

Convex Optimization with First Order Schemes

Gabriel Peyré



www.numerical-tours.com

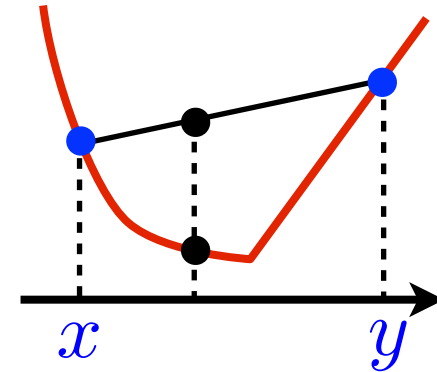


Convex Optimization

Setting: $G : \mathcal{H} \rightarrow \mathbb{R} \cup \{+\infty\}$

\mathcal{H} : Hilbert space. Here: $\mathcal{H} = \mathbb{R}^N$.

$$\text{Problem: } \min_{x \in \mathcal{H}} G(x)$$

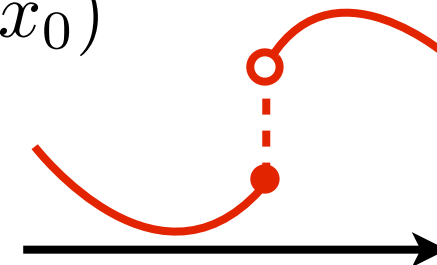


Class of functions:

$$\text{Convex: } G(tx + (1-t)y) \leq tG(x) + (1-t)G(y) \quad t \in [0, 1]$$

$$\text{Lower semi-continuous: } \liminf_{x \rightarrow x_0} G(x) \geq G(x_0)$$

$$\text{Proper: } \{x \in \mathcal{H} \mid G(x) \neq +\infty\} \neq \emptyset$$



$$\text{Indicator: } \iota_{\mathcal{C}}(x) = \begin{cases} 0 & \text{if } x \in \mathcal{C}, \\ +\infty & \text{otherwise.} \end{cases}$$

(\mathcal{C} closed and convex)

Example: ℓ^1 Regularization

Inverse problem: measurements $y = \mathcal{K}f_0 + w$

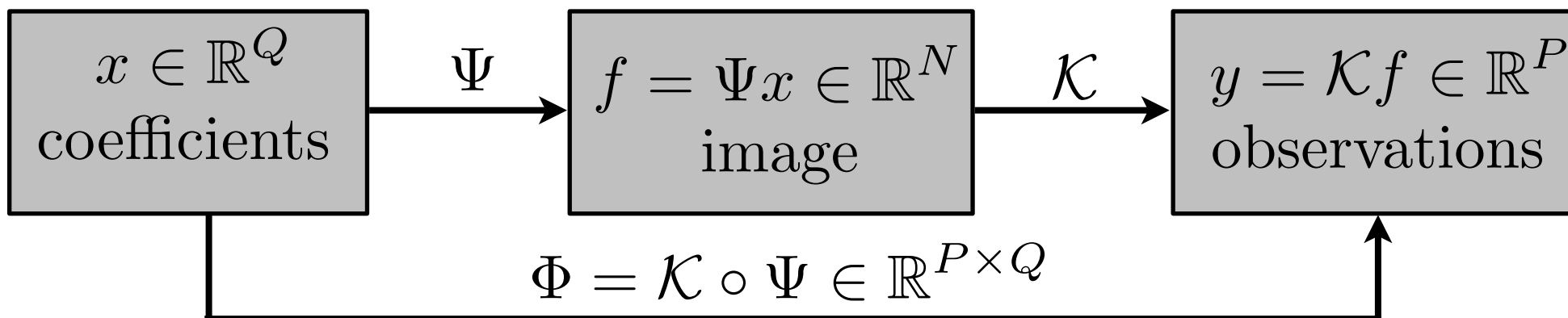


\mathcal{K}



$$\mathcal{K} : \mathbb{R}^N \rightarrow \mathbb{R}^P, \quad P \leq N$$

Model: $f_0 = \Psi x_0$ sparse in dictionary $\Psi \in \mathbb{R}^{N \times Q}$, $Q \geq N$.



Sparse recovery: $f^* = \Psi x^*$ where x^* solves

$$\min_{x \in \mathbb{R}^N} \frac{1}{2} \|y - \Phi x\|^2 + \lambda \|x\|_1$$

Fidelity

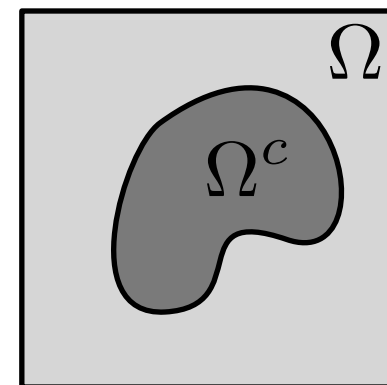
Regularization

Example: ℓ^1 Regularization

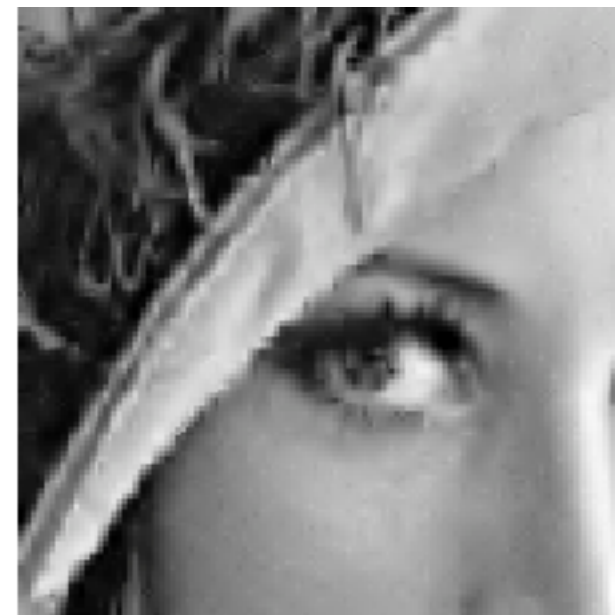
Inpainting: masking operator \mathcal{K}

$$(\mathcal{K}f)_i = \begin{cases} f_i & \text{if } i \in \Omega, \\ 0 & \text{otherwise.} \end{cases}$$

$$\mathcal{K} : \mathbb{R}^N \rightarrow \mathbb{R}^P \quad P = |\Omega|$$



$\Psi \in \mathbb{R}^{N \times Q}$ translation invariant wavelet frame.



