

DTU





January 17-18, 2023

02946 Scientific Computing for X-Ray Computed Tomography

Optimization Methods for Tomography

About me

Martin S. Andersen

Associate Professor
Section for Scientific Computing
DTU Compute

Head of Studies, MSc Eng MMC

PhD from UCLA (EE)
MSc from Aalborg University (EE)

Research in optimization and its applications

June 2023: 02953 Convex Optimization (5 ECTS)



Tentative schedule

Chapter 13 in textbook

Tuesday, January 17

Unconstrained optimization
Lipschitz continuity
Majorization minimization
Convexity
Step size rules & stopping criteria
Power iteration

Wednesday, January 18

Constrained optimization
Convex sets
Proximal gradient method
Optimality conditions
Accelerated proximal gradient method
Smoothing techniques

Optimization for tomography

Maximum likelihood (ML) estimation

$$\hat{\mathbf{x}} = \operatorname{argmax}_{\mathbf{x}} \{\pi(\mathbf{b} | \mathbf{x})\} = \operatorname{argmin}_{\mathbf{x}} \{-\ln \pi(\mathbf{b} | \mathbf{x})\}$$

Maximum a posteriori (MAP) estimation

$$\begin{aligned}\hat{\mathbf{x}} &= \operatorname{argmax}_{\mathbf{x}} \{\pi(\mathbf{b} | \mathbf{x})\pi(\mathbf{x})\} \\ &= \operatorname{argmin}_{\mathbf{x}} \{-\ln \pi(\mathbf{b} | \mathbf{x}) - \ln \pi(\mathbf{x})\}\end{aligned}$$

Example

$$\text{minimize } \frac{1}{2} \|\mathbf{b} - \mathbf{Ax}\|_2^2 + \gamma R(\mathbf{x}) + \text{const.}$$

Unconstrained optimization

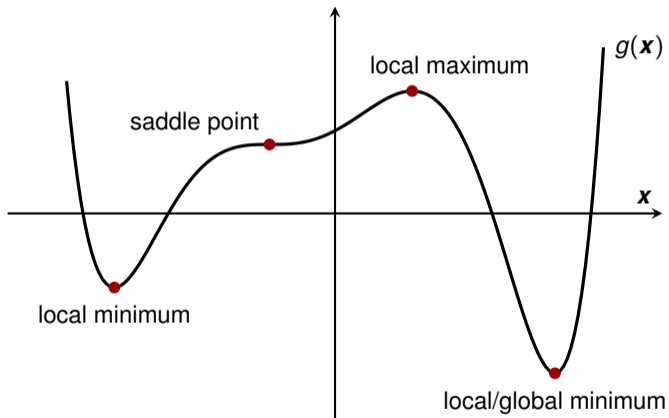
minimize $g(\mathbf{x})$

- variable $\mathbf{x} \in \mathbb{R}^n$
- objective function $g : \mathbb{R}^n \rightarrow \mathbb{R}$ continuously differentiable
- global minimum at \mathbf{x}^* if $g(\mathbf{y}) \geq g(\mathbf{x}^*)$ for all $\mathbf{y} \in \mathbb{R}^n$
- \mathbf{x} is a stationary point of g if

$$\nabla g(\mathbf{x}) = \begin{bmatrix} \frac{\partial g(\mathbf{x})}{\partial x_1} \\ \vdots \\ \frac{\partial g(\mathbf{x})}{\partial x_n} \end{bmatrix} = \mathbf{0}$$

- stationarity is a necessary condition for global optimality

Stationary points



Gradient method

Iterative update of image

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - t_k \nabla g(\mathbf{x}^{(k)}), \quad k = 0, 1, 2, \dots$$

- step size $t_k > 0$
- directional derivative of g at $\mathbf{x}^{(k)}$ in the direction $-\nabla g(\mathbf{x}^{(k)})$ is

$$-\nabla g(\mathbf{x}^{(k)})^T \nabla g(\mathbf{x}^{(k)}) = -\|\nabla g(\mathbf{x}^{(k)})\|_2^2$$

- directional derivative is negative unless $\mathbf{x}^{(k)}$ is a stationary point
- implies that $-\nabla g(\mathbf{x}^{(k)})$ is a descent direction if $\mathbf{x}^{(k)}$ is not stationary
- descent is guaranteed if we choose t_k such that $g(\mathbf{x}^{(k+1)}) < g(\mathbf{x}^{(k)})$

Exact line search

Cauchy's step size rule: minimize g along the current search direction

$$t_k = \operatorname{argmin}_{t>0} \left\{ g(\mathbf{x}^{(k)} - t\nabla g(\mathbf{x}^{(k)})) \right\}$$

- “greedy” heuristic
- may be as expensive to solve as original problem

Example: least-squares objective

$$g(\mathbf{x}) = \frac{1}{2} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2^2$$

- gradient $\nabla g(\mathbf{x}) = \mathbf{A}^T(\mathbf{A}\mathbf{x} - \mathbf{b})$
- exact line search

$$t_k = \operatorname{argmin}_{t>0} \left\{ g(\mathbf{x}^{(k)} - t\nabla g(\mathbf{x}^{(k)})) \right\} = \frac{\|\nabla g(\mathbf{x}^{(k)})\|_2^2}{\|\mathbf{A}\nabla g(\mathbf{x}^{(k)})\|_2^2} = \frac{\|\mathbf{A}^T \boldsymbol{\rho}^{(k)}\|_2^2}{\|\mathbf{A}\mathbf{A}^T \boldsymbol{\rho}^{(k)}\|_2^2}$$

which follows from $\boldsymbol{\rho}^{(k)} = \mathbf{b} - \mathbf{A}\mathbf{x}^{(k)}$ and

$$\frac{d}{dt} g(\mathbf{x}^{(k)} - t\nabla g(\mathbf{x}^{(k)})) = t\|\mathbf{A}\nabla g(\mathbf{x}^{(k)})\|_2^2 - \|\nabla g(\mathbf{x}^{(k)})\|_2^2 = 0$$

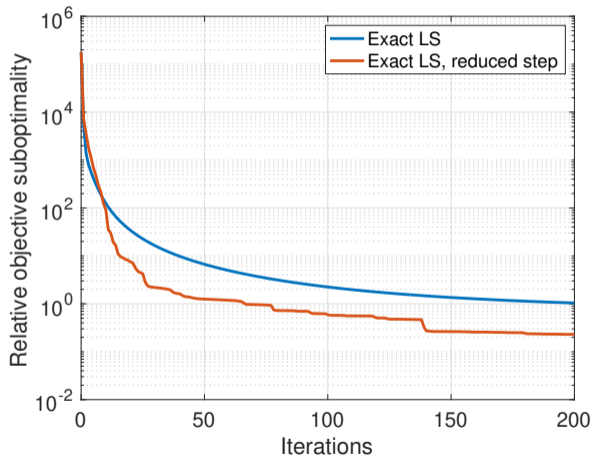
Example: least-squares objective (cont.)

Relative suboptimality

$$\frac{|g(\mathbf{x}^{(k)}) - g(\mathbf{x}^*)|}{|g(\mathbf{x}^*)|}$$

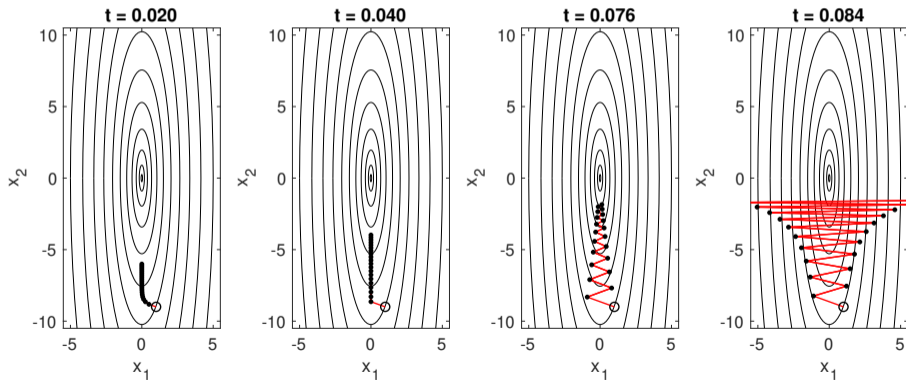
Adjusted step size

$$\gamma t_k, \quad \gamma = 0.9$$



Example: gradient method with fixed step size (first 20 iterations)

$$g(\mathbf{x}) = \frac{1}{2}(25x_1^2 + x_2^2)$$



Lipschitz continuity

Gradient ∇g is Lipschitz continuous if there exists a constant L such that

$$\|\nabla g(\mathbf{x}) - \nabla g(\mathbf{y})\|_2 \leq L\|\mathbf{x} - \mathbf{y}\|_2, \quad \text{for all } \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$$

Quadratic upper bound

If ∇g is Lipschitz continuous with constant L , then

$$g(\mathbf{y}) \leq g(\mathbf{x}) + \nabla g(\mathbf{x})^T(\mathbf{y} - \mathbf{x}) + \frac{L}{2}\|\mathbf{y} - \mathbf{x}\|_2^2, \quad \text{for all } \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$$

