Proximal Gradient Method. Proximal operator

Seminar

Optimization for ML. Faculty of Computer Science. HSE University



Regularized / Composite Objectives

Many nonsmooth problems take the form

$$\min_{x\in\mathbb{R}^n}\varphi(x)=f(x)+r(x)$$

• Lasso, L1-LS, compressed sensing

$$f(x) = \frac{1}{2} ||Ax - b||_2^2, r(x) = \lambda ||x||_1$$

• L1-Logistic regression, sparse LR

$$f(x) = -y \log h(x) - (1-y) \log(1-h(x)), r(x) = \lambda ||x||_1$$



Non-smooth convex optimization lower bounds

convex (non-smooth)	strongly convex (non-smooth)
$f(x_k) - f^* \sim \mathcal{O}\left(rac{1}{\sqrt{k}} ight)$	$f(x_k) - f^* \sim \mathcal{O}\left(rac{1}{k} ight)$
$k_arepsilon \sim \mathcal{O}\left(rac{1}{arepsilon^2} ight)$	$k_arepsilon \sim \mathcal{O}\left(rac{1}{arepsilon} ight)$

Non-smooth convex optimization lower bounds

convex (non-smooth)	strongly convex (non-smooth)
$f(x_k) - f^* \sim \mathcal{O}\left(\frac{1}{\sqrt{k}}\right)$	$f(x_k) - f^* \sim \mathcal{O}\left(\frac{1}{k}\right)$ $k \to \mathcal{O}\left(\frac{1}{k}\right)$
$\kappa_{arepsilon} ightarrow \mathcal{O}\left(rac{1}{arepsilon^2} ight)$	$\kappa_{\varepsilon} \sim O\left(\frac{-}{\varepsilon}\right)$

- Subgradient method is optimal for the problems above.
- One can use Mirror Descent (a generalization of the subgradient method to a possiby non-Euclidian distance) with the same convergence rate to better fit the geometry of the problem.
- However, we can achieve standard gradient descent rate $O\left(\frac{1}{k}\right)$ (and even accelerated version $O\left(\frac{1}{k^2}\right)$) if we will exploit the structure of the problem.



Proximal operator

i Proximal operator

For a convex set $E \in \mathbb{R}^n$ and a convex function $f: E \to \mathbb{R}$ operator $\text{prox}_f(x)$ s.t.

$$\operatorname{prox}_{f}(x) = \operatorname*{argmin}_{y \in E} \left[f(y) + \frac{1}{2} ||y - x||_{2}^{2} \right]$$

is called **proximal operator** for function f at point x



Let \mathbb{I}_S be the indicator function for closed, convex S. Recall orthogonal projection $\pi_S(y)$

Let \mathbb{I}_S be the indicator function for closed, convex S. Recall orthogonal projection $\pi_S(y)$

$$\pi_S(y) := \arg\min_{x \in S} \frac{1}{2} ||x - y||_2^2.$$



Let \mathbb{I}_S be the indicator function for closed, convex S. Recall orthogonal projection $\pi_S(y)$

$$\pi_S(y) := \arg\min_{x \in S} \frac{1}{2} \|x - y\|_2^2.$$

With the following notation of indicator function

$$\mathbb{I}_S(x) = egin{cases} 0, & x \in S, \ \infty, & x
otin S, \end{cases}$$

Rewrite orthogonal projection $\pi_S(y)$ as

$$\pi_S(y) := \arg\min_{x \in \mathbb{R}^n} \frac{1}{2} \|x - y\|^2 + \mathbb{I}_S(x).$$



Let \mathbb{I}_S be the indicator function for closed, convex S. Recall orthogonal projection $\pi_S(y)$

$$\pi_S(y) := \arg\min_{x \in S} \frac{1}{2} \|x - y\|_2^2.$$

With the following notation of indicator function

$$\mathbb{I}_{S}(x) = \begin{cases} 0, & x \in S, \\ \infty, & x \notin S, \end{cases}$$

Rewrite orthogonal projection $\pi_S(y)$ as

$$\pi_S(y) := \arg\min_{x \in \mathbb{R}^n} \frac{1}{2} \|x - y\|^2 + \mathbb{I}_S(x).$$

Proximity: Replace \mathbb{I}_S by some convex function!

$${\rm prox}_r(y) = {\rm prox}_{r,1}(y) := \arg\min\frac{1}{2}\|x-y\|^2 + r(x)$$



Proximal Gradient Method

Proximal Gradient Method Theorem

Consider the proximal gradient method

$$x_{k+1} = \operatorname{prox}_{\alpha r} \left(x_k - \alpha \nabla f(x_k) \right)$$

for the criterion $\phi(x) = f(x) + r(x)$ s.t.:

- *f* is convex, differentiable with Lipschitz gradients;
- r is convex and prox-friendly.

