

An Introduction to Non-smooth Optimization

Lecture 02 - Intro of the Intro

Jingwei LIANG

Institute of Natural Sciences, Shanghai Jiao Tong University

Email: optimization.sjtu@gmail.com

Office: Room 355, No. 6 Science Building

Outline

- ① Motivating example
- ② Regularization
- ③ Applications
- ④ First-order optimization methods



Least square regression



Let $m \in \mathbb{N}_{++}$. For $i = 1, \dots, m$, given each $x_i \in \mathbb{R}$,

$$y_i = ax_i + b + \epsilon_i$$

with ϵ_i being random noise.



Least square regression



Let $m \in \mathbb{N}_{++}$. For $i = 1, \dots, m$, given each $x_i \in \mathbb{R}$,

$$y_i = ax_i + b + \epsilon_i$$

with ϵ_i being random noise.

Matrix-vector representation,

$$\mathbf{A} = \begin{bmatrix} x_1 & 1 \\ x_2 & 1 \\ \vdots & \\ x_m & 1 \end{bmatrix}, \quad \mathbf{x} = \begin{bmatrix} a \\ b \end{bmatrix}, \quad \mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{bmatrix} \quad \text{and} \quad \boldsymbol{\epsilon} = \begin{bmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_m \end{bmatrix}.$$

The system of equations reads

$$\begin{cases} ax_1 + b + \epsilon_1 = y_1, \\ \vdots \\ ax_m + b + \epsilon_m = y_m. \end{cases} \iff \mathbf{y} = \mathbf{Ax} + \boldsymbol{\epsilon}.$$

Least square regression



Least square regression: estimating \mathbf{x} from \mathbf{y}

$$\min_{\mathbf{x} \in \mathbb{R}^2} \|\mathbf{Ax} - \mathbf{y}\|^2.$$

Assume that \mathbf{A} has full column rank

$$\begin{aligned} \min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{y}\|^2 &\iff \mathbf{0} = \mathbf{A}^T(\mathbf{Ax} - \mathbf{y}) \\ &\iff \mathbf{A}^T \mathbf{Ax} = \mathbf{A}^T \mathbf{y} \\ &\iff \mathbf{x} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{y} \end{aligned}$$

If $\mathbf{A}^T \mathbf{A}$ is not invertible,

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \gamma_k \mathbf{A}^T(\mathbf{Ax}^{(k)} - \mathbf{y}) \rightarrow \mathbf{x}^*.$$

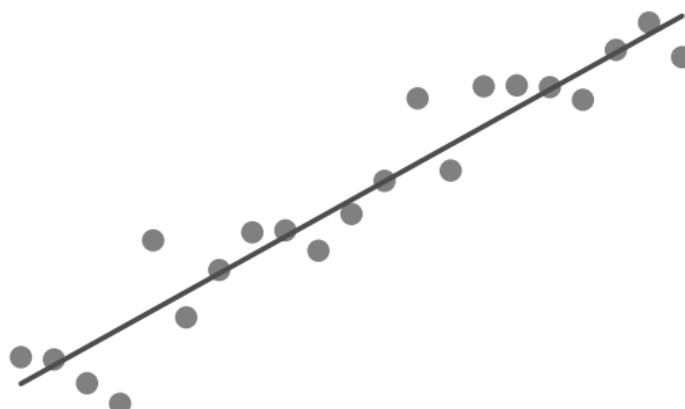
Least square regression



Least square regression: estimating \mathbf{x} from \mathbf{y}

$$\min_{\mathbf{x} \in \mathbb{R}^2} \|\mathbf{Ax} - \mathbf{y}\|^2.$$

— $y = 0.23x - 0.08$

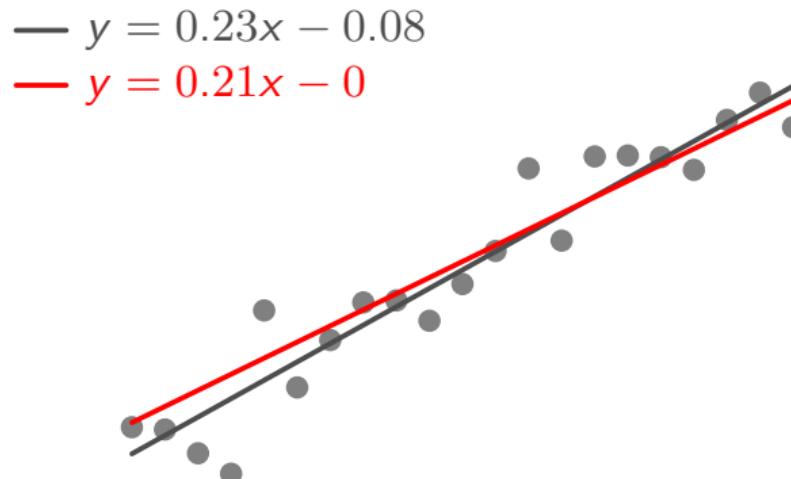


Non-negative least square



Non-negative least square regression:

$$\min_{\mathbf{x} \in \mathbb{R}^2} \|\mathbf{Ax} - \mathbf{y}\|^2 \quad \text{such that} \quad x_i \geq 0, i = 1, 2.$$