Quadratic upper bound: derivation

- Define restriction of g to line through \mathbf{x} and \mathbf{y} : $\phi(\tau) = g(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x}))$
- Newton–Leibniz integral rule: $\phi(1) - \phi(0) = \int_0^1 \phi'(\tau) d\tau$

$$\begin{aligned}g(\mathbf{y}) - g(\mathbf{x}) &= \int_0^1 \nabla g(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x}))^T (\mathbf{y} - \mathbf{x}) d\tau \\&= \nabla g(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \int_0^1 (\nabla g(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x})) - \nabla g(\mathbf{x}))^T (\mathbf{y} - \mathbf{x}) d\tau \\&\leq \nabla g(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \int_0^1 \|\nabla g(\mathbf{x} + \tau(\mathbf{y} - \mathbf{x})) - \nabla g(\mathbf{x})\|_2 \|\mathbf{y} - \mathbf{x}\|_2 d\tau \\&\leq \nabla g(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \int_0^1 \tau L \|\mathbf{y} - \mathbf{x}\|_2^2 d\tau\end{aligned}$$

Example

Least-squares objective

$$g(\mathbf{x}) = \frac{1}{2} \|\mathbf{b} - \mathbf{Ax}\|_2^2$$

with gradient $\nabla g(\mathbf{x}) = \mathbf{A}^T(\mathbf{Ax} - \mathbf{b})$

Implies that

$$\|\nabla g(\mathbf{y}) - \nabla g(\mathbf{x})\|_2 = \|\mathbf{A}^T \mathbf{A}(\mathbf{y} - \mathbf{x})\|_2 \leq \|\mathbf{A}^T \mathbf{A}\|_2 \|\mathbf{y} - \mathbf{x}\|_2$$

and hence ∇g is Lipschitz continuous with constant $L = \|\mathbf{A}\|_2^2$

Twice continuously differentiable functions

Suppose g is twice continuously differentiable with Hessian matrix

$$\nabla^2 g(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 g(\mathbf{x})}{\partial x_1^2} & \frac{\partial^2 g(\mathbf{x})}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 g(\mathbf{x})}{\partial x_1 \partial x_n} \\ \frac{\partial^2 g(\mathbf{x})}{\partial x_2 \partial x_1} & \frac{\partial^2 g(\mathbf{x})}{\partial x_2^2} & \cdots & \frac{\partial^2 g(\mathbf{x})}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 g(\mathbf{x})}{\partial x_n \partial x_1} & \frac{\partial^2 g(\mathbf{x})}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 g(\mathbf{x})}{\partial x_n^2} \end{bmatrix}$$

Bounded Hessian

∇g is Lipschitz continuous with constant L iff $\|\nabla^2 g(\mathbf{x})\|_2 \leq L$ for all \mathbf{x}

Exercise 13.1: Step size rules for least-squares problems

Consider the gradient method applied to the least-squares objective function $g(\mathbf{x}) = \frac{1}{2} \|\mathbf{b} - \mathbf{Ax}\|_2^2$, i.e.,

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - t_k \mathbf{A}^T (\mathbf{Ax}^{(k)} - \mathbf{b}), \quad k = 0, 1, 2, \dots$$

where $\mathbf{x}^{(0)}$ is an initial guess. For each of the following step size rules, show that the gradient iteration can be implemented such that each iteration only requires a single matrix-vector multiplication with \mathbf{A} and one with \mathbf{A}^T .

- 1 The step size t_k is constant, i.e., $t_k = t > 0$ for all k .
- 2 The step size t_k is found by means of the exact line search.

Exercise 13.1: Step size rules for least-squares problems (solution)

- 1 Constant step size: $t_k = t$

$$\begin{aligned}\nabla g(\mathbf{x}^{(k)}) &= \mathbf{A}^T(\mathbf{A}\mathbf{x}^{(k)} - \mathbf{b}) = -\mathbf{A}^T \boldsymbol{\rho}^{(k)} \\ \boldsymbol{\rho}^{(k+1)} &= \mathbf{b} - \mathbf{A}(\mathbf{x}^{(k)} - t\nabla g(\mathbf{x}^{(k)})) = \boldsymbol{\rho}^{(k)} + t\mathbf{A}\nabla g(\mathbf{x}^{(k)})\end{aligned}$$

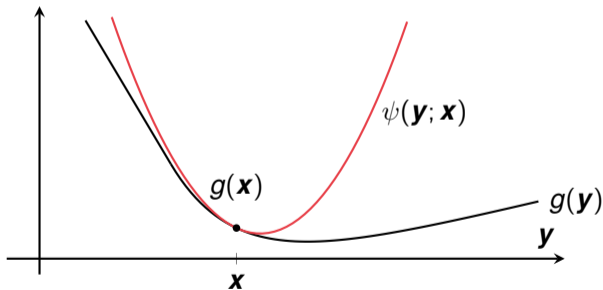
- 2 Exact line search: $t_k = \frac{\|\nabla g(\mathbf{x}^{(k)})\|_2^2}{\|\mathbf{A}\nabla g(\mathbf{x}^{(k)})\|_2^2} = \frac{\|\mathbf{A}^T \boldsymbol{\rho}^{(k)}\|_2^2}{\|\mathbf{A}\nabla g(\mathbf{x}^{(k)})\|_2^2}$

$$\begin{aligned}\nabla g(\mathbf{x}^{(k)}) &= -\mathbf{A}^T \boldsymbol{\rho}^{(k)} \\ \mathbf{y} &= \mathbf{A}\nabla g(\mathbf{x}^{(k)}) \\ \boldsymbol{\rho}^{(k+1)} &= \boldsymbol{\rho}^{(k)} + \frac{\|\nabla g(\mathbf{x}^{(k)})\|_2^2}{\|\mathbf{y}\|_2^2} \mathbf{y}\end{aligned}$$

Majorization

A function $\psi(\mathbf{y}; \mathbf{x})$ is said to be a majorization of g at \mathbf{x} if

$$\psi(\mathbf{x}; \mathbf{x}) = g(\mathbf{x}) \quad \text{and} \quad \psi(\mathbf{y}; \mathbf{x}) \geq g(\mathbf{y}), \quad \text{for all } \mathbf{y}$$



Majorization minimization

Iterative update based on majorization

$$\mathbf{x}^{(k+1)} = \operatorname{argmin}_{\mathbf{y}} \left\{ \psi(\mathbf{y}; \mathbf{x}^{(k)}) \right\}, \quad k = 0, 1, 2, \dots$$

- $\mathbf{x}^{(k+1)}$ minimizes the majorization $\psi(\mathbf{y}; \mathbf{x}^{(k)})$ so

$$\psi(\mathbf{x}^{(k+1)}; \mathbf{x}^{(k)}) \leq \psi(\mathbf{x}^{(k)}; \mathbf{x}^{(k)})$$

- properties of majorization imply that

$$g(\mathbf{x}^{(k+1)}) \leq \psi(\mathbf{x}^{(k+1)}; \mathbf{x}^{(k)}) \leq \psi(\mathbf{x}^{(k)}; \mathbf{x}^{(k)}) = g(\mathbf{x}^{(k)})$$

Majorization minimization: quadratic majorization

Use quadratic upper bound to construct majorization

$$\psi(\mathbf{y}; \mathbf{x}) = g(\mathbf{x}) + \nabla g(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|_2^2$$

with ∇g Lipschitz continuous with constant L

Gradient method with constant step size

$$\mathbf{x}^{(k+1)} = \operatorname{argmin}_{\mathbf{y}} \left\{ \psi(\mathbf{y}; \mathbf{x}^{(k)}) \right\} = \mathbf{x}^{(k)} - t_k \nabla g(\mathbf{x}^{(k)}), \quad t_k = \frac{1}{L}$$

Analysis of gradient method with constant step size

Majorization property $g(\mathbf{x}^{(k+1)}) \leq \psi(\mathbf{x}^{(k+1)}; \mathbf{x}^{(k)})$ implies that

$$g(\mathbf{x}^{(k+1)}) \leq g(\mathbf{x}^{(k)}) - \frac{1}{2L} \|\nabla g(\mathbf{x}^{(k)})\|_2^2$$

- summing inequality for $k = 0, \dots, N$,

$$\frac{1}{2L} \sum_{k=0}^N \|\nabla g(\mathbf{x}^{(k)})\|_2^2 \leq g(\mathbf{x}^{(0)}) - g(\mathbf{x}^{(N+1)}) \leq g(\mathbf{x}^{(0)}) - g(\mathbf{x}^*)$$

- converges to stationary point if $g(\mathbf{x}^{(0)}) - g(\mathbf{x}^*)$ is finite
- step size $t_k \in (0, 2/L)$ yields a descent unless $\nabla g(\mathbf{x}^{(k)}) = 0$

Exercise 13.2: Lipschitz continuous gradients

Suppose $g_1 : \mathbb{R}^n \rightarrow \mathbb{R}$ and $g_2 : \mathbb{R}^n \rightarrow \mathbb{R}$ are continuously differentiable.

Show that if ∇g_1 and ∇g_2 are Lipschitz continuous with constants L_1 and L_2 , respectively, then $\nabla g(\mathbf{x}) = \nabla g_1(\mathbf{x}) + \nabla g_2(\mathbf{x})$ is Lipschitz continuous with constant $L = L_1 + L_2$.

Exercise 13.2: Lipschitz continuous gradients (solution)

We have that $\nabla(g_1 + g_2) = \nabla g_1 + \nabla g_2$, and hence (triangle inequality)

$$\begin{aligned} \|\nabla g_1(\mathbf{y}) + \nabla g_2(\mathbf{y}) - \nabla g_1(\mathbf{x}) - \nabla g_2(\mathbf{x})\|_2 &\leq \\ &\|\nabla g_1(\mathbf{y}) - \nabla g_1(\mathbf{x})\|_2 + \|\nabla g_2(\mathbf{y}) - \nabla g_2(\mathbf{x})\|_2. \end{aligned}$$

Thus, if ∇g_1 is L_1 -Lipschitz and ∇g_2 is L_2 -Lipschitz, then

$$\|\nabla g_1(\mathbf{y}) + \nabla g_2(\mathbf{y}) - \nabla g_1(\mathbf{x}) - \nabla g_2(\mathbf{x})\|_2 \leq (L_1 + L_2)\|\mathbf{y} - \mathbf{x}\|_2$$

which shows that $\nabla g_1 + \nabla g_2$ is $(L_1 + L_2)$ -Lipschitz.

