

Sparsity Methods for Systems and Control

Algorithms for Convex Optimization

Masaaki Nagahara^{1,2}

¹The University of Kitakyushu
nagahara@kitakyu-u.ac.jp

²IIT Bombay (Visiting Faculty)

Table of Contents

- 1 Basics of convex optimization
- 2 Proximal Operator
- 3 Proximal splitting methods for ℓ^1 optimization
- 4 Proximal gradient method for ℓ^1 regularization
- 5 Generalized LASSO and ADMM

Table of Contents

- 1 Basics of convex optimization
- 2 Proximal Operator
- 3 Proximal splitting methods for ℓ^1 optimization
- 4 Proximal gradient method for ℓ^1 regularization
- 5 Generalized LASSO and ADMM

ℓ^1 optimization by CVX

ℓ^1 Optimization

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \|x\|_1 \quad \text{subject to} \quad \Phi x = y.$$

The MATLAB CVX code

```
cvx_begin
    variable x(n)
    minimize( norm(x, 1) )
    subject to
        y == Phi * x
cvx_end
```

- Useful for a small or middle scale problems
- Not that useful for
 - **large-scale** problems like image processing
 - **real-time** applications for control system
- You need to build an **efficient algorithm** by yourself for your **specific** problem.

Convex set

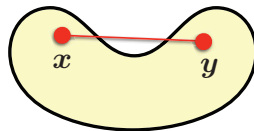
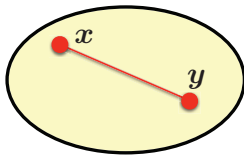
Convex set

Let C be a subset of \mathbb{R}^n . C is said to be a **convex set** if the following inclusion

$$tx + (1 - t)y \in C$$

holds for any vectors $x, y \in C$ and for any real number $t \in [0, 1]$.

- Convex and non-convex sets



Effective domain

- The **effective domain** of a function $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ is defined by

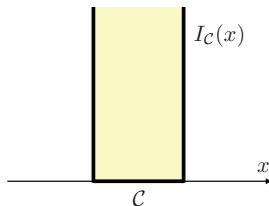
$$\text{dom}(f) \triangleq \{x \in \mathbb{R}^n : f(x) < \infty\}.$$

- Indicator function**

$$f(x) = \begin{cases} 0, & \text{if } \|x\|_2 \leq 1, \\ \infty, & \text{if } \|x\|_2 > 1. \end{cases}$$

- The effective domain of the indicator function is

$$\text{dom}(f) = \{x \in \mathbb{R}^n : \|x\|_2 \leq 1\}.$$



Convex function

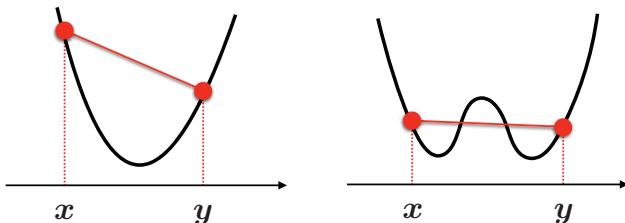
Convex function

Let $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ be a proper function. The function f is said to be a **convex function** if the following inequality

$$f(tx + (1 - t)y) \leq tf(x) + (1 - t)f(y)$$

holds for any vectors $x, y \in \text{dom}(f)$ and for any real number $t \in [0, 1]$.

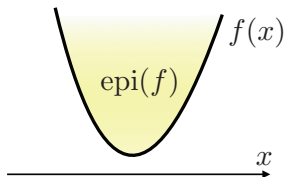
- Convex and non-convex functions



Epigraph

- The **epigraph** $\text{epi}(f)$ of function f is defined by

$$\text{epi}(f) \triangleq \{(x, t) \in \mathbb{R}^n \times \mathbb{R} : x \in \text{dom}(f), f(x) \leq t\}.$$



- Proper, convex, and closed function

function f	epigraph $\text{epi}(f)$
convex	convex set
closed	closed set
proper	non-empty

Convex optimization problem

Convex optimization problem

Let $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ be a proper, closed, and convex function, and $C \subset \mathbb{R}^n$ be a closed convex set. Then, a **convex optimization problem** is a problem to find a vector $\mathbf{x}^* \in \mathbb{R}^n$ that minimizes the function $f(\mathbf{x})$ over the set $C \subset \mathbb{R}^n$. The problem is briefly written as

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad f(\mathbf{x}) \quad \text{subject to} \quad \mathbf{x} \in C.$$

- The function $f(\mathbf{x})$ is called a **cost function** or an **objective function**.
- The set C is called a **constraint set** or a **feasible set**.
- The entries of C is called **feasible solutions**.
- The inclusion $\mathbf{x} \in C$ is called a **constraint**.