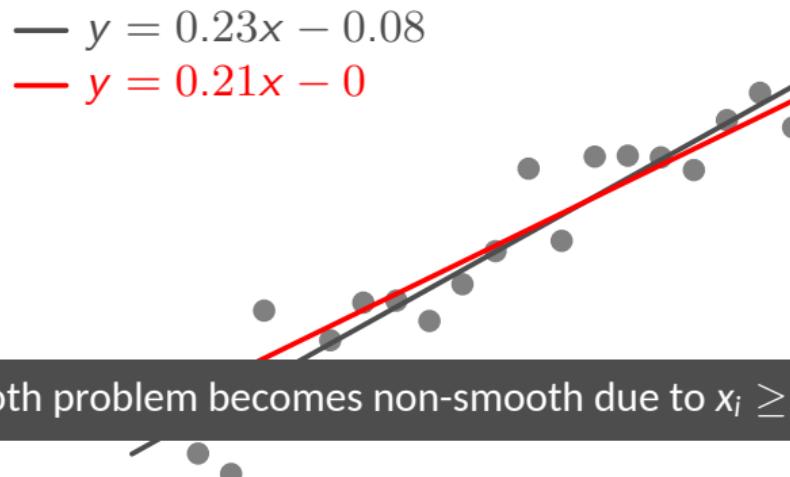


Non-negative least square



Non-negative least square regression:

$$\min_{\mathbf{x} \in \mathbb{R}^2} \|\mathbf{Ax} - \mathbf{y}\|^2 \quad \text{such that} \quad x_i \geq 0, i = 1, 2.$$



Smooth problem becomes non-smooth due to $x_i \geq 0, i = 1, 2$.

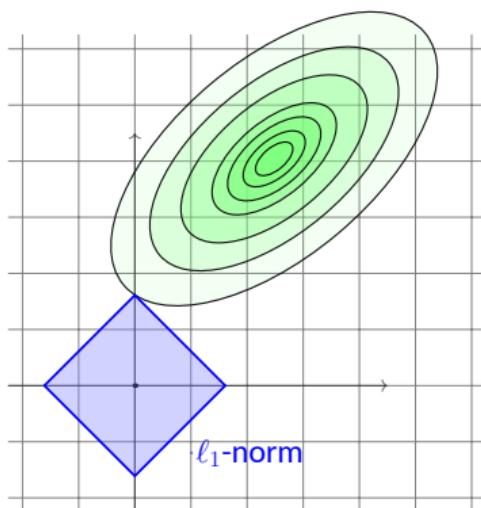
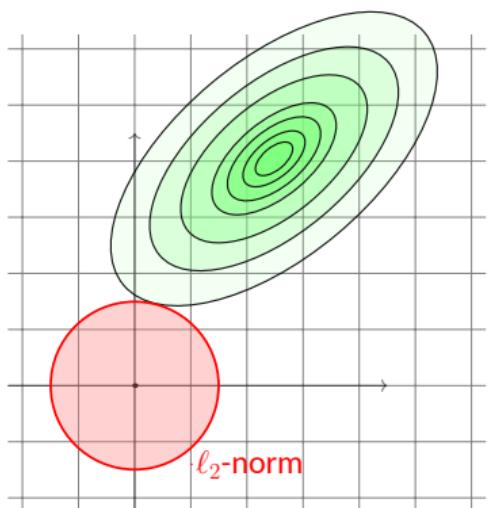
Other constraints on x



Besides non-negativity, we can consider the following requirement on the solution: let $p \in \{1, 2\}$ and $\delta > 0$,

$$\min_{\mathbf{x}} \|\mathbf{Ax} - \mathbf{y}\|^2$$

such that $\|\mathbf{x}\|_p \leq \delta$.



Regularization

Sparsity, low-rank, non-negativity



Let $R(\mathbf{x})$ be a function promoting prior information, e.g. non-negativity or norm constraint...

Regularized least square

$$\min_{\mathbf{x} \in \mathbb{R}^n} R(\mathbf{x}) + \frac{1}{2\mu} \|\mathbf{Ax} - \mathbf{y}\|^2.$$

- $\mu > 0$ provides a balance between diffusion and fidelity.
- The choices of $R(\mathbf{x})$ depends on the prior information.

Let $R(\mathbf{x})$ be a function promoting prior information, e.g. non-negativity or norm constraint...

