

音響システム特論 第1回

Advanced Topics of Acoustic Systems Day 1

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. 5, 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 Inverse Problem and Statistical Estimation

- Modeling in inverse problem

- Linear discrete model and maximum likelihood estimation

- Bayesian inference

- Wiener filter with application to speech enhancement

① Inverse Problem and Statistical Estimation

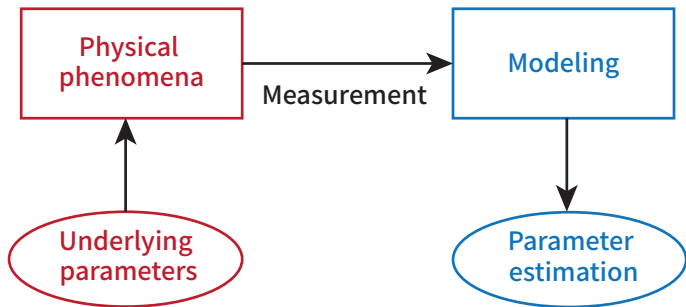
Modeling in inverse problem

Linear discrete model and maximum likelihood estimation

Bayesian inference

Wiener filter with application to speech enhancement

Inverse Problem



- Many engineering problems can be considered as parameter estimation from physical phenomena.
- This type of problems is generally called **Inverse Problem** (逆問題).

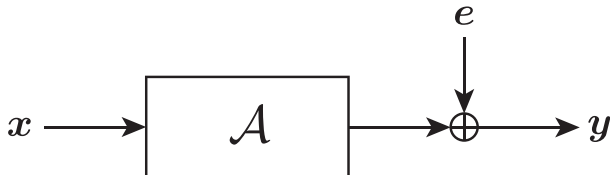
General measurement model

- Measurement $\mathbf{y} \in \mathbb{R}^M$ (or \mathbb{C}^M) is modeled as

$$\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{e}$$

where $\mathbf{x} \in \mathbb{R}^N$ (or \mathbb{C}^N) is parameter to be estimated (**model parameter**), \mathcal{A} is forward measurement operator, and $\mathbf{e} \in \mathbb{R}^M$ (or \mathbb{C}^M) is additive noise.

- Goal is to estimate \mathbf{x} from \mathbf{y} . \mathcal{A} is assumed to be given unless otherwise stated.



Examples of measurement model

■ Denoising

- x is clean speech, and y is noisy speech.
- Measurement operator is finite-dimensional identity matrix I

$$y = x + e$$

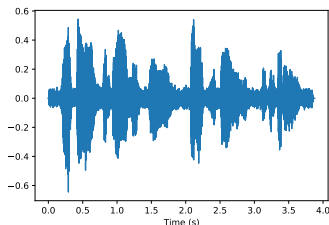
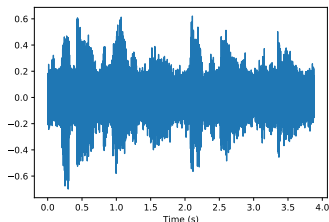


Figure: Left: Noisy speech, Right: Denoised speech

Examples of measurement model

■ Fourier analysis

- Commonly-used (nonparametric) spectral estimation

$$y_n = \sum_{k=1}^K a_k \cos\left(\frac{2\pi kn}{N}\right) + \sum_{k=1}^K b_k \sin\left(\frac{2\pi kn}{N}\right) + e_n$$

$(n \in \mathbb{Z}([0, N - 1]))$

Then, Fourier coefficients $\{a_k\}_{k=1}^K$ and $\{b_k\}_{k=1}^K$ are model parameters, i.e., $\mathbf{x} = [a_1, \dots, a_K, b_1, \dots, b_K]^T$.

- Measurement operator becomes finite-dimensional matrix

$$\mathbf{A} = \begin{bmatrix} 1 & \cdots & 1 & 0 & \cdots & 0 \\ \cos\left(\frac{2\pi}{N}\right) & \cdots & \cos\left(\frac{2\pi K}{N}\right) & \sin\left(\frac{2\pi}{N}\right) & \cdots & \sin\left(\frac{2\pi K}{N}\right) \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ \cos\left(\frac{2\pi(N-1)}{N}\right) & \cdots & \cos\left(\frac{2\pi K(N-1)}{N}\right) & \sin\left(\frac{2\pi(N-1)}{N}\right) & \cdots & \sin\left(\frac{2\pi K(N-1)}{N}\right) \end{bmatrix}$$

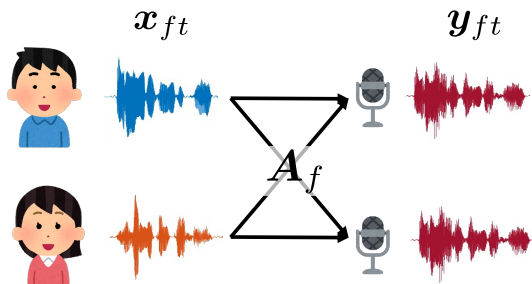
Examples of measurement model

■ Source separation

- Time-frequency-domain source signals x is measured by microphones y .
- Measurement operator becomes finite-dimensional transfer function matrix (mixture matrix).

$$y = Ax + e$$

In **blind source separation**, A is assumed to be unknown.



