

Chapter 4

Majorize-minimize methods

Contents (class version)

4.0 Introduction	4.2
4.1 Majorize-minimize (MM) principles	4.5
Properties of MM methods	4.6
Quadratic majorizer for smooth cost functions	4.7
4.2 Applications / MM examples	4.10
MM methods for LRMC	4.10
LRMC by iterative low-rank approximation	4.12
LASSO / sparse regression / compressed sensing	4.13
Convexity approach to separability	4.16
Poisson measurements (MLEM)	4.21
Line search using Huber's majorizer	4.23
Huber hinge function	4.32
Acceleration methods	4.35
4.3 Summary	4.36

4.0 Introduction

The **majorize-minimize (MM)** “algorithm” is really a family of methods for deriving algorithms for various applications. The applications are so numerous that entire books have been written on MM methods [1]. These notes focus on a couple of the most important *approaches* for designing majorizers for MM algorithms, in the context of several SIPML *applications* where **gradient descent** is inapplicable.

Approach overview

There are numerous techniques for designing MM algorithms; see [2, Ch. 14] at <http://web.eecs.umich.edu/~fessler/book/c-ox.pdf>

These notes focus on two main techniques that are used especially widely:

- **quadratic** majorizers, and
- majorizers built on **convexity**.

Perhaps surprisingly, both of these design methods can be useful even for **nonconvex** problems.

Applications

These notes focus on the following concrete motivating applications. The common thread of most of these applications is that the standard **gradient descent** algorithm is inapplicable due to terms that are not **differentiable**, or that are differentiable but do not have a **Lipschitz continuous** gradient.

- More efficient **line search** for certain cost functions

$$\arg \min_{\alpha} \Psi(\mathbf{x} + \alpha \mathbf{d}).$$

- Low-rank **matrix completion** (**LRMC**) (see EECS 551 notes) using the nondifferentiable rank constraint:

$$\hat{\mathbf{X}} = \arg \min_{\mathbf{X} \in \mathbb{F}^{M \times N}} \frac{1}{2} \|\mathbf{M} \odot (\mathbf{X} - \mathbf{Y})\|_{\text{F}}^2 + \chi_{\{\text{rank}(\mathbf{X}) \leq K\}}.$$

- **LASSO** / **sparse regression**

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \Psi(\mathbf{x}), \quad \Psi(\mathbf{x}) = \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2 + \beta \|\mathbf{x}\|_1 = \sum_{i=1}^M h_i([\mathbf{A}\mathbf{x}]_i) + \beta \|\mathbf{x}\|_1$$

$$h_i(t) = \frac{1}{2} |t - y_i|^2.$$

The sparsity promoting term $\|\cdot\|_1$ is not differentiable, yet it is essential for problems with numerous features (machine learning) or under-sampled measurements (**compressed sensing**).

- Binary classifier design using convex loss functions and sparsity regularization:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \Psi(\mathbf{x}), \quad \Psi(\mathbf{x}) = \sum_{i=1}^M h([\mathbf{A}\mathbf{x}]_i) + \beta \|\mathbf{x}\|_1,$$

where the i th row of \mathbf{A} is the product of the binary label $y_i = \pm 1$ and the i th feature vector. We will consider any convex loss function h including **hinge** loss (not differentiable), **exponential** loss (does not have a **Lipschitz continuous** derivative), and **logistic** loss. Although the logistic loss is **convex** with **Lipschitz continuous** derivative, the $\|\cdot\|_1$ term is not differentiable.

- Measurements having Poisson distributions (gamma rays, X-rays, low light level optical imaging) where $\mathbf{y} \sim \text{Poisson}\{\mathbf{A}\mathbf{x}\}$ for which the natural regularized estimator is:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \Psi(\mathbf{x}), \quad \Psi(\mathbf{x}) = \sum_{i=1}^M h_i([\mathbf{A}\mathbf{x}]_i) + \beta R(\mathbf{x}), \quad h_i(\lambda) = \lambda - y_i \log(\lambda),$$

where $h_i(\cdot)$ is the negative log-likelihood of the **Poisson distribution**.

This $h_i(\cdot)$ has a Lipschitz continuous derivative on $(0, \infty)$. (?)

A: True

B: False

??

The latter three applications all have data terms of the form $\sum_{i=1}^M h_i([\mathbf{A}\mathbf{x}]_i)$ for different convex h_i functions.

4.1 Majorize-minimize (MM) principles

To solve an optimization problem like

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x} \in \mathcal{X}} \Psi(\mathbf{x})$$

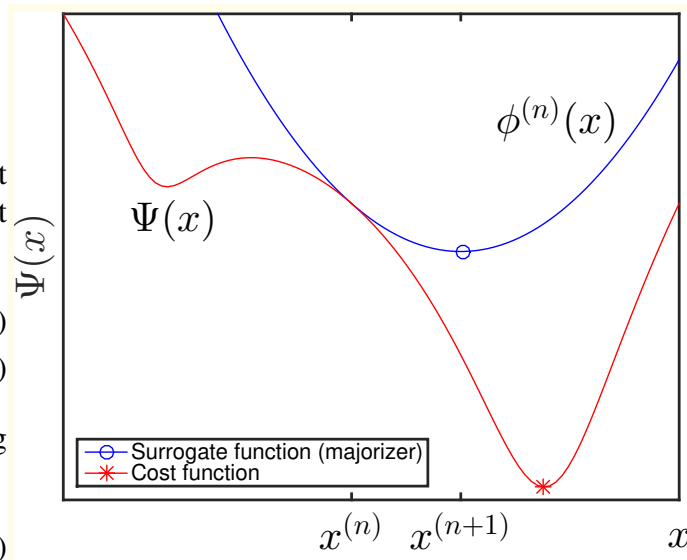
using a **majorize-minimize (MM)** approach, we first design a **majorizer** or **surrogate function** $\phi_k(\mathbf{x})$ that satisfies the following two conditions:

$$\Psi(\mathbf{x}_k) = \phi_k(\mathbf{x}_k) \quad (4.1)$$

$$\Psi(\mathbf{x}) \leq \phi_k(\mathbf{x}), \quad \forall \mathbf{x} \in \mathcal{X}. \quad (4.2)$$

Then the MM algorithm update is simply the following minimization step:

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x} \in \mathcal{X}} \phi_k(\mathbf{x}). \quad (4.3)$$



Properties of MM methods

Descent property

Any MM algorithm will monotonically decrease the cost function because of the **sandwich inequality**:

$$\Psi(\mathbf{x}_{k+1}) \leq \phi_k(\mathbf{x}_{k+1}) \leq \phi_k(\mathbf{x}_k) = \Psi(\mathbf{x}_k).$$

Algebraic properties of majorizers

Affine invariance:

If ϕ_k is a majorizer for Ψ , then $a\phi_k + b$ is a majorizer for $a\Psi + b$ for any $a > 0$.

Linearity:

If $f_k(\mathbf{x})$ is a majorizer for $f(\mathbf{x})$ and $g_k(\mathbf{x})$ is a majorizer for $g(\mathbf{x})$, then $af_k(\mathbf{x}) + bg_k(\mathbf{x})$ is a majorizer for $af(\mathbf{x}) + bg(\mathbf{x})$ for any $a, b > 0$.

