

Should go to Convex Optimization Self study Proximal Alg

8 **Proximal Algorithms Introduction**



Page 1

* Introduction:

Newton's method: base operations: low level: (linear algebra operations, compute $\nabla f(x), \nabla^2 f(x)$)

Proximal algorithms: analogous tool to Newton's method (but has a higher level of abstraction) solves non-smooth, constrained, large scale, distributed convex optimization problems well suited for large data set

- base operation: evaluate proximal operator of a function = solve a small optimization problem, often admit closed form solution v simple specialized methods can be used to solve very quickly
- has interesting interpretations
 - connected to many different topics in optimization applied math

Page 2

$f: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$, cpc, extended value # $f: \text{cpc} \mapsto \text{epi } f = \{(x, s) \in \mathbb{R}^n \times \mathbb{R} \mid f(x) \leq s\} : (x, s) \in \text{cpc}$ # Closed convex proper cpc function

closed proper convex

$\text{dom } f = \{x \in \mathbb{R}^n \mid f(x) < \infty\}$ # set of points for which f takes on finite values

f proximal operator of f at Q : $\text{prox}_f(Q) : \mathbb{R}^n \rightarrow \mathbb{R}^n = \arg \min_{x \in \mathbb{R}^n} (f(x) + \frac{1}{2} \|x - Q\|_2^2)$

$f(x) + \frac{1}{2} \|x - Q\|_2^2 = \underbrace{f(x)}_D + \underbrace{\frac{1}{2} \|x - Q\|_2^2}_{\text{Euclidean norm squared}}$

\rightarrow this is strongly convex # $f: \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\} \leftarrow \exists_m (f(x) - \frac{\mu}{2} \|x - Q\|_2^2) : \mathbb{R}^n$

\rightarrow the argmin of a strongly convex function will be unique $\forall Q \in \mathbb{R}^n$ (even when $\text{dom } f \subseteq \mathbb{R}^n$)

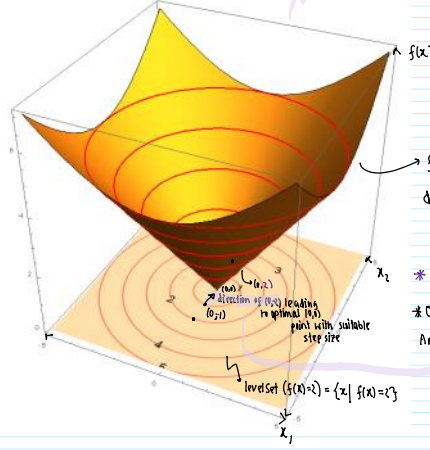
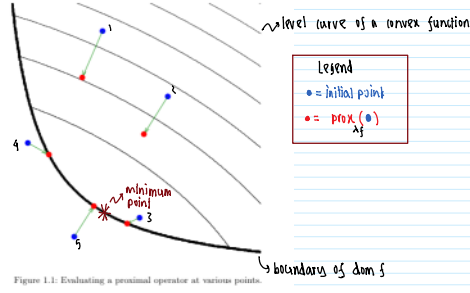
Proximal operator of scaled function λf : $\lambda > 0 \Rightarrow \text{prox}_{\lambda f}(Q) = \arg \min_x (\lambda f(x) + \frac{1}{2} \|x - Q\|_2^2) = \arg \min_x (\lambda (f(x) + \frac{1}{2\lambda} \|x - Q\|_2^2)) = \arg \min_x (f(x) + \frac{1}{2\lambda} \|x - Q\|_2^2)$

$\lambda > 0$ the argmin will be the same

code to generate function:

Page 3

Interpretations:



- * Some characteristics:
- points in the domain will stay in the domain, and move towards the optimal point $(e.g., \text{point } (1, 2))$
- points outside the domain (q, s) will move to the boundary of the domain and towards the minimum of the function.
- λ controls the extent to which proximal operator maps to the optimal point
 - \rightarrow large \rightarrow mapped points near the minimum
 - \rightarrow smaller \rightarrow smaller movement near the minimum
- λ : Analogous to step size in Newton's method

Basic interpretations of prox operator

$\text{prox}_{\lambda f}(v) = \arg \min_x (f(x) + \frac{1}{2\lambda} \|x - v\|_2^2)$

• can picture as minimizing $f(x)$ and $(x$ being near to $v)$

$\lambda \rightarrow 0 \rightarrow \text{prox}_{\lambda f}(v) \rightarrow \arg \min_x (f(x))$

$\lambda \rightarrow \infty \rightarrow \text{prox}_{\lambda f}(v) \rightarrow \arg \min_x (\frac{1}{2\lambda} \|x - v\|_2^2)$

Page 4

When f is Indicator function $I_C(x) = \begin{cases} 0, & x \in C \\ \infty, & x \notin C \end{cases}$

$C: [1, 2]$

$\text{prox}_{\lambda I_C}(v) = \text{prox}_{\lambda f}(v) = \arg \min_x (f(x) + \frac{1}{2\lambda} \|x - v\|_2^2) = \arg \min_x (I_C(x) + \frac{1}{2\lambda} \|x - v\|_2^2)$

