\|^2.$$

**NB:** proximal gradient is a descent method.

Consider  $\gamma_k \equiv \gamma \in ]0, 1/L]$

- For each step

$$\Phi(x_k) - \Phi(x_{k+1}) \geq \frac{\gamma}{2} \|x_k - x_{k+1}\|^2.$$

- Regarding  $\Phi(x^*)$

$$\Phi(x_{k+1}) - \Phi(x^*) \leq \frac{\gamma}{2} (\|x_k - x^*\|^2 - \|x_{k+1} - x^*\|^2).$$

- Summing up

$$\begin{aligned} k(\Phi(x_k) - \Phi(x^*)) &\leq \sum_{i=1}^k (\Phi(x_i) - \Phi(x^*)) \\ &\leq \frac{\gamma}{2} \sum_{i=1}^k (\|x_{i-1} - x^*\|^2 - \|x_i - x^*\|^2) \leq \frac{\gamma}{2} \|x_0 - x^*\|^2 \end{aligned}$$

- $O(1/k)$  rate

$$\Phi(x_k) - \Phi(x^*) \leq \frac{\gamma \|x_0 - x^*\|^2}{2k}.$$

**NB:** not optimal and can be accelerated.

- 1 Subgradient descent
- 2 Proximal gradient descent
- 3 Proximal mapping**
- 4 Inertial proximal gradient
- 5 Fast iterative shrinkage-thresholding algorithm (FISTA)
- 6 Restarting FISTA
- 7 Numerical experiments



## Proximal mapping

The proximal mapping (proximity operator) of a function  $R \in \Gamma_0$  is defined by

$$\text{prox}_{\gamma R}(y) \stackrel{\text{def}}{=} \operatorname{argmin}_{x \in \mathbb{R}^n} \gamma R(x) + \frac{1}{2} \|x - y\|^2.$$

**Optimality condition** denote  $y_+ \stackrel{\text{def}}{=} \text{prox}_{\gamma R}(y)$ ,

$$\begin{aligned} 0 \in \gamma \partial R(y_+) + y_+ - y &\iff y \in (\text{Id} + \gamma \partial R)(y_+) \\ &\iff y_+ = (\text{Id} + \gamma \partial R)^{-1}(y). \end{aligned}$$

**Projection**  $R(x) = \iota_\Omega(x)$ ,  $\partial \iota_\Omega(x) = \mathcal{N}_\Omega(x) = \{g : \langle g, v - x \rangle \leq 0\}$   
 $\mathcal{P}_\Omega = (\text{Id} + \mathcal{N}_\Omega)^{-1}.$

## Examples

- Hyperplane:  $\Omega = \{x : a^T x = b\}$ ,  $a \neq 0$

$$\mathcal{P}_\Omega = x + \frac{b - a^T x}{\|a\|^2} a.$$

- Affine subspace:  $\Omega = \{x : Ax = b\}$  with  $A \in \mathbb{R}^{m \times n}$ ,  $\text{rank}(A) = m < n$

$$\mathcal{P}_\Omega = x + A^T (AA^T)^{-1} (b - Ax).$$

- Half space:  $\Omega = \{x : a^T x \leq b\}$ ,  $a \neq 0$

$$\mathcal{P}_\Omega = x + \frac{b - a^T x}{\|a\|^2} a \text{ if } a^T x > b \quad \text{and} \quad x \text{ if } a^T x \leq b.$$

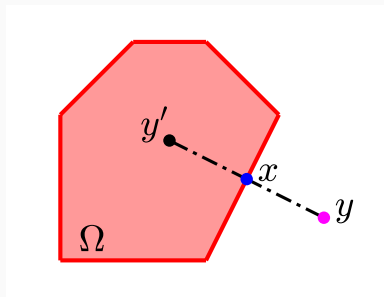
- Nonnegative orthant:  $\Omega = \mathbb{R}_+^n$

$$\mathcal{P}_\Omega = (\max\{0, x_i\})_i.$$

**Projection**  $R(x) = \iota_{\Omega}(x)$ ,  $\partial \iota_{\Omega}(x) = \mathcal{N}_{\Omega}(x) = \{g : \langle g, v - x \rangle \leq 0\}$   
 $\mathcal{P}_{\Omega} = (\text{Id} + \mathcal{N}_{\Omega})^{-1}.$