Adding a suitable nonnegative function:

If $\zeta(\mathbf{x}, \mathbf{x}_k) \geq 0$ is a (nonnegative) function for which $\zeta(\mathbf{x}_k, \mathbf{x}_k) = 0$, then the following function defines a **majorizer** for Ψ :

$$\phi_k(\mathbf{x}) \triangleq \Psi(\mathbf{x}) + \zeta(\mathbf{x}, \mathbf{x}_k). \quad (4.4)$$

These properties allow us to design majorizers for individual pieces of a composite cost function and then combine them together to form a majorizer for the entire cost function.

Quadratic majorizer for smooth cost functions

Suppose $\Psi(\mathbf{x}) : \mathbb{F}^N \rightarrow \mathbb{R}$ is smooth, meaning its gradient is \mathbf{S} -Lipschitz continuous per (3.14):

$$\|\mathbf{S}^{-1}(\nabla \Psi(\mathbf{x}) - \nabla \Psi(\mathbf{z}))\|_2 \leq \|\mathbf{S}'(\mathbf{x} - \mathbf{z})\|_2, \quad \forall \mathbf{x}, \mathbf{z} \in \mathbb{F}^N.$$

Then a majorizer for $\Psi(\mathbf{x})$ (even for a **nonconvex** Ψ) is

$$\phi_k(\mathbf{x}) \triangleq \Psi(\mathbf{x}_k) + \text{real}\{\langle \nabla \Psi(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle\} + \frac{1}{2} \|\mathbf{S}'(\mathbf{x} - \mathbf{x}_k)\|_2^2. \quad (4.5)$$

Proof (generalizing [3, p. 22]). Clearly $\phi_k(\mathbf{x}_k) = \Psi(\mathbf{x}_k)$. By **Taylor's theorem** with remainder:

$$\begin{aligned} \Psi(\mathbf{x}) &= \Psi(\mathbf{z}) + \text{real}\left\{\int_0^1 \langle \nabla \Psi(\mathbf{z} + \tau(\mathbf{x} - \mathbf{z})), \mathbf{x} - \mathbf{z} \rangle d\tau\right\} \\ &= \Psi(\mathbf{z}) + \text{real}\{\langle \nabla \Psi(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle\} + \text{real}\left\{\int_0^1 \langle \nabla \Psi(\mathbf{z} + \tau(\mathbf{x} - \mathbf{z})) - \nabla \Psi(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle d\tau\right\} \\ &= \Psi(\mathbf{z}) + \text{real}\{\langle \nabla \Psi(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle\} \\ &\quad + \text{real}\left\{\int_0^1 \langle \mathbf{S}^{-1}(\nabla \Psi(\mathbf{z} + \tau(\mathbf{x} - \mathbf{z})) - \nabla \Psi(\mathbf{z})), \mathbf{S}'(\mathbf{x} - \mathbf{z}) \rangle d\tau\right\}. \end{aligned}$$

$$\begin{aligned}
&\implies |\Psi(\mathbf{x}) - \Psi(\mathbf{z}) - \text{real}\{\langle \nabla \Psi(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle\}| \\
&= \left| \text{real} \left\{ \int_0^1 \langle \mathbf{S}^{-1}(\nabla \Psi(\mathbf{z} + \tau(\mathbf{x} - \mathbf{z})) - \nabla \Psi(\mathbf{z})), \mathbf{S}'(\mathbf{x} - \mathbf{z}) \rangle d\tau \right\} \right| \\
&\leq \left| \int_0^1 \langle \mathbf{S}^{-1}(\nabla \Psi(\mathbf{z} + \tau(\mathbf{x} - \mathbf{z})) - \nabla \Psi(\mathbf{z})), \mathbf{S}'(\mathbf{x} - \mathbf{z}) \rangle d\tau \right| \\
&\leq \int_0^1 \|\mathbf{S}^{-1}(\nabla \Psi(\mathbf{z} + \tau(\mathbf{x} - \mathbf{z})) - \nabla \Psi(\mathbf{z}))\|_2 \|\mathbf{S}'(\mathbf{x} - \mathbf{z})\|_2 d\tau \\
&= \|\mathbf{S}'(\mathbf{x} - \mathbf{z})\|_2 \int_0^1 \|\mathbf{S}^{-1}(\nabla \Psi(\mathbf{z} + \tau(\mathbf{x} - \mathbf{z})) - \nabla \Psi(\mathbf{z}))\|_2 d\tau \\
&\leq \|\mathbf{S}'(\mathbf{x} - \mathbf{z})\|_2 \int_0^1 \|\tau \mathbf{S}'(\mathbf{x} - \mathbf{z})\|_2 d\tau = \|\mathbf{S}'(\mathbf{x} - \mathbf{z})\|_2^2 \int_0^1 \tau d\tau = \frac{1}{2} \|\mathbf{S}'(\mathbf{x} - \mathbf{z})\|_2^2,
\end{aligned}$$

using the \mathbf{S} -Lipschitz continuity and the **Cauchy-Schwarz inequality**.

Thus

$$\begin{aligned}
\Psi(\mathbf{x}) - \Psi(\mathbf{z}) - \text{real}\{\langle \nabla \Psi(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle\} &\leq \frac{1}{2} \|\mathbf{S}'(\mathbf{x} - \mathbf{z})\|_2^2 \\
\implies \Psi(\mathbf{x}) &\leq \Psi(\mathbf{z}) + \text{real}\{\langle \nabla \Psi(\mathbf{z}), \mathbf{x} - \mathbf{z} \rangle\} + \frac{1}{2} \|\mathbf{S}'(\mathbf{x} - \mathbf{z})\|_2^2.
\end{aligned}$$

This holds for all $\mathbf{x}, \mathbf{z} \in \mathbb{F}^N$ and now take $\mathbf{z} = \mathbf{x}_k$ to establish (4.5). □

A possible converse is explored in a **HW** problem.

The MM algorithm corresponding to the majorizer (4.5) has the following minimization step:

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \phi_k(\mathbf{x}) \text{ (from (4.3))}$$

$$\phi_k(\mathbf{x}) \triangleq \Psi(\mathbf{x}_k) + \text{real}\{\langle \nabla \Psi(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle\} + \frac{1}{2} \|\mathbf{S}'(\mathbf{x} - \mathbf{x}_k)\|_2^2 \text{ (from (4.5))}$$

$$\nabla \phi_k(\mathbf{x}) =$$

$$\mathbf{x}_{k+1} =$$

This MM algorithm minimization step is the same as:

A: GD

B: PGD

C: PSD

D: PCG

E: None of these

??

