

Convex Optimization

Part 4: Proximal point method

Namhoon Lee

POSTECH

16 Nov 2022

Admin

Assignment 3

- ▶ grading still in progress

Assignment 4

- ▶ will be posted on PLMS this week

Proximal point method

an algorithm for minimizing a closed convex function f :

$$\begin{aligned}x_{k+1} &= \text{prox}_{t_k f}(x_k) \\ &= \arg \min_u \left(f(u) + \frac{1}{2t_k} \|u - x_k\|_2^2 \right)\end{aligned}$$

- ▶ can be viewed as proximal gradient method with $g(x) = 0$
- ▶ of interest if prox evaluations are much easier than minimizing f directly
- ▶ in practice, inexact prox evaluations may be sufficient
- ▶ step size $t_k > 0$ affects number of iterations, cost of prox evaluations
- ▶ basis of the augmented Lagrangian method

Convergence

Assumptions

- ▶ f is closed and convex (hence, $\text{prox}_{t f}(x)$ is uniquely defined for all x)
- ▶ optimal value f^* is finite and attained at x^*

Result

$$f(x_k) - f^* \leq \frac{\|x_0 - x^*\|_2^2}{2 \sum_{i=0}^{k-1} t_i} \quad \text{for } k \geq 1$$

- ▶ implies convergence if $\sum_i t_i \rightarrow \infty$
- ▶ rate is $1/k$ if t_i is fixed, or variable but bounded away from zero
- ▶ t_i is arbitrary; however cost of prox evaluations will depend on t_i

Proof.

apply analysis of proximal gradient method with $g(x) = 0$; find the lemma for the bound on proximal gradient update



Accelerated proximal point algorithms

- ▶ we take $g(x) = 0$ in FISTA:

$$x_1 = \text{prox}_{t_0 f}(x_0)$$
$$x_{k+1} = \text{prox}_{t_k f} \left(x_k + \theta_k \left(\frac{1}{\theta_{k-1}} - 1 \right) (x_k - x_{k-1}) \right) \quad \text{for } k \geq 1$$

- ▶ choose any $t_k > 0$, determine θ_k from equation

$$\frac{\theta_k^2}{t_k} = (1 - \theta_k) \frac{\theta_{k-1}^2}{t_{k-1}}$$

- ▶ converges if $\sum_i \sqrt{t_i} \rightarrow \infty$
- ▶ rate is $1/k^2$ if t_i is fixed or variable but bounded away from zero

Standard problem format

Primal and dual problem

$$\begin{array}{ll} \text{primal:} & \text{minimize } f(x) + g(Ax) \\ \text{dual:} & \text{maximize } -g^*(z) - f^*(-A^\top z) \end{array}$$

Examples

- ▶ set constraints ($g(y) = \delta_C(y)$):

$$\begin{array}{ll} \text{minimize} & f(x) \\ \text{subject to} & Ax \in C \end{array}$$

- ▶ regularized norm approximation ($g(y) = \|y - b\|$):

$$\text{minimize } f(x) + \|Ax - b\|$$

Augmented Lagrangian method: proximal point method applied to the dual

Proximal mapping of dual function

Definition: proximal mapping of $h(z) = g^*(z) + f^*(-A^\top z)$ is defined as

$$\text{prox}_{th}(z) = \arg \min_u \left(g^*(u) + f^*(-A^\top u) + \frac{1}{2t} \|u - z\|_2^2 \right)$$

Dual expression: $\text{prox}_{th}(z) = z + t(A\hat{x} - \hat{y})$ where

$$(\hat{x}, \hat{y}) = \arg \min_{x,y} \left(f(x) + g(y) + z^\top (Ax - y) + \frac{t}{2} \|Ax - y\|_2^2 \right)$$

- ▶ \hat{x}, \hat{y} minimize the augmented Lagrangian (Lagrangian + quadratic penalty)
- ▶ $f(x) + g(y) + z^\top (Ax - y)$ is Lagrangian of primal problem reformulated as

$$\begin{array}{ll} \text{minimize} & f(x) + g(y) \\ \text{subject to} & Ax - y = 0 \end{array}$$

Proof.

- ▶ write augmented Lagrangian minimization as

$$\begin{aligned} & \text{minimize (over } x, y, w) && f(x) + g(y) + \frac{t}{2} \|w\|_2^2 \\ & \text{subject to} && Ax - y + z/t = w \end{aligned}$$

- ▶ optimality conditions (u is the multiplier for the equality constraint):

$$Ax - y + \frac{1}{t}z = w, \quad -A^\top u \in \partial f(x), \quad u \in \partial g(y), \quad tw = u$$

- ▶ eliminating w gives

$$u = z + t(Ax - y), \quad -A^\top u \in \partial f(x), \quad u \in \partial g(y)$$

- ▶ eliminating x, y gives

$$0 \in \partial g^*(u) - A\partial f^*(-A^\top u) + \frac{1}{t}(u - z)$$

this is the optimality condition for the problem in the definition of $u = \text{prox}_{th}(z)$

Augmented Lagrangian method

choose initial z_0 and repeat:

1. minimize augmented Lagrangian

$$(\hat{x}, \hat{y}) = \arg \min_{x,y} \left(f(x) + g(y) + \frac{t_k}{2} \|Ax - y + z_k/t_k\|_2^2 \right)$$

2. dual update

$$z_{k+1} = z_k + t_k(A\hat{x} - \hat{y})$$

- ▶ also known as method of multipliers
- ▶ this is the proximal point method applied to the dual problem
- ▶ as variants, can apply the accelerated proximal point methods to the dual
- ▶ usually implemented with inexact minimization step 1