$= \arg \min_x \left(\begin{cases} 0, & x \in C \\ \infty, & x \notin C \end{cases} + \frac{1}{2\lambda} \|x - v\|_2^2 \right) = \arg \min_{x \in C} (\frac{1}{2\lambda} \|x - v\|_2^2) = \Pi_C(v)$

$\therefore \text{prox}_{\lambda I_C}(v) = \Pi_C(v)$ ($\text{prox}_{\lambda I_C}(v) = \Pi_C(v)$)

• So, proximal operator is sort of generalized projections. We can hope that various properties of projection that will be obeyed by proximal operator

* Under some assumptions, $\text{prox}_{\lambda f}(v) \approx v - \lambda \nabla f(v)$ # remainder in steepest descent, $x^{(k+1)} = x^{(k)} - \lambda \nabla f(x^{(k)})$

\downarrow small λ - differentiable

This suggests: • close connections between proximal operators and gradient methods

Proximal map of λ times LC is the projection on C

- proximal operator may be useful in optimization
- λ : similar to stepsize

* $x^* = \text{prox}_{\lambda f}(x^*) \Leftrightarrow x^* = \underset{x}{\text{argmin}} (f(x))$

this suggests close connection between fixed point theory and proximal algorithm

- proximal algorithm can be interpreted as solving optimization problems by finding fixed points of appropriate operators

- 13: * Proximal Algorithms:
- Proximal gradient algorithm
 - ADMM
 - Parallel, distributed algorithm
 - Evaluating it easily \rightarrow most useful when this happens

* Why study proximal algorithms:

- extremely general (non-smooth, extended real values)
- fast
- amenable to distributed optimization very large scale *
- conceptually, mathematically simple

Proximal mapping = Resolvent of the subdifferential operator:

$\forall x \in \text{dom} f$ $x \xrightarrow{\partial f} \partial f(x)$
 a point \rightarrow a set

So, ∂f is a mapping or relation

$\partial f = \{ \} (x, \partial f(x)) \}$
 $x \in \text{dom} f$

now: $\text{prox}_{\lambda f} = (I + \lambda \partial f)^{-1}$
 resolvent of the operator ∂f with operator

$\forall x \in \text{dom} f$ $\text{prox}_{\lambda f}(x) = (I + \lambda \partial f)^{-1} x$ then proximal operator is the resolvent of subdifferential operator

Operations on relations: however $(I + \lambda f)^{-1} \text{dom}[\cdot] = \mathbb{R}^n$, single-valued f

relation, sum, inverse

Proof: f is subdifferentiable on its domain

2017-04-06 10:38 PM

$\forall x \in \text{dom} f$
 $(I + \lambda \partial f)^{-1}(x) = \{ \}$

let $z \in (I + \lambda \partial f)^{-1}(x) \Rightarrow \{ \}$
 $\Leftrightarrow x \in (I + \lambda \partial f)z = z + \lambda \partial f(z) = \{ z + \lambda \square_i \}$

$\Leftrightarrow \exists \square_i \in \partial f(z)$ $x = z + \lambda \square_i$
 $\cancel{x} \square_i = -\frac{1}{\lambda}(z-x) \quad [\cdot: \lambda > 0]$
 $\square_i + \frac{1}{\lambda}(z-x) = 0$

$\Leftrightarrow \exists \square_i \in \partial f(z) \quad \square_i + \frac{1}{\lambda}(z-x) = 0$

$\Leftrightarrow 0 \in \partial f(z) + \frac{1}{\lambda}(z-x) = \partial_z \left(f(z) + \frac{1}{2\lambda} \|z-x\|_2^2 \right) = \left[\partial_{f_A} \left(f(f_A) + \frac{1}{2\lambda} \|f_A - x\|_2^2 \right) \right]_{f_A=z}$

remember: $f(z) + \frac{1}{2\lambda} \|z-x\|_2^2$ is a strongly convex as $\square + \square_{\text{strongly}} = \square_{\text{strongly}} \Rightarrow \square_{\text{strict}}$

$\square_{\text{convex}} + \square_{\text{strongly convex}} = \square_{\text{strongly convex}}$
 # so (\cdot) has unique minimizer at z , as (\cdot) is strictly convex

$\Leftrightarrow z = \underset{f_A}{\text{argmin}} \left[f(f_A) + \frac{1}{2\lambda} \|f_A - x\|_2^2 \right]$
 $= \underset{f_A}{\text{argmin}} \frac{1}{\lambda} \left[\lambda f(f_A) + \frac{1}{2} \|f_A - x\|_2^2 \right]$
 (constant so can be dropped)
 $= \underset{f_A}{\text{argmin}} \left[\lambda f(f_A) + \frac{1}{2} \|f_A - x\|_2^2 \right]$
 $= \text{prox}_{\lambda f}(x)$

$\therefore z \in (I + \lambda \partial f)^{-1}(x) \Leftrightarrow z = \text{prox}_{\lambda f}(x)$ which is a singleton (so, $(I + \lambda \partial f)^{-1}$ is a function)

$$\therefore \operatorname{prox}_{\lambda S}(b) = (1 + \lambda \sigma^2)^{-1} b \quad \{\text{single valued}\}$$

(PROVED)