Notation

- Minimum value:

$$\min_{x \in C} f(x).$$

- Minimizer (set):

$$\arg \min_{x \in C} f(x) \triangleq \{x^* \in C : f(x^*) \leq f(x), \forall x \in C \cap \text{dom}(f)\}.$$

$$\begin{array}{ccc} \minimize_{x \in \mathbb{R}^n} & \underbrace{f(x)} & \text{subject to } \underbrace{x \in C} \\ & \text{cost function} & \text{constraint} \end{array}$$

$$\min_{x \in C} f(x) \quad \text{minimum value}$$

$$\arg \min_{x \in C} f(x) \quad \text{minimizer (set)}$$

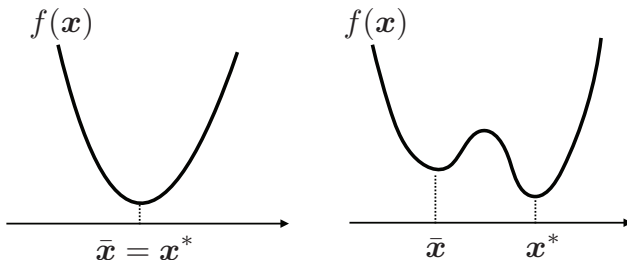
Global/local minimizers

- **Local minimizer**: there exists an open set \mathcal{B} that contains a feasible solution $\bar{x} \in C \cap \text{dom}(f)$ such that

$$f(x) \geq f(\bar{x}), \quad \forall x \in \mathcal{B} \cap C.$$

- **Global minimizer**: a feasible solution $x^* \in C$ that satisfies

$$f(x) \geq f(x^*), \quad \forall x \in C.$$



Theorem

For the convex optimization problem,

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \ f(x) \quad \text{subject to} \ x \in C.$$

any local minimizer is (if it exists) a global minimizer, and the set of global minimizers is a convex set.

- For convex optimization problems, you just need to **find a local minimizer**, which is **consequently a global minimizer**.

Strictly and strongly convex functions

- Let $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ be a proper function.
- The function f is said to be a **strictly convex function** if for any $x, y \in \text{dom}(f) \subset \mathbb{R}^n$ with $x \neq y$ and any $t \in (0, 1)$,

$$f(tx + (1 - t)y) < tf(x) + (1 - t)f(y)$$

- The function f is said to be a **strongly convex function** if there exists $\beta > 0$ such that for any $x, y \in \text{dom}(f) \subset \mathbb{R}^n$ and any $t \in [0, 1]$,

$$f(tx + (1 - t)y) \leq tf(x) + (1 - t)f(y) - t(1 - t)\frac{\beta}{2}\|x - y\|_2^2$$

The constant β is called a *modulus*.

Strongly convex functions

Theorem

Assume $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ is a proper, closed, and **strongly convex** function with modulus $\beta > 0$. Then f has the **unique minimizer** $\mathbf{x}^* \in \text{dom}(f)$. That is, for all $\mathbf{x} \in \text{dom}(f)$ such that $\mathbf{x} \neq \mathbf{x}^*$,

$$f(\mathbf{x}) > f(\mathbf{x}^*).$$

Moreover, for any $\mathbf{x} \in \text{dom}(f)$, we have

$$f(\mathbf{x}) \geq f(\mathbf{x}^*) + \frac{\beta}{2} \|\mathbf{x} - \mathbf{x}^*\|_2^2.$$

- This is an important property of strongly convex functions.
- This is used to define the **proximal operator** (see next Section).

Table of Contents

- 1 Basics of convex optimization
- 2 Proximal Operator**
- 3 Proximal splitting methods for ℓ^1 optimization
- 4 Proximal gradient method for ℓ^1 regularization
- 5 Generalized LASSO and ADMM

Proximal operator

Proximal operator

Let $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ be a proper, closed, and convex function. The **proximal operator** $\text{prox}_{\gamma f}$ with parameter $\gamma > 0$ is defined by

$$\text{prox}_{\gamma f}(v) \triangleq \arg \min_{x \in \text{dom}(f)} \left\{ f(x) + \frac{1}{2\gamma} \|x - v\|_2^2 \right\}.$$

- $\gamma = \infty$: Minimizer of $f(z)$:

$$\text{prox}_{\gamma f}(v) = \arg \min_{x \in \text{dom}(f)} f(x)$$

- $\gamma = 0$: Projection onto $\text{dom}(f)$:

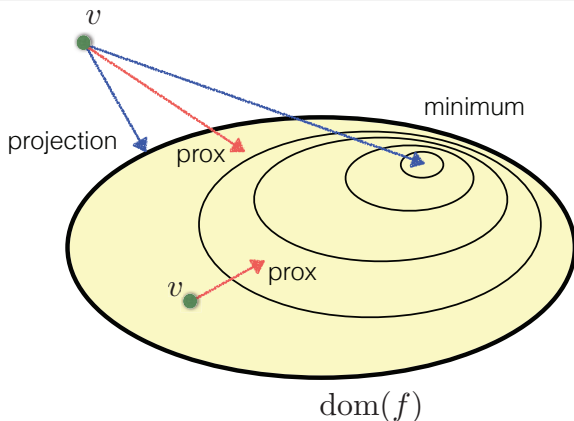
$$\text{prox}_{\gamma f}(v) = \arg \min_{x \in \text{dom}(f)} \frac{1}{2\gamma} \|x - v\|_2^2$$

- $\gamma \in (0, \infty)$: a mixture of those.

Proximal operator

Proximal operator

$$\text{prox}_{\gamma f}(v) \triangleq \arg \min_{x \in \text{dom}(f)} \left\{ f(x) + \frac{1}{2\gamma} \|x - v\|_2^2 \right\}.$$



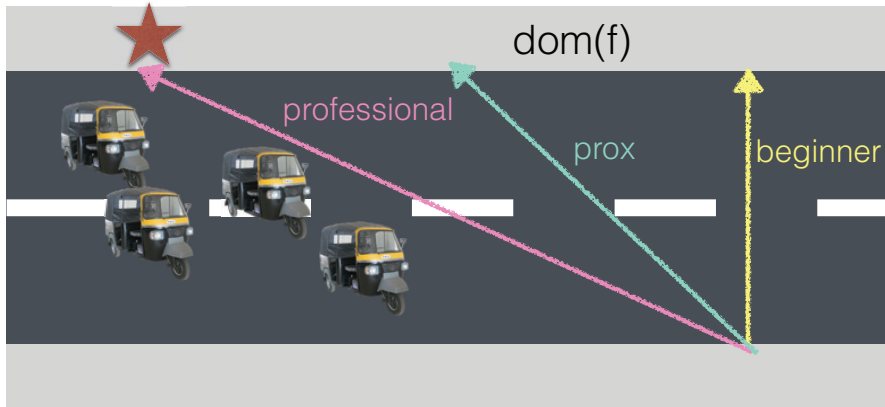
Proximal operator

The "crossing the street" problem.



Proximal operator

The "crossing the street" problem.



Proximal algorithm

Proximal algorithm

Initialization: give an initial vector $\mathbf{x}[0]$ and positive numbers $\gamma_0, \gamma_1, \gamma_2, \dots$

Iteration: for $k = 0, 1, 2, \dots$, do

$$\mathbf{x}[k+1] = \text{prox}_{\gamma_k f}(\mathbf{x}[k]) = \arg \min_{\mathbf{x} \in \text{dom}(f)} \left\{ f(\mathbf{x}) + \frac{1}{2\gamma_k} \|\mathbf{x} - \mathbf{x}[k]\|_2^2 \right\}.$$

- The algorithm minimizes the **strongly convex** function

$$g_k(\mathbf{x}) \triangleq f(\mathbf{x}) + \frac{1}{2\gamma_k} \|\mathbf{x} - \mathbf{x}[k]\|_2^2$$

at each step k , which is an approximation of f that may not be strongly convex.

Convergence theorem of proximal algorithm

Theorem

Suppose that the parameter sequence $\{\gamma_k\}$ satisfies $\gamma_k > 0$ for all k and

$$\sum_{k=0}^{\infty} \gamma_k = \infty.$$

Then, the vector sequence $\{x[k]\}$ generated by the proximal algorithm

$$x[k+1] = \text{prox}_{\gamma_k f}(x[k])$$

converges to one of the minimizers of f for any initial vector $x[0]$.

Proxiable functions

- Proximal operator

$$\text{prox}_{\gamma f}(v) \triangleq \arg \min_{x \in \text{dom}(f)} \left\{ f(x) + \frac{1}{2\gamma} \|x - v\|_2^2 \right\}.$$

needs to be obtained in a **closed form** to derive an efficient algorithm.