Regularized least square

$$\min_{\mathbf{x} \in \mathbb{R}^n} R(\mathbf{x}) + \frac{1}{2\mu} \|\mathbf{Ax} - \mathbf{y}\|^2.$$

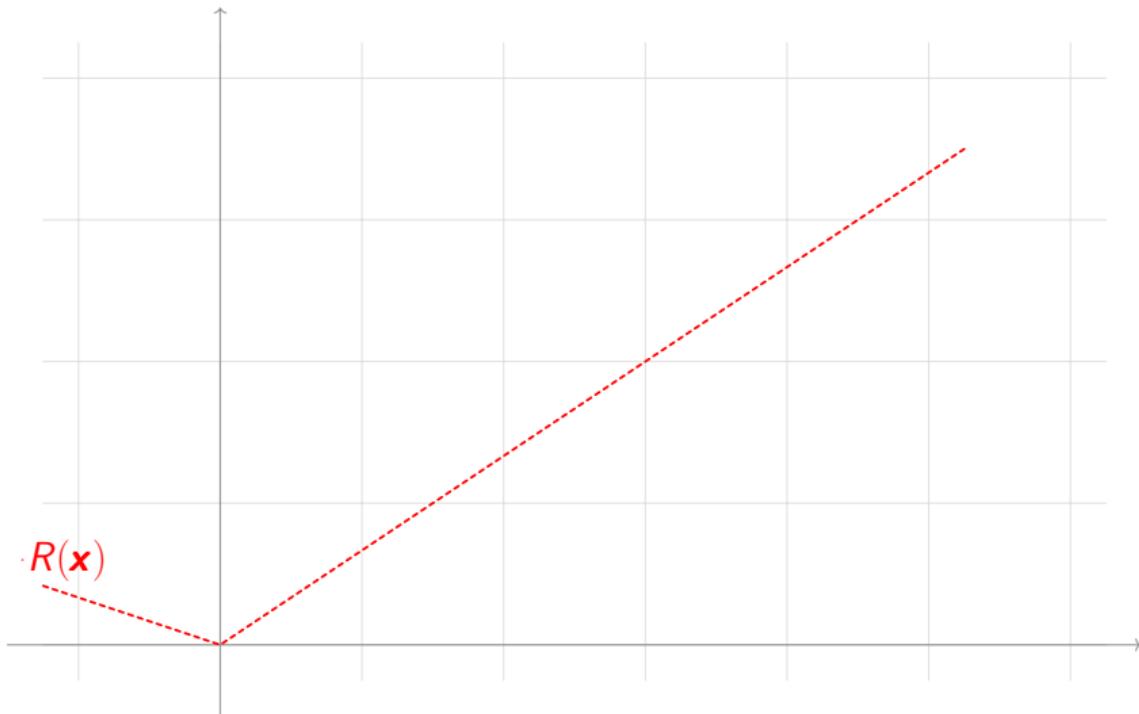
- $\mu > 0$ provides a balance between diffusion and fidelity.
- The choices of $R(\mathbf{x})$ depends on the prior information.

Regularization

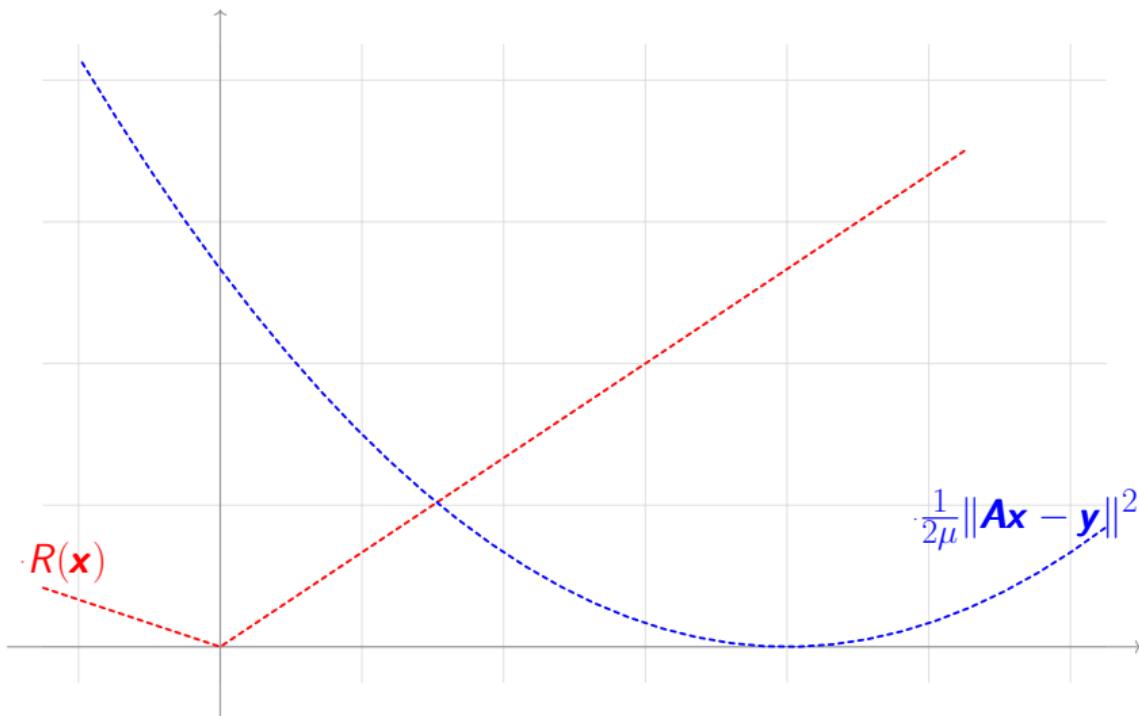
Function $R(\mathbf{x})$ is called regularization, it is a process that forces the solution to be “simpler”,

- obtain results for ill-posed problems (e.g. image processing).
- prevent overfitting (e.g. machine learning).