Special cases

Intuition for twice differentiable cost functions with bounded curvature. Eqn. (4.5) holds if:

$$\nabla^2 \Psi(\mathbf{x}) \preceq \mathbf{S}\mathbf{S}', \quad \forall \mathbf{x} \in \mathbb{R}^N.$$

If $\mathbf{S} = \sqrt{L}\mathbf{I}$ then

$$\Psi(\mathbf{x}) \leq \phi_k(\mathbf{x}) = \Psi(\mathbf{x}_k) + \text{real}\{\langle \nabla \Psi(\mathbf{x}_k), \mathbf{x} - \mathbf{x}_k \rangle\} + \frac{L}{2} \|\mathbf{x} - \mathbf{x}_k\|_2^2.$$

4.2 Applications / MM examples

MM methods for LRMC

Problem statement (for noisy data, see EECS 551 notes):

$$Y_{ij} = \begin{cases} X_{ij} + \varepsilon_{ij}, & (i, j) \in \Omega \\ 0, & \text{otherwise,} \end{cases} \quad \underbrace{M_{ij} \triangleq \begin{cases} 1, & (i, j) \in \Omega, \\ 0, & \text{otherwise,} \end{cases}}_{\text{mask}} \quad \underbrace{\tilde{M}_{i,j} = \begin{cases} 0, & (i, j) \in \Omega \\ 1, & (i, j) \notin \Omega, \end{cases}}_{\text{complement}}$$

for a sampling set $\Omega \subset \{1, \dots, M\} \times \{1, \dots, N\}$. Equivalently: $\tilde{M} \triangleq \mathbf{1}_M \mathbf{1}'_N - M$.

One possible LRMC formulation uses a (non-convex) **rank constraint**:

$$\hat{\mathbf{X}} = \arg \min_{\mathbf{X} : \text{rank}(\mathbf{X}) \leq K} q(\mathbf{X}), \quad q(\mathbf{X}) \triangleq \|\mathbf{M} \odot (\mathbf{X} - \mathbf{Y})\|_{\text{F}}^2. \quad (4.6)$$

Now define the following function:

$$Q(\mathbf{X}; \mathbf{Z}) = \|\mathbf{X} - \mathbf{Z} + \mathbf{M} \odot (\mathbf{Z} - \mathbf{Y})\|_{\text{F}}^2 = \left\| \mathbf{X} - (\tilde{\mathbf{M}} \odot \mathbf{Z} + \mathbf{M} \odot \mathbf{Y}) \right\|_{\text{F}}^2. \quad (4.7)$$

This function is a **majorizer** of $q(\mathbf{X})$ because

$$\begin{aligned} q(\mathbf{X}) &= Q(\mathbf{X}; \mathbf{X}) \\ q(\mathbf{X}) &\leq Q(\mathbf{X}; \mathbf{Z}), \quad \forall \mathbf{Z} \in \mathbb{F}^{M \times N}. \end{aligned}$$

Proof: Verifying $Q(\mathbf{X}; \mathbf{X}) = q(\mathbf{X})$ is easy.

(Read)

$$\begin{aligned}
 Q(\mathbf{X}; \mathbf{Z}) &= \|\mathbf{X} - \mathbf{Z} + \mathbf{M} \odot (\mathbf{Z} - \mathbf{Y})\|_{\mathbf{F}}^2 = \left\| (\tilde{\mathbf{M}} + \mathbf{M}) \odot (\mathbf{X} - \mathbf{Z}) + \mathbf{M} \odot (\mathbf{Z} - \mathbf{Y}) \right\|_{\mathbf{F}}^2 \\
 &= \left\| \tilde{\mathbf{M}} \odot (\mathbf{X} - \mathbf{Z}) + \mathbf{M} \odot (\mathbf{X} - \mathbf{Y}) \right\|_{\mathbf{F}}^2 = \left\| \tilde{\mathbf{M}} \odot (\mathbf{X} - \mathbf{Z}) \right\|_{\mathbf{F}}^2 + \left\| \mathbf{M} \odot (\mathbf{X} - \mathbf{Y}) \right\|_{\mathbf{F}}^2 \\
 &\geq \left\| \mathbf{M} \odot (\mathbf{X} - \mathbf{Y}) \right\|_{\mathbf{F}}^2 = q(\mathbf{X}),
 \end{aligned}$$

where the 4th equality holds because $\tilde{\mathbf{M}} \odot \mathbf{M} = \mathbf{0}$. □

In the earlier notation, define $q_k(\mathbf{X}) \triangleq Q(\mathbf{X}, \mathbf{X}_k)$ at the k th iteration. Then the MM requirements hold:

$$\begin{aligned}
 q_k(\mathbf{X}_k) &= Q(\mathbf{X}_k, \mathbf{X}_k) = q(\mathbf{X}_k) \\
 q_k(\mathbf{X}) &= Q(\mathbf{X}, \mathbf{X}_k) \geq q(\mathbf{X}), \quad \forall \mathbf{X}.
 \end{aligned}$$

I originally designed this majorizer by making a 2nd-order Taylor expansion of $q(\mathbf{X})$ about \mathbf{Z} and then using the fact that the elements of \mathbf{M} are all 0 or 1. We'll see another application of that approach later in this chapter. Here, the proof above illustrates a simpler way to design the majorizer, using the following definition:

$$Q(\mathbf{X}, \mathbf{Z}) \triangleq q(\mathbf{X}) + \left\| \tilde{\mathbf{M}} \odot (\mathbf{X} - \mathbf{Z}) \right\|_{\mathbf{F}}^2, \quad (4.8)$$

because the 2nd term is nonnegative and is zero when $\mathbf{X} = \mathbf{Z}$, per (4.4).

Now we use the **quadratic majorizer** (4.7) = (4.8) to develop a MM algorithm for LRMC that is practical for problems that are not too large.

LRMC by iterative low-rank approximation

Consider LRMC using the rank constraint (4.6):

$$\hat{\mathbf{X}} = \arg \min_{\mathbf{X} : \text{rank}(\mathbf{X}) \leq K} \|\mathbf{M} \odot (\mathbf{X} - \mathbf{Y})\|_{\text{F}}^2.$$

The MM algorithm update for this formulation is simply

$$\begin{aligned} \mathbf{X}_{k+1} &= \arg \min_{\mathbf{X} : \text{rank}(\mathbf{X}) \leq K} Q(\mathbf{X}; \mathbf{X}_k) \\ &= \arg \min_{\mathbf{X} : \text{rank}(\mathbf{X}) \leq K} \left\| \mathbf{X} - (\tilde{\mathbf{M}} \odot \mathbf{X}_k + \mathbf{M} \odot \mathbf{Y}) \right\|_{\text{F}}^2 \\ &= \arg \min_{\mathbf{X} : \text{rank}(\mathbf{X}) \leq K} \left\| \mathbf{X} - \tilde{\mathbf{X}}_k \right\|_{\text{F}}^2, \quad \tilde{\mathbf{X}}_k \triangleq \tilde{\mathbf{M}} \odot \mathbf{X}_k + \mathbf{M} \odot \mathbf{Y} = \begin{cases} Y_{i,j}, & (i,j) \in \Omega \\ [\mathbf{X}_k]_{i,j}, & (i,j) \notin \Omega. \end{cases} \end{aligned}$$

This algorithm alternates between two steps:

- Take the current guess \mathbf{X}_k and replace all the values at sampled locations with the measurements from \mathbf{Y} to get $\tilde{\mathbf{X}}_k$. (As mentioned earlier, this can be done “in place” using $\mathbf{X}[\mathbf{M}] = \mathbf{Y}[\mathbf{M}]$.)
- Perform low-rank (rank at most K) approximation (using **SVD**) to $\tilde{\mathbf{X}}_k$ to get the next iterate \mathbf{X}_{k+1} .

Because this is a MM method, we know it decreases the Frobenius norm cost function monotonically.