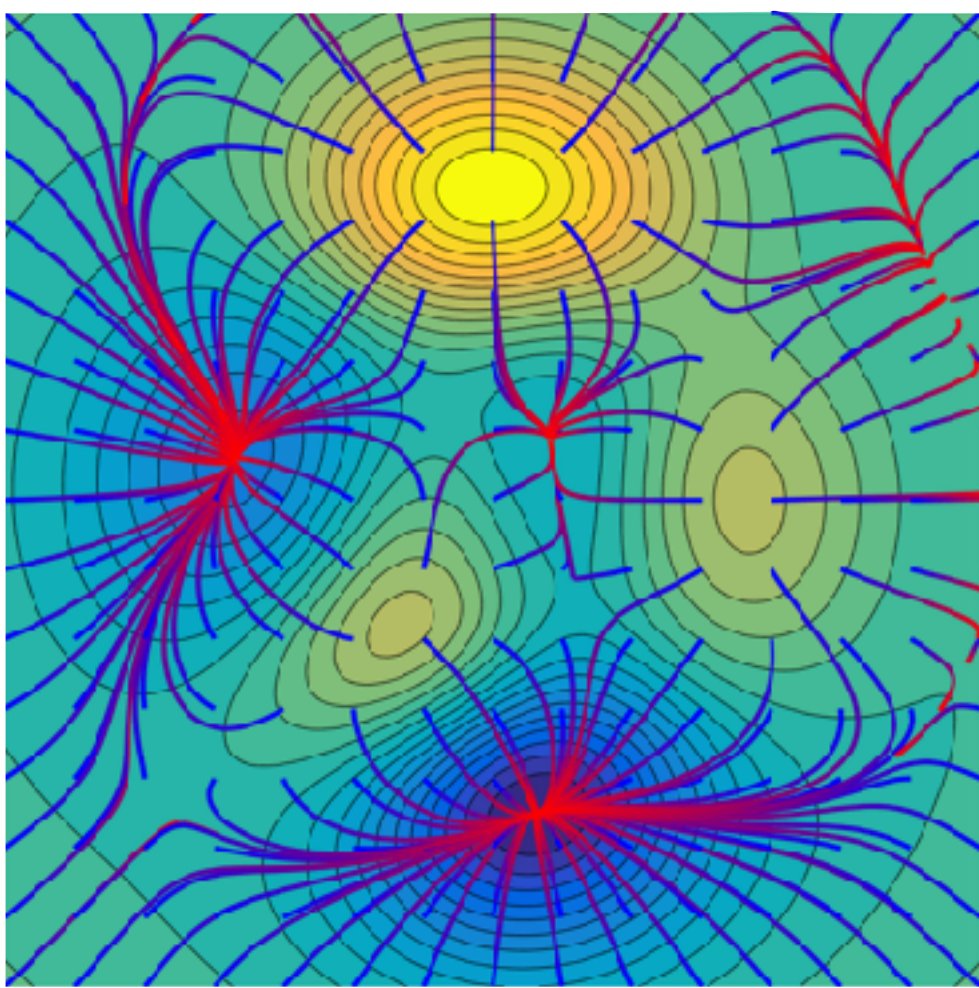
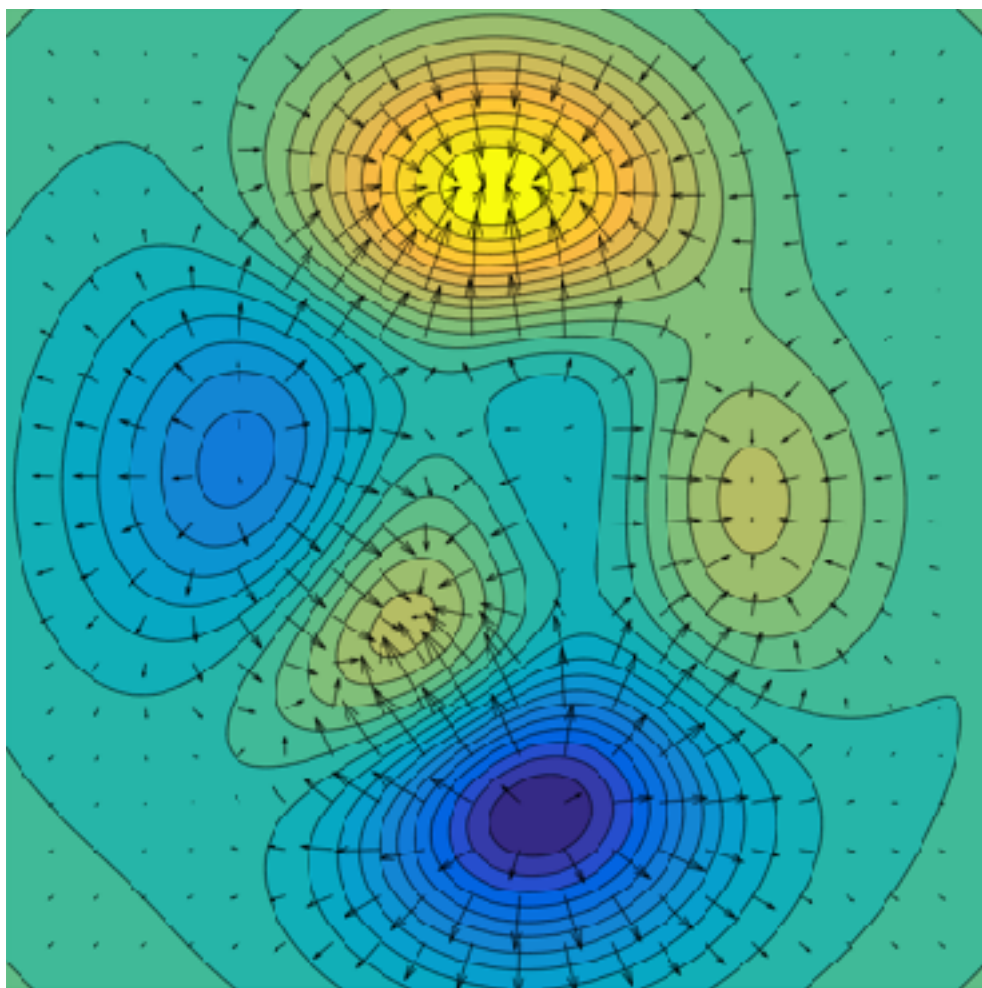
Original $f_0 = \Psi x_0$

$$y = \Phi x_0 + w$$

Recovery Ψx^*

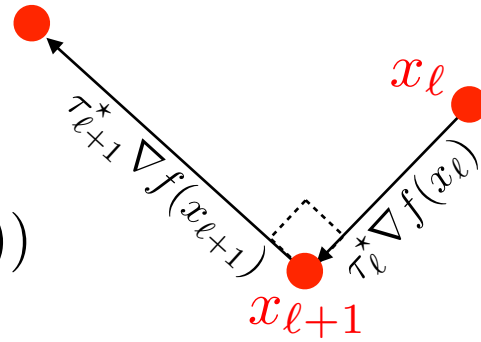
Overview

- **Smooth optimization**
- Subdifferential Calculus
- Proximal Calculus
- Forward Backward