Examples of measurement model

■ Interpolation / Inpainting

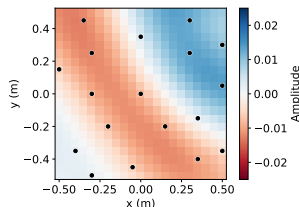
- True model parameter $x : \Omega \rightarrow \mathbb{R}$ is observed at fixed measurable set $\Omega_0 \subset \mathbb{R}^M$.

$$\mathbf{y} = x|_{\Omega_0} + \mathbf{e}$$

- When model parameter is finite-dimensional vector \mathbf{x} , measurement operator becomes finite-dimensional matrix

$$\mathbf{y} = \mathbf{S}\mathbf{x} + \mathbf{e}$$

Here, \mathbf{S} is diagonal matrix where $S_{n,n} = 1$ for sampled element and $S_{n,n} = 0$ for not sampled element.

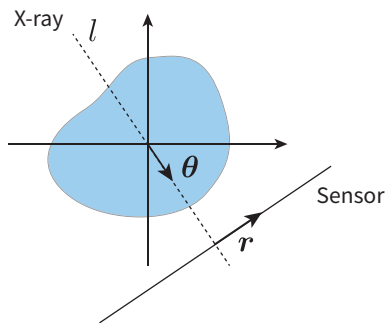


Examples of measurement model

■ Computational tomography (CT)

- X-ray traveling along line $l : s \rightarrow \mathbf{r} + s\boldsymbol{\theta}$ is detected by sensor, which is attenuated by material on l . Here, $\boldsymbol{\theta} \in \mathbb{S}^1$ and $\mathbf{r} \in \boldsymbol{\theta}^\perp$.
- Model parameter is density of object $x : \mathbb{R}^2 \rightarrow \mathbb{R}$. Then, measurement operator is represented as

$$\mathcal{A}(x)(\mathbf{r}, \boldsymbol{\theta}) = \int_{-\infty}^{\infty} x(\mathbf{r} + s\boldsymbol{\theta}) ds$$



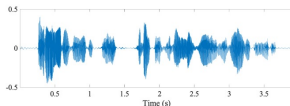
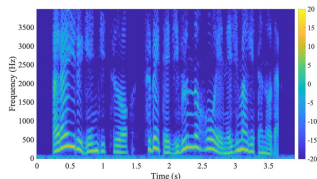
Examples of measurement model

■ Phase retrieval

- Model parameter is time-domain signal vector x .
- Goal is to estimate x only from magnitude spectrogram. Measurement operator becomes

$$\mathcal{A}(x) = |\mathbf{A}x|$$

Here, \mathbf{A} is STFT matrix, and $|\cdot|$ represents element-wise absolute value.



How to solve inverse problems?

- Well-posed problems must have i) *existence*, ii) *uniqueness*, and iii) *stability* of solution. Problems not having these properties are **ill-posed problems**.
- To estimate x from data y with given \mathcal{A} in a stable manner, a measure of data discrepancy penalized using regularizer is generally minimized.

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

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

- Various optimization algorithms can be applied, but *closed-form solution* and *linear estimator* are preferable in signal processing applications.

■ Variables

- x : scalar, \mathbf{x} : vector, \mathbf{X} : matrix (or tensor)
- Subscripts indicate their indexes, e.g., $x_{i,j}$ is (i, j) th entry of matrix \mathbf{X} .

■ Sets of real and complex numbers: \mathbb{R} and \mathbb{C}

■ Imaginary unit: $j = \sqrt{-1}$, complex conjugate: $(\cdot)^*$

■ Transpose: $(\cdot)^T$, conjugate transpose: $(\cdot)^H$

① Inverse Problem and Statistical Estimation

Modeling in inverse problem

Linear discrete model and maximum likelihood estimation

Bayesian inference

Wiener filter with application to speech enhancement

Linear discrete model

- Suppose that N model parameters $\mathbf{x} = [x_1, \dots, x_N]^T \in \mathbb{R}^N$ and M observation data $\mathbf{y} = [y_1, \dots, y_M]^T \in \mathbb{R}^M$ are related by finite-dimensional measurement matrix $\mathbf{A} \in \mathbb{R}^{M \times N}$ as

$$\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{n}$$

where $\mathbf{n} \in \mathbb{R}^M$ is vector of additive noise.

- This model is called *linear discrete model*.
- The problem to be solved is to estimate \mathbf{x} from \mathbf{y} with given \mathbf{A} .

Linear discrete model

- Assume additive noise \mathbf{n} follows multivariate Gaussian distribution of mean $\mathbf{0}$ and covariance $\sigma^2 \mathbf{I}$, i.e.,
 $\mathbf{n} \sim \mathcal{N}(\mathbf{n}|\mathbf{0}, \sigma^2 \mathbf{I})$

$$p(\mathbf{n}) = \frac{1}{(2\pi\sigma^2)^{M/2}} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{n}\|^2\right)$$

- Since linear transformation of random variable from Gaussian distribution also normally distributed,

$$p(\mathbf{y}) = \frac{1}{(2\pi\sigma^2)^{M/2}} \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2\right)$$

Maximum likelihood estimation

- Suppose $M > N$: *over-determined* (優決定) case
- **Maximum likelihood (ML) estimation (最尤推定)** of x from y : estimate x so that the observed data y is most probable.
- **Likelihood function (尤度関数)** of x , $\mathcal{L}(x)$, is defined as

$$\mathcal{L}(x) = p(y)$$

- ML estimate \hat{x}_{ML} is x that maximizes $\mathcal{L}(x)$:

$$\hat{x}_{\text{ML}} = \arg \max_x \mathcal{L}(x)$$

- For ease in computation, negative log likelihood is minimized instead of direct maximization of $\mathcal{L}(x)$ in general.

$$\hat{x}_{\text{ML}} = \arg \min_x -\ln \mathcal{L}(x)$$

Note that log function is monotonically increasing function.