SIRT-like methods

Recall the SIRT method

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \lambda_k \mathbf{D} \mathbf{A}^T \mathbf{M} (\mathbf{A} \mathbf{x}^{(k)} - \mathbf{b})$$

where $\lambda_k \in (0, 2)$ and \mathbf{M} and \mathbf{D} are positive diagonal matrices

May be viewed as *scaled* gradient method for minimizing

$$g(\mathbf{x}) = \frac{1}{2} (\mathbf{b} - \mathbf{A} \mathbf{x})^T \mathbf{M} (\mathbf{b} - \mathbf{A} \mathbf{x}) = \frac{1}{2} \|\mathbf{b} - \mathbf{A} \mathbf{x}\|_{\mathbf{M}}^2$$

with gradient $\nabla g(\mathbf{x}) = \mathbf{A}^T \mathbf{M} (\mathbf{A} \mathbf{x} - \mathbf{b})$

SIRT-like methods (cont.)

Gradient satisfies

$$\begin{aligned}\|\nabla g(\mathbf{y}) - \nabla g(\mathbf{x})\|_{\mathbf{D}} &= \|\mathbf{D}^{1/2} \mathbf{A}^T \mathbf{M} \mathbf{A} \mathbf{D}^{1/2} \mathbf{D}^{-1/2} (\mathbf{y} - \mathbf{x})\|_2 \\ &\leq \|\mathbf{M}^{1/2} \mathbf{A} \mathbf{D}^{1/2}\|_2^2 \|\mathbf{D}^{-1/2} (\mathbf{y} - \mathbf{x})\|_2\end{aligned}$$

Assuming that \mathbf{D} and \mathbf{M} satisfy $\|\mathbf{M}^{1/2} \mathbf{A} \mathbf{D}^{1/2}\|_2 \leq 1$,

$$\|\nabla g(\mathbf{y}) - \nabla g(\mathbf{x})\|_{\mathbf{D}} \leq \|\mathbf{y} - \mathbf{x}\|_{\mathbf{D}^{-1}}$$

SIRT-like methods (cont.)

Condition $\|\nabla g(\mathbf{y}) - \nabla g(\mathbf{x})\|_{\mathbf{D}} \leq \|\mathbf{y} - \mathbf{x}\|_{\mathbf{D}^{-1}}$ implies quadratic upper bound

$$g(\mathbf{y}) \leq g(\mathbf{x}) + \nabla g(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{1}{2} (\mathbf{y} - \mathbf{x})^T \mathbf{D}^{-1} (\mathbf{y} - \mathbf{x})$$

Majorization minimization method for minimizing g

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \mathbf{x}^{(k)} - \lambda_k \mathbf{D} \nabla g(\mathbf{x}^{(k)}) \\ &= \mathbf{x}^{(k)} - \lambda_k \mathbf{D} \mathbf{A}^T \mathbf{M} (\mathbf{A} \mathbf{x} - \mathbf{b}) \end{aligned}$$

with $\lambda_k \in (0, 2)$

SIRT-like methods (cont.)

$\|M^{1/2}AD^{1/2}\|_2 \leq 1$ is satisfied for D and M defined as

$$D_{jj}^{-1} = \sum_{i=1}^m |\mathbf{A}_{ij}|^\alpha, \quad M_{ii}^{-1} = \sum_{j=1}^n |\mathbf{A}_{ij}|^{2-\alpha}, \quad \alpha \in [0, 2]$$

- we define $|\mathbf{A}_{ij}|^0 = 1$ when $\mathbf{A}_{ij} = 0$
- objective function $g(\mathbf{x}) = \frac{1}{2} \|\mathbf{b} - \mathbf{Ax}\|_M^2$ depends on α
- Cimmino's method: $\alpha = 0$
- SIRT: $\alpha = 1$
- "parallel" coordinate descent: $\alpha = 2$

Exercise 13.3: SIRT-like methods

Recall that the SIRT iteration solves a weighted least-squares problem of the form

$$\text{minimize } \frac{1}{2} \|\mathbf{b} - \mathbf{Ax}\|_{\mathbf{M}}^2, \quad \mathbf{M} \text{ diag. positive definite.}$$

- 1 Show that $\|\mathbf{M}^{1/2}\mathbf{AD}^{1/2}\|_2 \leq 1$ if \mathbf{M} and \mathbf{D} are diagonal matrices and

$$\mathbf{D}_{jj}^{-1} = \sum_{i=1}^m |\mathbf{A}_{ij}|^\alpha, \quad \mathbf{M}_{ii}^{-1} = \sum_{j=1}^n |\mathbf{A}_{ij}|^{2-\alpha}, \quad \alpha \in [0, 2].$$

Hint: Show that $\|\mathbf{M}^{1/2}\mathbf{AD}^{1/2}\mathbf{x}\|_2^2 \leq \|\mathbf{x}\|_2^2$ when $\alpha \in [0, 2]$.

- 2 Implement the SIRT iteration in MATLAB with α as an input parameter.
- 3 Compute reconstructions for different α (see textbook for details).

Exercise 13.3: SIRT-like methods (solution, question 1)

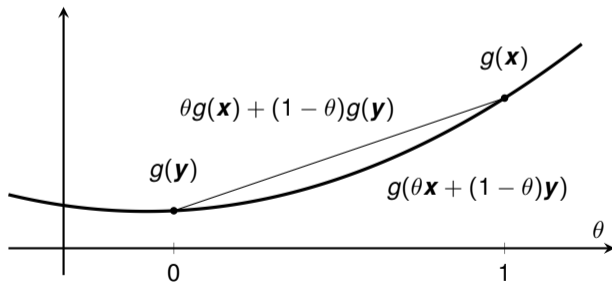
$$\begin{aligned}
 \|\mathbf{M}^{1/2} \mathbf{A} \mathbf{D}^{1/2} \mathbf{x}\|_2^2 &= \sum_{i=1}^m \left(\mathbf{M}_{ii}^{1/2} \sum_{j=1}^n \mathbf{A}_{ij} \mathbf{D}_{jj}^{1/2} x_j \right)^2 \\
 &\leq \sum_{i=1}^m \mathbf{M}_{ii} \left(\sum_{j=1}^n |\mathbf{A}_{ij}| |x_j| \mathbf{D}_{jj}^{1/2} \right)^2 \\
 &= \sum_{i=1}^m \mathbf{M}_{ii} \left(\sum_{j=1}^n |\mathbf{A}_{ij}|^{1-\alpha/2} |\mathbf{A}_{ij}|^{\alpha/2} |x_j| \mathbf{D}_{jj}^{1/2} \right)^2 \\
 \text{(Cauchy-Schwarz)} \quad &\leq \sum_{i=1}^m \mathbf{M}_{ii} \left(\sum_{j=1}^n |\mathbf{A}_{ij}|^{2-\alpha} \right) \left(\sum_{j=1}^n |\mathbf{A}_{ij}|^\alpha \mathbf{D}_{jj} x_j^2 \right) \\
 &= \sum_{i=1}^m \left(\sum_{j=1}^n |\mathbf{A}_{ij}|^\alpha \mathbf{D}_{jj} x_j^2 \right) = \sum_{j=1}^n \mathbf{D}_{jj} \left(\sum_{i=1}^m |\mathbf{A}_{ij}|^\alpha \right) x_j^2 = \|\mathbf{x}\|_2^2.
 \end{aligned}$$

Convexity

$g : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if and only if

$$g(\theta \mathbf{x} + (1 - \theta)\mathbf{y}) \leq \theta g(\mathbf{x}) + (1 - \theta)g(\mathbf{y}), \quad \theta \in [0, 1]$$

for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ (g is concave if $-g$ is convex)

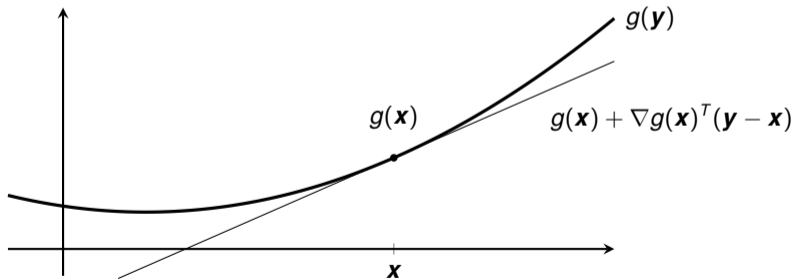


Convexity: first-order condition

Continuously differentiable g is convex if and only if

$$g(\mathbf{y}) \geq g(\mathbf{x}) + \nabla g(\mathbf{x})^T (\mathbf{y} - \mathbf{x})$$

for all $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$



Convexity: implications

- stationary points are global minimizers

$$g(\mathbf{y}) \geq g(\mathbf{x}^*) + \nabla g(\mathbf{x}^*)^T (\mathbf{y} - \mathbf{x}^*) = g(\mathbf{x}^*), \quad \text{for all } \mathbf{y} \in \mathbb{R}^n$$