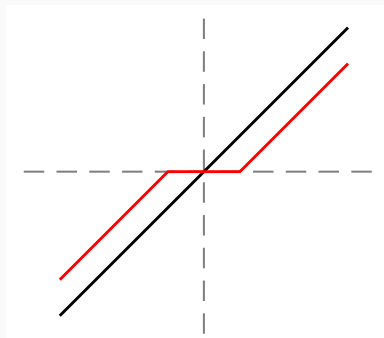
**Reflection**

$$\mathcal{R}_{\mathcal{N}_{\Omega}} = 2\mathcal{P}_{\Omega} - \text{Id} = \mathcal{P}_{\Omega} + (\mathcal{P}_{\Omega} - \text{Id}).$$



**Soft-threshold**:  $R(x) = |x|$ ,

$$\text{prox}_{\gamma R}(y) = \mathcal{T}_\gamma(y) = \begin{cases} y - \gamma : y > \gamma, \\ 0 : y \in [-\gamma, \gamma], \\ y + \gamma : y < -\gamma. \end{cases}$$



**Quadratic function**  $R(x) = \frac{1}{2}x^T Ax + b^T x + c, A \succeq 0$

$$\text{prox}_{\gamma R}(y) = (\text{Id} + \gamma A)^{-1}(y - \gamma b).$$

**Euclidean norm**  $R(x) = \|x\|$

$$\text{prox}_{\gamma R}(y) = \begin{cases} (1 - \frac{\gamma}{\|y\|})y : \|y\| > \gamma, \\ 0 : \text{o.w.} \end{cases}$$

**Logarithmic barrier**  $R(x) = -\sum_i \log(x_i)$

$$(\text{prox}_{\gamma R}(y))_i = \frac{y_i + \sqrt{y_i^2 + 4\gamma}}{2}, \quad i = 1, \dots, n.$$

**Nuclear norm**  $R(x) = \sum_i \sigma_i$

$$\text{prox}_{\gamma R}(y) = U\mathcal{T}_{\gamma}(\Sigma)V^T.$$

**Quadratic perturbation**  $H(x) = R(x) + \frac{\alpha}{2} \|x\|^2 + \langle x, u \rangle + \beta$ ,  $\alpha \geq 0$

$$\text{prox}_H = \text{prox}_{R/(\alpha+1)}\left(\frac{x-u}{\alpha+1}\right).$$

**Translation**  $H(x) = R(x - z)$

$$\text{prox}_H = z + \text{prox}_R(x - z).$$

**Scaling**  $H(x) = R(x/\rho)$

$$\text{prox}_H = \rho \text{prox}_{R/\rho^2}\left(\frac{x}{\rho}\right).$$

**Reflection**  $H(x) = R(-x)$

$$\text{prox}_H = -\text{prox}_R(-x).$$

**Composition**  $H = R \circ L$  with  $L$  being bijective bounded linear mapping such that  $L^{-1} = L^*$ ,

$$\text{prox}_H = L^* \circ \text{prox}_R \circ L.$$

- 1 Subgradient descent
- 2 Proximal gradient descent
- 3 Proximal mapping
- 4 Inertial proximal gradient**
- 5 Fast iterative shrinkage-thresholding algorithm (FISTA)
- 6 Restarting FISTA
- 7 Numerical experiments

## An inertial proximal gradient

**Initial** :  $x_0 \in \mathbb{R}^n$  and  $\gamma \in ]0, 2/L[$ ;

$$\begin{aligned}y_k &= x_k + a_k(x_k - x_{k-1}), \quad a_k \in [0, 1], \\x_{k+1} &= \text{prox}_{\gamma R}(y_k - \gamma \nabla F(x_k)).\end{aligned}$$

- Recovers inertial PPA when  $F = 0$ , and heavy-ball method when  $R = 0$ .
- Convergence via studying the Lyapunov function

$$\mathcal{E}(x_k) \stackrel{\text{def}}{=} \Phi(x_k) + \frac{a_k}{2\gamma} \|x_k - x_{k-1}\|^2.$$

- In general, no convergence rate.