(There is still no guarantee here of convergence of $\{\mathbf{X}_k\}$ to a global minimizer because the rank constraint set is nonconvex. For convex formulations see [4].)

LASSO / sparse regression / compressed sensing

These applications all involve the following **composite** cost function:

$$\hat{\mathbf{x}} = \arg \min_{\mathbf{x}} \Psi(\mathbf{x}), \quad \Psi(\mathbf{x}) = f(\mathbf{x}) + \beta \|\mathbf{x}\|_1, \quad f(\mathbf{x}) = \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{y}\|_2^2. \quad (4.9)$$

There are simple non-iterative solutions when \mathbf{A} is identity, diagonal, or unitary, but not otherwise.

First draw inspiration from (4.4) to define a **quadric majorizer** for $f(\mathbf{x})$:

$$f_k(\mathbf{x}) \triangleq \quad (4.10)$$

$$\begin{aligned} &= \\ &= \\ &= \end{aligned} \quad (4.11)$$

by **completing the square**, where c_1, c_2 are constants independent of \mathbf{x} .

Here we require $\mathbf{H} \succeq \mathbf{0}$ and $\mathbf{A}'\mathbf{A} + \mathbf{H} \succ \mathbf{0}$, but we do not insist that $\mathbf{H} \succ \mathbf{0}$.

ISTA: iterative soft thresholding algorithm

$$f_k(\mathbf{x}) =$$

Thus a majorizer for the original cost function is

$$\phi_k(\mathbf{x}) =$$

$$=$$

(4.12)

The minimization step is the **iterative soft thresholding algorithm (ISTA)** [5]:

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \phi_k(\mathbf{x}) =$$

(4.13)

It is also called the **proximal gradient method (PGM)**, especially for more general regularizers. This MM algorithm is remarkably simple, especially given that the 1-norm is non-differentiable.

A serious drawback of ISTA is its slow convergence $O(1/k)$ and it needs $\|\mathbf{A}\|_2$.

The threshold above is:

A: β

B: β/L

C: βL

D: βL^2

E: None

??

PGM with diagonal majorizers

Recall from HW1#10 :

(4.14)

With this choice of \mathbf{H} , the majorizer becomes

$$\phi_k(\mathbf{x}) =$$

and the MM algorithm becomes

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \phi_k(\mathbf{x}) =$$

(4.15)

where $\mathbf{D} = \text{diag}\{\mathbf{d}\}$ and \oslash denotes $\cdot /$ element-wise division. So the j th element of the gradient step is thresholded by β/d_j . No Lipschitz constant or operator norm is needed.

It still has slow $O(1/k)$ convergence rate, but there is a fast version called FISTA [6] and an even faster version called the **proximal optimized gradient method (POGM)** [7] that both have $O(1/k^2)$ bounds.

A HW problem will compare ISTA / FISTA / POGM with the `lasso_cls` constrained least-squares formulation from Ch. 1.

Convexity approach to separability

All of the preceding examples used **quadratic majorizers** for the data term. Now we turn to non-quadratic examples.

Many of cost functions have data-fit terms of this form:

$$f(\mathbf{x}) =$$

where each function h_i is **convex** (but not necessarily smooth).

Examples:

- $h_i(t) = \frac{1}{2} |t - y_i|^2$
- $h_i(t) = \max(1 - t, 0)$ (**hinge loss**, the tightest convex upper bound on the 0-1 loss)
- $h_i(t) = e^{-t}$ (**exponential loss**)
- $h_i(t) = t - y_i \log(t)$ (Poisson negative log-likelihood)

The commonality of all of these cost function components is convexity (only), *i.e.*,

(4.16)

The importance of this inequality is that the RHS is **additively separable**.

Here is the key technique for designing a **separable** majorizer using convexity:

$$[\mathbf{A}\mathbf{x}]_i =$$

$$h_i([\mathbf{A}\mathbf{x}]_i) =$$

$$\leq$$

$$\implies \sum_{i=1}^M h_i([\mathbf{A}\mathbf{x}]_i) \leq$$

Note that $\sum_{j=1}^N f_j^{(k)}(x_j^{(k)}) = \sum_{j=1}^N \sum_{i=1}^M \alpha_{ij} h_i([\mathbf{A}\mathbf{x}_k]_i) = \sum_{i=1}^M h_i([\mathbf{A}\mathbf{x}_k]_i)$, so (4.1) holds.

Now consider a cost function like

$$\Psi(\mathbf{x}) = \sum_{i=1}^M h_i([\mathbf{A}\mathbf{x}]_i) + \beta \|\mathbf{x}\|_1 = \sum_{i=1}^M h_i([\mathbf{A}\mathbf{x}]_i) + \beta \sum_{j=1}^N |x_j|.$$

The above inequalities establish the following separable majorizer:

$$\phi_k(\mathbf{x}) = \quad \quad \quad (4.17)$$

The (parallelizable!) minimization step of the MM algorithm for this majorizer is

$$\mathbf{x}_{k+1} = \arg \min_{\mathbf{x}} \phi_k(\mathbf{x}), \quad x_j^{(k+1)} =$$

Example. Consider again the LASSO problem (4.9), where $h_i(t) = \frac{1}{2} |t - y_i|^2$. The separable majorizer (4.17) is:

$$\begin{aligned}
 \phi_j^{(k)}(x_j) &= f_j^{(k)}(x_j) + \beta |x_j| = \sum_{i=1}^M \alpha_{ij} h_i \left(\frac{a_{ij}}{\alpha_{ij}} (x_j - x_j^{(k)}) + [\mathbf{A}\mathbf{x}_k]_i \right) + \beta |x_j| \\
 &= \sum_{i=1}^M \alpha_{ij} \frac{1}{2} \left| \frac{a_{ij}}{\alpha_{ij}} (x_j - x_j^{(k)}) + [\mathbf{A}\mathbf{x}_k]_i - y_i \right|^2 + \beta |x_j| \\
 &= \sum_{i=1}^M \frac{1}{\alpha_{ij}} \frac{1}{2} \left| a_{ij} (x_j - x_j^{(k)}) + \alpha_{ij} ([\mathbf{A}\mathbf{x}_k]_i - y_i) \right|^2 + \beta |x_j| \\
 &= \frac{1}{2} d_j \left| x_j - [\mathbf{x}_k - \mathbf{D}^{-1} \mathbf{A}'(\mathbf{A}\mathbf{x}_k - \mathbf{y})]_j \right|^2 + \beta |x_j| + c_3, \quad \mathbf{D} = \text{diag}\{d_j\}, \quad d_j \triangleq \sum_{i=1}^M |a_{ij}|^2 / \alpha_{ij},
 \end{aligned}$$

where c_3 is a constant independent of \mathbf{x} .

