

Proximal gradient method

Acknowledgement: this slides is based on Prof. Lieven Vandenberghes lecture notes

- motivation
- proximal mapping
- proximal gradient method with fixed step size
- proximal gradient method with line search

Proximal mapping

the proximal mapping (prox-operator) of a convex function h is defined as

$$\text{prox}_h(x) = \underset{u}{\operatorname{argmin}} \left(h(u) + \frac{1}{2} \|u - x\|_2^2 \right)$$

examples

- $h(x) = 0$: $\text{prox}_h(x) = x$
- $h(x) = I_C(x)$ (indicator function of C): prox_h is projection on C

$$\text{prox}_h(x) = \underset{u \in C}{\operatorname{argmin}} \|u - x\|_2^2 = P_C(x)$$

- $h(x) = \|x\|_1$: prox_h is the 'soft-threshold' (shrinkage) operation

$$\text{prox}_h(x)_i = \begin{cases} x_i - 1, & x_i \geq 1 \\ 0, & |x_i| \leq 1 \\ x_i + 1, & x_i \leq -1 \end{cases}$$

Proximal gradient method

unconstrained optimization with objective split in two components

$$\min_x f(x) = g(x) + h(x)$$

- g convex, differentiable, $\text{dom } g = \mathbf{R}^n$
- h convex with inexpensive prox-operator (many examples in the lecture on "proximal mapping")

proximal gradient algorithm

$$x^{(k)} = \text{prox}_{t_k h} \left(x^{(k-1)} - t_k \nabla g(x^{(k-1)}) \right)$$

$t_k > 0$ is step size, constant or determined by line search

Interpretation

$$x^+ = \text{prox}_{th}(x - t\nabla g(x))$$

from definition of proximal mapping:

$$\begin{aligned}x^+ &= \underset{u}{\operatorname{argmin}} \left(h(u) + \frac{1}{2t} \|u - x + t\nabla g(x)\|_2^2 \right) \\ &= \underset{u}{\operatorname{argmin}} \left(h(u) + g(x) + \nabla g(x)^\top (u - x) + \frac{1}{2t} \|u - x\|_2^2 \right)\end{aligned}$$

x^+ minimizes $h(u)$ plus a simple quadratic local model of $g(u)$ around x

Examples

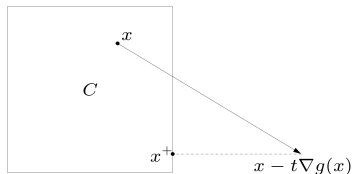
$$\min \quad g(x) + h(x)$$

gradient method: special case with $h(x) = 0$

$$x^+ = x - t\nabla g(x)$$

gradient projection method: special case with $h(x) = I_C(x)$

$$x^+ = P_C(x - t\nabla g(x))$$



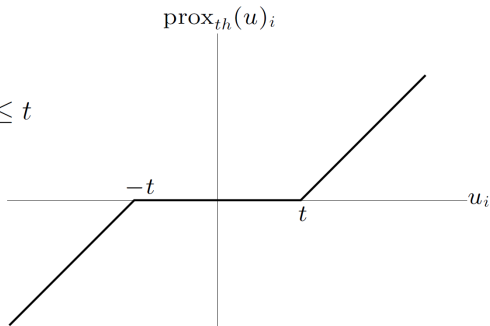
soft-thresholding: special case with $h(x) = \|x\|_1$

$$x^+ = \text{prox}_{th}(x - t\nabla g(x))$$

where

where

$$\text{prox}_{th}(u)_i = \begin{cases} u_i - t & u_i \geq t \\ 0 & -t \leq u_i \leq t \\ u_i + t & u_i \leq -t \end{cases}$$



Outline

- motivation
- **proximal mapping**
- proximal gradient method with fixed step size
- proximal gradient method with line search

Proximal mapping

if h is convex and closed (has a closed epigraph), then

$$\text{prox}_h(x) = \underset{u}{\operatorname{argmin}} \left(h(u) + \frac{1}{2} \|u - x\|_2^2 \right)$$

exists and is unique for all x

- from optimality conditions of minimization in the definition:

$$\begin{aligned} u = \text{prox}_h(x) &\Leftrightarrow x - u \in \partial h(u) \\ &\Leftrightarrow h(z) \geq h(u) + (x - u)^\top (z - u) \quad \forall z \end{aligned}$$