## A general inertial proximal gradient

**Initial :**  $x_0 \in \mathbb{R}^n$ ,  $x_{-1} = x_0$  and  $\gamma \in ]0, 2/L[$ ;

$$y_k = x_k + a_k(x_k - x_{k-1}), \quad a_k \in [0, 1],$$

$$z_k = x_k + b_k(x_k - x_{k-1}), \quad b_k \in [0, 1],$$

$$x_{k+1} = \text{prox}_{\gamma R}(y_k - \gamma \nabla F(z_k)).$$

- Convergence via studying the Lyapunov function

$$\mathcal{E}(x_k) \stackrel{\text{def}}{=} \Phi(x_k) + \frac{a_k}{2\gamma} \|x_k - x_{k-1}\|^2.$$

- In general, no convergence rate.
- Can be extend to multiple steps, e.g.

$$y_k = x_k + a_{0,k}(x_k - x_{k-1}) + a_{1,k}(x_{k-1} - x_{k-2}) + \cdots.$$

**Assumption**  $R = 0$ ,  $F = \frac{1}{2}\|Ax - f\|^2$  and  $(a_k, b_k) \equiv (a, b)$ .

- $A^T A$  is symmetric positive definite with  $A^T A \succeq \alpha \text{Id}$ .

- Taylor expansion

$$x_{k+1} = y_k - \gamma \nabla^2 F(x^*)(z_k - x^*).$$

- Let  $d_k = (x_k - x^*, x_{k-1} - x^*)^T$  and  $H = \nabla^2 F$ ,  $G = \text{Id} - \gamma H$ , then

$$d_{k+1} = \underbrace{\begin{bmatrix} (a-b)\text{Id} + (1+b)G, & -(a-b)\text{Id} - bG \\ \text{Id}, & 0 \end{bmatrix}}_M d_k.$$

- Spectral radius:  $\eta = \rho(G) = 1 - \gamma\alpha$  and  $\rho = \rho(M)$ ...

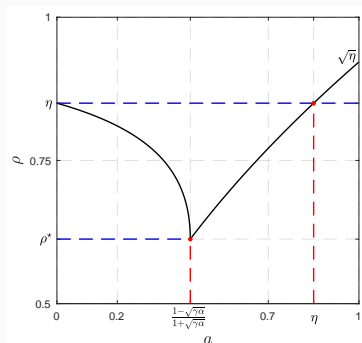
## Spectral radius $\rho$

Between  $\eta$  and  $\rho$ ,

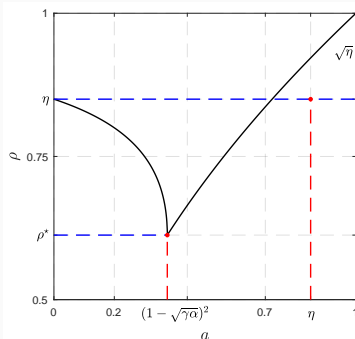
- $\eta$  and  $\rho$  satisfy the relation

$$\rho^2 - ((a - b) + (1 + b)\eta)\rho + (a - b) + b\eta = 0.$$

- Given any  $(a, b) \in [0, 1]^2$ , then  $\rho(M) < 1$  if, and only if  $\frac{2(b-a)-1}{1+2b} < \eta$ .



$b = a$



$b = 0$

- Given  $b \in [0, 1]$ , there exists optimal choice of  $a \in [0, 1]$  such that

$$\rho = 1 - \sqrt{\gamma\alpha}$$

can be obtained.

- Take  $b = a$ , for

$$a \in \left] \frac{1 - \sqrt{\gamma\alpha}}{1 + \sqrt{\gamma\alpha}}, 1 \right],$$

the leading eigenvalue of  $M$  is complex.

- Continue  $b = a$ , for

$$a \in ]\eta, 1],$$

the inertial scheme is actually **slower** than the original scheme.

