

音響システム特論 第2回

Advanced Topics of Acoustic Systems Day 2

小山 翔一 / Shoichi Koyama

東京大学 大学院情報理工学系研究科
Graduate School of Information Science and Technology, The University of Tokyo

Oct. 12, 2021

■ 目的 / Goals

- 音響信号を対象としたモデリングや最適化の方法論を学ぶことで、信号処理における発展的な数的手法を、応用を通して理解することを目指す。
- Gain understanding of advanced mathematical techniques in signal processing through applications by learning modeling and optimization methods for acoustics signals.

■ 概要 / Summary

- 前半は、音響信号を対象とした統計的推定法やアレイ信号処理、逆問題に関して講義する。具体的には、指向性制御、音源位置推定、空間音響、音場計測・制御などである。後半は、より発展的な内容に関して、主にゲスト講師による講義を行う。
- The first half is lectures about statistical estimation, array signal processing, and inverse problems for acoustic signals. More specifically, beamforming, source localization, spatial audio, and sound field analysis and control. The second half is mainly talks by guest lecturers about more advanced topics.

- 火曜 2 限 オンライン講義 / Tuesday 2nd period, Online lecture
- キーワード / Keywords
 - 音響信号処理, 統計的信号処理, アレイ信号処理, スパースモデリング, 逆問題, 指向性制御, 音源位置推定, 空間音響, 音場計測, 音場制御
 - Acoustic signal processing, Statistical signal processing, Array signal processing, Sparse modeling, Inverse problems, Beamforming, Source localization, Spatial audio, Sound field analysis, Sound field control
- 講義は日本語で行いますが, スライドはできるだけ英語にします。 / Lectures are in Japanese, but slides are in English as much as possible.

講義日程 / Schedule

- 10/5 逆問題と統計的推定法 / Inverse Problem and Statistical Estimation
- 10/12 信号処理における最適化法 / Optimization in Signal Processing
- 10/19 休講 / Canceled
- 10/26 スパースモデリング / Sparse Modeling
- 11/2 アレイ信号処理 / Array Signal Processing
- 11/9 波動場のモデリング / Wavefield Modeling
- 11/16 空間音響の基礎 / Basics of Spatial Audio
- 11/30 [Guest] 高道慎之介先生 / Prof. Shinnosuke Takamichi
- 12/7 [Guest] 伊藤信貴先生 (新領域) / Prof. Nobutaka Ito (GSFS)
- 12/14 休講 / Canceled
- 12/21 [Guest] 井本桂右先生 (同志社大) / Prof. Keisuke Imoto (Doshisha)
- 1/4 音場の分析と合成 / Sound Field Analysis and Synthesis
- 1/11 [Guest] 猿渡洋先生 / Prof. Hiroshi Saruwatari

- 講義資料 / Course materials
 - Downloadable at ITC-LMS or <http://www.sh01.org/ja/teaching/>
- 成績評価 / Grading system
 - レポート課題 (2回) / Reporting assignments (twice)

① Optimization in signal processing

Gradient method and adaptive filter

Convex optimization

Majorization-minimization algorithm

Optimization in signal processing

- Recall typical optimization problem appears in signal processing
- Model parameters x are estimated from observation values y as

$$\hat{x} = \arg \min_x \mathcal{D}(\mathcal{A}(x), y) + \mathcal{R}_\theta(x)$$

Here, \mathcal{A} is forward measurement operator, \mathcal{D} is data fidelity term, and \mathcal{R}_θ is regularization term with parameters θ .

- Closed-form solution is preferable, but the world is not that simple, unfortunately...
 - Cost function is nonlinear, nonconvex, or indifferentiable
 - Online algorithm is necessary for realtime application

Optimization in signal processing

- When there is no closed-form solution, **iterative method** is necessary.
 - Optimization problems appear in signal processing are *continuous optimization* in most cases.
 - *Discrete optimization* rarely appears with some exceptions (e.g., sensor placement problem)

Developing efficient optimization algorithm is one of important research topics in signal processing

Examples

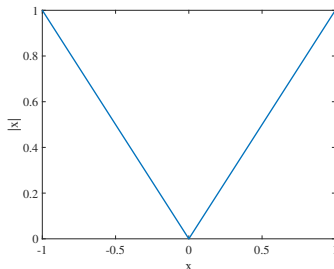
■ ℓ_1 -regularization problem

- Least-squares error + Sparsity-promoting regularization

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \frac{1}{2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|_2^2 + \lambda \|\mathbf{x}\|_1$$

Here, $\mathbf{x} \in \mathbb{R}^N$, $\mathbf{y} \in \mathbb{R}^M$, $\mathbf{A} \in \mathbb{R}^{M \times N}$, and $N > M$.

- Equivalent to MAP estimation with Laplace distribution for prior
- **Cost function is convex, but indifferentiable.**



Examples

■ Phase retrieval problem

- Optimization problem is generally formulated as

$$\underset{\mathbf{x} \in \mathbb{C}^N}{\text{minimize}} \|\mathbf{y} - |\mathbf{A}\mathbf{x}|\|_2^2$$

Here, $|\cdot|$ represents element-wise absolute value, $\mathbf{A} \in \mathbb{C}^{M \times N}$ is STFT matrix, $\mathbf{y} \in \mathbb{R}_+^M$ is given vectorized magnitude spectrogram, and $\mathbf{x} \in \mathbb{C}^N$ is time-domain signal to be recovered.

- Cost function is nonconvex and indifferentiable.

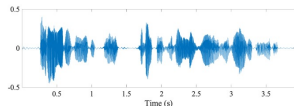
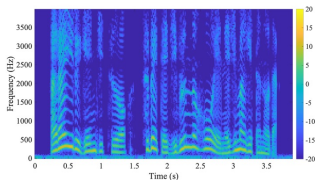


Table of Contents

① Optimization in signal processing

Gradient method and adaptive filter

Convex optimization

Majorization-minimization algorithm

Gradient method

- **Gradient method (勾配法)** is basic algorithm using gradient of objective function for unconstrained optimization problem

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} f(\mathbf{x})$$

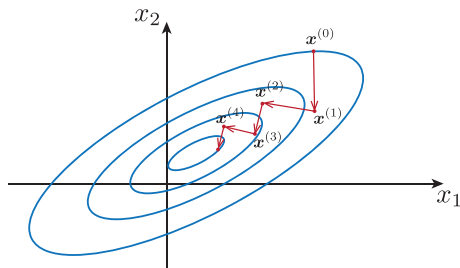
$f: \mathbb{R}^N \rightarrow \mathbb{R}$ is assumed to be differentiable.