- A **proxiable** function is a function that has a closed-form proximal operator.
- The following functions are proxiable:
 - quadratic functions (including the ℓ^2 norm)
 - indicator functions
 - the ℓ^1 norm

Quadratic function

- The **quadratic function**

$$f(x) = \frac{1}{2}x^\top \Phi x - y^\top x,$$

where Φ is a **positive-definite** matrix.

- The proximal operator is given by

$$\begin{aligned}\text{prox}_{\gamma f}(v) &= \arg \min_{x \in \mathbb{R}^n} \left\{ \frac{1}{2}x^\top \Phi x - y^\top x + \frac{1}{2\gamma}(x - v)^\top (x - v) \right\} \\ &= \left(\Phi + \frac{1}{\gamma}I \right)^{-1} \left(y + \frac{1}{\gamma}v \right).\end{aligned}$$

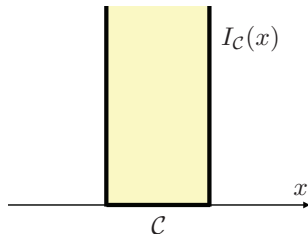
Indicator function

Indicator function

For a subset C in \mathbb{R}^n , the **indicator function** is defined by

$$I_C(x) = \begin{cases} 0, & x \in C \\ \infty, & x \notin C \end{cases}$$

- C : non-empty, closed, and convex $\Rightarrow I_C(x)$: a proper, closed, and convex function
- Draw the epigraph of I_C .



- The proximal operator of $I_C(\mathbf{x})$ is given by

$$\begin{aligned}\text{prox}_{\gamma I_C}(\mathbf{v}) &= \arg \min_{\mathbf{x} \in \mathbb{R}^n} \left\{ I_C(\mathbf{x}) + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{v}\|_2^2 \right\} \\ &= \arg \min_{\mathbf{x} \in C} \|\mathbf{x} - \mathbf{v}\|_2^2 \\ &= \Pi_C(\mathbf{v}).\end{aligned}$$

where Π_C is the **projection operator** onto the nonempty, closed, and convex set C .

- The proximal operator of the ℓ^1 norm $\|x\|_1$ has a closed form

$$\text{prox}_{\gamma\|\cdot\|_1}(v) = S_\gamma(v),$$

where $S_\gamma : \mathbb{R}^n \rightarrow \mathbb{R}^n$ is the **soft-thresholding operator** defined by

$$[S_\gamma(v)]_i = \begin{cases} v_i - \gamma, & v_i \geq \gamma, \\ 0, & -\gamma < v_i < \gamma, \\ v_i + \gamma, & v_i \leq -\gamma. \end{cases}$$

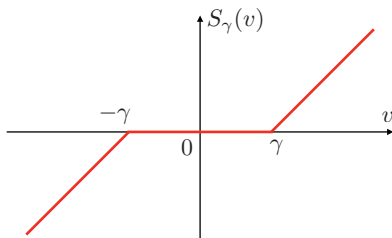


Table of Contents

- 1 Basics of convex optimization
- 2 Proximal Operator
- 3 Proximal splitting methods for ℓ^1 optimization**
- 4 Proximal gradient method for ℓ^1 regularization
- 5 Generalized LASSO and ADMM

ℓ^1 optimization

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \|x\|_1 \quad \text{subject to} \quad \Phi x = y,$$

- $\Phi \in \mathbb{R}^{m \times n}$ and $y \in \mathbb{R}^m$ are given
- $m < n$
- Φ has full row rank, that is, $\text{rank}(\Phi) = m$.

ℓ^1 optimization

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \|x\|_1 \quad \text{subject to} \quad \Phi x = y,$$

- Constraint set

$$C \triangleq \{x \in \mathbb{R}^n : \Phi x = y\}.$$

- Indicator function

$$I_C(x) = \begin{cases} 0, & \text{if } \Phi x = y, \\ \infty, & \text{if } \Phi x \neq y. \end{cases}$$

- Equivalent unconstrained optimization

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \|x\|_1 + I_C(x).$$

Splitting method

- Equivalent unconstrained optimization

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad \|x\|_1 + I_C(x).$$

- $\|x\|_1 + I_C(x)$ is proper, closed, and convex but **not proximal**.
 - The proximal algorithm cannot be directly applied.
 - Both functions,

$$f_1(x) \triangleq \|x\|_1, \quad f_2(x) \triangleq I_C(x)$$

are **proximal**

- We **split** the cost function as $f = f_1 + f_2$.

Douglas-Rachford splitting algorithm

- General optimization problem

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f_1(x) + f_2(x),$$

- f_1 and f_2 are proper, closed, and convex functions.
- f_1 and f_2 are **proximable**.