- gradient method with step size $t_k = \gamma/L$ and $\gamma \in (0, 2)$ satisfies

$$g(\mathbf{x}^{(k)}) - g(\mathbf{x}^*) \leq \frac{2L\|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2}{4 + \gamma(2 - \gamma)k}.$$

if ∇g is Lipschitz continuous with constant L

- suboptimality satisfies $g(\mathbf{x}^{(k)}) - g^* = O(1/k)$
- at most $O(1/\epsilon)$ iterations required before $g(\mathbf{x}^{(k)}) - g^* \leq \epsilon$

Strong convexity

g is strongly convex with parameter $\mu > 0$ if

$$g(\theta \mathbf{x} + (1 - \theta) \mathbf{y}) \leq \theta g(\mathbf{x}) + (1 - \theta) g(\mathbf{y}) - \frac{\theta(1 - \theta)\mu}{2} \|\mathbf{x} - \mathbf{y}\|_2^2$$

for all $\theta \in [0, 1]$ and $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$

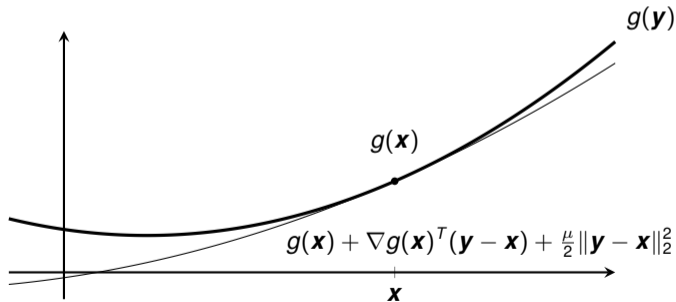
Interpretation: $\tilde{g}(\mathbf{x}) = g(\mathbf{x}) - \frac{\mu}{2} \|\mathbf{x}\|_2^2$ is convex

Strong convexity: first-order condition

Continuously differentiable g : strongly convex with parameter $\mu > 0$ if

$$g(\mathbf{y}) \geq g(\mathbf{x}) + \nabla g(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{\mu}{2} \|\mathbf{y} - \mathbf{x}\|_2^2, \quad \mathbf{x}, \mathbf{y} \in \mathbb{R}^n$$

Implies that minimizer is unique



Strong convexity: implications

- minimizing quadratic lower bound wrt. \mathbf{y} yields

$$g(\mathbf{y}) \geq g(\mathbf{x}) - \frac{1}{2\mu} \|\nabla g(\mathbf{x})\|_2^2 \implies g(\mathbf{x}) - g(\mathbf{x}^*) \leq \frac{1}{2\mu} \|\nabla g(\mathbf{x})\|_2^2$$

- substitute \mathbf{x}^* for \mathbf{y} in first-order condition

$$\begin{aligned} g(\mathbf{x}^*) &\geq g(\mathbf{x}) + \nabla g(\mathbf{x})^T (\mathbf{x}^* - \mathbf{x}) + \frac{\mu}{2} \|\mathbf{x}^* - \mathbf{x}\|_2^2 \\ &\geq g(\mathbf{x}) - \|\nabla g(\mathbf{x})\|_2 \|\mathbf{x}^* - \mathbf{x}\|_2 + \frac{\mu}{2} \|\mathbf{x}^* - \mathbf{x}\|_2^2 \end{aligned}$$

$g(\mathbf{x}^*) \leq g(\mathbf{x})$ implies that

$$\|\mathbf{x}^* - \mathbf{x}\|_2 \leq \frac{2}{\mu} \|\nabla g(\mathbf{x})\|_2$$

Strong convexity: implications (cont.)

- gradient method with step size $t_k = 2/(L + \mu)$ with satisfies

$$g(\mathbf{x}^{(k)}) - g(\mathbf{x}^*) \leq \frac{L}{2} \left(\frac{L - \mu}{L + \mu} \right)^{2k} \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2^2,$$

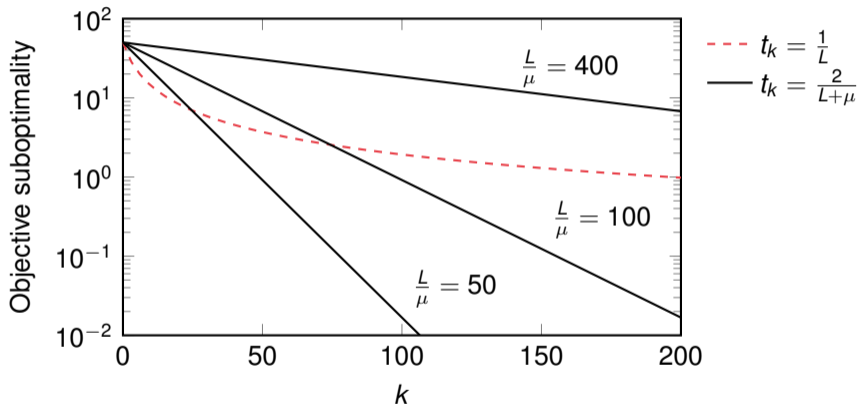
and

$$\|\mathbf{x}^{(k)} - \mathbf{x}^*\|_2 \leq \left(\frac{L - \mu}{L + \mu} \right)^k \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2$$

if g is μ -strongly convex with L -Lipschitz gradient

- implies that $\mathbf{x}^{(k)}$ converges linearly to \mathbf{x}^*

Comparison of worst-case suboptimality bounds



Example: least-squares problem

$$g(\mathbf{x}) = \frac{1}{2} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2^2$$

Linearization of gradient around \mathbf{x}^* yields

$$\nabla g(\mathbf{x}^{(k)}) = \nabla g(\mathbf{x}^*) + \nabla^2 g(\mathbf{x}^*)(\mathbf{x}^{(k)} - \mathbf{x}^*)$$

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - t\mathbf{A}^T\mathbf{A}(\mathbf{x}^{(k)} - \mathbf{x}^*),$$

Subtract \mathbf{x}^* from both sides, take norm, and use $\|\mathbf{M}\mathbf{x}\|_2 \leq \|\mathbf{M}\|_2\|\mathbf{x}\|_2$

$$\begin{aligned}\mathbf{x}^{(k+1)} - \mathbf{x}^* &= (\mathbf{I} - t\mathbf{A}^T\mathbf{A})(\mathbf{x}^{(k)} - \mathbf{x}^*) \\ &= (\mathbf{I} - t\mathbf{A}^T\mathbf{A})^{k+1}(\mathbf{x}^{(0)} - \mathbf{x}^*)\end{aligned}$$

$$\|\mathbf{x}^{(k)} - \mathbf{x}^*\|_2 \leq \|\mathbf{I} - t\mathbf{A}^T\mathbf{A}\|_2^k \|\mathbf{x}^{(0)} - \mathbf{x}^*\|_2$$

Example: least-squares problem (cont.)

Suppose eigenvalues of $\mathbf{A}^T \mathbf{A}$ belong to the interval $[\mu, L]$ where $L = \|\mathbf{A}\|_2^2$

Choose t such that it minimizes the spectral radius of $\mathbf{I} - t\mathbf{A}^T \mathbf{A}$

$$\begin{aligned} t^* &= \operatorname{argmin}_t \{\|\mathbf{I} - t\mathbf{A}^T \mathbf{A}\|_2\} \\ &= \operatorname{argmin}_t \left\{ \max_{\lambda \in [\mu, L]} |1 - t\lambda| \right\} \\ &= \operatorname{argmin}_t \left\{ \max\{1 - t\mu, 1 - tL, t\mu - 1, tL - 1\} \right\} \\ &= \frac{2}{L + \mu} \end{aligned}$$

Spectral radius of $\mathbf{I} - t^* \mathbf{A}^T \mathbf{A}$ is $(L - \mu)/(L + \mu)$

Exercises 13.4 and 13.5

13.4 Strong convexity. Suppose g is a twice continuously differentiable and strongly convex function with strong convexity parameter μ .

- 1 Show that the smallest eigenvalue of $\nabla^2 g(\mathbf{x})$ is bounded below by μ .
- 2 Consider the regularized least-squares objective function

$$g(\mathbf{x}) = \frac{1}{2} \|\mathbf{b} - \mathbf{Ax}\|_2^2 + \frac{\delta}{2} \|\mathbf{x}\|_2^2, \quad \delta > 0.$$

Derive the Lipschitz constant for ∇g and a lower bound on μ .

13.5 Poisson measurements. The negative log-likelihood function is

$$g(\mathbf{x}) = \mathbf{1}^T \exp(-\mathbf{Ax}) + \exp(-\mathbf{b})^T \mathbf{Ax} + \text{const.}$$

where $\mathbf{b} = -\log(I/I_0)$ and I is assumed to be positive.
(Refer to textbook for questions.)

Exercise 13.4: Strong convexity (solution)

Suppose g is a twice continuously differentiable and strongly convex function with strong convexity parameter μ .