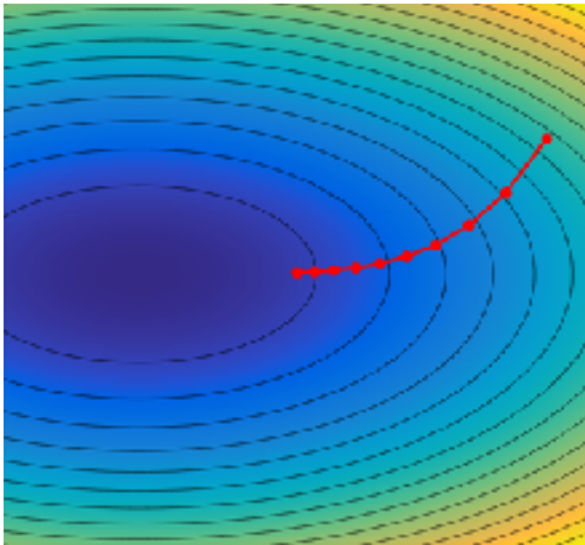


$$x_{l+1} = x_l - \tau_l \nabla f(x_l)$$

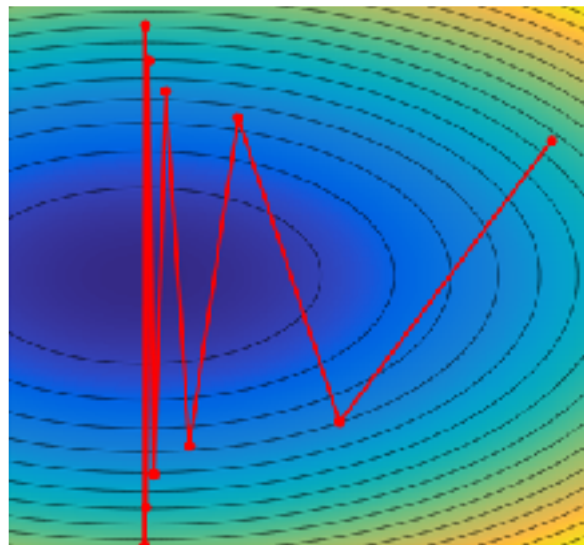
$$\tau_l^* = \operatorname{argmin}_{\tau} f(x_l - \tau \nabla f(x_l))$$



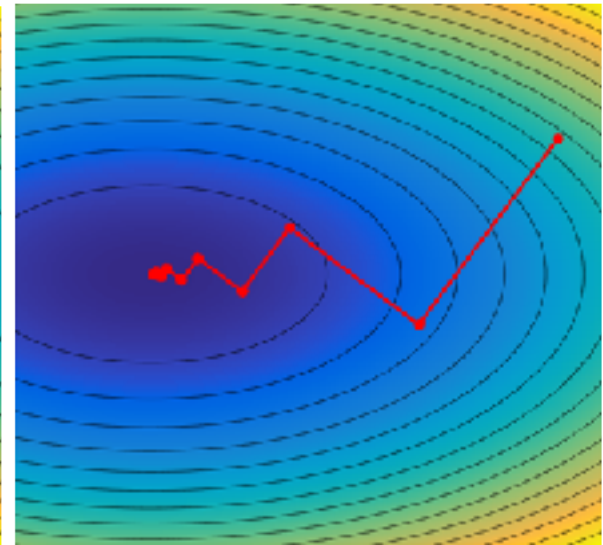
$$\nabla f(x_l) \perp \nabla f(x_{l+1})$$



Small τ_l

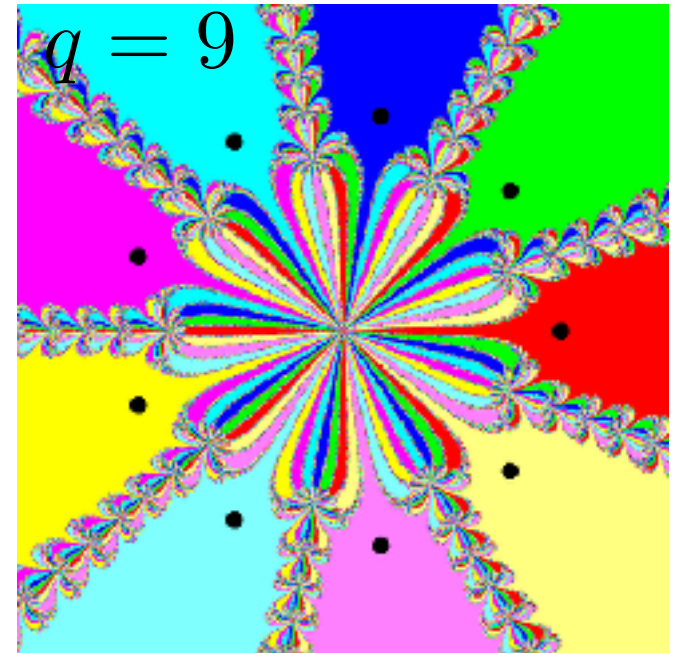
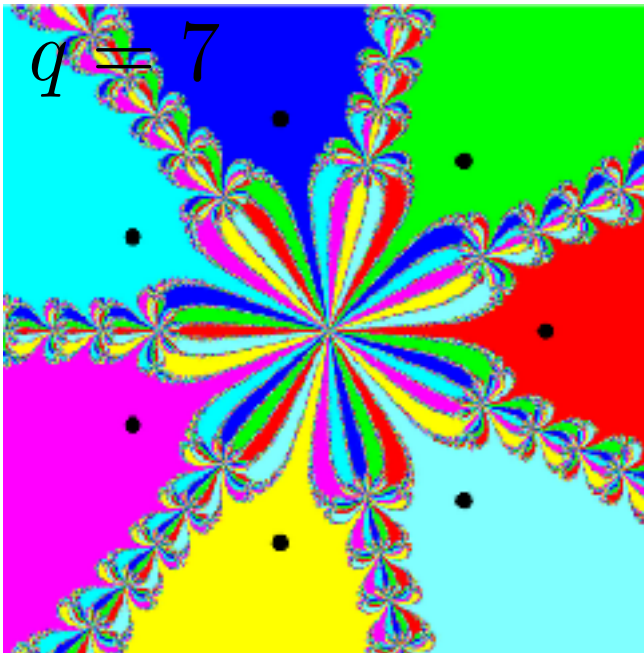
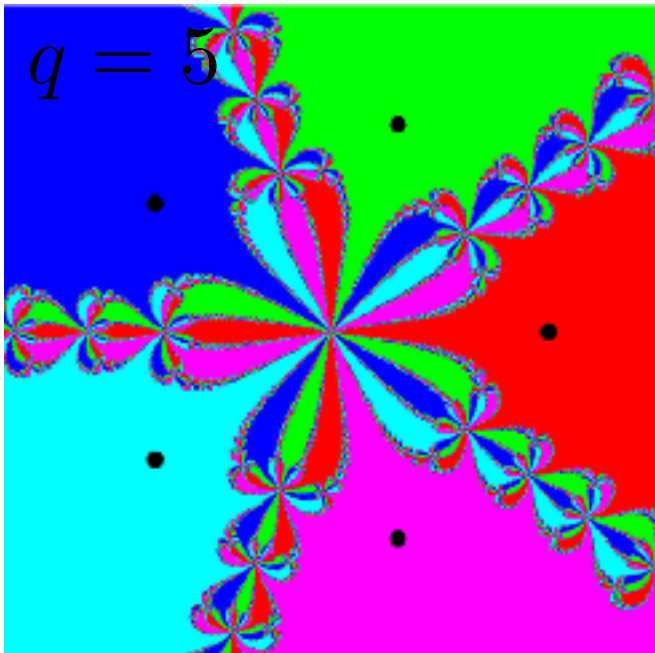