- **Steepest descent method (最急降下法)** is representative first-order method, which iteratively update \mathbf{x} in direction to $-\nabla f(\mathbf{x}^{(i)})$ (i is iteration index)
- First-order methods, e.g., steepest descent and conjugate gradient methods, are usually preferred to second-order methods, e.g., Newton's method, in signal processing because of computational efficiency.

Gradient method

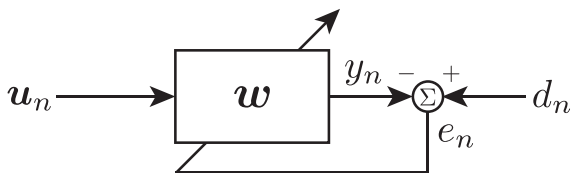
Algorithm Steepest descent method

- 1: Set initial value $\mathbf{x}^{(0)}$ and $i \leftarrow 0$
 - 2: **while** stopping condition is not satisfied **do**
 - 3: Compute descent direction $\mathbf{d}^{(i)} \leftarrow -\nabla f(\mathbf{x}^{(i)})$
 - 4: Line search of step size $\mu \geq 0$
 - 5: Update $\mathbf{x}^{(i+1)} \leftarrow \mathbf{x}^{(i)} + \mu \mathbf{d}^{(i)}$
 - 6: $i \leftarrow i + 1$
 - 7: **end while**
-



Adaptive filter

- Adaptive filter (適応フィルタ) is linear filter whose transfer function is adjusted based on observed data by adaptive process
- Gradient method can be applied to obtain adaptive filtering algorithm
- Application examples of adaptive filter are
 - Active noise control (ANC: 能動騒音制御)
 - Acoustic echo cancellation (AEC: 音響エコー除去)



Adaptive filter

- Recall **Wiener filter** to obtain filter that extract desired signal
- Input time-series signal \mathbf{u}_n and filter coefficients \mathbf{w} are

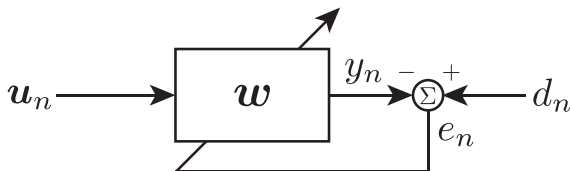
$$\mathbf{u}_n = [u_n, u_{n-1}, \dots, u_{n-K+1}]^T$$

$$\mathbf{w} = [w_1, w_2, \dots, w_K]^T$$

$\mathbb{E}[u_n] = 0$ is assumed.

- Then, output signal y_n is convolution of \mathbf{u}_k and \mathbf{w} as

$$y_n = \mathbf{w}^T \mathbf{u}_n = \sum_{j=1}^K w_j u_{n-j+1}$$



Adaptive filter

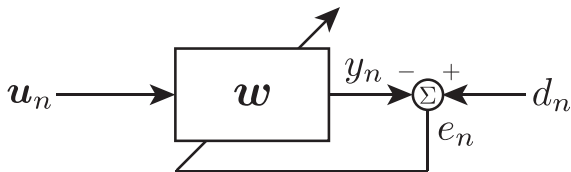
- Wiener filter is obtained by MMSE estimate of $e_n = d_n - y_n$

$$\underset{\mathbf{w} \in \mathbb{R}^K}{\text{minimize}} \mathcal{J} := \mathbb{E}[|e_n|^2]$$

- Closed-form solution of \mathbf{w} is obtained as

$$\hat{\mathbf{w}} = \mathbf{R}_u^{-1} \mathbf{r}_{ud}$$

where $\mathbf{R}_u = \mathbb{E}[\mathbf{u}_n \mathbf{u}_n^T]$ and $\mathbf{r}_{ud} = \mathbb{E}[\mathbf{u}_n d_n]$.



Steepest descent method for adaptive filter

- Online algorithm for adaptively update filter w can be obtained based **steepest descent method**
- Gradient of the cost function \mathcal{J} at $(n - 1)$ th iteration, w_{n-1} , is obtained as

$$\nabla \mathcal{J} = \left. \frac{\partial \mathcal{J}}{\partial w} \right|_{w=w_{n-1}} = -2r_{ud} + 2R_u w_{n-1}$$

- Filter w at n th iteration is obtained by updating w_{n-1} to direction of $-\nabla \mathcal{J}$ as

$$\begin{aligned} w_n &= w_{n-1} - \eta \nabla \mathcal{J} \\ &= w_{n-1} + \mu (r_{ud} - R_u w_{n-1}) \end{aligned}$$

where $\mu (= 2\eta)$ is step-size parameter.

Steepest descent method for adaptive filter

- Update formula is reformulated as

$$\begin{aligned}\mathbf{w}_n &= \mathbf{w}_{n-1} + \mu(\mathbb{E}[\mathbf{u}_n d_n] - \mathbb{E}[\mathbf{u}_n \mathbf{u}_n^T] \mathbf{w}_{n-1}) \\ &= \mathbf{w}_{n-1} + \mu \mathbb{E}[\mathbf{u}_n (d_n - \mathbf{u}_n^T \mathbf{w}_{n-1})] \\ &= \mathbf{w}_{n-1} + \mu \mathbb{E}[\mathbf{u}_n \xi_n]\end{aligned}$$

Here, *prior estimation error* ξ_n is defined as

$$\xi_n = d_n - \mathbf{u}_n^T \mathbf{w}_{n-1}$$

- Newton's method** is also applicable by calculating Hessian matrix $\nabla^2 \mathcal{J}$, but computationally expensive.

NLMS algorithm

- Update formula of steepest descent method still include expectation, which is difficult to calculate in practice.
- In *least-mean-square (LMS) algorithm*, the expectation is replaced with instantaneous value.

$$\mathbf{w}_n = \mathbf{w}_{n-1} + \mu \mathbf{u}_n \xi_n$$

- In **normalized least-mean-square (NLMS) algorithm**, gradient is normalized by power of input.

$$\mathbf{w}_n = \mathbf{w}_{n-1} + \frac{\mu}{\|\mathbf{u}_n\|^2} \mathbf{u}_n \xi_n$$

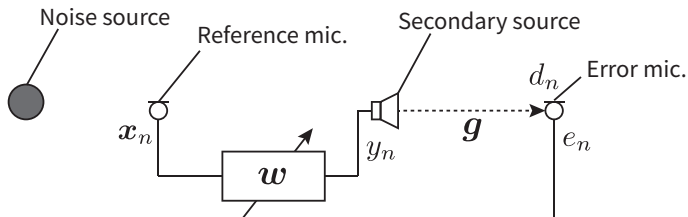
Convergence of NLMS algorithm is guaranteed by setting μ as

$$0 < \mu < 2$$

- There are a lot of variants, e.g., *affine projection* and *recursive least squares* algorithms.

