# First-order methods in optimization - Evaluation

Olivier Leblanc, Guillaume Thiran.

## 1 Part 1 - Slide 45

### 1.1 6. (with code)

$$f(\boldsymbol{x}) = 2x_{[1]} + x_{[2]} = \max_{\boldsymbol{y}} \left\{ \sum_i y_i x_i; \sum_i y_i = 3, 0 \leq y_i \leq 2 \right\} = \sigma_{\{\boldsymbol{y}|\mathbf{1}^\top \boldsymbol{y}=3, \mathbf{0} \leq \boldsymbol{y} \leq 2\mathbf{1}\}}.$$

Hence,

$$\text{prox}_f(\boldsymbol{x}) = \boldsymbol{x} - \mathcal{P}_{\{\boldsymbol{y}|\mathbf{1}^\top \boldsymbol{y}=3, \mathbf{0} \leq \boldsymbol{y} \leq 2\mathbf{1}\}}(\boldsymbol{x}).$$

Writing $C = \{\boldsymbol{y}|\mathbf{1}^\top \boldsymbol{y} = 3, \mathbf{0} \leq \boldsymbol{y} \leq 2\mathbf{1}\}$, it can be compared with $H_{\boldsymbol{a},b} \bigcap \text{Box}[\boldsymbol{l}, \boldsymbol{u}]$, with $\boldsymbol{a} = \mathbf{1}, b = 3, \boldsymbol{l} = \mathbf{0}, \boldsymbol{u} = \mathbf{2}$.

```python
import numpy as np

def f(x):
    xsorted = np.sort(x)
    return 2*xsorted[0] + xsorted[1]

def projbox(x, l, u):
    return np.minimum(np.maximum(x,l), u)

def projH_inter_box(x, a, b, l, u):
    mu = 1
    factor = 1
    val = 10

    while (np.abs(val)>1e-8):
        val = a@projbox(x-mu*a, l, u) - b
        mu *= (1+factor)**(np.sign(val))
        factor /= 1.2

    return projbox(x-mu*a, l, u)

def proxf(x):
    return x - projH_inter_box(x, np.ones(len(x)), 3, np.zeros(len(x)), 2*np.ones(len(x)))

print(proxf(np.array([2,1,4,1,2,1])))
```

The output is : $[1.49999994, 1., 2., 1., 1.49999994, 1.]$.

### 1.2 8.

$$f(t) = \begin{cases} 1/t, & t > 0, \\ \infty, & \text{else.} \end{cases}$$

$$\text{prox}_{\lambda f}(t) = \arg\min_u \begin{cases} 1/\lambda u, & u > 0, \\ \infty, & \text{else.} \end{cases} + \frac{1}{2} \|u - t\|_2^2$$

Hence,

$$\frac{-1}{\lambda u^2} + u - t = 0 \Leftrightarrow u^3 - \lambda t u^2 - \lambda = 0.$$

Finally,

$$\text{prox}_{\lambda f}(t) = \{u > 0 | u^3 - \lambda t u^2 - \lambda = 0\}. \tag{1}$$

## 1.3  9.

$$f(\boldsymbol{X}) = \begin{cases} \operatorname{tr} \boldsymbol{X}^{-1}, & \boldsymbol{X} \succ 0, \\ \infty, & \text{else.} \end{cases} = \begin{cases} \sum_{i=1}^{n} \frac{1}{\lambda_i}, & \boldsymbol{X} \succ 0, \\ \infty, & \text{else.} \end{cases}$$