Large τ_l



Optimal $\tau_l = \tau_l^*$

Newton method: $z_{k+1} = z_k - \frac{f(z_k)}{f'(z_k)}$



Attraction bassins for $f(z) = z^q - 1$

Overview

- Smooth optimization
- **Subdifferential Calculus**
- Proximal Calculus
- Forward Backward

Sub-differential

Sub-differential:

$$\partial G(x) = \{u \in \mathcal{H} \mid \forall z, G(z) \geq G(x) + \langle u, z - x \rangle\}$$

Smooth functions:

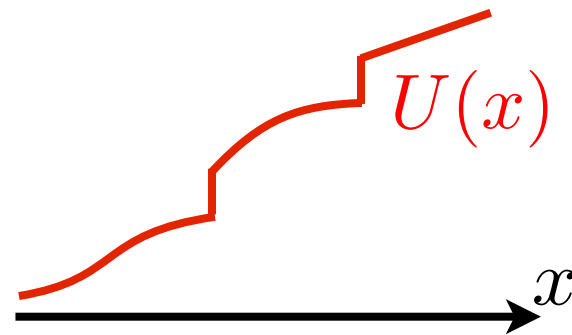
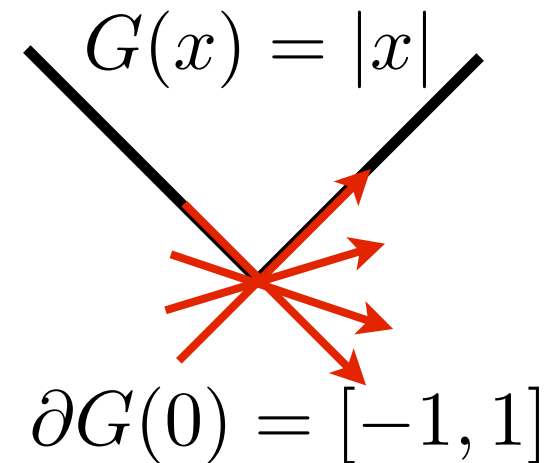
$$\text{If } F \text{ is } C^1, \partial F(x) = \{\nabla F(x)\}$$

First-order conditions:

$$x^* \in \operatorname{argmin}_{x \in \mathcal{H}} G(x) \iff 0 \in \partial G(x^*)$$

Monotone operator: $U(x) = \partial G(x)$

$$\forall (u, v) \in U(x) \times U(y), \quad \langle y - x, v - u \rangle \geq 0$$



Example: ℓ^1 Regularization

$$x^* \in \operatorname{argmin}_{x \in \mathbb{R}^Q} G(x) = \frac{1}{2} \|y - \Phi x\|^2 + \lambda \|x\|_1$$

$$\partial G(x) = \Phi^* (\Phi x - y) + \lambda \partial \|\cdot\|_1(x)$$

$$\partial \|\cdot\|_1(x)_i = \begin{cases} \operatorname{sign}(x_i) & \text{if } x_i \neq 0, \\ [-1, 1] & \text{if } x_i = 0. \end{cases}$$

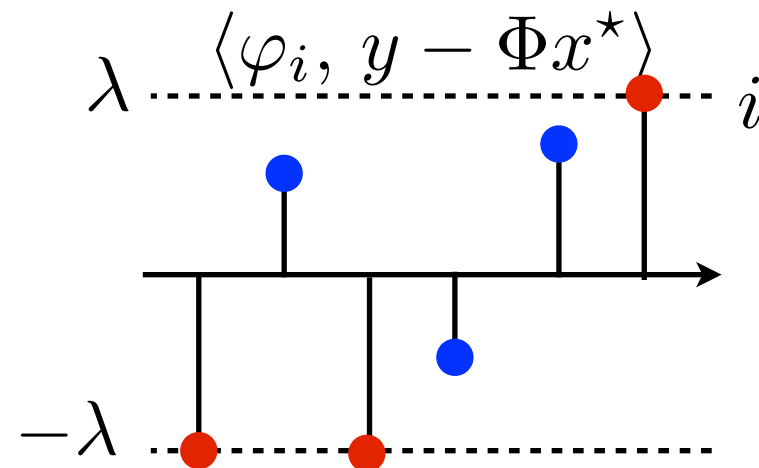
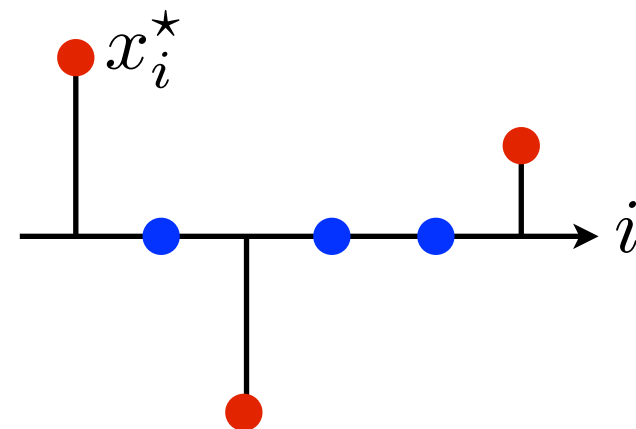
Support of the solution: ●

$$I = \{i \in \{0, \dots, N-1\} \mid x_i^* \neq 0\}$$

First-order conditions:

$$\exists s \in \mathbb{R}^N, \quad \Phi^* (\Phi x^* - y) + \lambda s = 0$$

$$\begin{cases} s_I = \operatorname{sign}(x_I), \\ \|s_{I^c}\|_\infty \leq 1. \end{cases}$$