- 1 Subgradient descent
- 2 Proximal gradient descent
- 3 Proximal mapping
- 4 Inertial proximal gradient
- 5 Fast iterative shrinkage-thresholding algorithm (FISTA)**
- 6 Restarting FISTA
- 7 Numerical experiments

## FISTA

**Initial** :  $x_0 \in \mathbb{R}^n$ ,  $x_{-1} = x_0$ ,  $\gamma = 1/L$  and  $t_0 = 1$ ;

$$t_k = \frac{1 + \sqrt{1 + 4t_{k-1}^2}}{2}, \quad a_k = \frac{t_{k-1} - 1}{t_k},$$

$$y_k = x_k + a_k(x_k - x_{k-1}),$$

$$x_{k+1} = \text{prox}_{\gamma R}(y_k - \gamma \nabla F(y_k)).$$

- A special case of inertial proximal gradient descent.
- Inertial parameters

$$t_k \approx \frac{k+1}{2} \quad \text{and} \quad a_k \rightarrow 1.$$

**Nesterov** compute  $\phi_k \in ]0, 1[$  from equation

$$\phi_k^2 = (1 - \phi_k)\phi_{k-1}^2$$

and  $a_k = \frac{\phi_{k-1}(1 - \phi_{k-1})}{\phi_{k-1}^2 + \phi_k}.$

- $\phi_k$  reads

$$\phi_k = \frac{-\phi_{k-1}^2 + \sqrt{\phi_{k-1}^4 + 4\phi_{k-1}^2}}{2} = \frac{2\phi_{k-1}^2}{\phi_{k-1}^2 + \sqrt{\phi_{k-1}^4 + 4\phi_{k-1}^2}}.$$

- Let  $t_k = 1/\phi_k$ ,

$$\frac{1}{t_k} = \frac{2}{1 + \sqrt{1 + 4t_{k-1}^2}}.$$

- Which leads to

$$t_k = \frac{1 + \sqrt{1 + 4t_{k-1}^2}}{2}.$$

- Moreover,  $a_k = \frac{t_{k-1} - 1}{t_k}.$

**NB:** FISTA is not a descent method.

- Denote  $f_k = \Phi(x_k) - \Phi(x^*)$  and  $u_k = t_k x_k - (t_k - 1)x_{k-1} - x^*$ , then

$$\frac{2}{L} t_k^2 f_k - \frac{2}{L} t_{k+1}^2 f_{k+1} \geq \|u_{k+1}\|^2 - \|u_k\|^2.$$

- Let  $c_k, d_k$  be positive sequences, if

$$c_k - c_{k+1} \geq d_{k+1} - d_k \forall k \geq 1, \text{ with } c_1 + d_1 < C, C > 0$$

then  $c_k < C$  for all  $k \geq 1$ .

- $\frac{2}{L} t_k^2 f_k \leq \|x_0 - x^*\|^2$

- $t_k \geq \frac{k+1}{2},$

$$\Phi(x_k) - \Phi(x^*) \leq \frac{L \|x_0 - x^*\|^2}{2(k+1)^2}.$$



- 1 Subgradient descent
- 2 Proximal gradient descent
- 3 Proximal mapping
- 4 Inertial proximal gradient
- 5 Fast iterative shrinkage-thresholding algorithm (FISTA)
- 6 Restarting FISTA**
- 7 Numerical experiments



Why FISTA oscillates

- for LSE, leading eigenvalue for the system is complex.
- over extrapolation, momentum beats gradient.