As $\boldsymbol{X} \in \mathbb{S}^n$, it is a *spectral function*. Hence one can write $f(\boldsymbol{X}) = g(\lambda_1(\boldsymbol{X}), \cdots, \lambda_n(\boldsymbol{X})) = \sum_{i=1}^{n} h(\lambda_i)$. With the Singular Value Decomposition (SVD) of $\boldsymbol{X}$ as $\boldsymbol{X} = \boldsymbol{U} \operatorname{diag} \lambda(\boldsymbol{X}) \boldsymbol{U}^T$, this gives

$$\begin{aligned} \operatorname{prox}_{\lambda f}(\boldsymbol{X}) &= \boldsymbol{U} \operatorname{diag}(\operatorname{prox}_{\lambda g}[\lambda_1, \cdots, \lambda_n]) \boldsymbol{U}^\top \\ &= \boldsymbol{U} \operatorname{diag}(\operatorname{prox}_{\lambda h}(\lambda_1), \cdots, \operatorname{prox}_{\lambda h}(\lambda_n)) \boldsymbol{U}^\top \\ &= \boldsymbol{U} \operatorname{diag}(\{u > 0 | u^3 - \lambda \lambda_i u^2 - \lambda = 0\}) \boldsymbol{U}^\top. \end{aligned}$$

## 1.4  10.

$$\lambda f(\boldsymbol{x}) = \lambda(\|\boldsymbol{x}\|_2 - 1)^2 = \lambda \|\boldsymbol{x}\|_2^2 - 2\lambda \|\boldsymbol{x}\|_2 + \lambda.$$

Using the provided tables, one identifies it with $g(\boldsymbol{x}) + \frac{c}{2} \|\boldsymbol{x}\|_2^2 + \langle \boldsymbol{a}, \boldsymbol{x} \rangle + \gamma$, with $g(\boldsymbol{x}) = -2\lambda \|\boldsymbol{x}\|_2$, $c = 2\lambda$, $\boldsymbol{a} = \boldsymbol{0}$, $\gamma = 0$. Hence,

$$\operatorname{prox}_{\lambda f}(\boldsymbol{x}) = \operatorname{prox}_{\frac{-2\lambda \|\cdot\|_2}{1+2\lambda}} \left( \frac{\boldsymbol{x}}{1+2\lambda} \right) = \left( 1 + \frac{2\lambda}{\|\boldsymbol{x}\|_2} \right) \frac{\boldsymbol{x}}{1+2\lambda}, \quad \boldsymbol{x} \neq \boldsymbol{0}.$$

# 2  Part 2 - Exercise 0 - p40

$$\min_{\boldsymbol{x} \in \mathbb{R}^n} \|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{b}\|_2^2 + \underbrace{\frac{\lambda_1}{2} \|\boldsymbol{x}\|_2^2 + \lambda 2 \|\boldsymbol{x}\|_1}_{\text{elastic net}},$$

where $\boldsymbol{A} \in \mathbb{R}^{m \times n}$, $\in \mathbb{R}^m$ and $\lambda_1, \lambda_2 > 0$. Choosing $f(\boldsymbol{x}) = \|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{b}\|_2^2 + \frac{\lambda_1}{2} \|\boldsymbol{x}\|_2^2$ as the $\sigma$-strongly convex and differentiable part and $g(\boldsymbol{x}) = \lambda_2 \|\boldsymbol{x}\|_1$ as the closed convex but non differentiable part, one has $\nabla f(\boldsymbol{x}) = \boldsymbol{A}^\top(\boldsymbol{A}\boldsymbol{x} - \boldsymbol{b}) + \lambda_1 \boldsymbol{I}\boldsymbol{x}$ and $\operatorname{prox}_{\lambda g}(\boldsymbol{x}) = \mathcal{T}_{\lambda_2}(\boldsymbol{x})$.