- 1 Show that the smallest eigenvalue of $\nabla^2 g(\mathbf{x})$ is bounded below by μ .

$$\tilde{g}(\mathbf{x}) = g(\mathbf{x}) - \frac{\mu}{2} \|\mathbf{x}\|_2^2, \quad \nabla^2 \tilde{g}(\mathbf{x}) = \nabla^2 g(\mathbf{x}) - \mu \mathbf{I} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T - \mu \mathbf{I}$$

- 2 Consider the regularized least-squares objective function

$$g(\mathbf{x}) = \frac{1}{2} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2^2 + \frac{\delta}{2} \|\mathbf{x}\|_2^2, \quad \delta > 0.$$

Derive the Lipschitz constant for ∇g and a lower bound on μ .

$$\nabla^2 g(\mathbf{x}) = \mathbf{A}^T \mathbf{A} + \delta \mathbf{I} \implies L = \|\mathbf{A}\|_2^2 + \delta, \quad \mu \geq \delta$$

Exercise 13.5: Poisson measurement model (solution)

- 1 Show that $g(\mathbf{x})$ is a convex function of \mathbf{x} .

$$\nabla g(\mathbf{x}) = \mathbf{A}^T (\exp(-\mathbf{b}) - \exp(-\mathbf{Ax})), \quad \nabla^2 g(\mathbf{x}) = \mathbf{A}^T \text{diag}(\exp(-\mathbf{Ax})) \mathbf{A}$$

$$\mathbf{y}^T \mathbf{A}^T \text{diag}(\exp(-\mathbf{Ax})) \mathbf{A} \mathbf{y} = \|\text{diag}(\exp(-\mathbf{Ax}))^{1/2} \mathbf{A} \mathbf{y}\|_2^2 \geq 0, \quad \forall \mathbf{y}$$

- 2 Derive the first-order optimality condition associated with

$$\hat{\mathbf{x}}_{\text{ml}} = \underset{\mathbf{x}}{\text{argmin}} \{g(\mathbf{x})\}.$$

- 3 Show that the gradient of $g(\mathbf{x})$ is Lipschitz continuous on \mathbb{R}_+^n .

$$\|\mathbf{A}^T \text{diag}(\exp(-\mathbf{Ax})) \mathbf{A}\|_2 \leq \|\mathbf{A}\|_2^2 \max_i (\exp(-\mathbf{r}_i^T \mathbf{x}))$$

- 4 Show that if $\mathbf{Ax} = \mathbf{b}$ is consistent, then \mathbf{x} satisfies the first-order optimality condition $\nabla g(\mathbf{x}) = 0$ if $\mathbf{Ax} = \mathbf{b}$.

Power iteration for matrix norm estimation

$$\|\mathbf{H}\|_2 = \sup_{\mathbf{x} \neq 0} \left\{ \frac{\mathbf{x}^T \mathbf{H} \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \right\} = \lambda_{\max}(\mathbf{H}) \quad \mathbf{H} \text{ symmetric}$$

Power iteration

$$\mathbf{x}^{(k+1)} = \mathbf{H} \mathbf{x}^{(k)} / \|\mathbf{H} \mathbf{x}^{(k)}\|_2, \quad k = 0, 1, 2, \dots, \quad \text{with } \mathbf{x}^{(0)} \text{ random}$$

$$\hat{\lambda}^{(k)} = \|\mathbf{H} \mathbf{x}^{(k)}\|_2 \xrightarrow{\text{a.s.}} \lambda_{\max}(\mathbf{H}) \quad \text{as } k \rightarrow \infty$$

Why it works: let $\mathbf{H} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T$ and $\boldsymbol{\alpha} = \mathbf{V}^T \mathbf{x}^{(0)}$

$$\mathbf{x}^{(k)} = \mathbf{H}^k \mathbf{x}^{(0)} / \|\mathbf{H}^k \mathbf{x}^{(0)}\|_2, \quad k = 1, 2, \dots$$

$$\mathbf{H}^k \mathbf{x}^{(0)} = \mathbf{V} \mathbf{\Lambda}^k \mathbf{V}^T \mathbf{x}^{(0)} = \sum_{i=1}^n \alpha_i \lambda_i^k \mathbf{v}_i, \quad k = 1, 2, \dots$$

Power iteration for matrix norm estimation

$$\|\mathbf{H}\|_2 = \sup_{\mathbf{x} \neq 0} \left\{ \frac{\mathbf{x}^T \mathbf{H} \mathbf{x}}{\mathbf{x}^T \mathbf{x}} \right\} = \lambda_{\max}(\mathbf{H}) \quad \mathbf{H} \text{ symmetric}$$

Power iteration

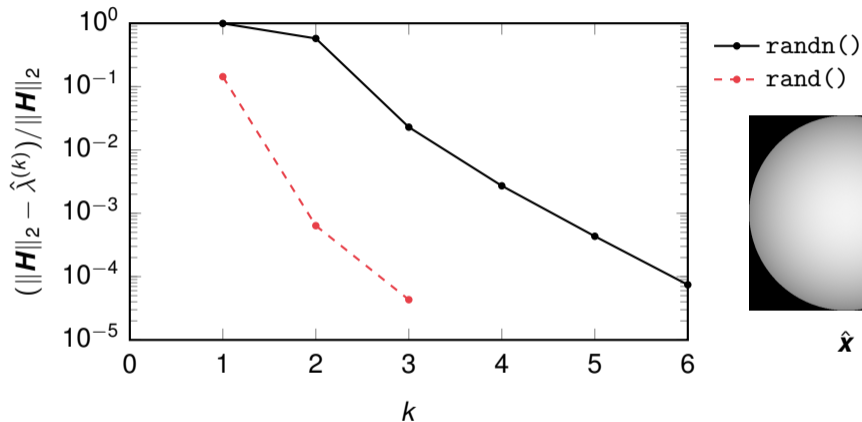
$$\mathbf{x}^{(k+1)} = \mathbf{H} \mathbf{x}^{(k)} / \|\mathbf{H} \mathbf{x}^{(k)}\|_2, \quad k = 0, 1, 2, \dots, \quad \text{with } \mathbf{x}^{(0)} \text{ random}$$

$$\hat{\lambda}^{(k)} = \|\mathbf{H} \mathbf{x}^{(k)}\|_2 \xrightarrow{\text{a.s.}} \lambda_{\max}(\mathbf{H}) \quad \text{as } k \rightarrow \infty$$

Why it works: let $\mathbf{H} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}^T$ and $\alpha = \mathbf{V}^T \mathbf{x}^{(0)}$

$$\mathbf{x}^{(k)} = \mathbf{H}^k \mathbf{x}^{(0)} / \|\mathbf{H}^k \mathbf{x}^{(0)}\|_2, \quad k = 1, 2, \dots$$

$$\lambda_1^{-k} \mathbf{H}^k \mathbf{x}^{(0)} = \alpha_1 \mathbf{v}_1 + \alpha_2 \left(\frac{\lambda_2}{\lambda_1} \right)^k \mathbf{v}_2 + \dots + \alpha_n \left(\frac{\lambda_n}{\lambda_1} \right)^k \mathbf{v}_n$$

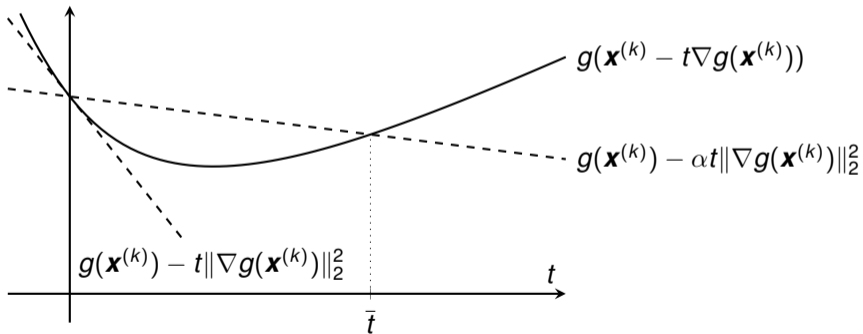
Example: $H = A^T A$ 

Remarks: (i) avoid forming H , (ii) similar to MATLAB's `normest()`

Backtracking line search

Armijo condition for gradient method

$$g(\mathbf{x}^{(k)} - t\nabla g(\mathbf{x}^{(k)})) \leq g(\mathbf{x}^{(k)}) - \alpha t \|\nabla g(\mathbf{x}^{(k)})\|_2^2$$



Backtracking line search (cont.)

Backtracking line search

Require: $\alpha \in (0, \frac{1}{2})$, $\beta \in (0, 1)$, and $t = t_0 > 0$
while $g(\mathbf{x}^{(k)} - t\nabla g(\mathbf{x}^{(k)})) > g(\mathbf{x}^{(k)}) - \alpha t \|\nabla g(\mathbf{x}^{(k)})\|_2^2$ **do**
 $t \leftarrow t\beta$
end while

- α controls a trade-off between max. step length and required decrease
- β controls backtracking “aggressiveness”
- typical values are $\alpha = 10^{-2}$ and $\beta = 0.7$

Barzilai–Borwein step size rules

Quadratic approximation

$$g(\mathbf{y}) \approx g(\mathbf{x}) + \nabla g(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{\alpha}{2} \|\mathbf{y} - \mathbf{x}\|_2^2$$

$$\nabla g(\mathbf{y}) - \nabla g(\mathbf{x}) \approx \alpha (\mathbf{y} - \mathbf{x})$$

Define $\Delta \mathbf{y} = \nabla g(\mathbf{x}^{(k)}) - \nabla g(\mathbf{x}^{(k-1)})$ and $\Delta \mathbf{s} = \mathbf{x}^{(k)} - \mathbf{x}^{(k-1)}$ ($k \geq 1$)

$$t_k^{\text{BB1}} = \alpha_k^{-1}, \quad \alpha_k = \operatorname{argmin}_{\alpha} \left\{ \|\Delta \mathbf{y} - \alpha \Delta \mathbf{s}\|_2^2 \right\} = \frac{\Delta \mathbf{s}^T \Delta \mathbf{y}}{\|\Delta \mathbf{s}\|_2^2}$$

$$t_k^{\text{BB2}} = \operatorname{argmin}_{\beta} \left\{ \|\beta \Delta \mathbf{y} - \Delta \mathbf{s}\|_2^2 \right\} = \frac{\Delta \mathbf{s}^T \Delta \mathbf{y}}{\|\Delta \mathbf{y}\|_2^2}$$