Adaptive filter in ANC

- Adaptive filter is used to obtain optimal filter w to reduce power of error microphone $\mathbb{E}[|e_n|^2]$ by secondary loudspeaker output



How to handle complex variables...?

- In acoustic signal processing, many optimization problems are formulated in *time-frequency domain*, which means cost function becomes function of complex variables.
- Let us consider derivative of single complex variable function $f(z)$ where $z = x + jy$ and $x, y \in \mathbb{R}$.
- Function $f(z)$ can be represented as

$$f(z) = u(x, y) + jv(x, y)$$

where $u, v: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ are real and imaginary parts of $f(z)$.

How to handle complex variables...?

- Simple strategy is to treat f as function of two independent variables x and y :

$$\frac{\partial}{\partial x} f(x, y) = \frac{\partial u}{\partial x} + j \frac{\partial v}{\partial x}, \quad \frac{\partial}{\partial y} f(x, y) = \frac{\partial u}{\partial y} + j \frac{\partial v}{\partial y}$$

However, this procedure is sometimes complicated...

- In **Wirtinger derivative**, $f(z)$ is considered as function of z and z^* :

$$f(z, z^*) = u \left(\frac{z + z^*}{2}, \frac{z - z^*}{2} \right) + jv \left(\frac{z + z^*}{2}, \frac{z - z^*}{2} \right)$$

$$\frac{\partial}{\partial z} f(z, z^*) = \frac{1}{2} \left(\frac{\partial f(z)}{\partial x} - j \frac{\partial f(z)}{\partial y} \right)$$

$$\frac{\partial}{\partial z^*} f(z, z^*) = \frac{1}{2} \left(\frac{\partial f(z)}{\partial x} + j \frac{\partial f(z)}{\partial y} \right)$$

How to handle complex variables...?

- $\partial f / \partial z^*$ is generally used in optimization problems:

$$\frac{\partial}{\partial z^*} f(z, z^*) = \frac{1}{2} \left\{ \frac{\partial u}{\partial x} - \frac{\partial v}{\partial y} + j \left(\frac{\partial u}{\partial x} + \frac{\partial v}{\partial y} \right) \right\}$$

- Example:

$$\frac{\partial |z|^2}{\partial z^*} = \frac{\partial z z^*}{\partial z^*} = z$$

Complex matrix differentiation

- Can be extended to vectors and matrices.
- Inner product

$$\frac{\partial}{\partial \mathbf{z}^*} \mathbf{z}^H \mathbf{a} = \mathbf{a}, \quad \frac{\partial}{\partial \mathbf{z}^*} \mathbf{a}^H \mathbf{z} = \mathbf{0}$$

- Quadratic form

$$\frac{\partial}{\partial \mathbf{z}^*} \mathbf{z}^H \mathbf{A} \mathbf{z} = \mathbf{A} \mathbf{z}$$

- Trace

$$\frac{\partial}{\partial \mathbf{A}^*} \text{trace}(\mathbf{A} \mathbf{B}) = \mathbf{B}^T, \quad \frac{\partial}{\partial \mathbf{A}^*} \text{trace}(\mathbf{A} \mathbf{B} \mathbf{A}^H) = \mathbf{A}^* \mathbf{B}^T$$

① Optimization in signal processing

Gradient method and adaptive filter

Convex optimization

Majorization-minimization algorithm

Why convex optimization?

- **Convex optimization** is the methods for minimizing *convex function over convex sets*, for example,

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} && f(\mathbf{x}) \\ & \text{subject to} && g_i(\mathbf{x}) \leq 0 \quad (i = 1, \dots, P) \end{aligned}$$

where f, g_1, \dots, g_P are convex functions.

- In particular, **first-order methods** in convex optimization have been attracted attention in signal processing and machine learning since about 2005.
- Three reasons:
 - Applicable to optimization problem including indifferentially (but convex) cost function (e.g., ℓ_1 -regularization problem)
 - Suitable for large-scale problems because of low time and space complexities
 - Global optimality is guaranteed by established algorithms

Preliminaries

- Set C is convex if for any $x, y \in C$ and any $\lambda \in [0, 1]$,

$$\lambda x + (1 - \lambda)y \in C$$

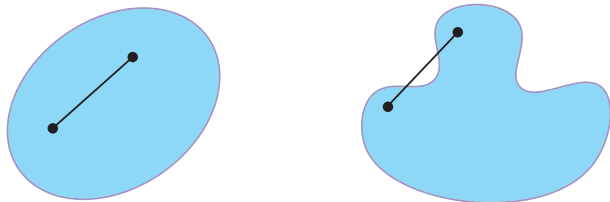


Figure: Left: Convex set, Right: Nonconvex set

Preliminaries

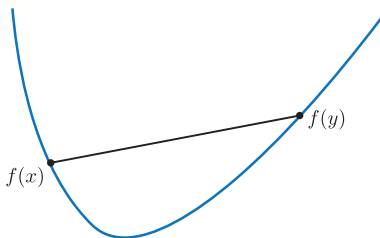
- Function $f: \mathbb{R}^N \rightarrow \mathbb{R} \cup \{+\infty\}$ is convex if for all $\mathbf{x}, \mathbf{y} \in \text{dom}(f)$ and $\lambda \in [0, 1]$,

$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y})$$

Here, $\text{dom}(f)$ is *effective domain* (実効定義域) of f defined as

$$\text{dom}(f) := \{\mathbf{x} \in \mathbb{R}^N \mid f(\mathbf{x}) < \infty\}$$

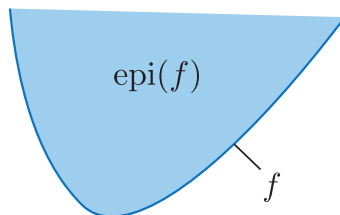
i.e., set of points for which f takes finite value.



- *Epigraph* (エピグラフ) of f is defined as

$$\text{epi}(f) := \{(\mathbf{x}, t) \in \mathbb{R}^N \times \mathbb{R} \mid f(\mathbf{x}) \leq t\}$$

- Epigraph of convex function is convex set.