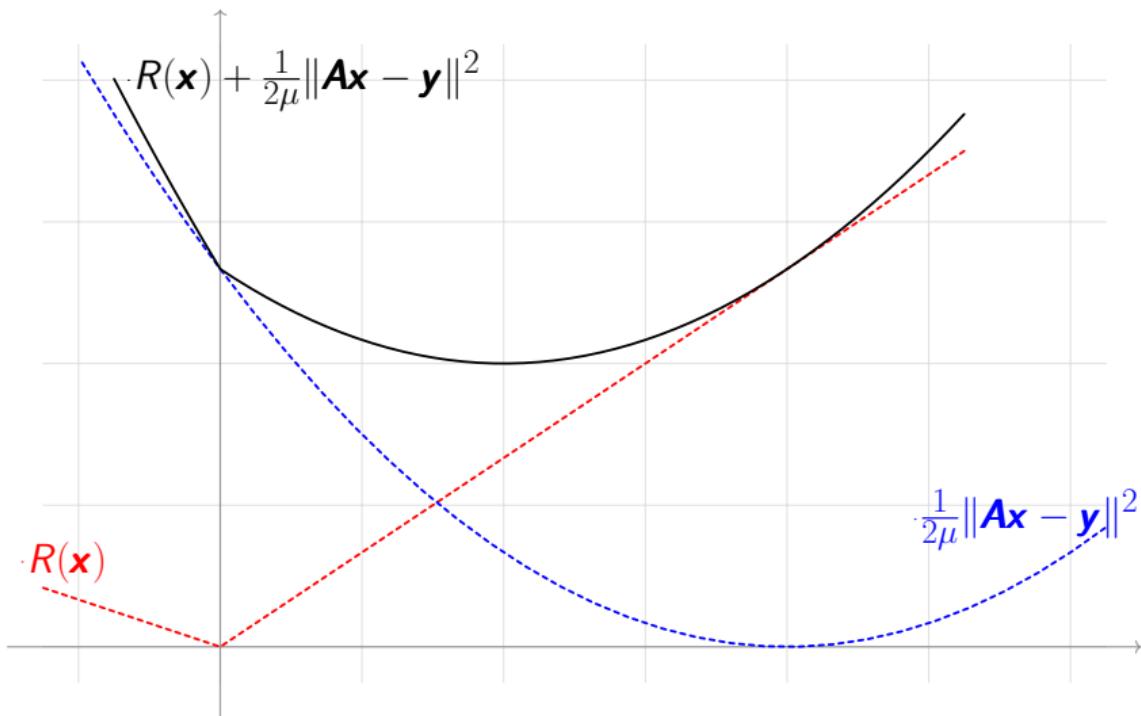
Regularized least square



Regularized least square



Regularized least square





A more general example

$$\min_{\mathbf{x} \in \mathbb{R}^n} R(\mathbf{x}) + F(\mathbf{x}).$$

Choices of $F(\mathbf{x})$: let $\mathbf{a} \in \mathbb{R}^n$ and $b \in \mathbb{R}$

- **Quadratic loss**

$$F(\mathbf{x}) = \frac{1}{2}(\mathbf{a}^T \mathbf{x} - b)^2.$$

- **Logistic loss**

$$F(\mathbf{x}) = \log(1 + e^{-b\mathbf{a}^T \mathbf{x}}).$$

- **Squared hinge loss**

$$F(\mathbf{x}) = \max \{1 - t(\mathbf{a}^T \mathbf{x} + b), 0\}, \quad t \in \{-1, 1\}.$$

Choices of $R(\mathbf{x}) \longrightarrow$

Examples: norms



Let $\mathbf{x} \in \mathbb{R}^n$

- **ℓ_1 -norm**

$$\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|.$$

- **ℓ_2 -norm**

$$\|\mathbf{x}\|_2 = \left(\sum_{i=1}^n (x_i)^2 \right)^{1/2}.$$

- **group ℓ_1 -norm** let $\mathcal{G} = \{g_1, g_2, \dots, g_\ell\}$ be a partition of $\{1, 2, \dots, n\}$,

$$\|\mathbf{x}\|_{1,2} = \sum_{g_j \in \mathcal{G}} \left(\sum_{i \in g_j} (x_i)^2 \right)^{1/2}.$$

Examples: total variation



Example - Total variation (TV) [Rudin, Osher & Fatemi '92]

Let ∇ be the discrete gradient operator

$$\|\nabla \mathbf{x}\|_1$$

In 1D case:

$$\nabla = \begin{bmatrix} -1 & 1 & & & \\ & -1 & 1 & & \\ & & \ddots & \ddots & \\ & & & -1 & 1 \\ & & & & 0 \end{bmatrix}.$$