So the minimization step of the MM algorithm is simply soft thresholding:

$$\mathbf{x}_{k+1} = \text{soft}(\mathbf{x}_k - \mathbf{D}^{-1} \mathbf{A}'(\mathbf{A}\mathbf{x}_k - \mathbf{y}), \beta \oslash \mathbf{d}).$$

To finalize the algorithm we must specify α_{ij} factors satisfying (4.16). A natural choice, cf. (4.14), is:

$$\alpha_{ij} = \frac{|a_{ij}|}{\sum_{k=1}^N |a_{ik}|} \implies d_j = \sum_{i=1}^M \frac{|a_{ij}|^2}{\alpha_{ij}} = \sum_{i=1}^M \frac{|a_{ij}|^2}{|a_{ij}| / \sum_{k=1}^N |a_{ik}|} = \sum_{i=1}^M |a_{ij}| \left(\sum_{k=1}^N |a_{ik}| \right) = [|\mathbf{A}|' |\mathbf{A}| \mathbf{1}]_j.$$

If $\beta = 0$, this MM algorithm is the same as:

A: GD

B: PGD

C: PSD

D: PCG

E: None of these

??

If $\beta > 0$, this MM algorithm is the same as:

A: GD

B: PGD

C: ISTA/PGM (4.13)

D: ISTA/PGM (4.15)

E: None of these

??

Example. Exponential loss for binary classifier design $h_i(t) = e^{-t}$.

(Read)

For this application the separable majorizer (4.17) is

$$\begin{aligned}\phi_j^{(k)}(x_j) &= f_j^{(k)}(x_j) + \beta |x_j| = \sum_{i=1}^M \alpha_{ij} h_i \left(\frac{a_{ij}}{\alpha_{ij}} (x_j - x_j^{(k)}) + [\mathbf{A}\mathbf{x}_k]_i \right) + \beta |x_j| \\ &= \sum_{i=1}^M \alpha_{ij} \exp \left(- \left(\frac{a_{ij}}{\alpha_{ij}} (x_j - x_j^{(k)}) + [\mathbf{A}\mathbf{x}_k]_i \right) \right) + \beta |x_j|.\end{aligned}$$

There is no closed-form expression for the MM algorithm update, but one can perform parallel 1D minimization for each coordinate $j = 1, \dots, N$:

$$x_j^{(k+1)} = \arg \min_{x_j} \phi_j^{(k)}(x_j) = \arg \min_{x_j} \sum_{i=1}^M \alpha_{ij} \exp \left(- \left(\frac{a_{ij}}{\alpha_{ij}} (x_j - x_j^{(k)}) + [\mathbf{A}\mathbf{x}_k]_i \right) \right) + \beta |x_j|.$$

The same principle applies to the **hinge** loss function. However, in both cases the 1D minimization itself seems challenging: the exponential does not have a Lipschitz derivative and the hinge is not differentiable.

Poisson measurements (MLEM)

For the Poisson noise problem with $h_i(t) = t - y_i \log(t)$, we must have $a_{ij} \geq 0$ and we use the convexity inequality with a slightly different approach:

$$\begin{aligned}
 [\mathbf{Ax}]_i &= \sum_{j=1}^N a_{ij} x_j = \\
 h_i([\mathbf{Ax}]_i) &= h_i\left(\sum_{j=1}^N \left(\frac{a_{ij} x_j^{(k)}}{[\mathbf{Ax}]_i}\right) \frac{x_j}{x_j^{(k)}} [\mathbf{Ax}]_i\right) \leq \sum_{j=1}^N \left(\frac{a_{ij} x_j^{(k)}}{[\mathbf{Ax}]_i}\right) h_i\left(\frac{x_j}{x_j^{(k)}} [\mathbf{Ax}]_i\right) \\
 \implies \sum_{i=1}^M h_i([\mathbf{Ax}]_i) &\leq \sum_{j=1}^N f_j^{(k)}(x_j), \quad f_j^{(k)}(x_j) \triangleq \sum_{i=1}^M \left(\frac{a_{ij} x_j^{(k)}}{[\mathbf{Ax}]_i}\right) h_i\left(\frac{x_j}{x_j^{(k)}} [\mathbf{Ax}]_i\right)
 \end{aligned}$$

Here we are essentially using iteration-dependent convexity parameters:

$$\alpha_{ij} =$$

Again considering 1-norm regularization, the separable majorizer is

$$\phi_k(\mathbf{x}) = \sum_{j=1}^N \phi_j^{(k)}(x_j), \quad \phi_j^{(k)}(x_j) = f_j^{(k)}(x_j) + \beta |x_j|.$$

Here we need $x_j \geq 0$ so $|x_j| = x_j$.

Differentiating to find the minimizer over x_j :

$$\begin{aligned} 0 = \frac{\partial}{\partial x_j} \phi_j^{(k)}(x_j) &= \sum_{i=1}^M a_{ij} \dot{h}_i \left(\frac{x_j}{x_j^{(k)}} [\mathbf{A}\mathbf{x}_k]_i \right) + \beta = \sum_{i=1}^M a_{ij} \left(1 - \frac{x_j^{(k)} y_i}{x_j [\mathbf{A}\mathbf{x}_k]_i} \right) + \beta \\ &= \sum_{i=1}^M a_{ij} - \frac{x_j^{(k)}}{x_j} \sum_{i=1}^M a_{ij} \frac{y_i}{[\mathbf{A}\mathbf{x}_k]_i} + \beta \\ \Rightarrow x_j^{(k+1)} &= \frac{x_j^{(k)}}{\sum_{i=1}^M a_{ij} + \beta} \sum_{i=1}^M a_{ij} \frac{y_i}{[\mathbf{A}\mathbf{x}_k]_i}. \end{aligned}$$

The version with $\beta = 0$ is used in clinical PET and SPECT systems daily.

The 1-norm seems to be of much less help in enforcing sparsity here compared to the case of the 2-norm data term [8].

Line search using Huber's majorizer

Recall that for “inverse problem” cost functions of the form $\Psi(\mathbf{x}) = \sum_{j=1}^J f_j(\mathbf{B}_j \mathbf{x})$ where each component function f_j is convex and smooth, we used GD to solve the 1D **line-search** minimization problem:

$$\alpha_k = \arg \min_{\alpha} h_k(\alpha), \quad h_k(\alpha) = \Psi(\mathbf{x}_k + \alpha \mathbf{d}_k) = \sum_{j=1}^J f_j(\mathbf{u}_j + \alpha \mathbf{v}_j), \quad \mathbf{u}_j \triangleq \mathbf{B}_j \mathbf{x}_k, \quad \mathbf{v}_j \triangleq \mathbf{B}_j \mathbf{d}_k.$$

To apply GD for the line search, in **HW** you found a Lipschitz constant of the derivative of $h_k(\alpha)$ to be

$$L_{h_k} = \sum_{j=1}^J L_{f_j} \|\mathbf{v}_j\|_2^2.$$

That Lipschitz constant leads to correct but undesirably small step sizes for many functions of interest, thereby possibly requiring excess numbers of inner line search iterations.

We focus on component functions of the form (note the JULIA dot):

$$f_j(\mathbf{v}) = \mathbf{1}' \psi_j \cdot (\mathbf{v}),$$

for which

$$L_{f_j} = L_{\psi}.$$

If ψ has a Lipschitz continuous derivative, then specializing (4.5):

$$\psi(t) \leq \psi(s) + L_{\dot{\psi}}(t-s) + \frac{L_{\ddot{\psi}}}{2}(t-s)^2. \quad (4.18)$$

The GD-based line-search used in HW was essentially a MM method using this quadratic majorizer. I call this the **maximum curvature** majorizer because if

$$\left| \dot{\psi}(t) - \dot{\psi}(s) \right| \leq L_{\dot{\psi}} |t - s|$$

and ψ is twice differentiable at t , then

$$\ddot{\psi}(t) =$$

So the Lipschitz constant for the derivative of a function is a (tight) bound for its maximum second derivative.