Maximum likelihood estimation

- Recall that the observation data is normally distributed:

$$p(\mathbf{y}) = \frac{1}{(2\pi\sigma^2)^{M/2}} \exp\left(-\frac{1}{2\sigma^2}\|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2\right)$$

- Then, the negative log likelihood becomes

$$-\ln \mathcal{L}(\mathbf{x}) = -\frac{1}{2\sigma^2}\|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 + C$$

C denotes constant term that does not include \mathbf{x} .

- Therefore, ML estimate becomes

$$\hat{\mathbf{x}}_{\text{ML}} = \arg \min_{\mathbf{x}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2$$

This problem is equivalent to **least-squares problem (最小二乗問題)**, i.e., minimization of square error between \mathbf{y} and $\mathbf{A}\mathbf{x}$.

Maximum likelihood estimation

- ML estimate is obtained by minimizing the following cost function.

$$\hat{\mathbf{x}}_{\text{ML}} = \arg \min_{\mathbf{x}} \mathcal{J}(\mathbf{x})$$
$$\mathcal{J}(\mathbf{x}) = \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2$$

- Since \mathcal{J} is quadratic of \mathbf{x} , desired \mathbf{x} satisfies $\partial\mathcal{J}/\partial\mathbf{x} = 0$.

$$\frac{\partial\mathcal{J}}{\partial\mathbf{x}} = \frac{\partial}{\partial\mathbf{x}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 = -2\mathbf{A}^{\text{T}}(\mathbf{y} - \mathbf{A}\mathbf{x}) = 0$$

- Therefore, ML estimate is

$$\hat{\mathbf{x}}_{\text{ML}} = \left(\mathbf{A}^{\text{T}}\mathbf{A}\right)^{-1} \mathbf{A}^{\text{T}}\mathbf{y}$$

Maximum likelihood estimation

- When the noise is correlated, covariance matrix is not diagonal. Here, assume $\mathbf{n} \sim \mathcal{N}(\mathbf{n}|\mathbf{0}, \mathbf{\Sigma})$.

$$p(\mathbf{n}) = \frac{1}{(2\pi)^{M/2} |\mathbf{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2} \mathbf{n}^T \mathbf{\Sigma}^{-1} \mathbf{n}\right)$$

- ML estimate becomes

$$\begin{aligned} \hat{\mathbf{x}}_{\text{ML}} &= \arg \min_{\mathbf{x}} (\mathbf{y} - \mathbf{A}\mathbf{x})^T \mathbf{\Sigma}^{-1} (\mathbf{y} - \mathbf{A}\mathbf{x}) \\ &= (\mathbf{A}^T \mathbf{\Sigma} \mathbf{A})^{-1} \mathbf{A}^T \mathbf{\Sigma}^{-1} \mathbf{y} \end{aligned}$$

Matrix differentiation

Inner product

$$\frac{\partial}{\partial \mathbf{x}} \mathbf{x}^\top \mathbf{a} = \mathbf{a}, \quad \frac{\partial}{\partial \mathbf{x}} \mathbf{a}^\top \mathbf{x} = \mathbf{a}$$

- Since $\mathbf{x}^\top \mathbf{a} = \sum_i x_i a_i$,

$$\frac{\partial}{\partial x_i} \mathbf{x}^\top \mathbf{a} = a_i$$

- Since $\mathbf{a}^\top \mathbf{x} = \sum_i a_i x_i$,

$$\frac{\partial}{\partial x_i} \mathbf{a}^\top \mathbf{x} = a_i$$

Matrix differentiation

Quadratic form

$$\frac{\partial}{\partial \mathbf{x}} \mathbf{x}^\top \mathbf{A} \mathbf{x} = (\mathbf{A} + \mathbf{A}^\top) \mathbf{x}$$

- Since $\mathbf{x}^\top \mathbf{A} \mathbf{x} = \sum_i \sum_j a_{ij} x_i x_j$,

$$\begin{aligned} \frac{\partial}{\partial x_k} \mathbf{x}^\top \mathbf{A} \mathbf{x} &= \frac{\partial}{\partial x_k} \sum_i \sum_j a_{ij} x_i x_j = 2a_{kk} x_k + \sum_{i \neq k} a_{ik} x_i + \sum_{j \neq k} a_{kj} x_j \\ &= \sum_i a_{ik} x_i + \sum_j a_{kj} x_j \end{aligned}$$

- Or,

$$\frac{\partial}{\partial \mathbf{x}} \mathbf{x}^\top \mathbf{A} \mathbf{x} = \frac{\partial (\mathbf{x}^\top \mathbf{A}) \mathbf{x}}{\partial \mathbf{x}} + \frac{\partial \mathbf{x}^\top (\mathbf{A} \mathbf{x})}{\partial \mathbf{x}} = \mathbf{A}^\top \mathbf{x} + \mathbf{A} \mathbf{x}$$

Matrix differentiation

- Trace of matrix also appears in signal processing / machine learning.

$$\text{trace}(\mathbf{A}) = \sum_i a_{ii}$$

- Trace allows permutation of multiplication order:

$$\text{trace}(\mathbf{AB}) = \text{trace}(\mathbf{BA})$$

$$\text{trace}(\mathbf{ABC}) = \text{trace}(\mathbf{CAB}) = \text{trace}(\mathbf{BCA})$$