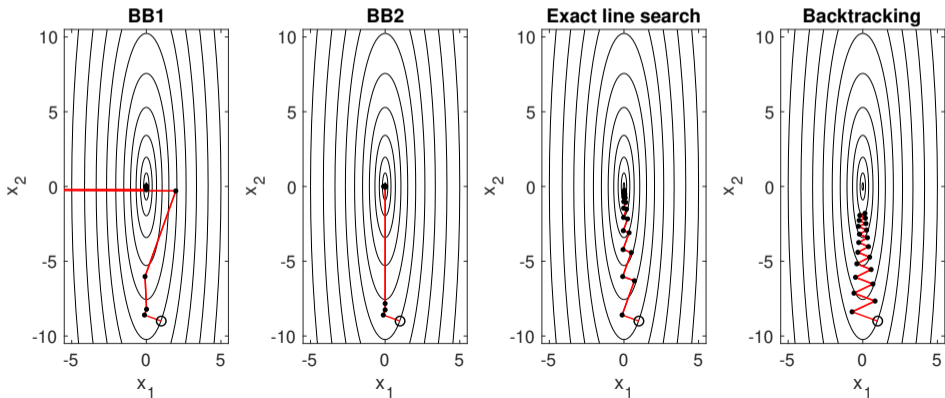
Barzilai–Borwein step size rules (cont.)

- first step size ($k = 0$) must be chosen using another method
- not a descent method ($g(\mathbf{x}^{(k+1)}) \leq g(\mathbf{x}^{(k)})$) not guaranteed
- convergence guaranteed if g is strongly convex and quadratic
- safe-guarding is generally required to ensure convergence

Example: least-squares problem

$$t_k^{\text{BB1}} = \frac{\|\nabla g(\mathbf{x}^{(k-1)})\|_2^2}{\|\mathbf{A}\nabla g(\mathbf{x}^{(k-1)})\|_2^2},$$

$$t_k^{\text{BB2}} = \frac{\|\mathbf{A}\nabla g(\mathbf{x}^{(k-1)})\|_2^2}{\|\mathbf{A}^T\mathbf{A}\nabla g(\mathbf{x}^{(k-1)})\|_2^2}$$



Stopping criteria

Approximate stationarity conditions

$$\|\nabla g(\mathbf{x}^{(k)})\|_2 \leq \epsilon, \quad \|\underbrace{\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)}}_{-t_k \nabla g(\mathbf{x}^{(k)})}\|_2 \leq \epsilon \|\mathbf{x}^{(k)}\|_2$$

not scale invariant; change of variables $\tilde{g}(\mathbf{y}) = g(\mathbf{C}\mathbf{y})$ yields

$$\|\nabla \tilde{g}(\mathbf{y}^{(k)})\|_2 = \|\mathbf{C}^T \nabla g(\mathbf{x}^{(k)})\|_2 \leq \epsilon, \quad \nabla \tilde{g}(\mathbf{y}) = \mathbf{C}^T \nabla g(\mathbf{C}\mathbf{y})$$

Strongly convex objective

$$\|\nabla g(\mathbf{x}^{(k)})\|_2 \leq \sqrt{2\mu\epsilon_{\text{obj}}}, \quad \|\nabla g(\mathbf{x}^{(k)})\|_2 \leq \frac{\mu\epsilon_{\text{dist}}}{2},$$

imply that $g(\mathbf{x}^{(k)}) - g(\mathbf{x}^*) \leq \epsilon_{\text{obj}}$ and $\|\mathbf{x}^{(k)} - \mathbf{x}^*\|_2 \leq \epsilon_{\text{dist}}$

Example: Tikhonov regularized least-squares

$$\text{minimize } g(\mathbf{x}) \equiv \frac{1}{2} \|\mathbf{b} - \mathbf{Ax}\|_2^2 + \frac{\alpha}{2} \|\mathbf{x}\|_2^2$$

$\alpha > 0$ is a lower bound on strong convexity parameter (Exercise 13.4)

Stopping criteria

$$\|\nabla g(\mathbf{x}^{(k)})\|_2 = \|\alpha \mathbf{x}^{(k)} - \mathbf{A}^T \boldsymbol{\varrho}^{(k)}\|_2 \leq \frac{\alpha \epsilon_{\text{dist}}}{2}, \quad \boldsymbol{\varrho}^{(k)} = \mathbf{b} - \mathbf{Ax}^{(k)}$$

ensures that $\|\mathbf{x}^{(k)} - \mathbf{x}^*\|_2 \leq \epsilon_{\text{dist}}$

Exercise 13.6: Step sizes

Apply the gradient method to the problem of minimizing

$$g(\mathbf{x}) = \frac{1}{2} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2^2$$

where \mathbf{A} and \mathbf{b} are generated as follows:

```
>> I0 = 1e6; n = 128;  
>> A = paralleltomo(n)*(2/n);  
>> x = reshape(phantomgallery('grains',n), [], 1);  
>> I = poissrnd(I0*exp(-A*x));  
>> b = -log(I/I0);
```

Plot (semi-log. y-axis) the obj. value for the first 200 iterations using:

- 1 Exact line search
- 2 Backtracking line search
- 3 BB1 step size
- 4 BB2 step size

Constrained optimization

$$\text{minimize } g(\mathbf{x}) + h(\mathbf{x})$$

- $g : \mathbb{R}^n \rightarrow \mathbb{R}$ convex and differentiable
- $h : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{-\infty, \infty\}$ closed convex
- h not necessarily continuously differentiable but “simple”

Special case

$$\begin{aligned} &\text{minimize } g(\mathbf{x}) \\ &\text{subject to } \mathbf{x} \in \mathcal{C}, \end{aligned}$$

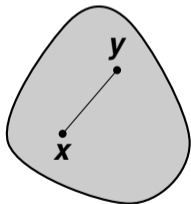
$$\text{corresponds to } h(\mathbf{x}) = I_{\mathcal{C}}(\mathbf{x}) = \begin{cases} 0, & \mathbf{x} \in \mathcal{C} \\ \infty, & \mathbf{x} \notin \mathcal{C} \end{cases}$$

Convex sets

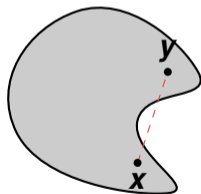
$\mathcal{C} \subseteq \mathbb{R}^n$ is a convex set if and only if

$$\theta \mathbf{x} + (1 - \theta) \mathbf{y} \in \mathcal{C}, \quad \theta \in [0, 1], \quad \text{for all } \mathbf{x}, \mathbf{y} \in \mathcal{C}$$

convex set



nonconvex set



Majorization minimization

Suppose ∇g is Lipschitz continuous with constant L

- majorization of $g + h$ at \mathbf{x}

$$\psi(\mathbf{y}; \mathbf{x}) = g(\mathbf{x}) + \nabla g(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{L}{2} \|\mathbf{y} - \mathbf{x}\|_2^2 + h(\mathbf{y})$$

- majorization minimization

$$\begin{aligned} \mathbf{x}^{(k+1)} &= \operatorname{argmin}_{\mathbf{y}} \left\{ \nabla g(\mathbf{x}^{(k)})^T \mathbf{y} + h(\mathbf{y}) + \frac{L}{2} \|\mathbf{y} - \mathbf{x}^{(k)}\|_2^2 \right\} \\ &= \operatorname{argmin}_{\mathbf{y}} \left\{ h(\mathbf{y}) + \frac{L}{2} \|\mathbf{y} - (\mathbf{x}^{(k)} - (1/L)\nabla g(\mathbf{x}^{(k)}))\|_2^2 \right\} \end{aligned}$$

Proximal gradient method

Proximal operator associated with h and $t > 0$

$$\text{prox}_{th}(\mathbf{x}) = \underset{\mathbf{y}}{\text{argmin}} \left\{ h(\mathbf{y}) + \frac{1}{2t} \|\mathbf{y} - \mathbf{x}\|_2^2 \right\}$$

- easy to evaluate if h is “simple”
- strong convexity implies that $\text{prox}_{th}(\mathbf{x})$ is unique

Proximal gradient method

$$\mathbf{x}^{(k+1)} = \text{prox}_{th}(\mathbf{x}^{(k)} - t\nabla g(\mathbf{x}^{(k)})), \quad t = \frac{1}{L}, \quad k = 0, 1, 2, \dots$$

Examples

- $h(\mathbf{x}) = I_{\mathcal{C}}(\mathbf{x})$ where \mathcal{C} is a closed, convex set

$$\text{prox}_{th}(\mathbf{x}) = P_{\mathcal{C}}(\mathbf{x}) = \underset{\mathbf{y} \in \mathcal{C}}{\text{argmin}} \{ \|\mathbf{y} - \mathbf{x}\|_2^2 \}$$

- $h(\mathbf{x}) = I_{\mathcal{C}}(\mathbf{x})$ with $\mathcal{C} = \{ \mathbf{x} \mid l_i \leq x_i \leq u_i, i = 1, \dots, n \}$

$$\text{prox}_{th}(\mathbf{x}) = \max(\mathbf{l}, \min(\mathbf{u}, \mathbf{x}))$$