Examples: total variation



Example - Total variation (TV) [Rudin, Osher & Fatemi '92]

Let ∇ be the discrete gradient operator

$$\|\nabla x\|_1$$



Original image



Horizontal gradient



Vertical gradient

Examples: wavelet frames



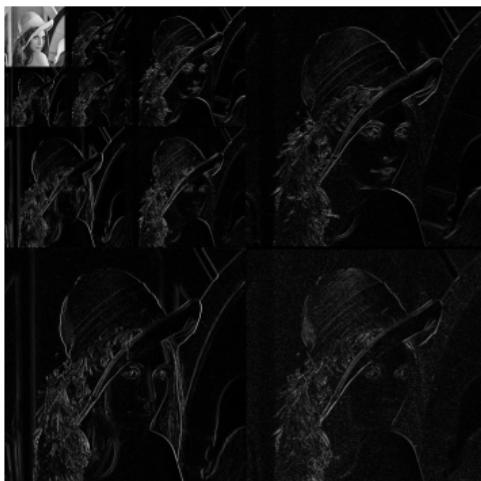
Example - Wavelet [Morlet, Meyer, Mallat, Daubechies, et al]

Family of filters $\{\psi_{j,k} : j, k \in \mathbb{Z}\}$

$$\psi_{j,k}(\cdot) = 2^{j/2} \psi(2^j \cdot -k).$$



Original image



Wavelet coefficients

Examples: low rank



Example - Nuclear norm [Recht, Fazel and Parrilo '10]

Let $\mathbf{A} \in \mathbb{R}^{m \times n}$, and $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^T$ be its singular value decomposition

$$\|\mathbf{A}\|_* = \sum_{i=1}^{\min\{m,n\}} s_{i,i}.$$



Rank 20



Rank 80



Rank 140

Examples: low rank



Example - Nuclear norm [Recht, Fazel and Parrilo '10]

Let $\mathbf{A} \in \mathbb{R}^{m \times n}$, and $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}^T$ be its singular value decomposition

$$\|\mathbf{A}\|_* = \sum_{i=1}^{\min\{m,n\}} s_{i,i}.$$



Rank 20



Rank 80



Rank 140

How to use regularization?

Applications

Image processing, computer vision



饮水思源 · 爱国荣校

Example: machine learning

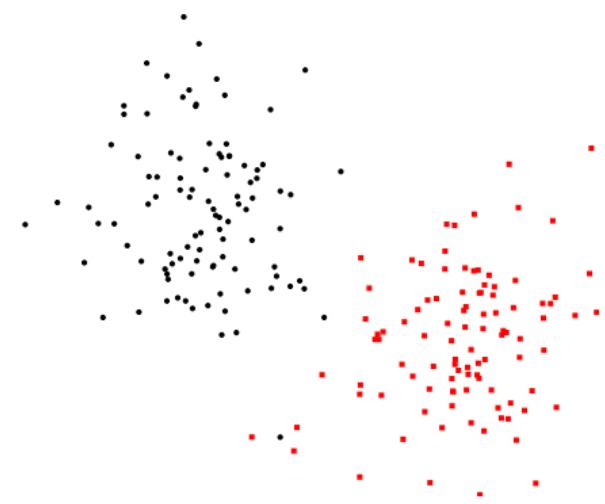


Example - Sparse logistic regression

Let $(\mathbf{a}_i, b_i) \in \mathbb{R}^n \times \{\pm 1\}$, $i = 1, \dots, m$,

$$\min_{(\mathbf{x}, y) \in \mathbb{R}^n \times \mathbb{R}} \mu \|\mathbf{x}\|_1 + \frac{1}{m} \sum_{i=1}^m f(\mathbf{x}^\top \mathbf{a}_i + y; b_i),$$

where $f(u_i; b_i) = \log(1 + e^{-u_i b_i})$.



Example: machine learning

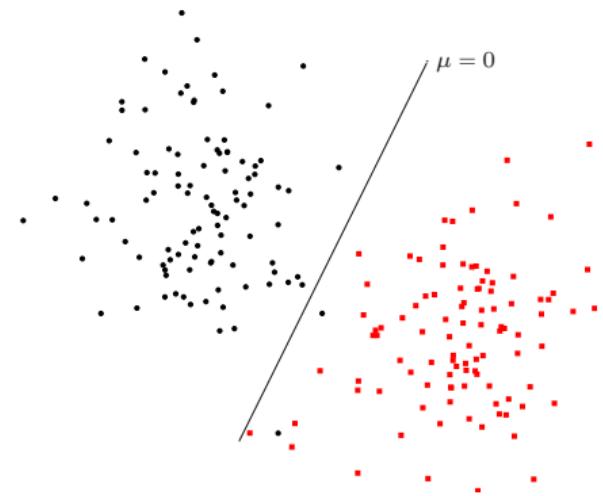


Example - Sparse logistic regression

Let $(\mathbf{a}_i, b_i) \in \mathbb{R}^n \times \{\pm 1\}$, $i = 1, \dots, m$,

$$\min_{(\mathbf{x}, y) \in \mathbb{R}^n \times \mathbb{R}} \mu \|\mathbf{x}\|_1 + \frac{1}{m} \sum_{i=1}^m f(\mathbf{x}^\top \mathbf{a}_i + y; b_i),$$

where $f(u_i; b_i) = \log(1 + e^{-u_i b_i})$.