Douglas-Rachford splitting algorithm

Initialization: give an initial vector $z[0]$ and a parameter $\gamma > 0$

Iteration: for $k = 0, 1, 2, \dots$ do

$$x[k+1] = \text{prox}_{\gamma f_1}(z[k])$$

$$z[k+1] = z[k] + \text{prox}_{\gamma f_2}(2x[k+1] - z[k]) - x[k+1]$$

Douglas-Rachford splitting for ℓ^1 optimization

- ℓ^1 optimization: minimize $_{x \in \mathbb{R}^n}$ $\|x\|_1 + I_C(x)$
 - $C = \{x \in \mathbb{R}^n : \Phi x = y\}$.
- $f_1(x) = \|x\|_1$ and $f_2(x) = I_C(x)$ are **proximable**.

$$\text{prox}_{\gamma f_1}(v) = S_\gamma(v),$$

$$\text{prox}_{\gamma f_2}(v) = \Pi_C(v) = v + \Phi^\top (\Phi \Phi^\top)^{-1} (y - \Phi v).$$

- Now we have got the algorithm!

Douglas-Rachford splitting algorithm for ℓ^1 optimization

Initialization: give an initial vector $z[0]$ and a parameter $\gamma > 0$

Iteration: for $k = 0, 1, 2, \dots$ do

$$x[k+1] = S_\gamma(z[k])$$

$$z[k+1] = z[k] + \Pi_C(2x[k+1] - z[k]) - x[k+1]$$

Douglas-Rachford splitting algorithm

Douglas-Rachford splitting algorithm for ℓ^1 optimization

Initialization: give an initial vector $z[0]$ and a parameter $\gamma > 0$

Iteration: for $k = 0, 1, 2, \dots$ do

$$x[k+1] = S_\gamma(z[k])$$

$$z[k+1] = z[k] + \Pi_C(2x[k+1] - z[k]) - x[k+1]$$

- Douglas-Rachford algorithm only requires
 - **simple continuous mapping** of the soft-thresholding function S_γ
 - **linear computation** of matrix-vector multiplication and vector addition
- Much **faster** and **easier to implement** than the standard interior-point method
 - The standard interior-point method requires **to solve linear equations at each step**.

Table of Contents

- 1 Basics of convex optimization
- 2 Proximal Operator
- 3 Proximal splitting methods for ℓ^1 optimization
- 4 Proximal gradient method for ℓ^1 regularization
- 5 Generalized LASSO and ADMM

ℓ^1 regularization (LASSO)

ℓ^1 regularization (LASSO)

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2} \|\Phi \mathbf{x} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{x}\|_1.$$

- Sum of two convex functions $f = f_1 + f_2$:

$$f_1(\mathbf{x}) = \frac{1}{2} \|\Phi \mathbf{x} - \mathbf{y}\|_2^2, \quad f_2(\mathbf{x}) = \lambda \|\mathbf{x}\|_1$$

- f_1 and f_2 are both **proximable** \rightarrow Douglas-Rachford splitting
- f_1 is also **differentiable**
- A yet faster algorithm exists for this type of optimization problem.

Proximal gradient algorithm

- Optimization problem:

$$\underset{x \in \mathbb{R}^n}{\text{minimize}} \quad f_1(x) + f_2(x),$$

- f_1 is differentiable and convex, satisfying $\text{dom}(f_1) = \mathbb{R}^n$
- f_2 is a proper, closed, and convex function.

Proximal gradient algorithm

Initialization: give an initial vector $x[0]$ and a real number $\gamma > 0$

Iteration: for $k = 0, 1, 2, \dots$ do

$$x[k+1] = \text{prox}_{\gamma f_2}(x[k] - \gamma \nabla f_1(x[k])).$$

- Proximal gradient algorithm

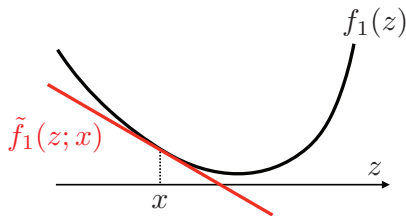
$$\mathbf{x}[k+1] = \underbrace{\text{prox}_{\gamma f_2}(\mathbf{x}[k] - \gamma \nabla f_1(\mathbf{x}[k]))}_{\triangleq \phi(\mathbf{x}[k])}.$$