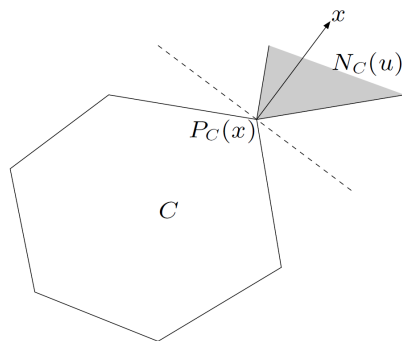
Projection on closed convex set

proximal mapping of indicator function I_C is Euclidean projection on C

$$\text{prox}_{I_C}(x) = \underset{u \in C}{\operatorname{argmin}} \|u - x\|_2^2 = P_C(x)$$

subgradient characterization

$$\begin{aligned} u = P_C(x) \\ \iff \\ (x - u)^\top (z - u) \leq 0 \quad \forall z \in C \end{aligned}$$



we will see that proximal mappings have many properties of projections

Nonexpansiveness

if $u = \text{prox}_h(x)$, $v = \text{prox}_h(y)$, then

$$(u - v)^\top (x - y) \geq \|u - v\|_2^2$$

prox_h is *firmly nonexpansive*, or *co-coercive* with constant 1

- follows from characterization of page 8 and monotonicity

$$x - u \in \partial h(u), y - v \in \partial h(v) \quad \Rightarrow \quad (x - u - y + v)^\top (u - v) \geq 0$$

- implies (from Cauchy-Schwarz inequality)

$$\|\text{prox}_h(x) - \text{prox}_h(y)\|_2 \leq \|x - y\|_2$$

prox_h is *nonexpansive*, or *Lipschitz continuous* with constant 1

Outline

- motivation
- proximal mapping
- **proximal gradient method with fixed step size**
- proximal gradient method with line search

Convergence of proximal gradient method

to minimize $g + h$, choose $x^{(0)}$ and repeat

$$x^{(k)} = \text{prox}_{t_k h} \left(x^{(k-1)} - t \nabla g(x^{(k-1)}) \right), \quad k \geq 1$$

assumptions

- g convex with $\text{dom } g = \mathbf{R}^n$; ∇g Lipschitz continuous with constant L :

$$\|\nabla g(x) - \nabla g(y)\|_2 \leq L\|x - y\|_2 \quad \forall x, y$$

- h is closed and convex (so that prox_{th} is well defined)
- optimal value f^* is finite and attained at x^* (not necessarily unique)

convergence result: $1/k$ rate convergence with fixed step size

$$t_k = 1/L$$

Gradient map

$$G_t(x) = \frac{1}{t}(x - \text{prox}_{th}(x - t\nabla g(x)))$$

$G_t(x)$ is the negative 'step' in the proximal gradient update

$$\begin{aligned}x^+ &= \text{prox}_{th}(x - t\nabla g(x)) \\ &= x - tG_t(x)\end{aligned}$$

- $G_t(x)$ is not a gradient or subgradient of $f = g + h$
- from subgradient definition of prox-operator (page 8),

$$G_t(x) \in \partial g(x) + \partial h(x - tG_t(x))$$

- $G_t(x) = 0$ if and only if x minimizes $f(x) = g(x) + h(x)$

Consequences of Lipschitz assumption

recall upper bound (lecture on "gradient method") for convex g with Lipschitz continuous gradient

$$g(y) \leq g(x) + \nabla g(x)^\top (y - x) + \frac{L}{2} \|y - x\|_2^2 \quad \forall x, y$$

- substitute $y = x - tG_t(x)$:

$$g(x - tG_t(x)) \leq g(x) - t\nabla g(x)^\top G_t(x) + \frac{t^2 L}{2} \|G_t(x)\|_2^2$$

- if $0 < t \leq 1/L$, then

$$g(x - tG_t(x)) \leq g(x) - t\nabla g(x)^\top G_t(x) + \frac{t}{2} \|G_t(x)\|_2^2 \quad (1)$$

A global inequality

if the inequality (1) holds, then for all z ,

$$f(x - tG_t(x)) \leq f(x) - G_t(x)^\top (x - z) + \frac{t}{2} \|G_t(x)\|_2^2 \quad (2)$$

proof: (define $v = G_t(x) - \nabla g(x)$)

$$\begin{aligned} f(x - tG_t(x)) &\leq g(x) - t\nabla g(x)^\top G_t(x) + \frac{t}{2} \|G_t(x)\|_2^2 + h(x - tG_t(x)) \\ &\leq g(z) - \nabla g(x)^\top (x - z) - t\nabla g(x)^\top G_t(x) + \frac{t}{2} \|G_t(x)\|_2^2 \\ &\quad + h(z) + v^\top (x - z - tG_t(x)) \\ &= g(z) + h(z) + G_t(x)^\top (x - z) - \frac{t}{2} \|G_t(x)\|_2^2 \end{aligned}$$