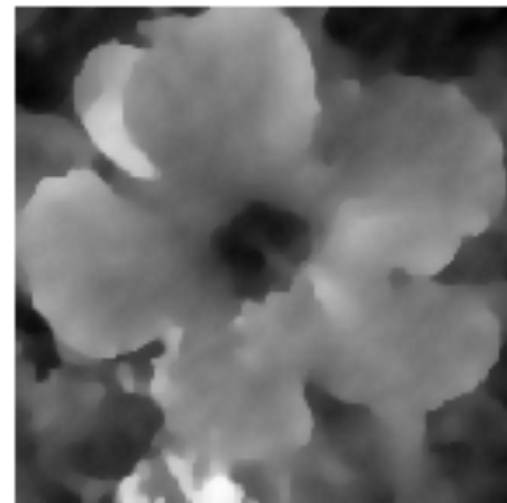
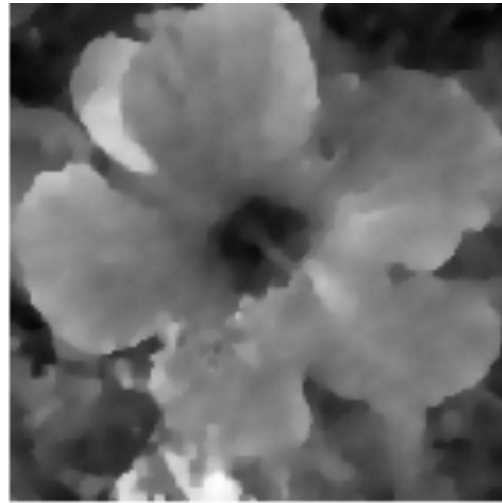
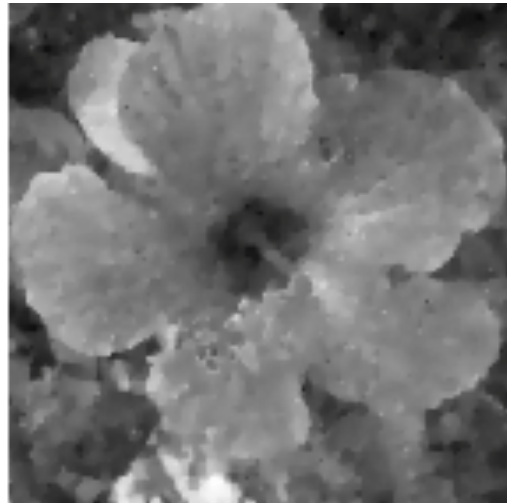
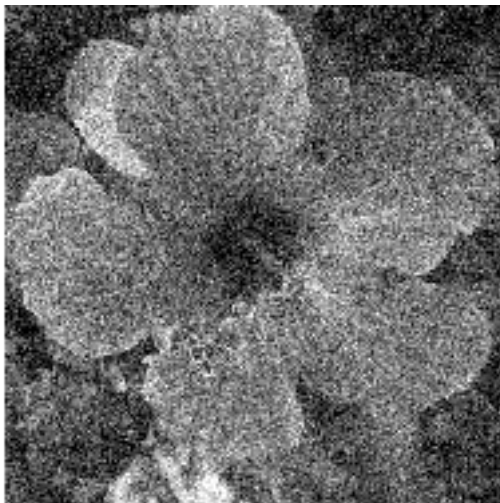
Example: Total Variation Denoising

Important: the optimization variable is f .

$$f^* \in \operatorname{argmin}_{f \in \mathbb{R}^N} \frac{1}{2} \|y - f\|^2 + \lambda J(f)$$

Finite difference gradient: $\nabla : \mathbb{R}^N \rightarrow \mathbb{R}^{N \times 2} \quad (\nabla f)_i \in \mathbb{R}^2$

Discrete TV norm: $J(f) = \sum_i \|(\nabla f)_i\|$



$\lambda = 0$ (noisy)



Example: Total Variation Denoising

$$f^* \in \operatorname{argmin}_{f \in \mathbb{R}^N} \frac{1}{2} \|y - f\|^2 + \lambda J(f)$$

$$J(f) = G(\nabla f) \quad G(u) = \sum_i \|u_i\|$$

Composition by linear maps: $\partial(J \circ A) = A^* \circ (\partial J) \circ A$

$$\partial J(f) = -\operatorname{div}(\partial G(\nabla f))$$

$$\partial G(u)_i = \begin{cases} \frac{u_i}{\|u_i\|} & \text{if } u_i \neq 0, \\ \{\eta \in \mathbb{R}^2 \mid \|\eta\| \leq 1\} & \text{if } u_i = 0. \end{cases}$$

First-order conditions: $\exists v \in \mathbb{R}^{N \times 2}, f^* = y + \lambda \operatorname{div}(v)$

$$\begin{cases} \forall i \in I, v_i = \frac{\nabla f_i^*}{\|\nabla f_i^*\|}, \\ \forall i \in I^c, \|v_i\| \leq 1 \end{cases} \quad I = \{i \mid (\nabla f^*)_i \neq 0\}$$

Overview

- Smooth optimization
- Subdifferential Calculus
- **Proximal Calculus**
- Forward Backward

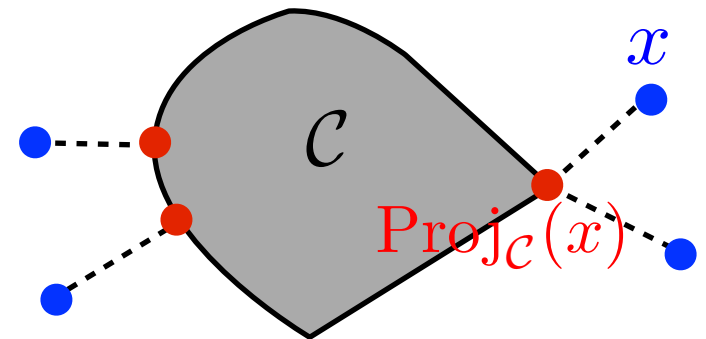
Proximal Operators

Proximal operator of G :

$$\text{Prox}_{\gamma G}(x) = \underset{z}{\operatorname{argmin}} \frac{1}{2} \|x - z\|^2 + \gamma G(z)$$

Indicators: $G(x) = \iota_{\mathcal{C}}(x)$

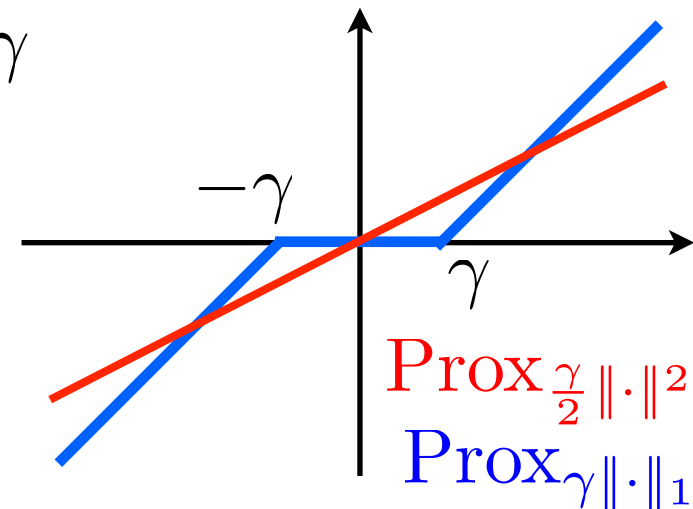
$$\begin{aligned} \text{Prox}_{\gamma G}(x) &= \text{Proj}_{\mathcal{C}}(x) \\ &= \underset{z \in \mathcal{C}}{\operatorname{argmin}} \|x - z\| \end{aligned}$$



ℓ^2 norm squared: $\text{Prox}_{\frac{\gamma}{2} \|\cdot\|^2}(x) = \frac{x}{1 + \gamma}$

ℓ^1 norm: $G(x) = \|x\|_1 = \sum_i |x_i|$

$$\text{Prox}_{\gamma G}(x)_i = \max \left(0, 1 - \frac{\gamma}{|x_i|} \right) x_i$$