- The function $\phi(\mathbf{x})$ is rewritten as

$$\begin{aligned}\phi(\mathbf{x}) &= \arg \min_{\mathbf{z} \in \mathbb{R}^n} \left\{ f_2(\mathbf{z}) + \frac{1}{2\gamma} \|\mathbf{z} - (\mathbf{x} - \gamma \nabla f_1(\mathbf{x}))\|_2^2 \right\} \\ &= \arg \min_{\mathbf{z} \in \mathbb{R}^n} \left\{ \underbrace{f_1(\mathbf{x}) + \nabla f_1(\mathbf{x})^\top (\mathbf{z} - \mathbf{x})}_{\triangleq \tilde{f}_1(\mathbf{z}; \mathbf{x})} + f_2(\mathbf{z}) + \frac{1}{2\gamma} \|\mathbf{z} - \mathbf{x}\|_2^2 \right\}.\end{aligned}$$

Geometrical interpretation

- The function $\tilde{f}_1(z; \mathbf{x})$ is a linear approximation of $f_1(z)$ around the point $\mathbf{x} \in \mathbb{R}^n$.



- The iteration becomes

$$\begin{aligned}\mathbf{x}[k+1] &= \arg \min_{\mathbf{z} \in \mathbb{R}^n} \left\{ \tilde{f}_1(\mathbf{z}; \mathbf{x}) + f_2(\mathbf{z}) + \frac{1}{2\gamma} \|\mathbf{z} - \mathbf{x}\|_2^2 \right\} \\ &= \text{prox}_{\gamma \tilde{f}}(\mathbf{x}[k])\end{aligned}$$

where $\tilde{f}(\mathbf{z}) = \tilde{f}_1(\mathbf{z}; \mathbf{x}) + f_2(\mathbf{z})$. This is a proximal algorithm for finding the minimizer of \tilde{f} .

Convergence analysis

Theorem

Assume the **gradient** ∇f_1 is **Lipschitz continuous** over \mathbb{R}^n with Lipschitz constant L . Assume also that the step size γ satisfies

$$\gamma \leq \frac{1}{L}.$$

Then the sequence $\{x[k]\}$ generated by the proximal gradient algorithm **converges** to an optimal solution x^* at the rate of $O(1/k)$.

- The gradient ∇f_1 is Lipschitz continuous over \mathbb{R}^n with Lipschitz constant L if

$$\|\nabla f_1(x) - \nabla f_1(y)\|_2 \leq L\|x - y\|_2, \quad \forall x, y \in \mathbb{R}^n.$$

ℓ^1 regularization (LASSO)

- ℓ^1 regularization problem (LASSO)

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2} \|\Phi \mathbf{x} - \mathbf{y}\|_2^2 + \lambda \|\mathbf{x}\|_1.$$

- Sum of two convex functions $f = f_1 + f_2$:

$$f_1(\mathbf{x}) = \frac{1}{2} \|\Phi \mathbf{x} - \mathbf{y}\|_2^2, \quad f_2(\mathbf{x}) = \lambda \|\mathbf{x}\|_1$$

with

$$\nabla f_1(\mathbf{x}) = \Phi^\top (\Phi \mathbf{x} - \mathbf{y}), \quad \text{prox}_{\gamma f_2}(\mathbf{v}) = S_{\gamma \lambda}(\mathbf{v}).$$

ISTA (Iterative Shrinkage Thresholding Algorithm)

Initialization: give an initial vector $\mathbf{x}[0]$ and parameter $\gamma > 0$

Iteration: for $k = 0, 1, 2, \dots$ do

$$\mathbf{x}[k+1] = S_{\gamma \lambda}(\mathbf{x}[k] - \gamma \Phi^\top (\Phi \mathbf{x}[k] - \mathbf{y})).$$

Convergence of ISTA

- Gradient

$$\nabla f_1(\mathbf{x}) = \Phi^\top(\Phi\mathbf{x} - \mathbf{y}).$$

- Lipschitz constant L of ∇f_1 :

$$\|\nabla f_1(\mathbf{x}_1) - \nabla f_1(\mathbf{x}_2)\|_2 = \|\Phi^\top\Phi(\mathbf{x}_1 - \mathbf{x}_2)\|_2 \leq \lambda_{\max}(\Phi^\top\Phi)\|\mathbf{x}_1 - \mathbf{x}_2\|_2$$

where $\lambda_{\max}(\Phi^\top\Phi)$ is the largest eigenvalue of $\Phi^\top\Phi$.

- We have $L = \lambda_{\max}(\Phi^\top\Phi) = \sigma_{\max}(\Phi)^2 = \|\Phi\|^2$, where $\sigma_{\max}(\Phi)$ is the largest singular value of Φ and $\|\Phi\|$ is a matrix norm.
- A sufficient condition for convergence:

$$\gamma \leq \frac{1}{L} = \frac{1}{\|\Phi\|^2}$$

- The convergence rate is $O(1/k)$.