line 2 follows from convexity of g and h , and $v \in \partial h(x - tG_t(x))$

Progress in one iteration

$$x^+ = x - tG_t(x)$$

- inequality (2) with $z = x$ shows the algorithm is a descent method:

$$f(x^+) \leq f(x) - \frac{t}{2} \|G_t(x)\|_2^2$$

- inequality (2) with $z = x^*$

$$\begin{aligned} f(x^+) - f^* &\leq G_t(x)^\top (x - x^*) - \frac{t}{2} \|G_t(x)\|_2^2 \\ &= \frac{1}{2t} (\|x - x^*\|_2^2 - \|x - x^* - tG_t(x)\|_2^2) \quad (3) \\ &= \frac{1}{2t} (\|x - x^*\|_2^2 - \|x^+ - x^*\|_2^2) \end{aligned}$$

(hence, $\|x^+ - x^*\|_2^2 \leq \|x - x^*\|_2^2$, i.e., distance to optimal set decreases)

Analysis for fixed step size

add inequalities (3) for $x = x^{(i-1)}, x^+ = x^{(i)}, t = t_i = 1/L$

$$\begin{aligned}\sum_{i=1}^k (f(x^{(i)}) - f^*) &\leq \frac{1}{2t} \sum_{i=1}^k \left(\|x^{(i-1)} - x^*\|_2^2 - \|x^{(i)} - x^*\|_2^2 \right) \\ &= \frac{1}{2t} \left(\|x^{(0)} - x^*\|_2^2 - \|x^{(k)} - x^*\|_2^2 \right) \\ &\leq \frac{1}{2t} \|x^{(0)} - x^*\|_2^2\end{aligned}$$

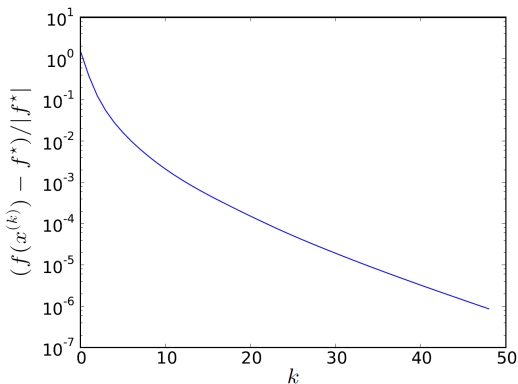
since $f(x^{(i)})$ is nonincreasing,

$$f(x^{(k)}) - f^* \leq \frac{1}{k} \sum_{i=1}^k (f(x^{(i)}) - f^*) \leq \frac{1}{2kt} \|x^{(0)} - x^*\|_2^2$$

conclusion: reaches $f(x^{(k)}) - f^* \leq \epsilon$ after $O(1/\epsilon)$ iterations

Quadratic program with box constraints

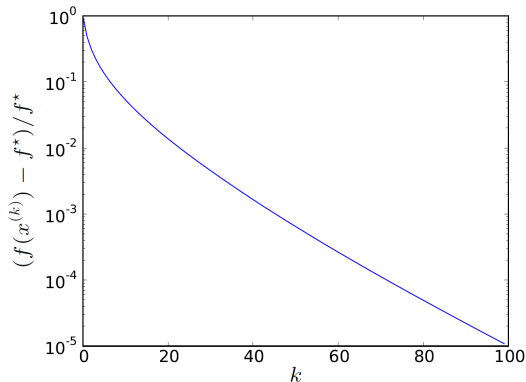
$$\begin{aligned} \min \quad & (1/2)x^\top Ax + b^\top x \\ \text{s.t.} \quad & 0 \preceq x \preceq 1 \end{aligned}$$



$n = 3000$; fixed step size $t = 1/\lambda_{\max}(A)$

1-norm regularized least-squares

$$\min \frac{1}{2} \|Ax - b\|_2^2 + \|x\|_1$$



randomly generated $A \in \mathbf{R}^{2000 \times 1000}$; step $t_k = 1/L$ with $L = \lambda_{\max}(A^T A)$

Outline

- motivation
- proximal mapping
- proximal gradient method with fixed step size
- **proximal gradient method with line search**

Line search

- the analysis for fixed step size starts with the inequality (1)

$$g(x - tG_t(x)) \leq g(x) - t\nabla g(x)^\top G_t(x) + \frac{t}{2}\|G_t(x)\|_2^2$$