- Element-wise representation:

$$\text{With } \mathbf{AB} = \sum_k a_{ik}b_{kj}, \quad \text{trace}(\mathbf{AB}) = \sum_i \sum_j a_{ij}b_{ji}$$

Matrix differentiation

Trace

- Since $\frac{\partial}{\partial a_{ij}} \text{trace}(\mathbf{A}\mathbf{B}) = b_{ji}$,

$$\frac{\partial}{\partial \mathbf{A}} \text{trace}(\mathbf{A}\mathbf{B}) = \mathbf{B}^\top$$

- Since $\frac{\partial}{\partial a_{ij}} \text{trace}(\mathbf{A}^\top \mathbf{B}) = b_{ij}$,

$$\frac{\partial}{\partial \mathbf{A}} \text{trace}(\mathbf{A}^\top \mathbf{B}) = \mathbf{B}$$

- Derivative of $\text{trace}(\mathbf{A}\mathbf{B}\mathbf{A}^\top)$

$$\begin{aligned} \frac{\partial}{\partial \mathbf{A}} \text{trace}(\mathbf{A}\mathbf{B}\mathbf{A}^\top) &= \frac{\partial}{\partial \mathbf{A}} \text{trace}((\mathbf{A}\mathbf{B})\mathbf{A}^\top) + \frac{\partial}{\partial \mathbf{A}} \text{trace}(\mathbf{A}(\mathbf{B}\mathbf{A}^\top)) \\ &= \mathbf{A}(\mathbf{B} + \mathbf{B}^\top) \end{aligned}$$

Least-squares method

- **Least-squares method** solves the following optimization problem to obtain the estimate \mathbf{x} :

$$\underset{\mathbf{x}}{\text{minimize}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2$$

- Least-squares solution is equivalent to ML estimate for linear discrete model when noise is Gaussian, i.e., $\mathbf{n} \sim \mathcal{N}(\mathbf{n}|\mathbf{0}, \sigma^2\mathbf{I})$.

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

Least-squares method

- To investigate robustness of least-squares solution \hat{x} against noise, consider singular value decomposition of \mathbf{A} as

$$\mathbf{A} = \sum_{j=1}^N \gamma_j \mathbf{u}_j \mathbf{v}_j^T$$

Singular vectors $\{\mathbf{u}_1, \dots, \mathbf{u}_M\}$ and $\{\mathbf{v}_1, \dots, \mathbf{v}_N\}$ form a set of orthonormal bases.

- **Linear estimator** $\mathbf{H} := (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T$ becomes

$$\mathbf{H} := (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T = \sum_{j=1}^N \frac{1}{\gamma_j} \mathbf{v}_j \mathbf{u}_j^T$$

Least-squares method

- Least-squares solution is rewritten as

$$\hat{\boldsymbol{x}} = \boldsymbol{H}\boldsymbol{y} = \boldsymbol{x} + \sum_{j=1}^N \frac{\boldsymbol{u}_j^{\top} \boldsymbol{n}}{\gamma_j} \boldsymbol{v}_j$$

The second term represents the effect of noise on the least-squares solution.

- $\hat{\boldsymbol{x}}$ is *unbiased estimator* (不偏推定量) because $\mathbb{E}[\hat{\boldsymbol{x}}] = \boldsymbol{x}$.
- When some singular values are close to 0, the terms of small singular values can amplify the effect of noise.

Regularized least-squares method

- One way to increase robustness against noise is **regularization (正則化)**.
- By using ℓ_2 -norm regularization (or Tikhonov regularization), the cost function becomes

$$\mathcal{J}(\mathbf{x}) = \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 + \lambda\|\mathbf{x}\|^2$$

The second term can be regarded as penalty term for the large amplitude of \mathbf{x} . λ is a constant for balancing data-fidelity and penalty terms.

- The regularized least-squares solution is obtained as

$$\hat{\mathbf{x}} = \left(\mathbf{A}^T \mathbf{A} + \lambda \mathbf{I}\right)^{-1} \mathbf{A}^T \mathbf{y}$$

Regularized least-squares method

- By using the singular value decomposition of \mathbf{A} ,

$$\hat{\mathbf{x}} = \left(\sum_{j=1}^N \frac{\gamma_j^2}{\gamma_j^2 + \lambda} \mathbf{v}_j \mathbf{v}_j^T \right) \mathbf{x} + \sum_{j=1}^N \frac{\gamma_j^2}{\gamma_j^2 + \lambda} (\mathbf{u}_j^T \mathbf{n}) \mathbf{v}_j$$

By observing the second term, i.e., the effect of noise, the amplification of the effect noise owing to small γ_j is suppressed by λ .

- Note that regularized least-squares solution $\hat{\mathbf{x}}$ is NOT unbiased estimator, i.e., $\mathbb{E}[\hat{\mathbf{x}}] \neq \mathbf{x}$.

Minimum-norm solution

- Next, suppose $M < N$: *underdetermined* (劣決定) case
- This problem has infinitely many solutions that minimize the cost function of least-squares error $\|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2$. (# of equations $<$ # of unknown variables)
- Preferable solution will minimize ℓ_2 -norm of \mathbf{x} , i.e., $\|\mathbf{x}\|^2$, with satisfying $\mathbf{y} = \mathbf{A}\mathbf{x}$.

$$\begin{aligned} & \underset{\mathbf{x}}{\text{minimize}} \|\mathbf{x}\|^2 \\ & \text{subject to } \mathbf{y} = \mathbf{A}\mathbf{x} \end{aligned}$$

Such a solution is referred to as **minimum-norm solution**.