Often we can find quadratic majorizers for ψ with lower curvature than $L_{\dot{\psi}}$.

Lower curvature majorizers lead to larger step sizes and hence convergence in fewer iterations.

Huber's majorizer

Theorem (Huber, 1981 [9, p. 184], [2, Ch. 14]). Suppose $\psi : \mathbb{R} \rightarrow \mathbb{R}$ satisfies the following conditions

$$\begin{aligned} \psi(t) &\text{ is differentiable,} \\ \psi(t) &= \psi(-t), \quad \forall t \text{ (symmetry),} \\ \omega_\psi(t) &\triangleq \dot{\psi}(t)/t \text{ is bounded and monotone nonincreasing for } t > 0. \end{aligned} \quad (4.19)$$

Then the parabola function defined by

$$\begin{aligned} q(t; s) &\triangleq \text{ } \\ &= \left(\psi(s) - \frac{\omega_\psi(s)}{2} s^2 \right) + \frac{\omega_\psi(s)}{2} t^2, \end{aligned} \quad (4.20)$$

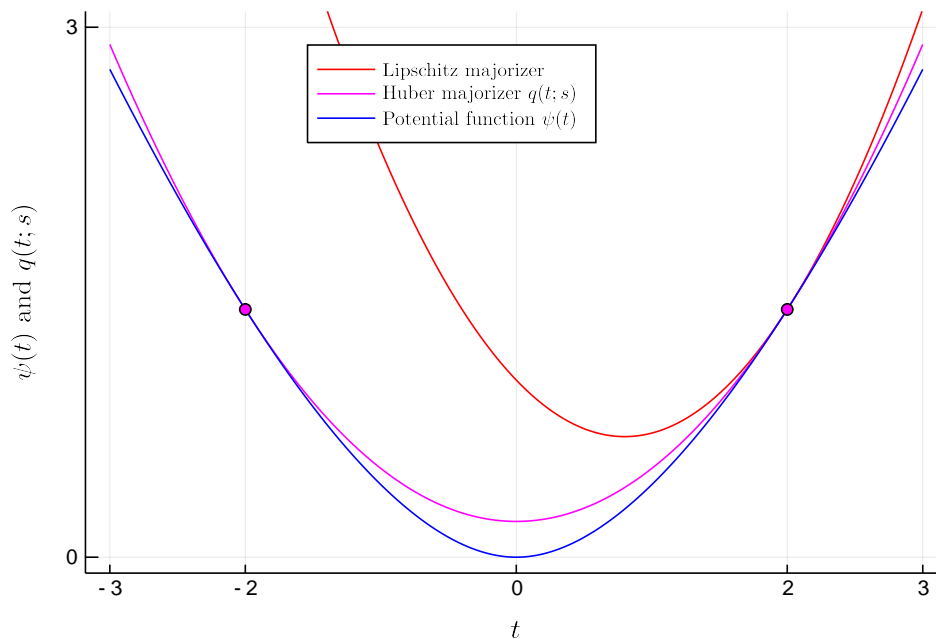
is a majorizer for ψ , *i.e.*, it satisfies conditions analogous to (4.1) and (4.2), namely

$$\begin{aligned} \psi(s) &= q(s; s) \\ \psi(t) &\leq q(t; s), \quad \forall s, t. \end{aligned}$$

The curvature ω_ψ is optimal in the sense of being the smallest value that ensures those requirements.

Challenge. Prove or disprove: Under Huber's conditions, $\dot{\psi}$ is Lipschitz continuous, and $L_{\dot{\psi}} = \omega_\psi(0)$.

Example. This figure illustrates the Huber quadratic majorizer (4.20) and the Lipschitz quadratic majorizer (4.18) for the **Fair potential** with $\delta = 3$ at $s = 2$. Clearly the Huber majorizer is tighter, allowing larger step sizes.



MM line search using Huber's majorizer

To realize these larger steps, one uses a MM approach instead of GD for the line search.

Dropping the iteration k subscript to simplify notation, the line-search problem is

$$\alpha_* = \arg \min_{\alpha} h(\alpha), \quad h(\alpha) = \Psi(\mathbf{x} + \alpha \mathbf{d}) = \sum_{j=1}^J f_j(\mathbf{u}_j + \alpha \mathbf{v}_j), \quad \begin{aligned} \mathbf{u}_j &\triangleq \mathbf{B}_j \mathbf{x}, \\ \mathbf{v}_j &\triangleq \mathbf{B}_j \mathbf{d}, \end{aligned} \quad f_j(\mathbf{v}) = \mathbf{1}' \psi_j \cdot (\mathbf{v}).$$

(Dropping k here also avoids confusion with the inner MM iteration for the line search.)

Assuming each ψ_j satisfies Huber's conditions, the following quadratic function is a majorizer for $h(\alpha)$ for the n th inner iteration:

$$\begin{aligned} \phi_n(\alpha) &\triangleq q(\alpha; \alpha_n) = \sum_{j=1}^J q_j(\alpha; \alpha_n) \\ q_j(\alpha; \alpha_n) &= \mathbf{1}' \left(\psi_j \cdot (\mathbf{u}_j + \alpha_n \mathbf{v}_j) + \dot{\psi}_j \cdot (\mathbf{u}_j + \alpha_n \mathbf{v}_j) \odot \mathbf{v}_j (\alpha - \alpha_n) + \frac{1}{2} \omega_j \cdot (\mathbf{u}_j + \alpha_n \mathbf{v}_j) \odot \mathbf{v}_j \cdot^2 (\alpha - \alpha_n)^2 \right) \\ &= c_0 + c_1(\alpha_n)(\alpha - \alpha_n) + \frac{1}{2} c_2(\alpha_n)(\alpha - \alpha_n)^2 \\ c_0 &= h(\alpha_n) \end{aligned}$$

$$c_1(\alpha_n) = \sum_{j=1}^J \operatorname{real} \left\{ \mathbf{v}_j' \dot{\psi}_j \cdot (\mathbf{u}_j + \alpha_n \mathbf{v}_j) \right\} = \dot{h}(\alpha_n)$$

$$c_2(\alpha_n) = \sum_{j=1}^J (\mathbf{v}_j \cdot \omega_j)^2 \omega_j \cdot (\mathbf{u}_j + \alpha_n \mathbf{v}_j) = \sum_{j=1}^J \mathbf{v}_j' \operatorname{diag}\{\omega_j \cdot (\mathbf{u}_j + \alpha_n \mathbf{v}_j)\} \mathbf{v}_j \leq \sum_{j=1}^J L_{\psi_j} \|\mathbf{v}_j\|_2^2.$$

Differentiating $q(\alpha; \alpha_n)$ w.r.t. α and setting to zero leads to the following MM update:

$$\alpha_{n+1} = \arg \min_{\alpha} \phi_n(\alpha) =$$

A HW problem will compare the speed of this to the GD-based line search used previously.