Proximal Calculus

Separability: $G(x) = G_1(x_1) + \dots + G_n(x_n)$

$$\text{Prox}_G(x) = (\text{Prox}_{G_1}(x_1), \dots, \text{Prox}_{G_n}(x_n))$$

Quadratic functionals: $G(x) = \frac{1}{2} \|\Phi x - y\|^2$

$$\begin{aligned} \text{Prox}_{\gamma G} &= (\text{Id} + \gamma \Phi^* \Phi)^{-1} \Phi^* \\ &= \Phi^* (\text{Id} + \gamma \Phi \Phi^*)^{-1} \end{aligned}$$

Composition by tight frame: $A \circ A^* = \text{Id}$

$$\text{Prox}_{G \circ A}(x) = A^* \circ \text{Prox}_G \circ A + \text{Id} - A^* \circ A$$

Ortho-basis A: $\text{Prox}_{G \circ A} = A^* \circ \text{Prox}_G \circ A$

Non-convex Proximal Operators

Proximal operator of G :

$$\text{Prox}_{\gamma G}(x) = \underset{z}{\operatorname{argmin}} \frac{1}{2} \|x - z\|^2 + \gamma G(z)$$

$$G(x) = \|x\|_1 = \sum_i |x_i|$$

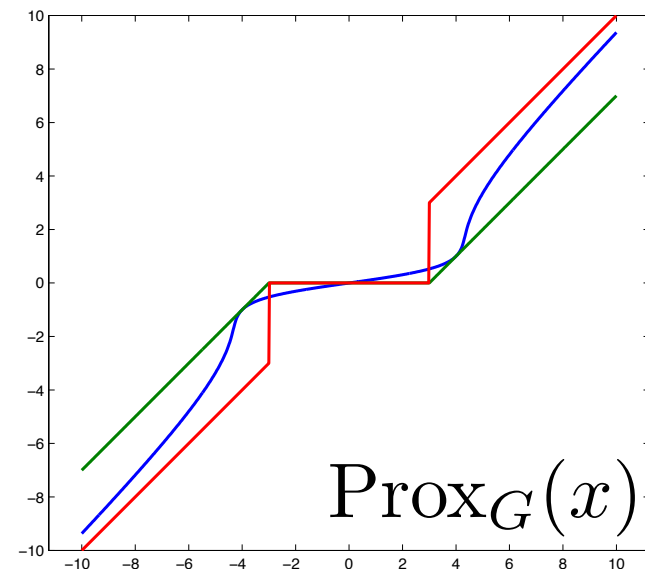
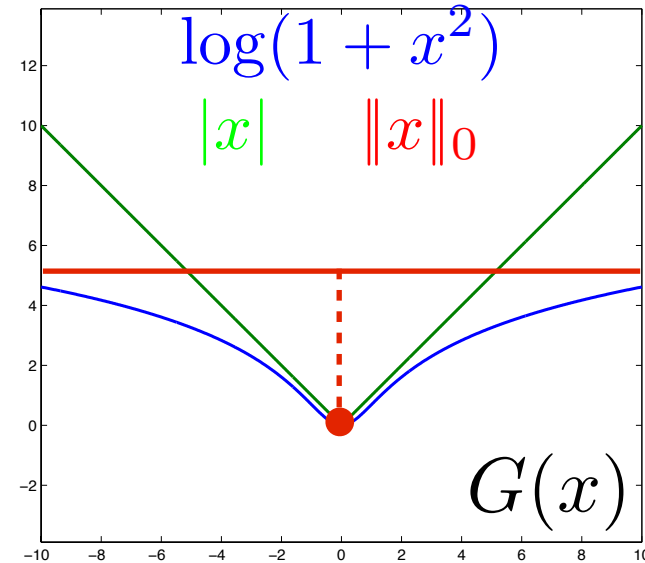
$$\text{Prox}_{\gamma G}(x)_i = \max \left(0, 1 - \frac{\gamma}{|x_i|} \right) x_i$$

$$G(x) = \|x\|_0 = |\{i \mid x_i \neq 0\}|$$

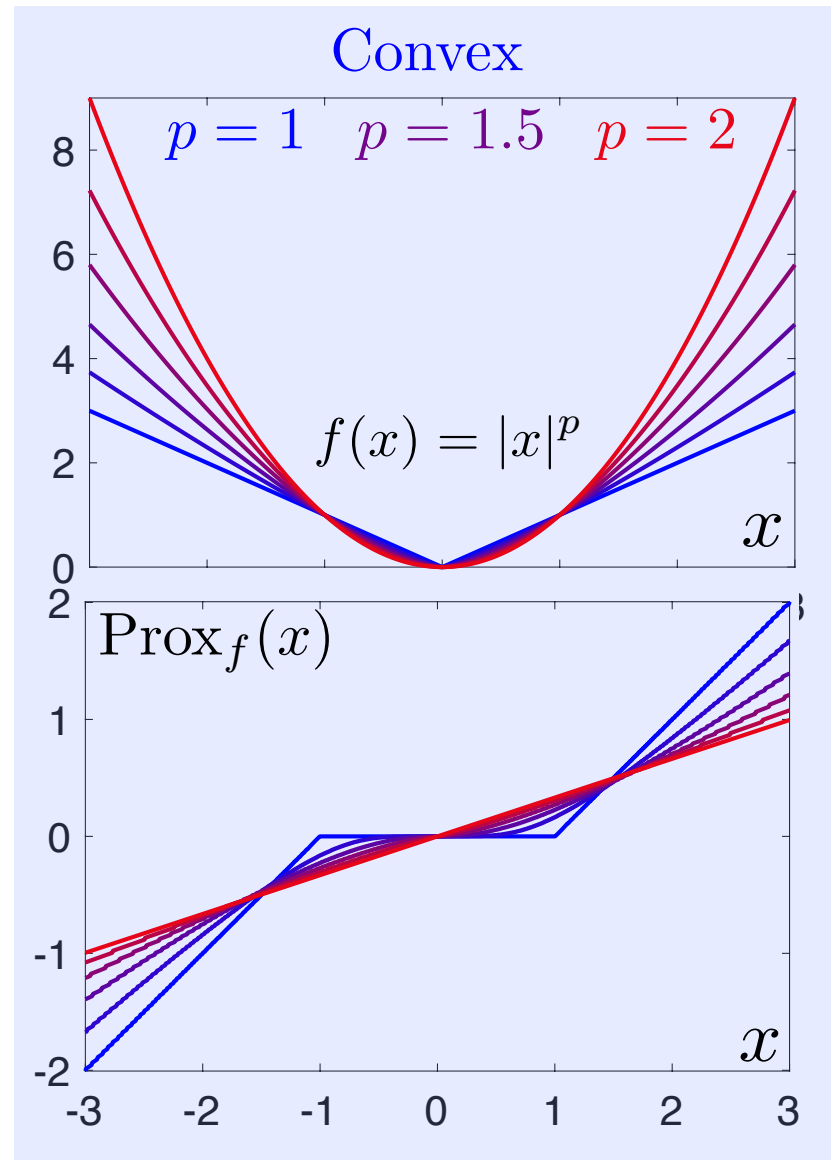
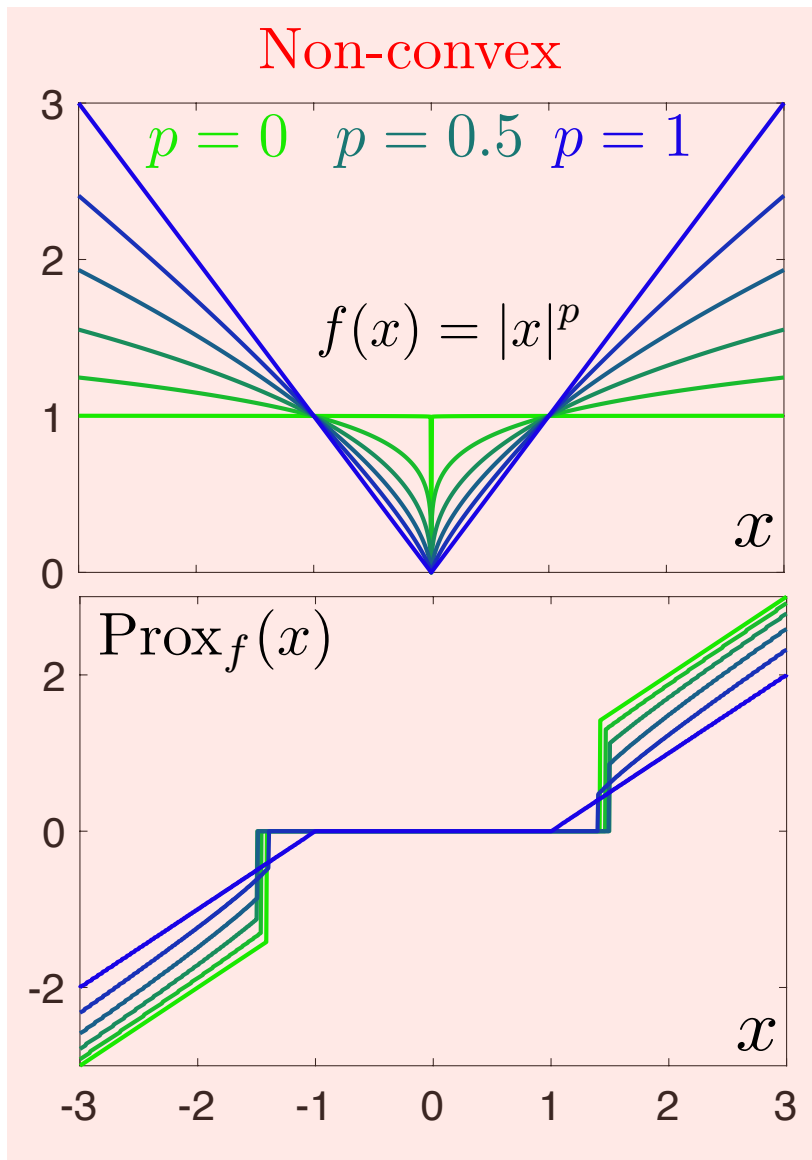
$$\text{Prox}_{\gamma G}(x)_i = \begin{cases} x_i & \text{if } |x_i| \geq \sqrt{2\gamma}, \\ 0 & \text{otherwise.} \end{cases}$$