Minimum-norm solution

- By using *Lagrange's method of underdetermined multipliers* (ラグランジュ未定乗数法), Lagrange function $L(\mathbf{x}, \boldsymbol{\theta})$ is defined with Lagrangian multipliers $\boldsymbol{\theta} \in \mathbb{R}^M$ as

$$L(\mathbf{x}, \boldsymbol{\theta}) = \|\mathbf{x}\|^2 + \boldsymbol{\theta}^\top (\mathbf{y} - \mathbf{A}\mathbf{x})$$

- Original constrained optimization problem is converted into unconstrained optimization problem of minimizing L w.r.t. \mathbf{x} and $\boldsymbol{\theta}$.

Minimum-norm solution

- By differentiating $L(\mathbf{x}, \boldsymbol{\theta})$ w.r.t. \mathbf{x} and $\boldsymbol{\theta}$,

$$\frac{\partial L(\mathbf{x}, \boldsymbol{\theta})}{\partial \mathbf{x}} = 2\mathbf{x} - \mathbf{A}^\top \boldsymbol{\theta}$$
$$\frac{\partial L(\mathbf{x}, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbf{y} - \mathbf{A}\mathbf{x}$$

- By setting $\partial L/\partial \mathbf{x} = \mathbf{0}$ and $\partial L/\partial \boldsymbol{\theta} = \mathbf{0}$,

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

Minimum-norm solution

- When noise is not negligible, data-fidelity constraint should be relaxed. Therefore, an alternative optimization problem is considered as

$$\underset{\mathbf{x}}{\text{minimize}} \|\mathbf{x}\|^2 \quad \text{subject to} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 \leq \eta$$

Here, η is a constant. However, this problem does not have closed-form solution.

- Instead, inequality constraint is replaced with equality constraint as

$$\underset{\mathbf{x}}{\text{minimize}} \|\mathbf{x}\|^2 \quad \text{subject to} \quad \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 = \eta$$

Minimum-norm solution

- The Lagrangian function is defined as

$$L(\mathbf{x}, \theta) = \|\mathbf{x}\|^2 + \theta (\|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 - \eta)$$

- By setting $\lambda = 1/\theta$, optimal solution is obtained as

$$\begin{aligned}\hat{\mathbf{x}} &= \arg \min_{\mathbf{x}} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 + \lambda \|\mathbf{x}\|^2 \\ &= \mathbf{A}^T (\mathbf{A}\mathbf{A}^T + \lambda \mathbf{I})^{-1} \mathbf{y}\end{aligned}$$

The following relation for inverse matrix is used to derived the last line.

$$(\mathbf{B}^{-1} + \mathbf{C}^T \mathbf{C})^{-1} \mathbf{C}^T = \mathbf{B}\mathbf{C}^T (\mathbf{C}\mathbf{B}\mathbf{C}^T + \mathbf{I})^{-1}$$

- λ is a balancing parameter between first and second terms of $L(\mathbf{x}, \theta)$.

① Inverse Problem and Statistical Estimation

Modeling in inverse problem

Linear discrete model and maximum likelihood estimation

Bayesian inference

Wiener filter with application to speech enhancement

Bayesian inference

- In Bayesian approach, the unknown variable (model parameter) x is also regarded as random variable, which is main difference from ML estimation.
- Based on **Bayes' theorem**, *posterior probability distribution* (事後確率分布) $p(x|y)$, i.e., probability distribution of x given y , is represented as

$$p(x|y) = \frac{p(y|x)p(x)}{\int p(y|x)p(x)dx} \propto p(y|x)p(x)$$

$p(x)$ is *prior probability distribution* (事前確率分布) and $p(y|x)$ is *likelihood* (尤度).

How to determine $p(\mathbf{x})$...?

- If no prior information on \mathbf{x} , $p(\mathbf{x})$ becomes constant (non-informative prior distribution: 無情報事前分布), and Bayes' estimation corresponds to ML estimation.
- *Conjugate prior* (共役事前分布) is often used for ease in computation. Then, posterior distribution results in the same probability distribution family of prior distribution.
- For example, Gaussian family is self-conjugate w.r.t. Gaussian likelihood function. **When likelihood and prior distributions are Gaussian, posterior distribution is also Gaussian.**

MAP and MMSE estimation

Two approaches to Bayes' estimation of x .

- **Maximum a posteriori (MAP) estimation (MAP 推定)** is to estimate x so that posterior distribution $p(\mathbf{x}|\mathbf{y})$ is maximized.

$$\hat{\mathbf{x}}_{\text{MAP}} = \arg \max_{\mathbf{x}} p(\mathbf{x}|\mathbf{y})$$

- **Minimum mean square error (MMSE) estimation (MMSE 推定)** is to estimate x so that mean square error between true and estimate is minimized.

$$\hat{\mathbf{x}}_{\text{MMSE}} = \arg \min_{\hat{\mathbf{x}}} \mathbb{E}[\|\hat{\mathbf{x}} - \mathbf{x}\|^2]$$