1D optimization using Huber's majorizer

Facts:

$$\arg \min_x \sum_{n=1}^N \frac{1}{2} |x - x_n|^2 = \frac{1}{N} \sum_{n=1}^N x_n \quad (\text{sample mean})$$

$$\arg \min_x \sum_{n=1}^N |x - x_n| = \text{median}(x_1, \dots, x_N) \quad (\text{sample median})$$

What about more general 1D case:

$$\hat{x} = \arg \min_x f(x), \quad f(x) = \sum_{n=1}^N \psi(x - x_n).$$

In general there is no closed-form solution for \hat{x} . If ψ satisfies Huber's conditions then we can form a quadratic majorizer:

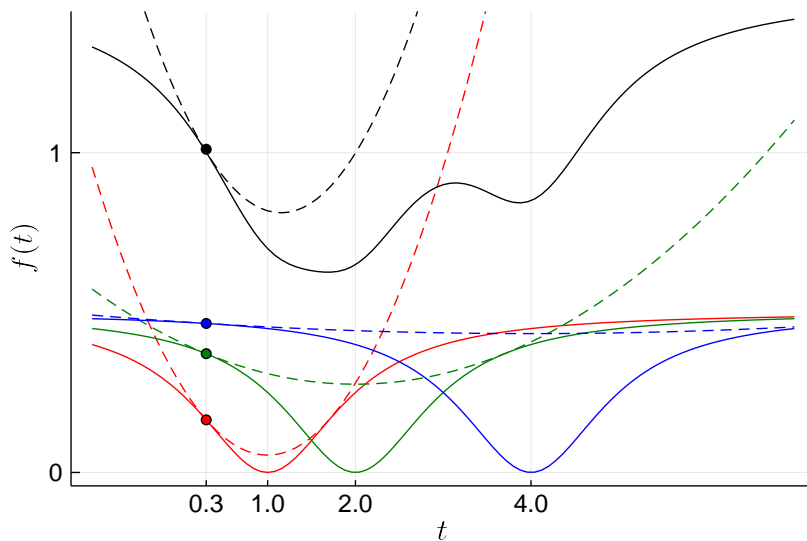
$$f(x) \leq \phi^{(t)}(x) = \sum_{n=1}^N \psi(x^{(t)} - x_n) + \dot{\psi}(x^{(t)} - x_n)(x - x^{(t)}) + \omega_{\psi}(x^{(t)} - x_n) \frac{1}{2}(x - x^{(t)})^2$$

and update by zeroing its derivative:

$$x^{(t+1)} = x^{(t)} - \frac{\dot{\phi}^{(t)}(x^{(t)})}{\ddot{\phi}^{(t)}(x^{(t)})} = x^{(t)} - \frac{\sum_{n=1}^N \dot{\psi}(x^{(t)} - x_n)}{\sum_{n=1}^N \omega_{\psi}(x^{(t)} - x_n)}.$$

This update will decrease $f(x)$ monotonically, whereas Newton's method has no such guarantee in general.

Example. Consider the Geman & McClure potential $\psi(t) = \frac{t^2/2}{1+t^2}$ for which $\omega_\psi(t) = \frac{1}{(1+t^2)^2}$, and the case where $x_1 = 1$, $x_2 = 2$, $x_3 = 4$. Here the Huber quadratic majorizer works as needed, whereas the Newton parabola would fail.



Huber's conditions and uniqueness

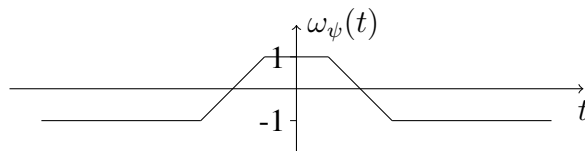
Do Huber's 3 conditions (4.19) ensure ψ has a unique minimizer?

ψ differentiable, symmetric and $\omega_\psi(t) = \dot{\psi}(t)/t$ is bounded and nonincreasing for $t > 0$.

Example.

Consider this weighting function and note that

$$\psi(t) = \int_0^t \dot{\psi}(s) ds = \int_0^t s \omega_\psi(s) ds.$$

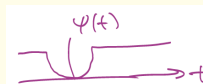
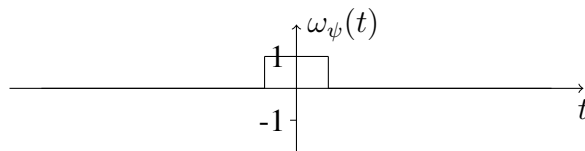


Now suppose we add a 4th condition:

$$\omega_\psi \geq 0$$

Now we have that $\dot{\psi}(t) \geq 0 \forall t \geq 0$,

so $\psi(t)$ is non-decreasing for $t \geq 0$.



Huber hinge function

The **hinge loss** function is not differentiable, making it unsuitable for gradient-based methods.

The (continuously differentiable) **Huber hinge** loss function for classifier design is defined, for $\delta > 0$, as [10]:

$$h(t; \delta) = \begin{cases} 1 - t - \delta/2, & t \leq 1 - \delta \\ \frac{1}{2\delta}(t - 1)^2, & 1 - \delta \leq t \leq 1 \\ 0, & 1 \leq t. \end{cases}$$

The Lipschitz constant of the derivative of this function is:

A: $1/\delta^2$

B: $1/\delta$

C: 1

D: δ

E: δ^2

??

A quadratic majorizer with this curvature is used in [10, eqn. (34)] to train a **support-vector machine** (SVM).

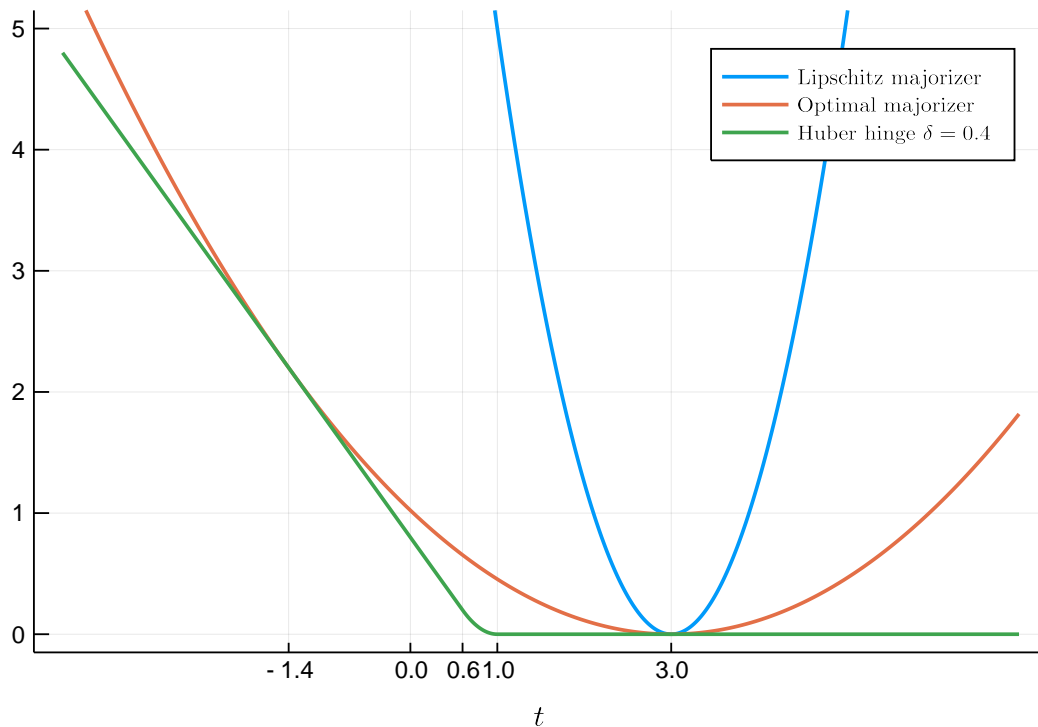
For $s > 1$, find the quadratic function having *optimal*, i.e., smallest possible curvature:

$$q(t; s) =$$

What is $c(s)$?