Preliminaries

- Suppose convex function $f: \mathbb{R}^N \rightarrow \mathbb{R} \cup \{+\infty\}$
- f is **proper**, if its effective domain

$$\text{dom}(f) := \{\mathbf{x} \in \mathbb{R}^N \mid f(\mathbf{x}) < \infty\}$$

is not empty.

- f is **closed** if its epigraph

$$\text{epi}(f) := \{(\mathbf{x}, t) \in \mathbb{R}^N \times \mathbb{R} \mid f(\mathbf{x}) \leq t\}$$

is closed set, which is equivalent to f is **lower semicontinuous**.

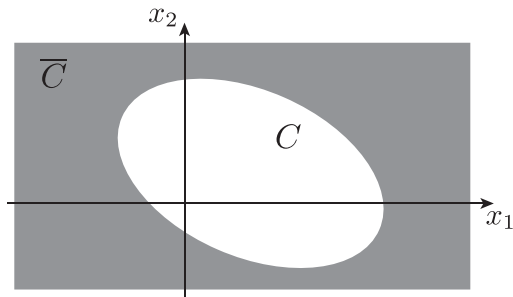
- Let $\Gamma_0(\mathbb{R}^N)$ be class of proper lower semicontinuous convex functions on \mathbb{R}^N .

Preliminaries

- For example, **indicator function** (指示関数) ι_C defined with closed nonempty convex set C as

$$\iota_C(\mathbf{x}) := \begin{cases} 0, & \mathbf{x} \in C \\ \infty, & \text{otherwise} \end{cases}$$

belongs to $\Gamma_0(\mathbb{R}^N)$.



- Single step of steepest descent method for f from \mathbf{v} to \mathbf{x} with step size γ can be rewritten as

$$\begin{aligned}\mathbf{x} &= \mathbf{v} - \gamma \nabla f(\mathbf{v}) \\ &= \arg \min_{\mathbf{x}} \left\{ f(\mathbf{v}) + \langle \nabla f(\mathbf{v}), \mathbf{x} - \mathbf{v} \rangle + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{v}\|_2^2 \right\}\end{aligned}$$

- Minimization problem of

- first-order approximation of f at \mathbf{v} : $f(\mathbf{v}) + \langle \nabla f(\mathbf{v}), \mathbf{x} - \mathbf{v} \rangle$
- proximity term with weight $1/2\gamma$: $\frac{1}{2\gamma} \|\mathbf{x} - \mathbf{v}\|_2^2$

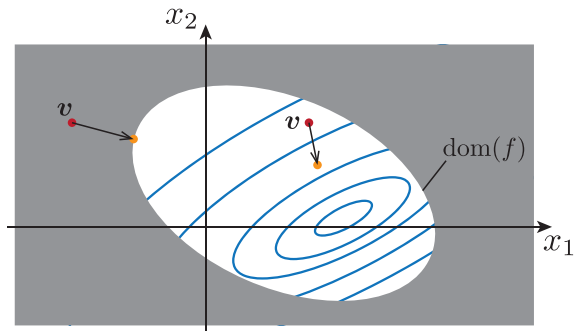
Proximal operator

- **Proximal operator/mapping** (近接作用素/写像) is powerful tool in convex optimization, which can be regarded as generalization of gradient step.
- Proximal operator $\text{prox}_{\gamma f}: \mathbb{R}^N \rightarrow \mathbb{R}^N$ for $f \in \Gamma_0(\mathbb{R}^N)$ with parameter $\gamma > 0$ is defined as

$$\text{prox}_{\gamma f}(\mathbf{v}) = \arg \min_{\mathbf{x} \in \mathbb{R}^N} \left\{ f(\mathbf{x}) + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{v}\|_2^2 \right\}$$

- $\text{prox}_{\gamma f}(\mathbf{v})$ is a point that compromises between minimizing f and being near to \mathbf{v} (*proximal point*), and γ is relative weight.

Proximal operator



- Under some assumptions, proximal operator is approximated as

$$\text{prox}_{\gamma f}(\mathbf{v}) \approx \mathbf{v} - \gamma \nabla f(\mathbf{v})$$

when γ is small and f is differentiable.

- One of the benefits of using proximal operator is f can be *nonsmooth*.

Proximal point algorithm

- Proximal point/minimization algorithm (近接点/最小化アルゴリズム) is the simplest algorithm based on proximal operator.
- Just repeatedly applying $\text{prox}_{\gamma f}$ as

Algorithm Proximal point algorithm

- 1: Set initial value $\mathbf{x}^{(0)}$ and $i \leftarrow 0$
 - 2: **while** stopping condition is not satisfied **do**
 - 3: $\mathbf{x}^{(i+1)} \leftarrow \text{prox}_{\gamma^{(i)} f}(\mathbf{x}^{(i)})$
 - 4: $i \leftarrow i + 1$
 - 5: **end while**
-

Proximal point algorithm

- Proximal point algorithm is equivalent to steepest descent method for smooth approximation function (more precisely, *Moreau envelope*) of f with step size $\gamma^{(i)}$.
- Useful when proximal operator can be quickly evaluated.

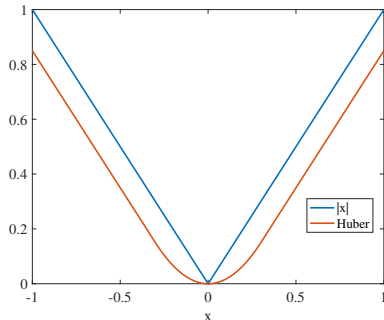


Figure: Moreau envelope of $|x|$ becomes Huber function.

How to evaluate proximal operator?

- Obviously, computability of proximal operator depends on f .
- Fortunately, proximal operators for various f useful in applications can be efficiently evaluated (*proximable*)
- Exhaustive list of proximal operators can be found at

<http://proximity-operator.net/>

(Sometimes, proximal-operator-like operator is formally defined for $f \notin \Gamma_0(\mathbb{R}^N)$.)

Examples of proximal operator

- Indicator function for closed nonempty convex set C

$$\iota_C(\mathbf{x}) := \begin{cases} 0, & \mathbf{x} \in C \\ \infty, & \text{otherwise} \end{cases}$$

- Proximal operator of $\iota_C \in \Gamma_0(\mathbb{R}^N)$ is

$$\begin{aligned} \text{prox}_{\iota_C}(\mathbf{v}) &= \arg \min_{\mathbf{x} \in \mathbb{R}^N} \left\{ \iota_C(\mathbf{x}) + \frac{1}{2} \|\mathbf{x} - \mathbf{v}\|_2^2 \right\} \\ &= \arg \min_{\mathbf{x} \in C} \|\mathbf{x} - \mathbf{v}\|_2^2 \\ &:= P_C(\mathbf{v}) \end{aligned}$$

prox_{ι_C} is projection operator onto C , i.e., P_C .