Fast ISTA (FISTA)

- Accelerated algorithm of ISTA = **FISTA (Fast ISTA)**
- The convergence rate is $O(1/k^2)$.

FISTA

Initialization: give initial vectors $\mathbf{x}[0]$, $\mathbf{z}[0]$, initial number $t[0]$, and parameter $\gamma > 0$

Iteration: for $k = 0, 1, 2, \dots$ do

$$\mathbf{x}[k+1] = S_{\gamma\lambda}(\mathbf{z}[k] - \gamma\Phi^\top(\Phi\mathbf{z}[k] - \mathbf{y})),$$

$$t[k+1] = \frac{1 + \sqrt{1 + 4t[k]^2}}{2},$$

$$\mathbf{z}[k+1] = \mathbf{x}[k+1] + \frac{t[k] - 1}{t[k+1]}(\mathbf{x}[k+1] - \mathbf{x}[k]).$$

Table of Contents

- 1 Basics of convex optimization
- 2 Proximal Operator
- 3 Proximal splitting methods for ℓ^1 optimization
- 4 Proximal gradient method for ℓ^1 regularization
- 5 Generalized LASSO and ADMM

- A generalized regularization problem:

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2} \|\Phi \mathbf{x} - \mathbf{y}\|_2^2 + \lambda \|\Psi \mathbf{x}\|_1,$$

- Ψ is a matrix.
- We call this the **generalized LASSO**.
- If $\Psi = I$, then this is LASSO.
- $\|\Psi \mathbf{x}\|_1$ is **not proximal** in general.
 - No closed-form expression for its proximal operator.
- We need yet another algorithm.

- Optimization problem

$$\underset{x \in \mathbb{R}^n, z \in \mathbb{R}^p}{\text{minimize}} \quad f_1(x) + f_2(z) \quad \text{subject to} \quad z = \Psi x,$$

- f_1, f_2 : proper, closed, and convex
- $\Psi \in \mathbb{R}^{p \times n}$
- Alternating Direction Method of Multipliers (ADMM)

ADMM

Initialization: give initial vectors $z[0], v[0] \in \mathbb{R}^p$, and real number $\gamma > 0$

Iteration: for $k = 0, 1, 2, \dots$ do

$$x[k+1] := \arg \min_{x \in \mathbb{R}^n} \left\{ f_1(x) + \frac{1}{2\gamma} \|\Psi x - z[k] + v[k]\|^2 \right\},$$

$$z[k+1] := \text{prox}_{\gamma f_2}(\Psi x[k+1] + v[k]),$$

$$v[k+1] := v[k] + \Psi x[k+1] - z[k+1].$$

Convergence of ADMM

- Assume f_1 and f_2 are proper, closed, and convex functions.
- Assume also that the Lagrangian

$$L(x, z, \lambda) = f_1(x) + f_2(z) + \lambda^\top (\Psi x - z).$$

has a saddle point, that is, there exist x^* , z^* , and λ^* such that

$$L(x^*, z^*, \lambda) \leq L(x^*, z^*, \lambda^*) \leq L(x, z, \lambda^*), \quad \forall x, z, \lambda.$$

- Then,
 - The residual

$$r[k] \triangleq \Psi x[k] - z[k] \rightarrow \mathbf{0} \quad \text{as } k \rightarrow \infty.$$

- The objective value

$$f_1(x[k]) + f_2(z[k]) \rightarrow f^* \triangleq \inf_{\substack{x \in \mathbb{R}^n, z \in \mathbb{R}^p \\ \Psi x = z}} \{f_1(x) + f_2(z)\} \quad \text{as } k \rightarrow \infty.$$

- If $\Psi^\top \Psi$ is invertible, then the sequence

$$(x[k], z[k]) \rightarrow (x^*, z^*) \quad \text{as } k \rightarrow \infty.$$

ADMM for generalized LASSO

- Generalized LASSO

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad \frac{1}{2} \|\Phi \mathbf{x} - \mathbf{y}\|_2^2 + \lambda \|\Psi \mathbf{x}\|_1,$$

- The first update in ADMM:

$$\begin{aligned} \arg \min_{\mathbf{x} \in \mathbb{R}^n} \left\{ \frac{1}{2} \|\Phi \mathbf{x} - \mathbf{y}\|_2^2 + \frac{1}{2\gamma} \|\Psi \mathbf{x} - \mathbf{z}[k] + \mathbf{v}[k]\|_2^2 \right\} \\ = (\Phi^\top \Phi + \gamma^{-1} \Psi^\top \Psi)^{-1} (\Phi^\top \mathbf{y} + \gamma^{-1} \Psi^\top (\mathbf{z}[k] - \mathbf{v}[k])). \end{aligned}$$