Examples

$$\text{minimize } f(x) + g(Ax)$$

Equality constraints: g is indicator of $\{b\}$

$$\hat{x} = \arg \min_x \left(f(x) + \frac{t}{2} \|Ax - b + z/t\|_2^2 \right)$$
$$z := z + t(A\hat{x} - b)$$

Set constraint: g indicator of convex set C

$$\hat{x} = \arg \min_x \left(f(x) + \frac{t}{2} d(Ax + z/t)^2 \right)$$
$$z := z + t(A\hat{x} - P_C(A\hat{x} + z/t))$$

- ▶ in step 1 on previous page, $\hat{y} = P_C(A\hat{x} + z/t)$ where P_C is projection on C
- ▶ $d(u) = \|u - P_C(u)\|_2$ is Euclidean distance of u to C

Moreau-Yosida smoothing

Definition: the Moreau-Yosida regularization of a closed convex function f is

$$\begin{aligned} f_{(t)}(x) &= \inf_u \left(f(u) + \frac{1}{2t} \|u - x\|_2^2 \right) \quad (\text{with } t > 0) \\ &= f\left(\text{prox}_{tf}(x)\right) + \frac{1}{2t} \left\| \text{prox}_{tf}(x) - x \right\|_2^2 \end{aligned}$$

this is also known as the Moreau envelope of f

Immediate properties

- ▶ $f_{(t)}$ is convex (infimum over u of a convex function of x, u)
- ▶ domain of $f_{(t)}$ is \mathbb{R}^n (recall that $\text{prox}_{tf}(x)$ is defined for all x)

Examples

Indicator function: smoothed f is squared Euclidean distance

$$f(x) = \delta_C(x), \quad f_{(t)}(x) = \frac{1}{2t}d(x)^2$$

1-norm: smoothed function is Huber penalty

$$f(x) = \|x\|_1, \quad f_{(t)}(x) = \sum_{k=1}^n \phi_t(x_k)$$

$$\phi_t(z) = \begin{cases} z^2/(2t) & |z| \leq t \\ |z| - t/2 & |z| \geq t \end{cases}$$

Conjugate of Moreau envelope

$$f_{(t)}(x) = \inf_u \left(f(u) + \frac{1}{2t} \|u - x\|_2^2 \right)$$

- ▶ $f_{(t)}$ is infimal convolution of $f(u)$ and $\|v\|_2^2/(2t)$:

$$f_{(t)}(x) = \inf_{u+v=x} \left(f(u) + \frac{1}{2t} \|v\|_2^2 \right)$$

- ▶ conjugate is sum of conjugates of $f(u)$ and $\|v\|_2^2/(2t)$:

$$(f_{(t)})^*(y) = f^*(y) + \frac{t}{2} \|y\|_2^2$$

- ▶ hence, conjugate is strongly convex with parameter t

Gradient of Moreau envelope

$$f_{(t)}(x) = \sup_y \left(x^\top y - f^*(y) - \frac{t}{2} \|y\|_2^2 \right)$$

- ▶ maximizer y in definition is unique and satisfies

$$\begin{aligned} x - ty \in \partial f^*(y) &\iff y \in \partial f(x - ty) \\ &\iff y = \frac{1}{t}(x - \text{prox}_{tf}(x)) \end{aligned}$$

- ▶ maximizer y is the gradient of $f_{(t)}$:

$$\nabla f_{(t)}(x) = \frac{1}{t}(x - \text{prox}_{tf}(x)) = \text{prox}_{t^{-1}f^*}(x/t)$$

we applied the Moreau decomposition

- ▶ gradient $\nabla f_{(t)}$ is Lipschitz continuous with constant $1/t$

Interpretation of proximal point algorithm

apply gradient method to minimize Moreau envelope

$$\text{minimize } f_{(t)}(x) = \inf_u \left(f(u) + \frac{1}{2t} \|u - x\|_2^2 \right)$$

this is an exact smooth reformulation of problem of minimizing $f(x)$:

- ▶ solution x is minimizer of f
- ▶ $f_{(t)}$ is differentiable with Lipschitz continuous gradient ($L = 1/t$)

Gradient update: with fixed $t_k = 1/L = t$

$$x_{k+1} = x_k - t \nabla f_{(t)}(x_k) = \text{prox}_{tf}(x_k)$$

... the proximal point update with constant step size $t_k = t$

Interpretation of augmented Lagrangian algorithm

$$\text{minimize } f(x) + g(Ax)$$

- ▶ augmented Lagrangian iteration is

$$(\hat{x}, \hat{y}) = \arg \min_{x,y} \left(f(x) + g(y) + \frac{t}{2} \|Ax - y + (1/t)z\|_2^2 \right)$$
$$z := z + t(A\hat{x} - \hat{y})$$

- ▶ with fixed t , dual update is gradient step applied to a smoothed dual
- ▶ after eliminating y , primal step can be written as

$$\hat{x} = \arg \min_x \left(f(x) + g_{(1/t)}(Ax + (1/t)z) \right)$$

- ▶ second term $g_{(1/t)}(Ax + (1/t)z)$ is a smooth approximation of $g(Ax)$
- ▶ adding the offset z/t allows us to use a fixed t

Example

$$\text{minimize } f(x) + \|Ax - b\|_1$$

- ▶ augmented Lagrangian iteration is

$$(\hat{x}, \hat{y}) = \arg \min_{x, y} \left(f(x) + \|y - b\|_1 + \frac{t}{2} \|Ax - y + (1/t)z\|_2^2 \right)$$
$$z := z + t(A\hat{x} - \hat{y})$$

- ▶ primal step after eliminating y : \hat{x} is the solution of

$$\text{minimize } f(x) + \phi_{1/t}(Ax - b + (1/t)z)$$

with $\phi_{1/t}$ the Huber penalty applied componentwise

Any questions?