- $h(\mathbf{x}) = \|\mathbf{x}\|_1$

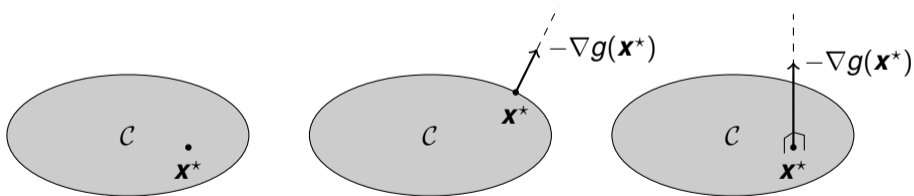
$$\text{prox}_{th}(\mathbf{x}) = \text{diag}(\text{sgn}(\mathbf{x})) \max(\text{abs}(\mathbf{x}) - t\mathbf{1}, \mathbf{0})$$

Optimality condition

\mathbf{x}^* is a minimizer of $g + h$ if and only if

$$\mathbf{x}^* = \text{prox}_{th}(\mathbf{x}^* - t\nabla g(\mathbf{x}^*)), \quad t > 0$$

Special case: $h(\mathbf{x}) = I_{\mathcal{C}}(\mathbf{x})$ where \mathcal{C} is closed, convex



Example: nonnegativity constraints

$$\begin{array}{ll} \text{minimize} & g(\mathbf{x}) \\ \text{subject to} & \mathbf{x} \in \mathcal{C} \end{array}$$

with $\mathcal{C} = \{\mathbf{x} \mid x_i \geq 0, i = 1, \dots, n\}$

Optimality condition

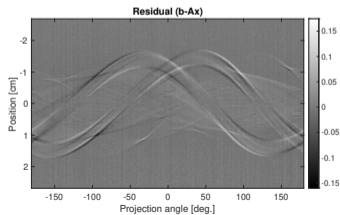
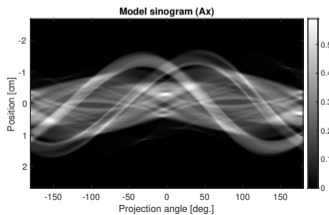
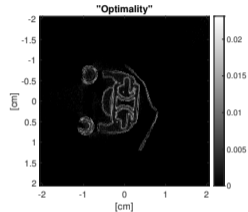
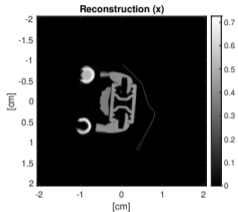
$$\mathbf{x}^* = \max(\mathbf{0}, \mathbf{x}^* - t \nabla g(\mathbf{x}^*)), \quad t > 0$$

or equivalently, for $i = 1, \dots, n$

$$(x_i^* = 0 \wedge [\nabla g(\mathbf{x}^*)]_i \geq 0) \quad \vee \quad (x_i^* > 0 \wedge [\nabla g(\mathbf{x}^*)]_i = 0)$$

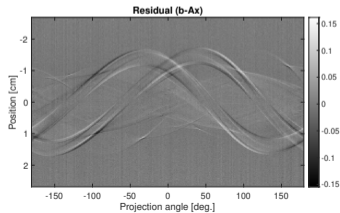
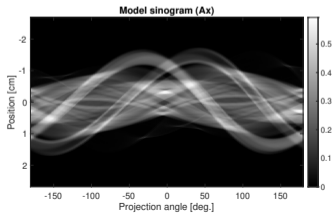
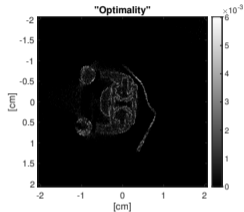
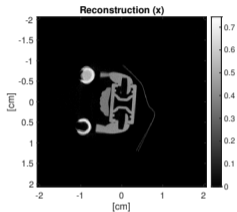
Example: reconstruction with nonnegativity constraints

PG: 200 iterations

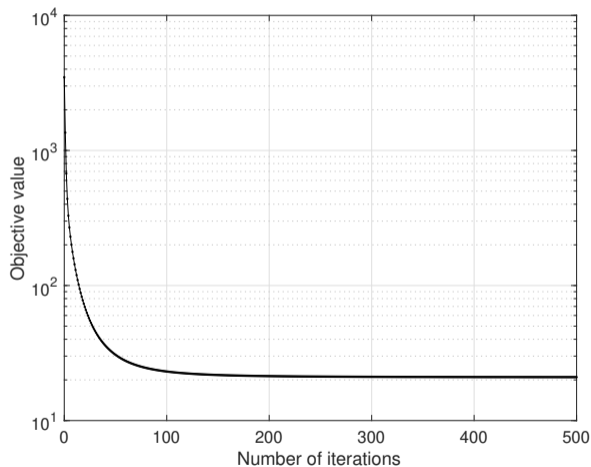


Example: reconstruction with nonnegativity constraints

PG: 500 iterations



Example: reconstruction with nonnegativity constraints



Accelerated proximal gradient method

Accelerated proximal gradient method

Require: initial vector $\mathbf{x}^{(0)}$, $\mathbf{y} = \mathbf{x}^{(0)}$, $t_0 = 1$

for $k = 0, 1, 2, \dots$ **do**

$$\mathbf{x}^{(k+1)} = \text{prox}_{(1/L)h}(\mathbf{y} - (1/L)\nabla g(\mathbf{y}))$$

$$t_{k+1} = \frac{1 + \sqrt{1 + 4t_k^2}}{2}$$

$$\mathbf{y} = \mathbf{x}^{(k+1)} + \frac{t_k - 1}{t_{k+1}}(\mathbf{x}^{(k+1)} - \mathbf{x}^{(k)})$$

end for

- improved worst-case bound: $f(\mathbf{x}^{(k)}) - f(\mathbf{x}^*) = O(1/k^2)$ where $f = g + h$
- not a descent method

Example

$$\text{minimize } g(\mathbf{x}) + h(\mathbf{x})$$

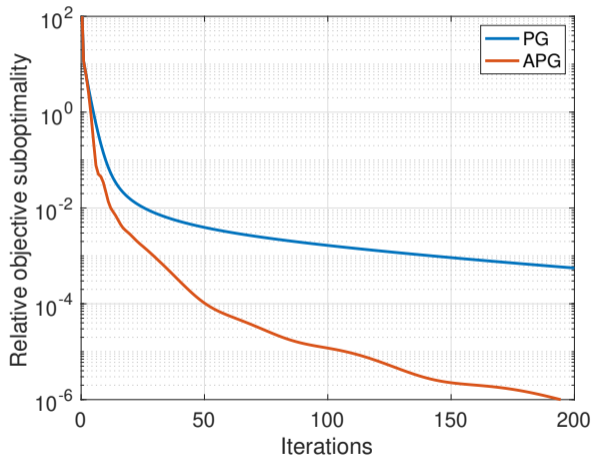
where

$$g(\mathbf{x}) = \frac{1}{2} \|\mathbf{b} - \mathbf{A}\mathbf{x}\|_2^2$$

$$h(\mathbf{x}) = \frac{\gamma}{2} \|\mathbf{x}\|_2 + l_c(\mathbf{x})$$

$$\mathcal{C} = \{\mathbf{x} \mid x_i \geq 0, i = 1, \dots, n\}$$

(several ways to “split” objective)



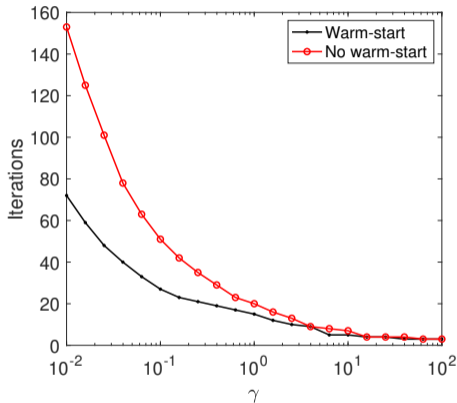
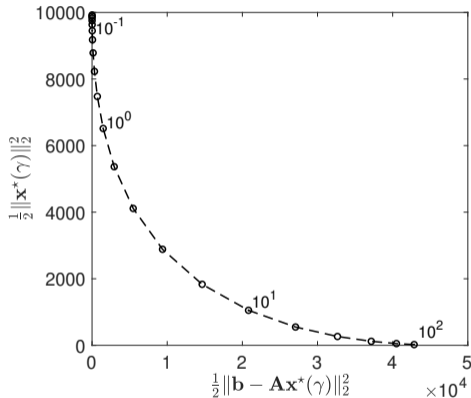
Regularized least-squares

$$\begin{aligned} & \text{minimize} && \frac{1}{2} \|\mathbf{b} - \mathbf{Ax}\|_2^2 + \gamma R(\mathbf{x}) \\ & \text{subject to} && \mathbf{x} \in \mathcal{C} \end{aligned}$$

- special cases: Tikhonov, generalized Tikhonov, and TV regularization
- trade-off between two objectives
- often of interest to solve problem with different values of γ
- trade-off curve (aka L-curve), parameterized by γ

$$\left(\frac{1}{2} \|\mathbf{b} - \mathbf{Ax}^*(\gamma)\|_2^2, R(\mathbf{x}^*(\gamma)) \right)$$

Tracing the trade-off curve: Tikhonov regularization



Warm-start: use $\mathbf{x}^*(\gamma)$ as initial guess when solving for $\mathbf{x}^*(\gamma')$

Total variation regularized least-squares

$$\text{TV}_a(\mathbf{x}) = \|\mathbf{D}\mathbf{x}\|_1 = \sum_{i=1}^{2n} |\mathbf{d}_i^T \mathbf{x}|, \quad \mathbf{D} = \begin{bmatrix} \mathbf{I}_N \otimes \mathbf{D}_M \\ \mathbf{D}_N \otimes \mathbf{I}_M \end{bmatrix}$$