Example: machine learning

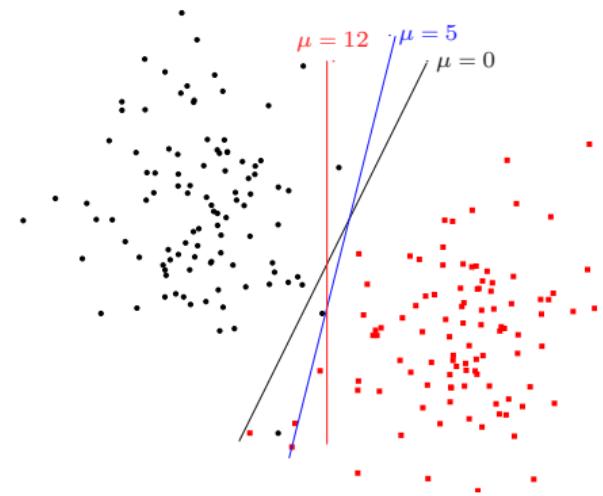


Example - Sparse logistic regression

Let $(\mathbf{a}_i, b_i) \in \mathbb{R}^n \times \{\pm 1\}, i = 1, \dots, m,$

$$\min_{(\mathbf{x}, y) \in \mathbb{R}^n \times \mathbb{R}} \mu \|\mathbf{x}\|_1 + \frac{1}{m} \sum_{i=1}^m f(\mathbf{x}^\top \mathbf{a}_i + y; b_i),$$

where $f(u_i; b_i) = \log(1 + e^{-u_i b_i}).$

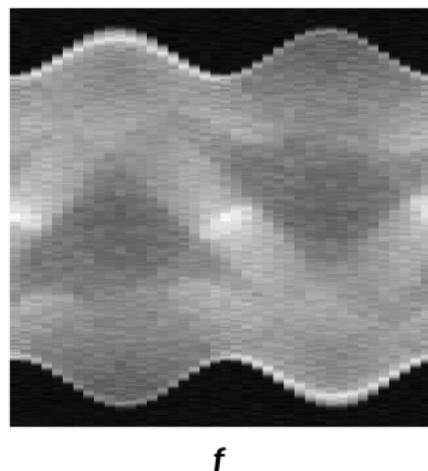


Mathematical formulation

$$\mathbf{f} = \mathcal{F}\bar{\mathbf{x}} + \varepsilon,$$

where

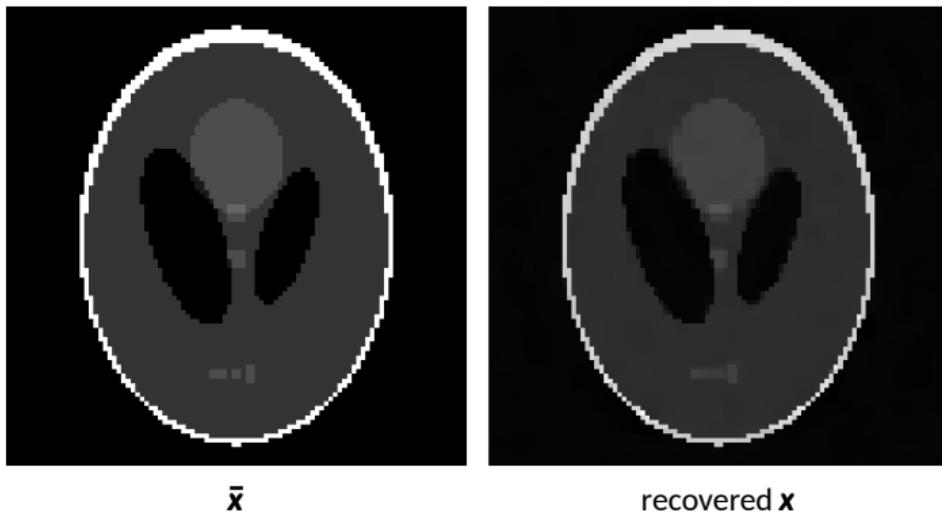
- $\bar{\mathbf{x}}$ is the true image which is **piecewise constant/smooth** — TV.
- \mathcal{F} is partial Fourier transform.
- ε is additive noise.



Example - TV based MRI reconstruction

Let $p \in \{1, 2\}$,

$$\min_{\mathbf{x} \in \mathbb{R}^{m \times n}} \mu \|\nabla \mathbf{x}\|_1 + \frac{1}{2} \|\mathbf{f} - \mathcal{F}\mathbf{x}\|_p^p.$$



Video decomposition



Mathematical formulation

$$\mathbf{f} = \bar{\mathbf{I}} + \bar{\mathbf{s}} + \varepsilon,$$

where

- $\bar{\mathbf{I}}$ is the background which is **low rank** — nuclear norm.
- $\bar{\mathbf{s}}$ is the foreground which is **sparse** — ℓ_1 -norm.
- ε is additive white Gaussian noise.