Examples of proximal operator

- ℓ_1 -norm of $\mathbf{x} \in \mathbb{R}^N$ defined as

$$f(\mathbf{x}) = \|\mathbf{x}\|_1 = \sum_{n=1}^N |x_n|$$

- Proximal operator of ℓ_1 -norm is derived as

$$\begin{aligned} \text{prox}_{\gamma\ell_1}(\mathbf{v}) &= \arg \min_{\mathbf{x} \in \mathbb{R}^N} \left\{ \|\mathbf{x}\|_1 + \frac{1}{2\gamma} \|\mathbf{x} - \mathbf{v}\|_2^2 \right\} \\ &= \arg \min_{\mathbf{x} \in \mathbb{R}^N} \sum_{n=1}^N \left\{ |x_n| + \frac{1}{2\gamma} (x_n - v_n)^2 \right\} \end{aligned}$$

x_n, v_n are n th element of \mathbf{x}, \mathbf{v} , respectively.

- Simplified to minimization problem of each element.

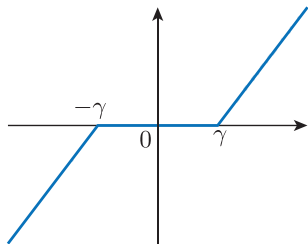
Examples of proximal operator

- Each element of $\text{prox}_{\gamma\ell_1}(\mathbf{v})$ is written as

$$\begin{aligned} [\text{prox}_{\gamma\ell_1}(\mathbf{v})]_n &= \begin{cases} v_n - \gamma, & v_n \geq \gamma \\ 0, & -\gamma < v_n < \gamma \\ v_n + \gamma, & v_n \leq -\gamma \end{cases} \\ &:= \mathcal{T}_\gamma(v_n) \end{aligned}$$

\mathcal{T}_γ is referred to as **soft thresholding operator**.

- For simplicity, $\text{prox}_{\gamma\ell_1}(\mathbf{v}) = \mathcal{T}_\gamma(\mathbf{v})$.



Examples of proximal operator

- Quadratic function defined with positive definite symmetric matrix $\mathbf{Q} \in \mathbb{R}^{N \times N}$ and vector $\mathbf{y} \in \mathbb{R}^N$

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

Note that this function is differentiable.

- Proximal operator of quadratic function is

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

Simply derived by evaluating gradient = 0.

- Applicable to stably compute $\mathbf{Q}^{-1} \mathbf{y}$.

Proximal splitting algorithms

- Generally, proximal point algorithm is not directly applicable to cost functions used in practical applications.
- Many cost functions in signal processing are written as

$$f(\mathbf{x}) = f_1(\mathbf{x}) + \cdots + f_K(\mathbf{x})$$

- Algorithms separately handling each (possibly proximable) cost function f_1, \dots, f_K are collectively called **proximal splitting algorithms** (近接分離アルゴリズム).
- Several convergence-guaranteed proximal splitting algorithms are introduced.

Proximal gradient method

- Suppose the following optimization problem:

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} f(\mathbf{x}) + g(\mathbf{x})$$

where $f, g \in \Gamma_0(\mathbb{R}^N)$, and f is differentiable.

- Proximal gradient / forward-backward splitting method (近接勾配法) for solving this problem:

Algorithm Proximal gradient method

- 1: Set initial value $\mathbf{x}^{(0)}$, $\gamma > 0$, and $i \leftarrow 0$
 - 2: **while** stopping condition is not satisfied **do**
 - 3: $\mathbf{x}^{(i+1)} \leftarrow \text{prox}_{\gamma g}(\mathbf{x}^{(i)} - \gamma \nabla f(\mathbf{x}^{(i)}))$
 - 4: $i \leftarrow i + 1$
 - 5: **end while**
-

Douglas–Rachford splitting

- Suppose the following optimization problem:

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad f(\mathbf{x}) + g(\mathbf{x})$$

where $f, g \in \Gamma_0(\mathbb{R}^N)$.

- **Douglas–Rachford splitting** for solving this problem:

Algorithm Douglas–Rachford splitting

- 1: Set initial value $\mathbf{z}^{(0)}$, $\gamma > 0$, and $i \leftarrow 0$
 - 2: **while** stopping condition is not satisfied **do**
 - 3: $\mathbf{x}^{(i+1)} \leftarrow \text{prox}_{\gamma f}(\mathbf{z}^{(i)})$
 - 4: $\mathbf{z}^{(i+1)} \leftarrow \mathbf{z}^{(i)} + \text{prox}_{\gamma g}(2\mathbf{x}^{(i+1)} - \mathbf{z}^{(i)}) - \mathbf{x}^{(i+1)}$
 - 5: $i \leftarrow i + 1$
 - 6: **end while**
-

Alternating direction method of multipliers

- Suppose the following optimization problem:

$$\underset{\mathbf{x} \in \mathbb{R}^N, \mathbf{z} \in \mathbb{R}^M}{\text{minimize}} \quad f(\mathbf{x}) + g(\mathbf{z}) \quad \text{subject to} \quad \mathbf{z} = \mathbf{A}\mathbf{x}$$

where $f \in \Gamma_0(\mathbb{R}^N)$, $g \in \Gamma_0(\mathbb{R}^M)$, and $\mathbf{A} \in \mathbb{R}^{M \times N}$.

- Alternating direction method of multipliers (ADMM: 交互方向乗数法) for solving this problem:

Algorithm ADMM

- 1: Set initial value $\mathbf{z}^{(0)}$, $\boldsymbol{\theta}^{(0)}$, $\gamma > 0$, and $i \leftarrow 0$
 - 2: **while** stopping condition is not satisfied **do**
 - 3: $\mathbf{x}^{(i+1)} \leftarrow \arg \min_{\mathbf{x} \in \mathbb{R}^N} \left\{ f(\mathbf{x}) + \frac{1}{2\gamma} \|\mathbf{A}\mathbf{x} - \mathbf{z}^{(i)} + \boldsymbol{\theta}^{(i)}\|_2^2 \right\}$
 - 4: $\mathbf{z}^{(i+1)} \leftarrow \text{prox}_{\gamma g}(\mathbf{A}\mathbf{x}^{(i+1)} + \boldsymbol{\theta}^{(i)})$
 - 5: $\boldsymbol{\theta}^{(i+1)} \leftarrow \boldsymbol{\theta}^{(i)} + \mathbf{A}\mathbf{x}^{(i+1)} - \mathbf{z}^{(i+1)}$
 - 6: $i \leftarrow i + 1$
 - 7: **end while**
-

Primal-dual splitting method

- Suppose the following optimization problem:

$$\underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} \quad f(\mathbf{x}) + g(\mathbf{x}) + h(\mathbf{A}\mathbf{x})$$

where $f, g \in \Gamma_0(\mathbb{R}^N)$, $h \in \Gamma_0(\mathbb{R}^M)$, $\mathbf{A} \in \mathbb{R}^{M \times N}$, and f is differentiable.