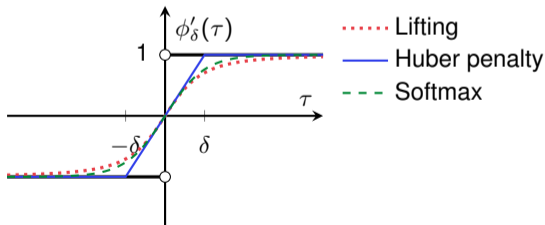
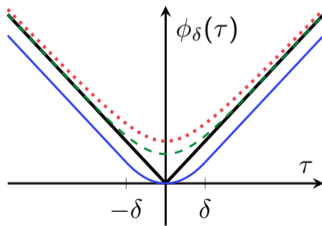
- $\text{TV}_a(\mathbf{x})$ is convex but not everywhere differentiable
- $\text{TV}_a(\mathbf{x})$ is not “simple” (proximal operator is not cheap to eval.)
- smooth approximation

$$\text{TV}_a^\delta(\mathbf{x}) = \sum_{i=1}^{2n} \phi_\delta(\mathbf{d}_i^T \mathbf{x}), \quad \nabla \text{TV}_a^\delta(\mathbf{x}) = \sum_{i=1}^{2n} \mathbf{d}_i \phi'_\delta(\mathbf{d}_i^T \mathbf{x})$$

- more advanced methods exist (splitting methods, etc.)

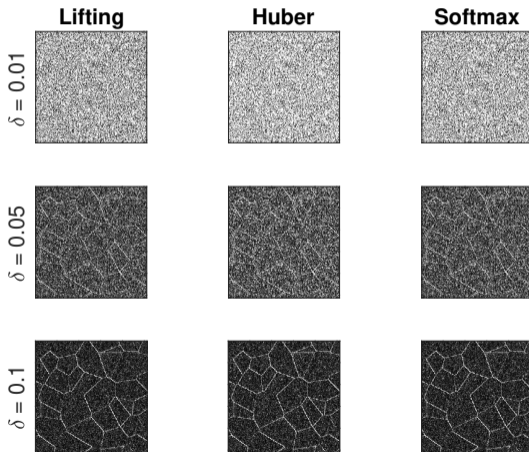
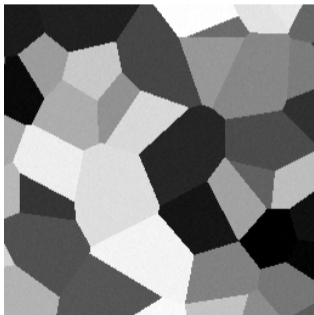
Smooth approximation to absolute value function

- Lifting: $\phi_\delta(\tau) = \left\| \begin{bmatrix} \tau \\ \delta \end{bmatrix} \right\|_2 = \sqrt{\tau^2 + \delta^2}$
- Huber penalty (scaled): $\phi_\delta(\tau) = \begin{cases} \frac{\tau^2}{2\delta}, & |\tau| \leq \delta \\ |\tau| - \frac{\delta}{2}, & |\tau| > \delta \end{cases}$
- Softmax: $\phi_\delta(\tau) = \delta \log(e^{\tau/\delta} + e^{-\tau/\delta})$



Example: gradient of $TV_a^\delta(\mathbf{x})$

```
>> N = 256;  
>> X = phantomgallery('grains',N) ...  
      + 1e-2*randn(N,N);
```



Extension to isotropic TV

$$\text{TV}_i^\delta(\mathbf{x}) = \sum_{i=1}^n \phi_\delta(\mathbf{D}_i \mathbf{x}), \quad \mathbf{D}_i = \begin{bmatrix} \mathbf{i}_i^T (\mathbf{I}_N \otimes \mathbf{D}_M) \\ \mathbf{i}_i^T (\mathbf{D}_N \otimes \mathbf{I}_M) \end{bmatrix}$$

$\phi_\delta: \mathbb{R}^2 \rightarrow \mathbb{R}$ approximates 2-norm of vector in \mathbb{R}^2

Approximation	$\phi_\delta(\mathbf{y})$	$\nabla \phi_\delta(\mathbf{y})$
Lifting	$\left\ \begin{bmatrix} \mathbf{y} \\ \delta \end{bmatrix} \right\ _2$	$\left\ \begin{bmatrix} \mathbf{y} \\ \delta \end{bmatrix} \right\ _2^{-1} \mathbf{y}$
Huber	$\begin{cases} \frac{\mathbf{y}^T \mathbf{y}}{2\delta}, & \ \mathbf{y}\ _2 \leq \delta \\ \ \mathbf{y}\ _2 - \frac{\delta}{2}, & \ \mathbf{y}\ _2 > \delta \end{cases}$	$\frac{1}{\max(\delta, \ \mathbf{y}\ _2)} \mathbf{y}$
Softmax	$\delta \log(e^{\ \mathbf{y}\ _2/\delta} + e^{-\ \mathbf{y}\ _2/\delta})$	$\begin{cases} \frac{\tanh(\ \mathbf{y}\ _2/\delta)}{\ \mathbf{y}\ _2} \mathbf{y}, & \mathbf{y} \neq \mathbf{0} \\ \mathbf{0}, & \mathbf{y} = \mathbf{0} \end{cases}$

Exercise 13.7: Smooth approximation of total variation penalty

Show that the three smooth approximations of the absolute value function all have a Lipschitz continuous derivative with Lipschitz constant $L = 1/\delta$.

- Lifting: $\phi_\delta(\tau) = \left\| \begin{bmatrix} \tau \\ \delta \end{bmatrix} \right\|_2 = \sqrt{\tau^2 + \delta^2}$
- Huber penalty (scaled): $\phi_\delta(\tau) = \begin{cases} \frac{\tau^2}{2\delta}, & |\tau| \leq \delta \\ |\tau| - \frac{\delta}{2}, & |\tau| > \delta \end{cases}$
- Softmax: $\phi_\delta(\tau) = \delta \log(e^{\tau/\delta} + e^{-\tau/\delta})$

Exercise 13.7: Smooth approximation of total variation penalty

- Lifting: $\phi_\delta(\tau) = \left\| \begin{bmatrix} \tau \\ \delta \end{bmatrix} \right\|_2 = \sqrt{\tau^2 + \delta^2}$

$$\phi'_\delta(\tau) = \frac{\tau}{\phi_\delta(\tau)}, \quad \phi''_\delta(\tau) = \frac{\phi_\delta(\tau) - \tau\phi'_\delta(\tau)}{\phi_\delta(\tau)^2} = \frac{\delta^2}{\phi_\delta(\tau)^3}$$

- Huber penalty (scaled): $\phi_\delta(\tau) = \begin{cases} \frac{\tau^2}{2\delta}, & |\tau| \leq \delta \\ |\tau| - \frac{\delta}{2}, & |\tau| > \delta \end{cases}$

$$\phi'_\delta(\tau) = \begin{cases} \frac{\tau}{\delta}, & |\tau| \leq \delta \\ \text{sgn}(\tau), & |\tau| > \delta \end{cases}, \quad \phi''_\delta(\tau) = \begin{cases} 1/\delta, & |\tau| < \delta \\ 0, & |\tau| > \delta \end{cases}$$

- Softmax: $\phi_\delta(\tau) = \delta \log(e^{\tau/\delta} + e^{-\tau/\delta})$

$$\phi'_\delta(\tau) = \tanh(\tau/\delta), \quad \phi''_\delta(\tau) = \frac{1}{\delta \cosh^2(\tau/\delta)}$$

Exercise 13.8: Regularized weighted least-squares problems

Consider the following weighted least-squares problems with two different regularization terms: (i) generalized Tikhonov regularization

$$\mathbf{x}_{\text{GTik}} = \operatorname{argmin}_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{b} - \mathbf{Ax}\|_{\mathbf{W}}^2 + \frac{\gamma}{2} \|\mathbf{Dx}\|_2^2 \right\} \quad (1)$$

and (ii) total variation regularization

$$\mathbf{x}_{\text{TV}} = \operatorname{argmin}_{\mathbf{x}} \left\{ \frac{1}{2} \|\mathbf{b} - \mathbf{Ax}\|_{\mathbf{W}}^2 + \gamma \|\mathbf{Dx}\|_1 \right\}. \quad (2)$$

The variable $\mathbf{x} \in \mathbb{R}^n$ represents an image of size $N \times N$ (i.e., $n = N^2$). (Refer to textbook for questions.)

Exercise 13.8 (solutions)

Use power iteration to estimate a Lipschitz constant for the gradient of $g(\mathbf{x})$ in the generalized Tikhonov problem. Plot the estimated Lipschitz constant for different values of γ .

Apply the power iteration method to the Hessian of $g(\mathbf{x})$ which is given by

$$\nabla^2 g(\mathbf{x}) = \mathbf{A}^T \mathbf{W} \mathbf{A} + \gamma \mathbf{D}^T \mathbf{D}.$$

Avoid forming the Hessian, e.g., by evaluating the matrix-vector product as $\mathbf{H}(\mathbf{x})$ where

$$\mathbf{H} = \text{@}(\mathbf{x}) \mathbf{A}' * (\mathbf{w} * (\mathbf{A} * \mathbf{x})) + \text{gamma} * (\mathbf{D}' * (\mathbf{D} * \mathbf{x}));$$

The vector \mathbf{w} contains the diagonal elements of \mathbf{W} .