Mathematical formulation

$$\mathbf{f} = \bar{\mathbf{I}} + \bar{\mathbf{s}} + \varepsilon,$$

where

- $\bar{\mathbf{I}}$ is the background which is **low rank** — nuclear norm.
- $\bar{\mathbf{s}}$ is the foreground which is **sparse** — ℓ_1 -norm.
- ε is additive white Gaussian noise.

Example - Principal component pursuit [Candès et al '11]

$$\min_{\mathbf{I}, \mathbf{s} \in \mathbb{R}^{m \times n}} \mu(\|\mathbf{I}\|_* + \nu \|\mathbf{s}\|_1) + \frac{1}{2} \|\mathbf{I} + \mathbf{s} - \mathbf{f}\|^2.$$

Video decomposition



$$\begin{matrix} \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} \\ \textcolor{red}{\square} & \textcolor{red}{\square} & \textcolor{red}{\square} \\ \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} \end{matrix} = \begin{matrix} \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} \end{matrix} + \begin{matrix} \textcolor{white}{\square} & \textcolor{white}{\square} & \textcolor{white}{\square} \\ \textcolor{red}{\square} & \textcolor{red}{\square} & \textcolor{red}{\square} \\ \textcolor{white}{\square} & \textcolor{white}{\square} & \textcolor{white}{\square} \end{matrix}$$

f \bar{I} \bar{s}

$$\begin{matrix} \textcolor{lightgreen}{\square} \\ \textcolor{red}{\square} \\ \textcolor{lightgreen}{\square} \\ \textcolor{red}{\square} \\ \textcolor{lightgreen}{\square} \\ \textcolor{red}{\square} \\ \textcolor{lightgreen}{\square} \end{matrix} = \begin{matrix} \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} \end{matrix} + \begin{matrix} \textcolor{white}{\square} \\ \textcolor{red}{\square} \\ \textcolor{white}{\square} \\ \textcolor{white}{\square} \\ \textcolor{white}{\square} \\ \textcolor{white}{\square} \\ \textcolor{white}{\square} \end{matrix}$$

f \bar{I} \bar{s}

Video decomposition



$$\begin{matrix} \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} \\ \textcolor{red}{\square} & \textcolor{red}{\square} & \textcolor{red}{\square} \\ \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} \end{matrix} = \begin{matrix} \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} & \textcolor{lightgreen}{\square} \end{matrix} + \begin{matrix} \textcolor{white}{\square} & \textcolor{white}{\square} & \textcolor{red}{\square} \\ \textcolor{white}{\square} & \textcolor{white}{\square} & \textcolor{red}{\square} \\ \textcolor{white}{\square} & \textcolor{white}{\square} & \textcolor{red}{\square} \end{matrix}$$

f \bar{I} \bar{s}

$$\begin{matrix} \textcolor{lightgreen}{\square} \\ \textcolor{red}{\square} \\ \textcolor{lightgreen}{\square} \\ \textcolor{red}{\square} \\ \textcolor{lightgreen}{\square} \\ \textcolor{red}{\square} \\ \textcolor{lightgreen}{\square} \end{matrix} = \begin{matrix} \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} \\ \textcolor{lightgreen}{\square} \end{matrix} + \begin{matrix} \textcolor{white}{\square} \\ \textcolor{red}{\square} \\ \textcolor{white}{\square} \\ \textcolor{red}{\square} \\ \textcolor{white}{\square} \\ \textcolor{red}{\square} \\ \textcolor{white}{\square} \end{matrix}$$

f \bar{I} \bar{s}

First-order optimization methods

Non-smooth optimization, first-order methods



饮水思源 · 爱国荣校

Problem - Non-smooth optimization problem

Let $r \in \mathbb{N}_{++}$

$$\min_{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_r} \left\{ \Phi(\mathbf{x}) \stackrel{\text{def}}{=} F(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_r) + \sum_{i=1}^r R_i(\mathbf{K}_i \mathbf{x}_i) \right\},$$

where

F : smooth data fidelity term...

R_i : non-smooth regularization terms...

\mathbf{K}_i : linear/nonlinear operators...

- Signal/imaging processing, compressed sensing, inverse problems
- Statistics, data science, machine learning
- Control theory, operation research, game theory
- ...

