I225E Statistical Signal Processing

10. Least Squares Estimation

MAWALIM and UNOKI

candylim@jaist.ac.jp and unoki@jaist.ac.jp

School of Information Science



Least Squares Estimation

What if the data distribution is unknown?

⇒ Least Squares Estimation

[Features]

- 1. No probabilistic assumption on data
- 2. Assumption only on signal model
- 3. No guaranty for optimality
- Performance cannot be assessed without assumption about probabilistic structure of data
- 5. Widely applied, due to its ease of implementation

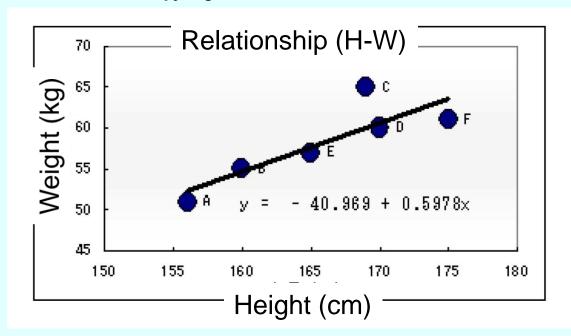
Least Squares Estimation

- Linear least squares
- Orthogonality principle
- Matrix inversion lemma
- Order-recursive least squares
- Sequential least squares

Idea is simple:

With respect to observed data s[n] and their estimation $\hat{s}[n;\theta]$, find coefficients θ of the estimation function in such a way that they minimize the squared error:

$$J(\theta) = \sum_{n=0}^{N-1} \{s[n] - \hat{s}[n; \theta]\}^2$$



3. Linear Least Squares

Estimation function constructed by linear summation of observed data:

$$\hat{s}[n;\theta] = \sum_{m=1}^{p} H_{nm} \theta_m$$

In vector formula, $\hat{s} = H\theta$

- Linear coefficients: $\theta = [\theta_1, \theta_2, \cdots, \theta_p]^T$
- Estimation: $\hat{s} = [\hat{s}[0], \hat{s}[1], \dots, \hat{s}[N-1]]^T$
- Observation: $s = [s[0], s[1], \dots, s[N-1]]^T$
- Observation matrix: $\mathbf{H} = |H_{nm}|_{n=0,\dots,N-1;m=1,\dots,p}$

Least squares error is

$$J(\theta) = \sum_{n=0}^{N-1} (s[n] - \hat{s}[n; \theta])^{2}$$

$$= (s - H\theta)^{T} (s - H\theta)$$

$$= s^{T} s - s^{T} H\theta - \theta^{T} H^{T} s + \theta^{T} H^{T} H\theta$$

$$= s^{T} s - 2\theta^{T} H^{T} s + \theta^{T} H^{T} H\theta$$

The least squares solution is

$$\frac{\partial J}{\partial \theta} = -2\mathbf{H}^T s + 2\mathbf{H}^T \mathbf{H} \theta = 0$$

$$\Rightarrow \quad \hat{\theta} = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{s}$$

If $H \in \Re^{N \times p}$ is of full rank p, $H^T H$ is invertible and the least squares solution can be obtained.

The least squares error is

$$J_{min} = J(\hat{\theta})$$

$$= (s - H\hat{\theta})^{T}(s - H\hat{\theta})$$

$$= (s - H(H^{T}H)^{-1}H^{T}s)^{T}(s - H(H^{T}H)^{-1}H^{T}s)$$

$$= s^{T}(I - H(H^{T}H)^{-1}H^{T})^{T}(I - H(H^{T}H)^{-1}H^{T})s$$

$$= s^{T}(I - H(H^{T}H)^{-1}H^{T})s$$

$$= (I - H(H^{T}H)^{-1}H^{T})^{2} = (I - H(H^{T}H)^{-1}H^{T})$$

This can be further simplified as

$$J_{min} = \mathbf{s}^T \mathbf{s} - \mathbf{s}^T \mathbf{H} (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{s}$$
$$= \mathbf{s}^T (\mathbf{s} - \mathbf{H} \hat{\theta}).$$

4. Orthogonality Principle

Decomposing the observation matrix H into column vectors h_1, h_2, \dots, h_p .

$$S = \begin{bmatrix} \boldsymbol{h}_1 & \boldsymbol{h}_2 & \cdots & \boldsymbol{h}_p \end{bmatrix} \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_p \end{bmatrix} = \sum_{i=1}^p \theta_i \boldsymbol{h}_i$$

Error vector is

$$\epsilon = \mathbf{s} - \mathbf{H}\boldsymbol{\theta} = s - \sum_{i=1}^{p} \theta_i \mathbf{h}_i$$

The least squares error and the local minimum condition are

$$J(\theta) = \|\epsilon\|^2 = \|\mathbf{s} - H\theta\|^2 = \|\mathbf{s} - \sum_{i=1}^p \theta_i \mathbf{h}_i\|^2$$
$$\frac{\partial J}{\partial \theta_i} = -2\mathbf{h}_i^T (\mathbf{s} - \sum_{i=1}^p \theta_i \mathbf{h}_i)$$
$$= -2\mathbf{h}_i^T \epsilon = 0$$
$$\to \epsilon \perp \mathbf{h}_i$$

If coefficients θ_i are chosen in such a way that error vector ϵ is orthogonal to observation vectors h_1, h_2, \dots, h_p , the least squares error J is minimum.