- (*Proximal Gradient*)

$$\begin{aligned} \boldsymbol{x}^{k+1} &= \operatorname{prox}_{\frac{1}{L} g}(\boldsymbol{x}^k - \tfrac{1}{L} \nabla f(\boldsymbol{x}^k)) \\ &= \mathcal{T}_{\frac{\lambda_2}{L}} \left( \boldsymbol{x}^k - \tfrac{1}{L}(\boldsymbol{A}^\top \boldsymbol{A} + \lambda_1 \boldsymbol{I})\boldsymbol{x}^k + \tfrac{1}{L} \boldsymbol{A}^\top \boldsymbol{b} \right), \end{aligned}$$

  with $L = \left\| \boldsymbol{A}^\top \boldsymbol{A} + \lambda_1 \boldsymbol{I} \right\| = \lambda_{\max}(\boldsymbol{A}^\top \boldsymbol{A}) + \lambda_1 = \|\boldsymbol{A}\|_2^2 + \lambda_1$.

- (*FISTA*)

$$\begin{cases} \boldsymbol{x}^{k+1} &= \mathcal{T}_{\frac{\lambda_2}{L}} \left( \boldsymbol{y}^k - \tfrac{1}{L}(\boldsymbol{A}^\top \boldsymbol{A} + \lambda_1 \boldsymbol{I})\boldsymbol{y}^k + \tfrac{1}{L} \boldsymbol{A}^\top \boldsymbol{b} \right) \\ t_{k+1} &= \frac{1+\sqrt{1+4t_k^2}}{2} \\ \boldsymbol{y}^{k+1} &= \boldsymbol{x}^{k+1} + \left( \frac{t_k - 1}{t_{k+1}} \right) (\boldsymbol{x}^{k+1} - \boldsymbol{x}^k) \end{cases}$$

- (*V-FISTA*)

$$\begin{cases} \boldsymbol{x}^{k+1} &= \mathcal{T}_{\frac{\lambda_2}{L}} \left( \boldsymbol{y}^k - \tfrac{1}{L}(\boldsymbol{A}^\top \boldsymbol{A} + \lambda_1 \boldsymbol{I})\boldsymbol{y}^k + \tfrac{1}{L} \boldsymbol{A}^\top \boldsymbol{b} \right) \\ \boldsymbol{y}^{k+1} &= \boldsymbol{x}^{k+1} + \left( \frac{\sqrt{\kappa} - 1}{\sqrt{\kappa} + 1} \right) (\boldsymbol{x}^{k+1} - \boldsymbol{x}^k), \end{cases}$$

  with $\kappa = L/\sigma$, and $\sigma = 1$.

Now, if one chooses $f(\boldsymbol{x}) = \|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{b}\|_2^2$ as the differentiable (but not strongly convex) part and $\alpha g(\boldsymbol{x}) = \frac{\alpha \lambda_1}{2} \|\boldsymbol{x}\|_2^2 + \alpha \lambda_2 \|\boldsymbol{x}\|_1$ as the closed convex but non differentiable part, one has $\nabla f(\boldsymbol{x}) = \boldsymbol{A}^\top(\boldsymbol{A}\boldsymbol{x} - \boldsymbol{b})$ and, by identification with $g'(\boldsymbol{x}) + \frac{c}{2} \|\boldsymbol{x}\|_2^2 + \langle \boldsymbol{a}, \boldsymbol{x} \rangle + \gamma$ with $g'(\boldsymbol{x}) = \alpha \lambda_2 \|\boldsymbol{x}\|_1$, $c = \alpha \lambda_1$, $\boldsymbol{a} = \boldsymbol{0}$, $\gamma = 0$, $\operatorname{prox}_{\alpha g}(\boldsymbol{x}) = \operatorname{prox}_{\frac{\alpha}{\alpha \lambda_1 + 1} \lambda_2 \|\cdot\|_1} \left( \frac{\boldsymbol{x}}{\alpha \lambda_1 + 1} \right) = \mathcal{T}_{\frac{\alpha \lambda_2}{\alpha \lambda_1 + 1}} \left( \frac{\boldsymbol{x}}{\alpha \lambda_1 + 1} \right)$.