Non-smooth, (non-convex), composite, high dimension

First-order methods: two basic ingredients

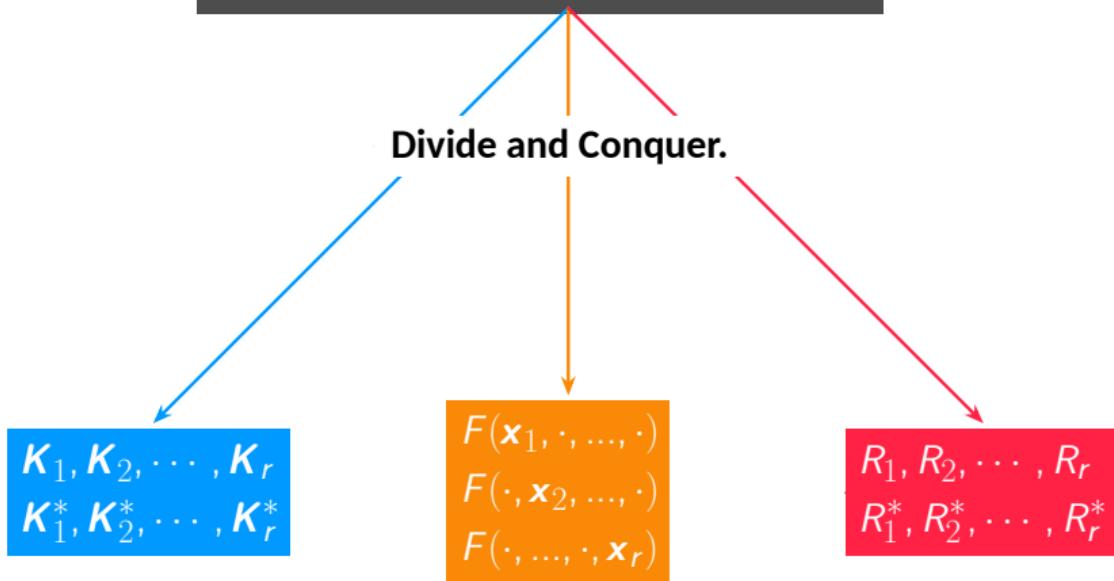


$$\min_{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_r} F(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_r) + \sum_{i=1}^r R_i(\mathbf{K}_i \mathbf{x}_i)$$

First-order methods: two basic ingredients



$$\min_{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_r} F(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_r) + \sum_{i=1}^r R_i(\mathbf{K}_i \mathbf{x}_i)$$



Algorithm - Gradient descent [Cauchy '1847]

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x})$$

where F is convex smooth differentiable with ∇F being L -Lipschitz.

Gradient descent (GD):

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \gamma_k \nabla F(\mathbf{x}^{(k)}), \quad \gamma_k \in]0, 2/L[.$$



Algorithm - Gradient descent [Cauchy '1847]

$$\min_{\mathbf{x} \in \mathbb{R}^n} F(\mathbf{x})$$

where F is convex smooth differentiable with ∇F being L -Lipschitz.

Gradient descent (GD):

$$\mathbf{x}^{(k+1)} = \mathbf{x}^{(k)} - \gamma_k \nabla F(\mathbf{x}^{(k)}), \quad \gamma_k \in]0, 2/L[.$$

Algorithm - Proximal point algorithm [Rockafellar '76]

$$\min_{\mathbf{x} \in \mathbb{R}^n} R(\mathbf{x})$$

with R being proper closed convex. Define “proximity operator” by

$$\text{prox}_{\gamma R}(\mathbf{v}) \stackrel{\text{def}}{=} \operatorname{argmin}_{\mathbf{x}} \gamma R(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{v}\|^2.$$

Proximal point algorithm (PPA):

$$\mathbf{x}^{(k+1)} = \text{prox}_{\gamma_k R}(\mathbf{x}^{(k)}), \quad \gamma_k > 0.$$

Definition - First-order methods

Numerical schemes that use *at most* the first-order differentiability, e.g. gradient or sub-gradient, of the objective.

$F + R$ Forward–Backward splitting [Lions & Mercier '79]

$R_1 + R_2$ Douglas–Rachford splitting [Douglas & Rachford '56; Lions & Mercier '79]

ADMM [Glowinski & Marrocco '75; Gabay & Mercier '76]...

$F + R(K\cdot)$ Primal–Dual splitting methods [Arrow, Hurwicz & Uzawa '58; Esser, Zhang & Chan '10; Chambolle & Pock '11]

$F + \sum_i R_i$ Generalized Forward–Backward splitting [Raguet, Fadili & Peyré '13]

– ...

Origins from numerical PDE back to 1950s, now ubiquitous in signal/image processing, inverse problems, data science, statistics, machine learning...

- Leonid I. Rudin, Stanley Osher, and Emad Fatemi. "Nonlinear total variation based noise removal algorithms." *Physica D: nonlinear phenomena* 60.1-4 (1992): 259-268.
- Ingrid Daubechies. "Ten lectures on wavelets." Society for industrial and applied mathematics, 1992.
- Stéphane Mallat. "A wavelet tour of signal processing." Elsevier, 1999.
- Benjamin Recht, Maryam Fazel, and Pablo A. Parrilo. "Guaranteed minimum-rank solutions of linear matrix equations via nuclear norm minimization." *SIAM review* 52.3 (2010): 471-501.
- Emmanuel J. Candès, Xiaodong Li, Yi Ma, John Wright. "Robust principal component analysis?." *Journal of the ACM (JACM)* 58.3 (2011): 1-37.