- Primal-dual splitting method (主-双対近接分離法) for solving this problem:

Algorithm Primal-dual splitting method

- 1: Set initial value $\mathbf{x}^{(0)}$, $\mathbf{z}^{(0)}$, $\gamma_1, \gamma_2 > 0$, and $i \leftarrow 0$
 - 2: **while** stopping condition is not satisfied **do**
 - 3: $\mathbf{x}^{(i+1)} \leftarrow \text{prox}_{\gamma_1 g}(\mathbf{x}^{(i)} - \gamma_1(\nabla f(\mathbf{x}^{(i)}) + \mathbf{A}^T \mathbf{z}^{(i)}))$
 - 4: $\mathbf{z}^{(i+1)} \leftarrow \text{prox}_{\gamma_2 h^*}(\mathbf{z}^{(i)} + \gamma_2(2\mathbf{x}^{(i+1)} - \mathbf{x}^{(i)}))$
 - 5: $i \leftarrow i + 1$
 - 6: **end while**
-

Proximal splitting in signal processing

- Typical optimization problem appears in signal processing

$$\hat{\boldsymbol{x}} = \arg \min_{\boldsymbol{x}} \mathcal{D}(\mathcal{A}(\boldsymbol{x}), \boldsymbol{y}) + \mathcal{R}_{\theta}(\boldsymbol{x})$$

- Proximal splitting algorithms are useful when \mathcal{D} and \mathcal{R}_{θ} are convex functions, and \mathcal{A} is linear operator.

⇒ Look up the list of proximal operators!

- Convexity is too strong constraint?

⇒ Effectiveness of proximal splitting algorithms are empirically shown in some nonconvex optimization problems.

① Optimization in signal processing

Gradient method and adaptive filter

Convex optimization

Majorization-minimization algorithm

Majorization-minimization algorithm

- Majorization-minimization (MM) algorithm (補助関数法 or 上界最小化アルゴリズム) is a general framework to develop optimization algorithms rather than specific algorithm.
- Many algorithms can be interpreted as special case of MM algorithm:
 - Steepest descent method, Newton's method, EM algorithm, proximal gradient method, multiplicative update rule in NMF, etc.
- Applicable to *nonconvex* optimization problem with guarantee of *monotonic non-increase of cost function*.
 - Surrogate/auxiliary function (代理関数/補助関数) of cost function is constructed
 - Parameter of surrogate function and parameter to be optimized are alternately updated

Majorization-minimization algorithm

- Objective function $\mathcal{L}(\mathbf{x})$ and its surrogate function $\mathcal{L}^+(\mathbf{x}, \boldsymbol{\xi})$ with additional variable $\boldsymbol{\xi}$.
- Surrogate function $\mathcal{L}^+(\mathbf{x}, \boldsymbol{\xi})$ is upper bound function of $\mathcal{L}(\mathbf{x})$, and is tangent to $\mathcal{L}(\mathbf{x})$ at $\boldsymbol{\xi} = \mathbf{x}$.

$$\mathcal{L}(\mathbf{x}) \leq \mathcal{L}^+(\mathbf{x}, \boldsymbol{\xi})$$

$$\mathcal{L}(\mathbf{x}) = \mathcal{L}^+(\mathbf{x}, \mathbf{x})$$

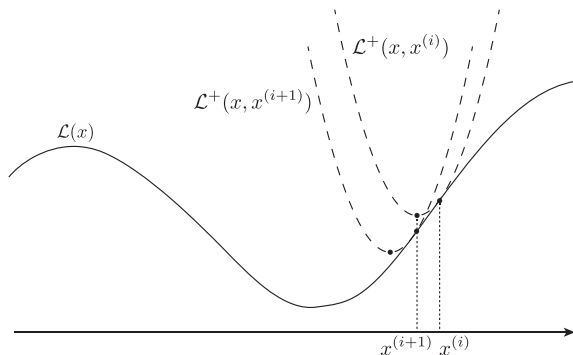
- Parameter to be optimized \mathbf{x} is iteratively updated as

$$\mathbf{x}^{(i+1)} = \arg \min_{\mathbf{x}} \mathcal{L}^+(\mathbf{x}, \mathbf{x}^{(i)})$$

Majorization-minimization algorithm

Algorithm MM algorithm

- 1: Set initial value $\mathbf{x}^{(0)}$, and $i \leftarrow 0$
- 2: **while** stopping condition is not satisfied **do**
- 3: $\mathbf{x}^{(i+1)} \leftarrow \arg \min_{\mathbf{x}} \mathcal{L}^+(\mathbf{x}, \mathbf{x}^{(i)})$
- 4: $i \leftarrow i + 1$
- 5: **end while**



Majorization-minimization algorithm

- Monotonic non-increase of objective function is guaranteed as

$$\mathcal{L}(\mathbf{x}^{(i+1)}) \leq \mathcal{L}^+(\mathbf{x}^{(i+1)}, \mathbf{x}^{(i)}) \leq \mathcal{L}^+(\mathbf{x}^{(i)}, \mathbf{x}^{(i)}) = \mathcal{L}(\mathbf{x}^{(i)})$$

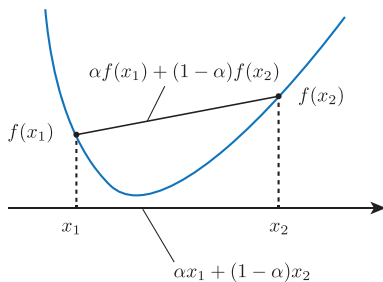
- Efficient update rule is obtained if minimization problem at each iteration can be simply solved, e.g., in closed form.
- $\mathcal{L}(\mathbf{x})$ can be nonconvex and/or indifferentiable.
- How to develop surrogate function is problem-dependent.

Tips on surrogate function design

- Jensen's inequality for convex function f

$$\sum_i \alpha_i f(\mathbf{x}_i) \geq f\left(\sum_i \alpha_i \mathbf{x}_i\right)$$

where $\sum_i \alpha_i = 1$ and $\alpha_i \geq 0 \forall i$. Equality holds if \mathbf{x}_i 's are equal.