MAP and MMSE estimation

- Mean square error is calculated as

$$\begin{aligned}\mathbb{E}[\|\hat{\mathbf{x}} - \mathbf{x}\|^2] &= \iint_{-\infty}^{\infty} (\hat{\mathbf{x}} - \mathbf{x})^T (\hat{\mathbf{x}} - \mathbf{x}) p(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} \\ &= \int_{-\infty}^{\infty} \left[\int_{-\infty}^{\infty} (\hat{\mathbf{x}} - \mathbf{x})^T (\hat{\mathbf{x}} - \mathbf{x}) p(\mathbf{x}|\mathbf{y}) d\mathbf{x} \right] p(\mathbf{y}) d\mathbf{y}\end{aligned}$$

Note that $p(\mathbf{x}, \mathbf{y}) = p(\mathbf{x}|\mathbf{y})p(\mathbf{y})$. (*multiplication theorem on probability*: 乗法定理)

- Therefore,

$$\hat{\mathbf{x}}_{\text{MMSE}} = \arg \min_{\hat{\mathbf{x}}} \int_{-\infty}^{\infty} (\hat{\mathbf{x}} - \mathbf{x})^T (\hat{\mathbf{x}} - \mathbf{x}) p(\mathbf{x}|\mathbf{y}) d\mathbf{x}$$

MAP and MMSE estimation

- By setting the derivative of the cost function w.r.t. \hat{x} to 0,

$$\frac{\partial}{\partial \hat{x}} \int_{-\infty}^{\infty} (\hat{x} - \mathbf{x})^T (\hat{x} - \mathbf{x}) p(\mathbf{x}|\mathbf{y}) d\mathbf{x} = 2 \int_{-\infty}^{\infty} (\hat{x} - \mathbf{x}) p(\mathbf{x}|\mathbf{y}) d\mathbf{x} = 0$$

Thus,

$$\hat{\mathbf{x}}_{\text{MMSE}} = \int_{-\infty}^{\infty} \mathbf{x} p(\mathbf{x}|\mathbf{y}) d\mathbf{x}$$

MMSE estimate corresponds to expectation of the posterior distribution.

- Generally complicated, but simply derived in the following two cases:
 - Both \mathbf{x} and \mathbf{y} are Gaussian. (Equivalent to MAP estimate)
 - \hat{x} is obtained by *linear estimator*, i.e., $\hat{x} = \mathbf{H}\mathbf{y} + \mathbf{g}$.

MAP estimation for Gaussian prior

- Assume $p(\mathbf{x}) \sim \mathcal{N}(\mathbf{x}|\mathbf{0}, \sigma_x^2 \mathbf{I})$

$$p(\mathbf{x}) = \frac{1}{(2\pi\sigma_x^2)^{N/2}} \exp\left(-\frac{1}{2\sigma_x^2}\|\mathbf{x}\|^2\right)$$

- MAP estimate is obtained as

$$\begin{aligned}\hat{\mathbf{x}}_{\text{MAP}} &= \arg \min_{\mathbf{x}} -\ln p(\mathbf{x}|\mathbf{y}) \\ &= \arg \min_{\mathbf{x}} [-\ln p(\mathbf{x}|\mathbf{y}) - \ln p(\mathbf{x})] \\ &= \arg \min_{\mathbf{x}} \left[\frac{1}{\sigma^2} \|\mathbf{y} - \mathbf{A}\mathbf{x}\|^2 + \frac{1}{\sigma_x^2} \|\mathbf{x}\|^2 \right] \\ &= \left(\mathbf{A}^\top \mathbf{A} + \frac{\sigma^2}{\sigma_x^2} \mathbf{I} \right)^{-1} \mathbf{A}^\top \mathbf{y}\end{aligned}$$

- Correspond to regularized least squares solution when $\lambda = \sigma^2/\sigma_x^2$, i.e., signal-to-noise ratio.

MAP estimation for Gaussian prior

Brief summary of more general case:

- Noise \mathbf{n} follows Gaussian distribution of mean $\mathbf{0}$ and precision (inverse of covariance) $\mathbf{\Lambda}$ ($:= \mathbf{\Sigma}^{-1}$).

$$\mathbf{n} \sim \mathcal{N}(\mathbf{n}|\mathbf{0}, \mathbf{\Lambda}^{-1})$$

Then, the likelihood is

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x}, \mathbf{\Lambda}^{-1})$$

- Prior distribution \mathbf{x} is Gaussian distribution of mean $\mathbf{0}$ and precision $\mathbf{\Lambda}_x$ ($:= \mathbf{\Sigma}_x^{-1}$).

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mathbf{0}, \mathbf{\Lambda}_x^{-1})$$

MAP estimation for Gaussian prior

- Since posterior distribution is also Gaussian, $p(\mathbf{x}|\mathbf{y})$ can be represented with its mean $\bar{\mathbf{x}}$ and precision $\mathbf{\Gamma}$ as

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\bar{\mathbf{x}}, \mathbf{\Gamma}^{-1})$$

Exponential part of $p(\mathbf{x}|\mathbf{y})$ becomes

$$-\frac{1}{2}(\mathbf{x} - \bar{\mathbf{x}})^{\top} \mathbf{\Gamma} (\mathbf{x} - \bar{\mathbf{x}}) = -\frac{1}{2} \mathbf{x}^{\top} \mathbf{\Gamma} \mathbf{x} + \mathbf{x}^{\top} \mathbf{\Gamma} \bar{\mathbf{x}} + C$$

C is constant term not including \mathbf{x} .