Then Proximal Gradient Method with fixed step size $\alpha = \frac{1}{L}$ converges with rate $O(\frac{1}{k})$

Quadratic Function ($A \ge 0$)

$$f(x) = \frac{1}{2}x^T A x + b^T x + c,$$



Quadratic Function ($A \ge 0$)

$$f(x) = \frac{1}{2}x^T A x + b^T x + c$$

$$prox_{tf}(x) = (I + tA)^{-1}(x - tb)$$

Euclidean Norm

 $f(x) = \|x\|_2,$



Quadratic Function ($A \ge 0$)

$$f(x) = \frac{1}{2}x^T A x + b^T x + c,$$

$$prox_{tf}(x) = (I + tA)^{-1}(x - tb)$$

Euclidean Norm

 $f(x) = \|x\|_2,$

$$\operatorname{prox}_{tf}(x) = \begin{cases} (1 - t/\|x\|_2)x & \text{if } \|x\|_2 \geq t \\ 0 & \text{otherwise} \end{cases}$$

Logarithmic Barrier

$$f(x) = -\sum_{i=1}^{n} \log x_i,$$



Quadratic Function ($A \ge 0$)

$$f(x) = \frac{1}{2}x^T A x + b^T x + c,$$

$$prox_{tf}(x) = (I + tA)^{-1}(x - tb)$$

Euclidean Norm

 $f(x) = \|x\|_2,$

$$\operatorname{prox}_{tf}(x) = \begin{cases} (1 - t/\|x\|_2)x & \text{if } \|x\|_2 \geq t \\ 0 & \text{otherwise} \end{cases}$$

Logarithmic Barrier

$$f(x) = -\sum_{i=1}^{n} \log x_i,$$

$$\operatorname{prox}_{tf}(x)_i = \frac{x_i + \sqrt{x_i^2 + 4t}}{2}, \quad i = 1, \dots, n$$



Simple Calculus Rules for Proximal Mappings

Separable Sum

$$f\left(\begin{bmatrix} x \\ y \end{bmatrix} \right) = g(x) + h(y), \quad \operatorname{prox}_f\left(\begin{bmatrix} x \\ y \end{bmatrix} \right) = \begin{bmatrix} \operatorname{prox}_g(x) \\ \operatorname{prox}_h(y) \end{bmatrix}$$

Scaling and Translation of Argument (for $a \neq 0$)

$$f(x) = g(ax+b), \quad \operatorname{prox}_f(x) = \frac{1}{a} \left(\operatorname{prox}_{a^2g}(ax+b) - b \right)$$

Right Scalar Multiplication (with $\lambda > 0$ **)**

$$f(x) = \lambda g(x/\lambda), \quad \operatorname{prox}_f(x) = \lambda \operatorname{prox}_{\lambda^{-1}g}(x/\lambda)$$



ISTA and FISTA

Methods for solving problems involving L1 regularization (e.g. Lasso).

ISTA (Iterative Shrinkage-Thresholding Algorithm)

• Step:

$$x_{k+1} = \operatorname{prox}_{\alpha\lambda||\cdot||_1} \left(x_k - \alpha \nabla f(x_k) \right)$$

• Convergence: $O(\frac{1}{k})$

FISTA (Fast Iterative Shrinkage-Thresholding Algorithm)

• Step:

$$x_{k+1} = \operatorname{prox}_{\alpha\lambda||\cdot||_1} \left(y_k - \alpha \nabla f(y_k) \right),$$





♥ O Ø 9

Problem 1. ReLU in prox

i Question

Find the $\operatorname{prox}_f(x)$ for $f(x) = \lambda \max(0, x)$:

$$\operatorname{prox}_{\lambda \max(0,\cdot)}(x) = \operatorname*{argmin}_{y \in \mathbb{R}} \left[\frac{1}{2} ||y - x||^2 + \lambda \max(0, y) \right]$$



Problem 2. Grouped l_1 -regularizer

i Question

$$\begin{aligned} \text{Find the } \operatorname{prox}_{f}(x) \text{ for } f(x) &= ||x||_{1/2} = \sum_{g=0}^{G} ||x_g||_2 \text{ where } x \in \mathbb{R}^n = \underbrace{[x_1, x_2, \dots, \underbrace{\dots, g}_{1}, \dots, \underbrace{x_{n-2}, x_{n-1}, x_n}_{G}]: \\ & \operatorname{prox}_{||x||_{1/2}}(x) = \operatornamewithlimits{argmin}_{y \in \mathbb{R}} \left[\frac{1}{2} ||y - x||_2^2 + \sum_{g=0}^{G} ||y_g||_2 \right] \end{aligned}$$



Linear Least Squares with *L*₁-regularizer

Proximal Methods Comparison for Linear Least Squares with L_1 -regularizer **\clubsuit**Open in Colab.



Image Denoising using ISTA

Problem: Recover clean image x_{true} from noisy $y = x_{\mathsf{true}} + \mathbf{n}$

Key Idea Clean images have sparse wavelet coefficients $\alpha = Wx_{true}$ (mostly zeros), where W is the Wavelet Transform operator. Noise coefficients are not sparse.

Optimization Approach: Find coefficients α minimizing a composite objective:



where \boldsymbol{W}^T is the inverse Wavelet Transform and $\boldsymbol{\lambda}$ balances the terms.

Result: The denoised image is obtained from the final coefficients α^* :

$$x_{\mathsf{denoised}} = \boldsymbol{W}^T \boldsymbol{\alpha}^*$$



Image Denoising using ISTA Popen Git.





Original with missing pixels

L1 - PSNR = 27.18 SSIM = 0.91 - Time: 1.92s



TV-PSRR = 13.22 SSH = 0.08 - Time 2.461

 $f \rightarrow \min_{x,y,z}$ Colab examples