- (*V-FISTA2*)

$$\begin{cases} \boldsymbol{x}^{k+1} & = \mathcal{T}_{\frac{\lambda_2/L_2}{\lambda_1/L_2+1}} \left( \boldsymbol{y}^k - \frac{1}{L_2}(\boldsymbol{A}^\top \boldsymbol{A})\boldsymbol{y}^k + \frac{1}{L}\boldsymbol{A}^\top \boldsymbol{b} \right) \\ \boldsymbol{y}^{k+1} & = \boldsymbol{x}^{k+1} + \left( \frac{\sqrt{\kappa_2}-1}{\sqrt{\kappa_2}+1} \right)(\boldsymbol{x}^{k+1} - \boldsymbol{x}^k), \end{cases}$$

with $L_2 = \|\boldsymbol{A}\|_2^2$, $\kappa_2 = L_2/\sigma$, and $\sigma = 1$.

One notices the only difference between *V-FISTA* and *V-FISTA2* occurs in the second line, with $\kappa_2 \neq \kappa$.
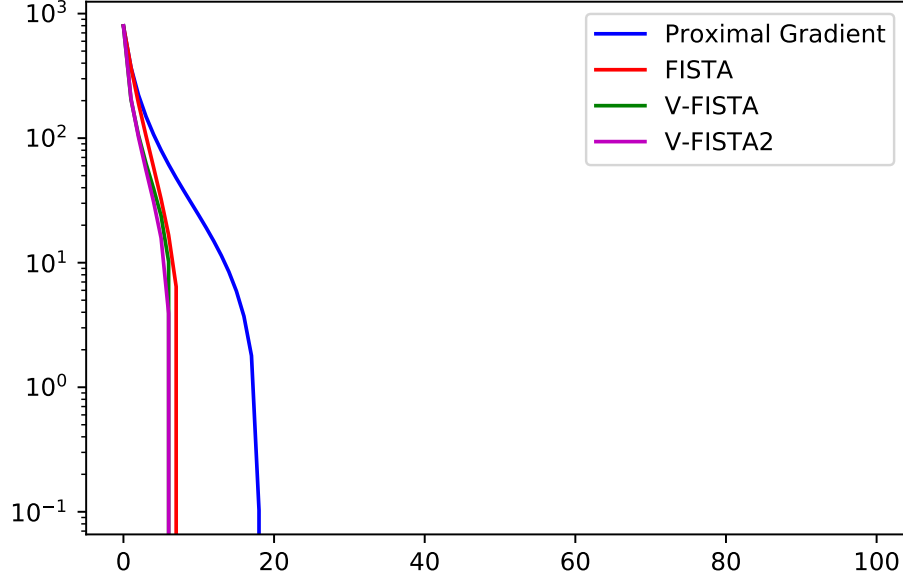


Figure 1: $F(\boldsymbol{x}^k) - F_{\mathrm{opt}}$ in log-scale along the y-axis for the first 100 iterations of each of the methods with all-zeros vectors and a constant stepsize.

As can be observed in Fig. 1, the V-FISTA methods worked the best, with a slight improvement in V-FISTA which does not follow the theory. The four first elements are:

- (-0.40403765 0.18475212 0.97264407 -0.99645397)

- (-0.43969331 0.01974521 1.42280231 -0.87819581)

- (-0.4319773 0.02881602 1.43373682 -0.9066518 )

- (-0.43207117 0.02954933 1.43438138 -0.9058778 )

- (-0.43210892 0.02959727 1.43437048 -0.9058291 )

for the Ground truth, and the four methods, PG, FISTA, V-FISTA, V-FISTA2, respectively.

# 3    Part 2 - Exercise 1 - p41

$$\min_{\boldsymbol{x}\in\mathbb{R}^{30}} \sqrt{\boldsymbol{x}^\top \boldsymbol{Q}\boldsymbol{x} + 2\boldsymbol{b}^\top \boldsymbol{x} + c} + 0.2\,\|\boldsymbol{D}\boldsymbol{x}\|_1$$