$$G(x) = \sum_i \log(1 + |x_i|^2)$$

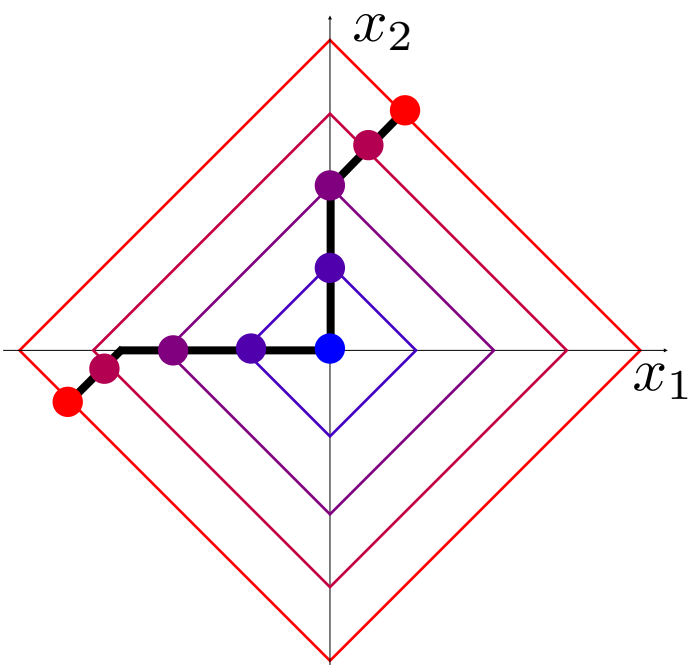
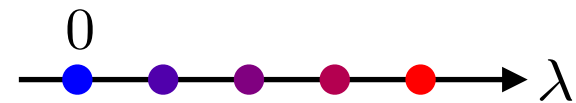
→ 3rd order polynomial root.



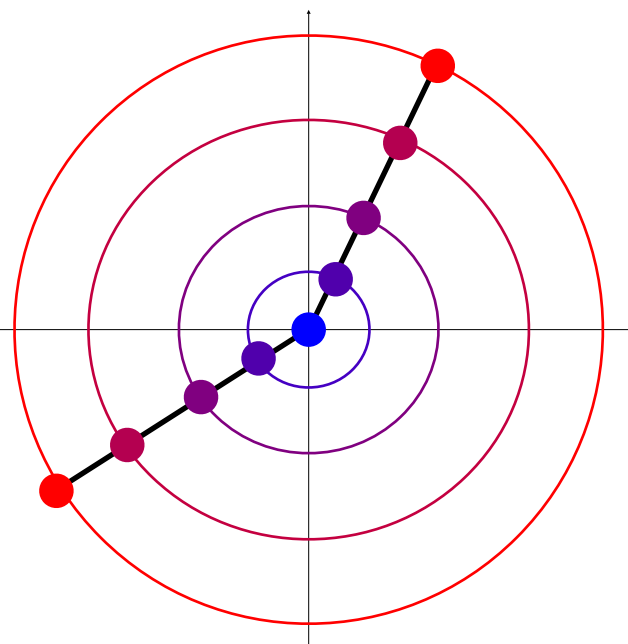
$$\text{Prox}_f(x) = \underset{x'}{\operatorname{argmin}} \frac{1}{2} \|x - x'\|^2 + f(x')$$