- The proximal operator in the second update is the soft-thresholding operator:

$$\text{prox}_{\gamma f_2}(\mathbf{v}) = S_{\gamma\lambda}(\mathbf{v})$$

ADMM for generalized LASSO

ADMM for generalized LASSO

Initialization: give initial vectors $z[0], v[0] \in \mathbb{R}^p$, and real number $\gamma > 0$

Iteration: for $k = 0, 1, 2, \dots$ do

$$x[k+1] = (\Phi^\top \Phi + \gamma^{-1} \Psi^\top \Psi)^{-1} (\Phi^\top y + \gamma^{-1} \Psi^\top (z[k] - v[k]))$$

$$z[k+1] = S_{\gamma\lambda}(\Psi x[k+1] + v[k])$$

$$v[k+1] = v[k] + \Psi x[k+1] - z[k+1].$$

- The inverse matrix $(\Phi^\top \Phi + \gamma^{-1} \Psi^\top \Psi)^{-1}$ is computed **offline**.
- If the matrix $\Phi^\top \Phi + \gamma^{-1} \Psi^\top \Psi$ is a **tridiagonal matrix**, the linear equation

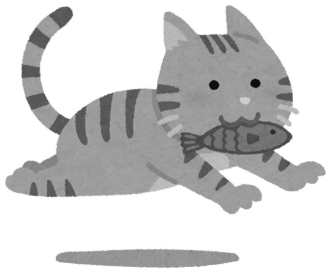
$$(\Phi^\top \Phi + \gamma^{-1} \Psi^\top \Psi)x = v$$

with unknown x can be solved in $O(n)$.

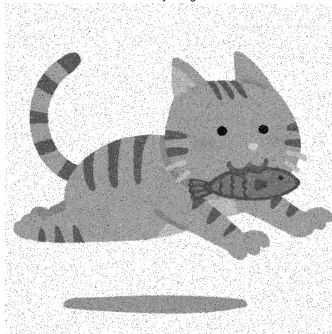
Application: Image denoising

- Remove noise from an image.
- **Preserve edges** at the same time.
 - Applying a low-pass filter does not work very well.

Original image



Noisy image



Total variation denoising

- $\mathbf{Y} \in \mathbb{R}^{n \times m}$: a noisy image
- Pull out each column vector, say $\mathbf{y} \in \mathbb{R}^n$, and solve the following optimization problem, one by one:

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{y}\|_2^2 + \lambda \sum_{i=1}^n |x_{i+1} - x_i|.$$

- $\|\mathbf{x} - \mathbf{y}\|_2^2$: proximity to the data
- $\sum_{i=1}^n |x_{i+1} - x_i|$: **total variation** to preserve edges
- $\lambda > 0$: weight

ADMM for total variation denoising

- Total variation denoising:

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad \|\mathbf{x} - \mathbf{y}\|_2^2 + \lambda \sum_{i=1}^n |x_{i+1} - x_i|.$$

- Define $\Phi = I$ and

$$\Psi = \begin{bmatrix} -1 & 1 & 0 & \dots & 0 \\ 0 & -1 & 1 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & 0 \\ 0 & \dots & 0 & -1 & 1 \\ 0 & \dots & 0 & 0 & -1 \end{bmatrix}.$$

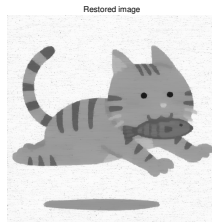
- Generalized LASSO

$$\underset{\mathbf{x} \in \mathbb{R}^n}{\text{minimize}} \quad \|\Phi \mathbf{x} - \mathbf{y}\|_2^2 + \lambda \|\Psi \mathbf{x}\|_1,$$

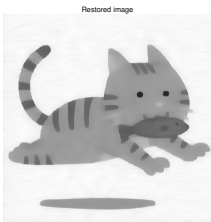
- We can use **ADMM**!

ADMM for total variation denoising

- The weight λ should be carefully chosen.



- $\lambda = 50$, $\lambda = 100$, and $\lambda = 200$



Conclusion

- In convex optimization, a local minimum is a global minimum.
- ℓ^1 optimization problems appeared in this course are convex optimization.
- Proximal operators are used to derive fast algorithms for convex optimization with non-differentiable ℓ^1 norm and constraints