## Restarting FISTA

**Initial :**  $x_0 \in \mathbb{R}^n$ ,  $x_{-1} = x_0$ ,  $\gamma = 1/L$  and  $t_0 = 1$ ;

**repeat :**

1. Run FISTA iteration
2. If  $\langle y_k - x_{k+1}, x_{k+1} - x_k \rangle > 0$ :  $t_k = 1$ ,  $y_k = x_k$ .

**until :** stopping criterion is satisfied.

- 1 Subgradient descent
- 2 Proximal gradient descent
- 3 Proximal mapping
- 4 Inertial proximal gradient
- 5 Fast iterative shrinkage-thresholding algorithm (FISTA)
- 6 Restarting FISTA
- 7 Numerical experiments**

## $\ell_1$ -regularised least square (LASSO)

$$\min_{x \in \mathbb{R}^n} \mu \|x\|_1 + \frac{1}{2} \|Ax - f\|^2.$$

## Sparse logistic regression

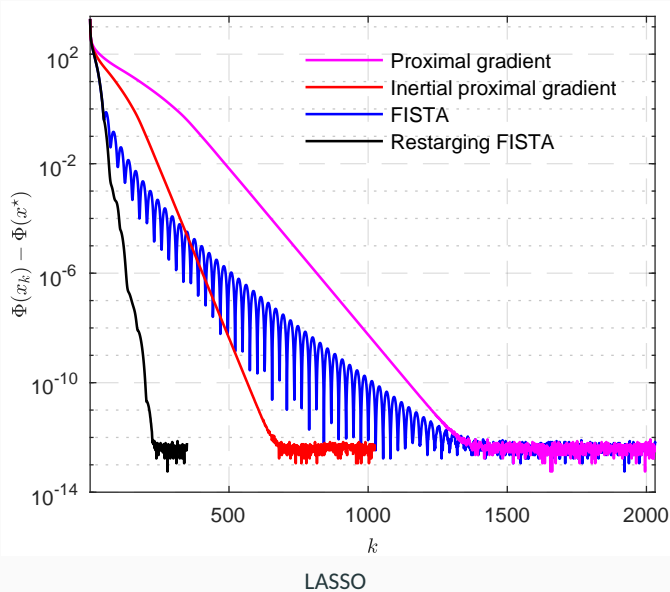
$$\min_{x \in \mathbb{R}^n} \mu \|x\|_1 + \frac{1}{m} \sum_{i=1}^m \log(1 + e^{-l_i h_i^T x}),$$

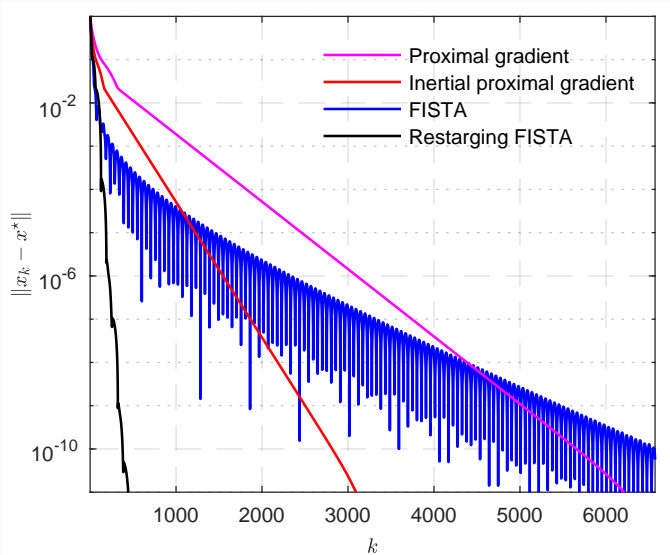
where  $\mu = 10^{-2}$ . The australian data set from LIBSVM<sup>1</sup> is considered.

---

<sup>1</sup><https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>

- Proximal gradient descent
- Inertial proximal gradient descent
- FISTA
- Restarting FISTA





Sparse logistic regression



- B. Polyak. "Introduction to optimization". Optimization Software, 1987.
- Y. Nesterov. "Introductory lectures on convex optimization: A basic course". Vol. 87. Springer Science & Business Media, 2013.
- A. Beck and M. Teboulle. "A fast iterative shrinkage-thresholding algorithm for linear inverse problems". SIAM Journal on Imaging Sciences, 2(1):183–202, 2009.
- H. Bauschke and P. L. Combettes. "Convex analysis and monotone operator theory in Hilbert spaces". Springer, 2011.
- B. O'Donoghue and E. J. Candès. "Adaptive restart for accelerated gradient schemes". Foundations of Computational Mathematics, pages 1–18, 2012.