- Based on Bayes' theorem and Gaussian assumption of $p(\mathbf{y}|\mathbf{x})$ and $p(\mathbf{x})$, the exponential part becomes

$$\begin{aligned} & -\frac{1}{2} \left[\mathbf{x}^{\top} \mathbf{\Lambda}_x \mathbf{x} + (\mathbf{y} - \mathbf{A}\mathbf{x})^{\top} \mathbf{\Lambda} (\mathbf{y} - \mathbf{A}\mathbf{x}) \right] \\ & = -\frac{1}{2} \mathbf{x}^{\top} (\mathbf{\Lambda}_x + \mathbf{A}^{\top} \mathbf{\Lambda} \mathbf{A}) \mathbf{x} + \mathbf{x}^{\top} \mathbf{A}^{\top} \mathbf{\Lambda} \mathbf{y} + C' \end{aligned}$$

MAP estimation for Gaussian prior

- By comparing coefficients,

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

- Another representation of $\bar{\mathbf{x}}$:

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

where $\Sigma_y := \mathbf{A} \Lambda_x^{-1} \mathbf{A}^\top + \Lambda^{-1}$.

- Thus, MAP estimate is obtained as

$$\hat{\mathbf{x}}_{\text{MAP}} = \bar{\mathbf{x}} = (\Lambda_x + \mathbf{A}^\top \Lambda \mathbf{A})^{-1} \mathbf{A}^\top \Lambda \mathbf{y}$$

Linear MMSE estimator for Gaussian prior

- When linear estimator $\hat{\mathbf{x}} = \mathbf{H}\mathbf{y} + \mathbf{g}$ is assumed, MMSE estimate is obtained as

$$\underset{\mathbf{H}, \mathbf{g}}{\text{minimize}} \mathcal{J}(\mathbf{H}, \mathbf{g}) := \mathbb{E}[\|\mathbf{x} - \hat{\mathbf{x}}\|^2] = \mathbb{E}[\|\mathbf{x} - (\mathbf{H}\mathbf{y} + \mathbf{g})\|^2]$$

- By setting $\partial\mathcal{J}/\partial\mathbf{H} = \mathbf{0}$ and $\partial\mathcal{J}/\partial\mathbf{g} = \mathbf{0}$,

$$\frac{\partial\mathcal{J}}{\partial\mathbf{H}} = -2\mathbb{E}[(\mathbf{x} - \mathbf{H}\mathbf{y} - \mathbf{g})\mathbf{y}^T] = \mathbf{0}$$

$$\frac{\partial\mathcal{J}}{\partial\mathbf{g}} = -2\mathbb{E}[\mathbf{x} - \mathbf{H}\mathbf{y} - \mathbf{g}] = \mathbf{0}$$

- When $\mathbb{E}[\mathbf{n}] = \mathbf{0}$ and $\mathbb{E}[\mathbf{x}] = \mathbf{0}$,

$$\mathbf{g} = \mathbb{E}[\mathbf{x}] - \mathbf{H}\mathbb{E}[\mathbf{y}] = \mathbf{0}$$

$$\mathbf{H} = \mathbb{E}[\mathbf{x}\mathbf{y}^T](\mathbb{E}[\mathbf{y}\mathbf{y}^T])^{-1} := \mathbf{R}_{xy}\mathbf{R}_y^{-1}$$

Linear MMSE estimator for Gaussian prior

- Finally, MMSE estimate is obtained as

$$\hat{\mathbf{x}}_{\text{MMSE}} = \mathbf{R}_{xy} \mathbf{R}_y^{-1} \mathbf{y}$$

- When $\mathbf{n} \sim \mathcal{N}(\mathbf{n}|\mathbf{0}, \Sigma)$, $\mathbf{x} \sim \mathcal{N}(\mathbf{x}|\mathbf{0}, \Sigma_x)$, and \mathbf{x} and \mathbf{n} are uncorrelated,

$$\mathbf{R}_{xy} = \mathbb{E}[\mathbf{x}\mathbf{y}^T] = \mathbb{E}[\mathbf{x}(\mathbf{A}\mathbf{x} + \mathbf{n})^T] = \Sigma_x \mathbf{A}^T$$

$$\mathbf{R}_y = \mathbb{E}[\mathbf{y}\mathbf{y}^T] = \mathbb{E}[(\mathbf{A}\mathbf{x} + \mathbf{n})(\mathbf{A}\mathbf{x} + \mathbf{n})^T] = \mathbf{A}\Sigma_x \mathbf{A}^T + \Sigma = \Sigma_y$$

Thus,

$$\hat{\mathbf{x}}_{\text{MMSE}} = \Sigma_x \mathbf{A}^T \Sigma_y^{-1} \mathbf{y}$$

⇒ Equivalent to MAP estimate for Gaussian prior.

Summary of statistical estimation

- Maximum likelihood (ML) estimation (最尤推定)

$$\hat{\mathbf{x}}_{\text{ML}} = \arg \min_{\mathbf{x}} -\ln p(\mathbf{y}|\mathbf{x})$$

- Maximum a posteriori (MAP) estimation (最大事後確率推定)

$$\hat{\mathbf{x}}_{\text{MAP}} = \arg \min_{\mathbf{x}} -\ln p(\mathbf{x}|\mathbf{y})$$

- Minimum mean-square error (MMSE) estimation (最小平均二乗誤差推定)

$$\hat{\mathbf{x}}_{\text{MMSE}} = \arg \min_{\hat{\mathbf{x}}} \mathbb{E}[\|\hat{\mathbf{x}} - \mathbf{x}\|^2]$$