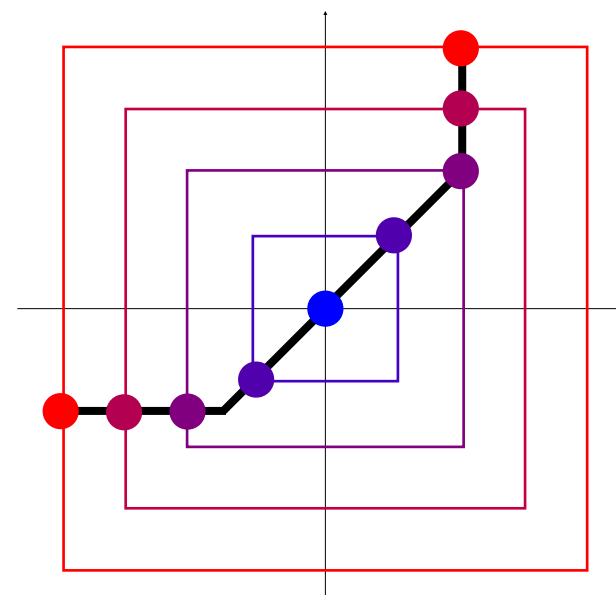
$$\text{Prox}_{\lambda f}(x) = \underset{x'}{\operatorname{argmin}} \frac{1}{2} \|x - x'\|^2 + \lambda f(x')$$



$$f(x) = |x_1| + |x_2|$$



$$f(x) = \sqrt{|x_1|^2 + |x_2|^2}$$



$$f(x) = \max(|x_1|, |x_2|)$$

Prox and Subdifferential

Resolvent of ∂G :

$$\begin{aligned} z = \text{Prox}_{\gamma G}(x) &\iff 0 \in z - x + \gamma \partial G(z) \\ \iff x \in (\text{Id} + \gamma \partial G)(z) &\iff z = (\text{Id} + \gamma \partial G)^{-1}(x) \end{aligned}$$

Inverse of a set-valued mapping:

$$\text{where } x \in U(y) \iff y \in U^{-1}(x)$$

$\text{Prox}_{\gamma G} = (\text{Id} + \gamma \partial G)^{-1}$ is a single-valued mapping

Fix point: $x^* \in \underset{x}{\text{argmin}} G(x)$

$$\iff 0 \in \partial G(x^*) \iff x^* \in (\text{Id} + \gamma \partial G)(x^*)$$

$$\iff x^* = (\text{Id} + \gamma \partial G)^{-1}(x^*) = \text{Prox}_{\gamma G}(x^*)$$

Gradient and Proximal Descents

Gradient descent: $x^{(\ell+1)} = x^{(\ell)} - \gamma_\ell \nabla G(x^{(\ell)})$ [explicit]

G is C^1 and ∇G is L -Lipschitz

Theorem: If $0 < \gamma_\ell < 2/L$, $x^{(\ell)} \rightarrow x^*$ a solution.

Sub-gradient descent: $x^{(\ell+1)} = x^{(\ell)} - \gamma_\ell v^{(\ell)}$, $v^{(\ell)} \in \partial G(x^{(\ell)})$

Theorem: If $\gamma_\ell \sim 1/\ell$, $x^{(\ell)} \rightarrow x^*$ a solution.

→ Problem: slow.

Proximal-point algorithm: $x^{(\ell+1)} = \text{Prox}_{\gamma_\ell G}(x^{(\ell)})$ [implicit]

Theorem: If $\gamma_\ell \geq c > 0$, $x^{(\ell)} \rightarrow x^*$ a solution.

→ $\text{Prox}_{\gamma G}$ hard to compute.

Overview

- Smooth optimization
- Subdifferential Calculus
- Proximal Calculus
- **Forward Backward**

Proximal Splitting Methods

Solve $\min_{x \in \mathcal{H}} E(x)$

Problem: $\text{Prox}_{\gamma E}$ is not available.

Splitting: $E(x) = \boxed{F(x)} + \sum_i \boxed{G_i(x)}$
Smooth Simple

Iterative algorithms using: $\begin{cases} \nabla F(x) \\ \text{Prox}_{\gamma G_i}(x) \end{cases}$

Forward-Backward: $\xrightarrow{\text{solves}} F + G$
Douglas-Rachford: $\longrightarrow \sum G_i$
Primal-Dual: $\longrightarrow \sum G_i \circ A$
Generalized FB: $\longrightarrow F + \sum G_i$

Smooth + Simple Splitting

Inverse problem: measurements $y = \mathcal{K}f_0 + w$



$$\mathcal{K} : \mathbb{R}^N \rightarrow \mathbb{R}^P, \quad P \leq N$$

Model: $f_0 = \Psi x_0$ sparse in dictionary Ψ .

Sparse recovery: $f^* = \Psi x^*$ where x^* solves

$$\min_{x \in \mathbb{R}^N} \boxed{F(x)} + \boxed{G(x)}$$

Smooth Simple

Data fidelity: $F(x) = \frac{1}{2} \|y - \Phi x\|^2$ $\Phi = \mathcal{K} \circ \Psi$

Regularization: $G(x) = \|x\|_1 = \sum_i |x_i|$

Forward-Backward

Fix point equation:

$$\begin{aligned}x^* \in \operatorname{argmin}_x F(x) + G(x) &\iff 0 \in \nabla F(x^*) + \partial G(x^*) \\ &\iff (x^* - \gamma \nabla F(x^*)) \in x^* + \gamma \partial G(x^*) \\ &\iff x^* = \operatorname{Prox}_{\gamma G}(x^* - \gamma \nabla F(x^*))\end{aligned}$$

Forward-backward:

$$x^{(\ell+1)} = \operatorname{Prox}_{\gamma G} \left(x^{(\ell)} - \gamma \nabla F(x^{(\ell)}) \right)$$

Projected gradient descent:

$$G = \iota_C$$

Theorem: Let ∇F be L -Lipschitz.

If $\gamma < 2/L$, $x^{(\ell)} \rightarrow x^*$ a solution of (\star)

Example: L1 Regularization

$$\min_x \frac{1}{2} \|\Phi x - y\|^2 + \lambda \|x\|_1 \iff \min_x F(x) + G(x)$$

$$F(x) = \frac{1}{2} \|\Phi x - y\|^2$$

$$\nabla F(x) = \Phi^* (\Phi x - y) \qquad L = \|\Phi^* \Phi\|$$

$$G(x) = \lambda \|x\|_1$$

$$\text{Prox}_{\gamma G}(x)_i = \max \left(0, 1 - \frac{\gamma \lambda}{|x_i|} \right) x_i$$

Forward-backward \iff Iterative soft thresholding

Convergence Speed

$$\min_x E(x) = F(x) + G(x)$$

∇F is L -Lipschitz.

G is simple.

Theorem: If $L > 0$, FB iterates $x^{(\ell)}$ satisfies

$$E(x^{(\ell)}) - E(x^*) = O(1/\ell)$$

C degrades with $L \rightarrow 0$.

Multi-steps Accelerations

Beck-Teboule accelerated FB: $t^{(0)} = 1$

$$\begin{aligned}x^{(\ell+1)} &= \text{Prox}_{1/L} \left(y^{(\ell)} - \frac{1}{L} \nabla F(y^{(\ell)}) \right) \\t^{(\ell+1)} &= \frac{1 + \sqrt{1 + 4(t^{(\ell)})^2}}{2} \\y^{(\ell+1)} &= x^{(\ell+1)} + \frac{t^{(\ell)} - 1}{t^{(\ell+1)}} (x^{(\ell+1)} - x^{(\ell)})\end{aligned}$$

(see also Nesterov method)

Theorem: If $L > 0$, $E(x^{(\ell)}) - E(x^*) = O(1/\ell^2)$

Complexity theory: optimal in a worse-case sense.