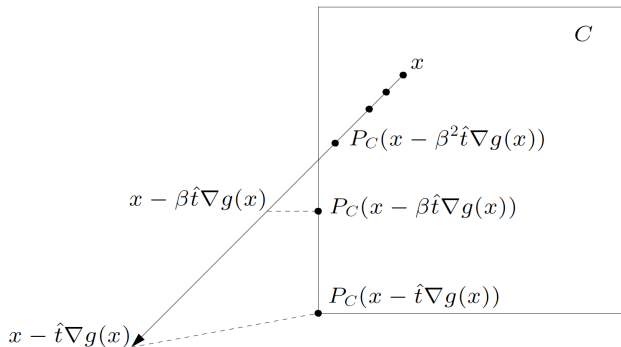
this inequality is known to hold for $0 < t \leq 1/L$

- if L is not known, we can satisfy (1) by a backtracking line search: start at some $t := \hat{t} > 0$ and backtrack ($t := \beta t$) until (1) holds
- step size t selected by the line search satisfies
 $t \geq t_{\min} = \min\{\hat{t}, \beta/L\}$
- requires one evaluation of g and proxth per line search iteration

several other types of line search work

example: line search for projected gradient method

$$x^+ = P_C(x - t\nabla g(x)) = x - tG_t(x)$$



backtrack until $x - tG_t(x)$ satisfies 'sufficient decrease' inequality (1)

Analysis with line search

from page 17, if (1) holds in iteration i , then $f(x^{(i)}) < f(x^{(i-1)})$ and

$$\begin{aligned} f(x^{(i)}) - f^* &\leq \frac{1}{2t_i} \left(\|x^{(i-1)} - x^*\|_2^2 - \|x^{(i)} - x^*\|_2^2 \right) \\ &\leq \frac{1}{2t_{\min}} \left(\|x^{(i-1)} - x^*\|_2^2 - \|x^{(i)} - x^*\|_2^2 \right) \end{aligned}$$

- adding inequalities for $i = 1$ to $i = k$ gives

$$\sum_{i=1}^k f(x^{(i)}) - f^* \leq \frac{1}{2t_{\min}} \left(\|x^{(0)} - x^*\|_2^2 \right)$$

- since $f(x^{(i)})$ is nonincreasing, obtain similar $1/k$ bound as for fixed t_i :

$$f(x^{(k)}) - f^* \leq \frac{1}{2kt_{\min}} \left(\|x^{(0)} - x^*\|_2^2 \right)$$

Another Backtracking Line Search Scheme

Let $h(x) = \mu\|x\|_1$, consider

$$\min_x f(x) = g(x) + h(x)$$

- Compute $\bar{x}^k = \text{prox}_{\tau^k h}(x - \tau^k \nabla g(x^k))$, $d^k = \bar{x}^k - x^k$,
 $\Delta^k = \nabla g(x^k)^\top d^k + (h(\bar{x}^k) - h(x^k))$
- Choose $\alpha_\rho > 0$ and $\sigma, \rho \in (0, 1)$. Choose $\alpha_k = \alpha_\rho \rho^h$ such that h is the smallest integer that satisfies

$$f(x^k + \alpha_\rho \rho^h) \leq f(x^k) + \sigma \alpha_\rho \rho^h \Delta^k$$

- Set $x^{k+1} = x^k + \alpha^k d^k$.

Improving Backtracking Line Search Scheme

Choosing τ^k : Barzilai-Borwein method

- $s^{k-1} = x^k - x^{k-1}$ and $y^{k-1} = \nabla g(x^k) - \nabla g(x^{k-1})$
- $\tau^{k,BB1} = \frac{(s^{k-1})^\top s^{k-1}}{(s^{k-1})^\top y^{k-1}}$ or $\tau^{k,BB2} = \frac{(s^{k-1})^\top y^{k-1}}{(y^{k-1})^\top y^{k-1}}$.

Choosing α^k : Nonmontone Armijo-like Line search



$$\psi_\mu(x^k + \alpha^k d^k) \leq C^k + \sigma \alpha^k \Delta^k$$

- $C^k = (\eta Q^{k-1} C^{k-1} + f(x^k))/Q^k$, $Q^k = \eta Q^{k-1} + 1$, $C^0 = f(x^0)$ and $Q^0 = 1$ (Zhang and Hager)

The continuation strategy

- Choose $\mu^0 > \mu^1 > \dots > \mu^l = \mu$
- Solve $z(\mu^i) = \arg \min_x g(x) + \mu^i \|x\|_1$ starting from $z(\mu^{i-1})$

convergence analysis of proximal gradient method

-  A. Beck and M. Teboulle, *A fast iterative shrinkage-thresholding algorithm for linear inverse problems*, SIAM Journal on Imaging Sciences (2009)
-  A. Beck and M. Teboulle, *Gradient-based algorithms with applications to signal recovery*, in: Y. Eldar and D. Palomar (Eds.), *Convex Optimization in Signal Processing and Communications* (2009)