① Inverse Problem and Statistical Estimation

Modeling in inverse problem

Linear discrete model and maximum likelihood estimation

Bayesian inference

Wiener filter with application to speech enhancement

Wiener filter

- **Wiener filter** is special case of MMSE estimator.
- Input time-series signal \mathbf{u}_n and filter coefficients \mathbf{w} are

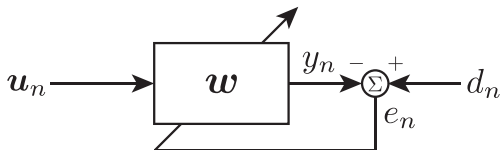
$$\mathbf{u}_n = [u_n, u_{n-1}, \dots, u_{n-K+1}]^T$$

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

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

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

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



Wiener filter

- Goal is to obtain \mathbf{w} such that output y_n corresponds to desired response d_n , i.e., error signal $e_n = d_n - y_n$ becomes 0.
- Cost function of wiener filter is defined as mean square error of e_n as

$$\begin{aligned}\mathcal{J} &= \mathbb{E}[|e_n|^2] \\ &= \mathbb{E}[(d_n - \mathbf{w}^\top \mathbf{u}_n)(d_n - \mathbf{w}^\top \mathbf{u}_n)^\top] \\ &= \mathbb{E}[|d_n|^2] - \mathbf{w}^\top \mathbb{E}[\mathbf{u}_n d_n] - \mathbb{E}[d_n \mathbf{u}_n] \mathbf{w} + \mathbf{w}^\top \mathbb{E}[\mathbf{u}_n \mathbf{u}_n^\top] \mathbf{w} \\ &= \sigma_d^2 - \mathbf{w}^\top \mathbf{r}_{ud} - \mathbf{r}_{du} \mathbf{w} + \mathbf{w}^\top \mathbf{R}_u \mathbf{w}\end{aligned}$$

Here, $\sigma_d = \mathbb{E}[|d_n|^2]$, $\mathbf{r}_{du} = \mathbb{E}[d_n \mathbf{u}_n]$, $\mathbf{r}_{ud} = \mathbb{E}[\mathbf{u}_n d_n]$, $\mathbf{R}_u = \mathbb{E}[\mathbf{u}_n \mathbf{u}_n^\top]$.

- By setting $\partial \mathcal{J} / \partial \mathbf{w} = 0$,

$$\mathbf{R}_u \mathbf{w} = \mathbf{r}_{ud}$$

This equation is called **Wiener-Hopf equation**.

- Therefore, optimal filter is

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

Noncausal Wiener filter

- By assuming input u_n is weak stationary, its autocorrelation matrix \mathbf{R}_u becomes Toeplitz.

$$\mathbf{R}_u = \begin{bmatrix} r_u(0) & r_u(1) & \cdots & r_u(K-1) \\ r_u(-1) & r_u(0) & \cdots & r_u(K-2) \\ \vdots & \vdots & \ddots & \vdots \\ r_u(1-K) & r_u(2-K) & \cdots & r_u(0) \end{bmatrix}$$

Here, $r_u(n') := \mathbb{E}[u_n u_{n-n'}]$.

- Wiener-Hopf equation is rewritten as

$$\sum_{k=1}^K w_k r_u(n-k+1) = r_{du}(n) \quad (n = 0, \dots, K-1)$$

where $r_{du}(n') := \mathbb{E}[d_n u_{n-n'}]$.

Noncausal Wiener filter

- By extending the range of n to $(-\infty, \infty)$, frequency-domain Wiener-Hopf equation is derived as

$$W(\omega)S_u(\omega) = S_{du}(\omega)$$

$W(\omega)$, $S_u(\omega)$, and $S_{du}(\omega)$ are Fourier transform of w_k , $r_u(n)$, and $r_{du}(n)$, respectively.

- Suppose that u_n is sum of desired signal d_n and noise v_n :

$$u_n = d_n + v_n$$

and d_n and v_n are uncorrelated.

- Then, noncausal Wiener filter is derived as

$$W(\omega) = \frac{S_{du}(\omega)}{S_u(\omega)} = \frac{S_d(\omega)}{S_d(\omega) + S_v(\omega)}$$

Wiener filter for denoising

■ Spectral subtraction method in speech enhancement:

- $S_v(\omega)$: Estimated from input signal in interval that speech activity is absent.
- $S_d(\omega)$: Approximated as subtraction of noise power from input power.

$$S_d(\omega) = \max(S_u(\omega) - S_v(\omega), 0)$$

- Example code:
 - https://colab.research.google.com/drive/1XE3F5h1fvy5jpfVWEK8iz_WV1vxZMdhF?usp=sharing

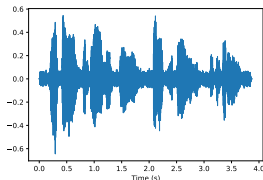
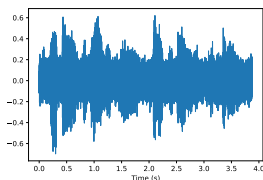


Figure: Left: Noisy speech, Right: Denoised speech

References



C. Bishop (2006)

Pattern Recognition and Machine Learning
[Springer-Verlag, New York.](#)



関原謙介 (2015)

ベイズ信号処理 – 信号・ノイズ・推定をベイズ的に考える
[共立出版, Tokyo.](#)



S. Boll (1979)

Suppression of acoustic noise in speech using spectral subtraction
[IEEE Trans. ASSP, vol. 27, no. 2, pp. 113–120.](#)