Tips on surrogate function design

■ Examples of Jensen's inequality

- Negative log function

$$-\sum_i \alpha_i \ln x_i \geq -\ln \left(\sum_i \alpha_i x_i \right)$$

- Inverse function

$$\sum_i \alpha_i \frac{1}{x_i} \geq \frac{1}{\sum_i \alpha_i x_i}$$

Again, $\sum_i \alpha_i = 1$ and $\alpha_i \geq 0 \forall i$. Equality holds if x_i 's are equal.

Tips on surrogate function design

- Tangent line of concave function

$$f(\boldsymbol{\xi}) + \langle \nabla f(\boldsymbol{\xi}), \boldsymbol{x} - \boldsymbol{\xi} \rangle \geq f(\boldsymbol{x})$$

- Example:
 - Log function

$$\ln \xi + \frac{x}{\xi} - 1 \geq \ln x$$

Equality holds for $x = \xi$.

- In general, several inequalities are combined to design surrogate function

Hyperparameter estimation

- Example of MM algorithm (actually, **EM algorithm**): hyperparameter estimation in linear discrete model
- MAP estimation in linear discrete model ($\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{n}$) with Gaussian likelihood $\mathcal{N}(\mathbf{n}|\mathbf{0}, \mathbf{\Lambda}^{-1})$ and prior $\mathcal{N}(\mathbf{x}|\mathbf{0}, \mathbf{\Lambda}_x^{-1})$

$$\hat{\mathbf{x}} = \left(\mathbf{\Lambda}_x + \mathbf{A}^\top \mathbf{\Lambda} \mathbf{A} \right)^{-1} \mathbf{A}^\top \mathbf{\Lambda} \mathbf{y}$$

requires to set **hyperparameters**, $\mathbf{\Lambda}_x$ and $\mathbf{\Lambda}$.

- Several techniques to determine hyperparameters (or regularization parameter) exist, such as (generalized) cross validation and L-curve method.
- A method to estimate $\mathbf{\Lambda}_x$ and $\mathbf{\Lambda}$ from observation data is introduced here.

Hyperparameter estimation

- Here, hyperparameters are denoted by one variable as $\theta = \{\Lambda_x, \Lambda\}$ for notational simplicity.
- Let us consider ML estimate of θ :

$$\hat{\theta} = \arg \max_{\theta} p(\mathbf{y}|\theta)$$

which is referred to as *marginal likelihood* or *data evidence*.

- $p(\mathbf{y}|\theta)$ is represented by *marginalizing* $p(\mathbf{y}, \mathbf{x}|\theta)$ w.r.t. \mathbf{x} as

$$p(\mathbf{y}|\theta) = \int p(\mathbf{y}, \mathbf{x}|\theta) d\mathbf{x} = \int p(\mathbf{y}|\mathbf{x})p(\mathbf{x}|\theta) d\mathbf{x}$$

However, the integral in right-hand side is hard to calculate.

- In **EM algorithm**, *average log likelihood* $\mathbb{E}[\ln p(\mathbf{y}, \mathbf{x}|\theta)]$ is used instead of *marginal log likelihood* $\ln p(\mathbf{y}|\theta)$.

Hyperparameter estimation

- Original cost function is marginal log likelihood:

$$\mathcal{L}(\boldsymbol{\theta}) = \ln p(\mathbf{y}|\boldsymbol{\theta}) = \int \ln p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta})d\mathbf{x}$$

- By using posterior distribution with $\boldsymbol{\theta}^{(i)}$ ($\boldsymbol{\theta}$ at i th iteration), $\mathcal{L}(\boldsymbol{\theta})$ is written as

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}) &= \ln \int p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta})d\mathbf{x} \\ &= \ln \int \frac{p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta})}{p(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}^{(i)})} p(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}^{(i)})d\mathbf{x} \\ &= \ln \mathbb{E} \left[\frac{p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta})}{p(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}^{(i)})} \right]\end{aligned}$$

Hyperparameter estimation

- Based on Jensen's inequality, $\ln \mathbb{E}[f(\mathbf{x})] \geq \mathbb{E}[\ln f(\mathbf{x})]$, the following inequality holds:

$$\begin{aligned}\mathcal{L}(\boldsymbol{\theta}) &= \ln \mathbb{E} \left[\frac{p(\mathbf{x}, \mathbf{y} | \boldsymbol{\theta})}{p(\mathbf{x} | \mathbf{y}, \boldsymbol{\theta}^{(i)})} \right] \\ &\geq \mathbb{E} \left[\ln \frac{p(\mathbf{x}, \mathbf{y} | \boldsymbol{\theta})}{p(\mathbf{x} | \mathbf{y}, \boldsymbol{\theta}^{(i)})} \right] \\ &= \int p(\mathbf{x} | \mathbf{y}, \boldsymbol{\theta}^{(i)}) \ln \frac{p(\mathbf{x}, \mathbf{y} | \boldsymbol{\theta})}{p(\mathbf{x} | \mathbf{y}, \boldsymbol{\theta}^{(i)})} d\mathbf{x} \\ &= \int p(\mathbf{x} | \mathbf{y}, \boldsymbol{\theta}^{(i)}) \ln p(\mathbf{x}, \mathbf{y} | \boldsymbol{\theta}) d\mathbf{x} \\ &\quad - \int p(\mathbf{x} | \mathbf{y}, \boldsymbol{\theta}^{(i)}) \ln p(\mathbf{x} | \mathbf{y}, \boldsymbol{\theta}^{(i)}) d\mathbf{x} \\ &:= \bar{\mathcal{L}}(\boldsymbol{\theta}, \boldsymbol{\theta}^{(i)})\end{aligned}$$

Equality holds for $\boldsymbol{\theta} = \boldsymbol{\theta}^{(i)}$.