(group)

Example. The following figure illustrates the case where $\delta = 2/5$ and $s = 3$.



If $1 - \delta \leq s \leq 1$ then the curvature of the optimal quadratic majorizer is?

A: $1/\delta^2$

B: $1/\delta$

C: 1

D: δ

E: δ^2

??

If $1 - \delta \leq s \leq 1$ then the vertex of the optimal quadratic majorizer is?

A: 0

B: s

C: 1

D: $1/s$

E: None of these

??

A HW problem considers the case $s < 1 - \delta$.

The “best quadratic majorizer” for the (nondifferentiable) case $\delta = 0$ (the ordinary hinge function) is given in [11, eqn. (9)]. They call it “sharp” quadratic majorization. However, the expression there is invalid at the hinge corner point, so the authors propose an ad hoc modification that destroys the guarantee of the monotone descent property of MM methods, so is quite undesirable. Using $\delta > 0$ gives us the Huber hinge for which our quadratic majorizer is optimal.

An alternative to the Huber hinge is the “squared hinge loss” function [12] given by $h(t) = (\max(1 - t, 0))^2$. This function is convex and continuously differentiable and has a Lipschitz continuous derivative. However it increases unnecessarily rapidly (quadratically) for negative values, whereas the Huber hinge increases linearly.

Acceleration methods

For unconstrained problems, one acceleration method that retains monotonicity is the over-relaxation approach of [13]. In the MM setting, one first makes a standard MM update:

$$\tilde{\mathbf{x}}_{t+1} \triangleq \arg \min_{\mathbf{x}} \phi_t(\mathbf{x})$$

and then takes an extrapolated step

$$\mathbf{x}_{t+1} \triangleq \tilde{\mathbf{x}}_{t+1} + \alpha_t(\tilde{\mathbf{x}}^{(t+1)} - \mathbf{x}_t)$$

where the step size α_t is chosen to ensure that

$$\phi_t(\mathbf{x}_{t+1}) \leq \phi_t(\mathbf{x}_t),$$

which in turn ensures that $\Psi(\mathbf{x}_t)$ decreases monotonically. This method is viable when evaluating the majorizer ϕ_t is less expensive than a line search on Ψ .

4.3 Summary

This chapter has touched on a few of the many methods for designing MM optimization algorithms, focusing on quadratic majorizers and majorizers based on convexity.

One can combine the methods. For example in transmission tomography the data-fit term involves $h_i(t) = (b_i \exp(-t) + r_i) - y_i \log(b_i \exp(-t) + r_i)$ which is nonconvex. Nevertheless, one can design a quadratic majorizer for it [14] and then make a separable quadratic majorizer of that [15] to derive a simple MM algorithm.

For generalizations of MM, see the “relatively smooth” approach of [16] and the generalized MM (G-MM) approach, see [17].

Bibliography

- [1] K. Lange. *MM optimization algorithms*. Soc. Indust. Appl. Math., 2016 (cit. on p. 4.2).
- [2] J. A. Fessler. *Image reconstruction: Algorithms and analysis*. Book in preparation. , 2006 (cit. on pp. 4.2, 4.25).
- [3] Y. Nesterov. *Introductory lectures on convex optimization: A basic course*. Springer, 2004 (cit. on p. 4.7).
- [4] C. Grussler and P. Giselsson. “Low-rank inducing norms with optimality interpretations”. In: *SIAM J. Optim.* 28.4 (Jan. 2018), 3057–78 (cit. on p. 4.12).
- [5] I. Daubechies, M. Defrise, and C. De Mol. “An iterative thresholding algorithm for linear inverse problems with a sparsity constraint”. In: *Comm. Pure Appl. Math.* 57.11 (Nov. 2004), 1413–57 (cit. on p. 4.14).
- [6] A. Beck and M. Teboulle. “A fast iterative shrinkage-thresholding algorithm for linear inverse problems”. In: *SIAM J. Imaging Sci.* 2.1 (2009), 183–202 (cit. on p. 4.15).
- [7] A. B. Taylor, J. M. Hendrickx, and Francois Glineur. “Exact worst-case performance of first-order methods for composite convex optimization”. In: *SIAM J. Optim.* 27.3 (Jan. 2017), 1283–313 (cit. on p. 4.15).
- [8] D. J. Lingensfelter, J. A. Fessler, and Z. He. “Sparsity regularization for image reconstruction with Poisson data”. In: *Proc. SPIE 7246 Computational Imaging VII*. 2009, 72460F (cit. on p. 4.22).
- [9] P. J. Huber. *Robust statistics*. New York: Wiley, 1981 (cit. on p. 4.25).
- [10] P. J. F. Groenen, G. Nalbantov, and J. C. Bioch. “SVM-Maj: a majorization approach to linear support vector machines with different hinge errors”. In: *Advances in Data Analysis and Classification* 2.1 (Apr. 2008), 17–43 (cit. on p. 4.32).
- [11] J. de Leeuw and K. Lange. “Sharp quadratic majorization in one dimension”. In: *Comp. Stat. Data Anal.* 53.7 (May 2009), 2471–84 (cit. on p. 4.34).
- [12] Z. Zhang, D. Liu, G. Dai, and M. I. Jordan. “Coherence functions with applications in large-margin classification methods”. In: *J. Mach. Learning Res.* 13 (Sept. 2012), 2705–34 (cit. on p. 4.34).
- [13] Y. Yu. “Monotonically overrelaxed EM algorithms”. In: *J. Computational and Graphical Stat.* 21.2 (2012), 518–37 (cit. on p. 4.35).
- [14] H. Erdogan and J. A. Fessler. “Monotonic algorithms for transmission tomography”. In: *IEEE Trans. Med. Imag.* 18.9 (Sept. 1999), 801–14 (cit. on p. 4.36).
- [15] H. Erdogan and J. A. Fessler. “Ordered subsets algorithms for transmission tomography”. In: *Phys. Med. Biol.* 44.11 (Nov. 1999), 2835–51 (cit. on p. 4.36).

- [16] H. Lu, R. M. Freund, and Y. Nesterov. “Relatively smooth convex optimization by first-order methods, and applications”. In: *SIAM J. Optim.* 28.1 (Jan. 2018), 333–54 (cit. on p. [4.36](#)).
- [17] S. N. Parizi, K. He, S. Sclaroff, and P. Felzenszwalb. “Generalized majorization-minimization”. In: *Proc. Intl. Conf. Mach. Learn.* Vol. 97. 2019, 5022–31 (cit. on p. [4.36](#)).