Hyperparameter estimation

- First term is average log likelihood. Second term is entropy of $p(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}^{(i)})$, which does not depend on $\boldsymbol{\theta}$.

$$\begin{aligned}\bar{\mathcal{L}}(\boldsymbol{\theta}, \boldsymbol{\theta}^{(i)}) &= \int p(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}^{(i)}) \ln p(\mathbf{x}, \mathbf{y}|\boldsymbol{\theta}) d\mathbf{x} \\ &\quad - \int p(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}^{(i)}) \ln p(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}^{(i)}) d\mathbf{x} \\ &= \Theta(\boldsymbol{\theta}|\boldsymbol{\theta}^{(i)}) + \mathcal{H}[p(\mathbf{x}|\mathbf{y}, \boldsymbol{\theta}^{(i)})]\end{aligned}$$

- By iteratively solving maximization problem of $\Theta(\boldsymbol{\theta}|\boldsymbol{\theta}^{(i)})$, monotonic non-decrease of $\mathcal{L}(\boldsymbol{\theta})$ is guaranteed.

$$\mathcal{L}(\boldsymbol{\theta}^{(i)}) = \bar{\mathcal{L}}(\boldsymbol{\theta}^{(i)}, \boldsymbol{\theta}^{(i)}) \leq \bar{\mathcal{L}}(\boldsymbol{\theta}^{(i+1)}, \boldsymbol{\theta}^{(i)}) = \mathcal{L}(\boldsymbol{\theta}^{(i+1)})$$

Hyperparameter estimation

- Specifically, average log likelihood is written as

$$\begin{aligned}\Theta(\Lambda_x, \Lambda) &= \frac{1}{2} \ln |\Lambda_x| - \frac{1}{2} \mathbb{E}[\mathbf{x}^\top \Lambda_x \mathbf{x}] + \frac{1}{2} \ln |\Lambda| - \frac{1}{2} \mathbb{E}[(\mathbf{y} - \mathbf{A}\mathbf{x})^\top \Lambda (\mathbf{y} - \mathbf{A}\mathbf{x})]\end{aligned}$$

Maximization problem of $\Theta(\Lambda_x, \Lambda)$ has closed-form solution.

- EM algorithm is summarized as follows:

E-step: Posterior average $\bar{\mathbf{x}}$ and precision matrix Γ are calculated as

$$\begin{aligned}\Gamma &= \Lambda_x + \mathbf{A}^\top \Lambda \mathbf{A} \\ \bar{\mathbf{x}} &= \Gamma^{-1} \mathbf{A}^\top \Lambda \mathbf{y}\end{aligned}$$

M-step: Update Λ_x and Λ as

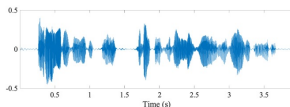
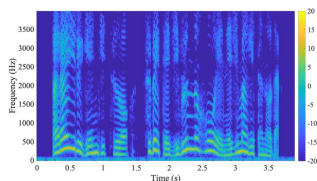
$$\begin{aligned}\hat{\Lambda}_x^{-1} &= \bar{\mathbf{x}} \bar{\mathbf{x}}^\top + \Gamma^{-1} \\ \hat{\Lambda}^{-1} &= (\mathbf{y} - \mathbf{A} \bar{\mathbf{x}})(\mathbf{y} - \mathbf{A} \bar{\mathbf{x}})^\top + \mathbf{A} \Gamma^{-1} \mathbf{A}^\top\end{aligned}$$

Phase retrieval

- Another example of MM algorithm: phase retrieval

$$\underset{\mathbf{x} \in \mathbb{C}^N}{\text{minimize}} \mathcal{L}(\mathbf{x}) := \|\mathbf{y} - |\mathbf{A}\mathbf{x}|\|_2^2$$

Here, $|\cdot|$ represents element-wise absolute value, $\mathbf{A} \in \mathbb{C}^{M \times N}$ is STFT matrix, $\mathbf{y} \in \mathbb{R}_+^M$ is given vectorized magnitude spectrogram, and $\mathbf{x} \in \mathbb{C}^N$ is time-domain signal to be recovered.



Phase retrieval

- By using $\mathbf{x}^{(i)} \in \mathbb{C}^N$, $\mathbf{v}^{(i)}$ is defined as

$$\mathbf{v}^{(i)} = \mathbf{y} \odot \exp\left(\mathbf{j} \arg(\mathbf{A}\mathbf{x}^{(i)})\right)$$

where $\arg(\cdot)$ is argument of complex variable. Then,

$$\| |\mathbf{A}\mathbf{x}| - \mathbf{y} \|^2 \leq \| \mathbf{A}\mathbf{x} - \mathbf{v}^{(i)} \|^2$$

holds. Equality holds for

$$\mathbf{x} = \mathbf{x}^{(i)}$$

- Proof is omitted...

Phase retrieval

- Surrogate function for MM algorithm is constructed as

$$\mathcal{L}^+(\mathbf{x}, \mathbf{v}^{(i)}) := \|\mathbf{A}\mathbf{x} - \mathbf{v}^{(i)}\|^2$$

which is simply minimized by

$$\arg \min_{\mathbf{x}} \mathcal{L}^+(\mathbf{x}, \mathbf{v}^{(i)}) = (\mathbf{A}^H \mathbf{A})^{-1} \mathbf{A}^H \mathbf{v}^{(i)}$$

- MM algorithm for phase retrieval is as follows:

- 1 Compute $\mathbf{v}^{(i)}$

$$\mathbf{v}^{(i)} = \mathbf{y} \odot \exp\left(\mathbf{j} \arg(\mathbf{A}\mathbf{x}^{(i)})\right)$$

- 2 Update \mathbf{x}

$$\mathbf{x}^{(i+1)} = (\mathbf{A}^H \mathbf{A})^{-1} \mathbf{A}^H \mathbf{v}^{(i)}$$

- MM algorithm for phase retrieval repeats
 - 1 Generation of time-frequency-domain signal with given magnitude and current phase
 - 2 Update estimate by inverse STFT
- Equivalent to **Griffin-Lim algorithm**.

References I



S. Haykin (2014)

Adaptive Filter Theory

[Pearson Education Limited, Harlow.](#)



S. Boyd and L. Vandenberghe (2004)

Convex Optimization

[Cambridge University Press, Cambridge.](#)



N. Parikh and S. Boyd (2013)

Proximal Algorithms

[Foundations and Trends in Optimization, vol. 1, no. 3, pp. 123–231.](#)



P. Combettes and J.-C. Pesquet (2011)

Proximal splitting methods in signal processing

[Fixed-Point Algorithms for Inverse Problems in Science and Engineering, pp. 185–212.](#)

References II



金森敬文, 鈴木大慈, 竹内一郎, 佐藤一誠 (2016)

機械学習のための連続最適化
[講談社, Tokyo.](#)



D. R. Hunter and K. Lange (2012)

A Tutorial on MM Algorithm
[The American Statistician](#), vol. 58, no. 1, pp. 30–37.



Y. Sun, P. Babu, and D. P. Palomar (2017)

Majorization-Minimization Algorithms in Signal Processing,
Communications, and Machine Learning
[IEEE Trans. Signal Process.](#), vol. 65, no. 3, pp. 794–816.