Inversion of partitioned matrix

 $(p+q) \times (p+q)$ matrix **M** partitioned as:

$$\mathbf{M} = \begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}, \mathbf{A} \in \mathbb{R}^{p \times p}, \mathbf{B} \in \mathbb{R}^{p \times q}, \mathbf{C} \in \mathbb{R}^{q \times p}, \mathbf{D} \in \mathbb{R}^{q \times q}$$

Inversion of partitioned matrix:

$$\mathbf{M}^{-1} = \begin{pmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{pmatrix}^{-1} = \begin{pmatrix} \left(\mathbf{A} - \mathbf{B} \mathbf{D}^{-1} \mathbf{C} \right)^{-1} & -\mathbf{A}^{-1} \mathbf{B} \left(\mathbf{D} - \mathbf{C} \mathbf{A}^{-1} \mathbf{B} \right)^{-1} \\ -\mathbf{D}^{-1} \mathbf{C} \left(\mathbf{A} - \mathbf{B} \mathbf{D}^{-1} \mathbf{C} \right)^{-1} & \left(\mathbf{D} - \mathbf{C} \mathbf{A}^{-1} \mathbf{B} \right)^{-1} \end{pmatrix}^{-1}$$

(Exercise) Prove by:

$$\mathbf{M}^{-1}\mathbf{M} = \mathbf{M}\mathbf{M}^{-1} = \mathbf{I}_{p+q}$$

Matrix inversion lemma

Suppose matrices:

$$\mathbf{A} \in \mathbb{R}^{p \times p}, \mathbf{B} \in \mathbb{R}^{p \times q}, \mathbf{C} \in \mathbb{R}^{q \times q}, \mathbf{D} \in \mathbb{R}^{q \times p}$$

Matrix inversion lemma

$$(\mathbf{A} + \mathbf{B}\mathbf{C}\mathbf{D})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{B}(\mathbf{C}^{-1} + \mathbf{D}\mathbf{A}^{-1}\mathbf{B})^{-1}\mathbf{D}\mathbf{A}^{-1}$$

(Exercise): Prove the lemma by showing:

$$(\mathbf{A} + \mathbf{BCD})(\mathbf{A} + \mathbf{BCD})^{-1} = (\mathbf{A} + \mathbf{BCD})^{-1}(\mathbf{A} + \mathbf{BCD}) = \mathbf{I}_{p}$$

5. Order-Recursive Least Squares

Denote observation matrix as H(k) ($\in R^{N \times k}$) and observed data as s ($\in R^N$). kth order least square estimation based on H(k) is given by

$$\hat{s}[n] = \sum_{m=1}^{k} H_{nm}(k) \theta_m(k).$$

Suppose a column is added to the observation matrix

$$H(k+1) = [H(k) \quad h_{k+1}] = [N \times k \quad N \times 1]$$

Then, (k + 1)th order least square estimation

$$\hat{s}[n] = \sum_{m=1}^{k+1} H_{nm}(k+1)\theta_m(k+1)$$

can be updated via the following recursive formula

$$\hat{\theta}(k+1) = \begin{bmatrix} \hat{\theta}(k) - \frac{\{H^{T}(k)H(k)\}^{-1}H^{T}(k)h_{k+1}h_{k+1}^{T}P_{k}^{\perp}s}{h_{k+1}^{T}P_{k}^{\perp}h_{k+1}} \\ \frac{h_{k+1}^{T}P_{k}^{\perp}s}{h_{k+1}^{T}P_{k}^{\perp}h_{k+1}} \end{bmatrix} = \begin{bmatrix} k \times 1 \\ 1 \times 1 \end{bmatrix}$$

Here, $P_k^{\perp} = I - H(k) \{H^T(k)H(k)\}^{-1}H^T(k)$. The minimum least square error is updated as

$$J_{min}(k+1) = J_{min}(k) - \frac{(h_{k+1}^T P_k^{\perp} s)^2}{h_{k+1}^T P_k^{\perp} h_{k+1}}$$

Moreover, inverse matrix, denoted as $D_k = (H^T(k)H(k))^{-1}$, D_k can be updated as follows

$$D_{k+1} = \begin{bmatrix} D_k + \frac{D_k H^T(k) h_{k+1} h_{k+1}^T H(k) D_k}{h_{k+1}^T P_k^{\perp} h_{k+1}} & -\frac{D_k H^T(k) h_{k+1}}{h_{k+1}^T P_k^{\perp} h_{k+1}} \\ -\frac{h_{k+1}^T H_k D_k}{h_{k+1}^T P_k^{\perp} h_{k+1}} & \frac{1}{h_{k+1}^T P_k^{\perp} h_{k+1}} \end{bmatrix}$$

[Procedure]

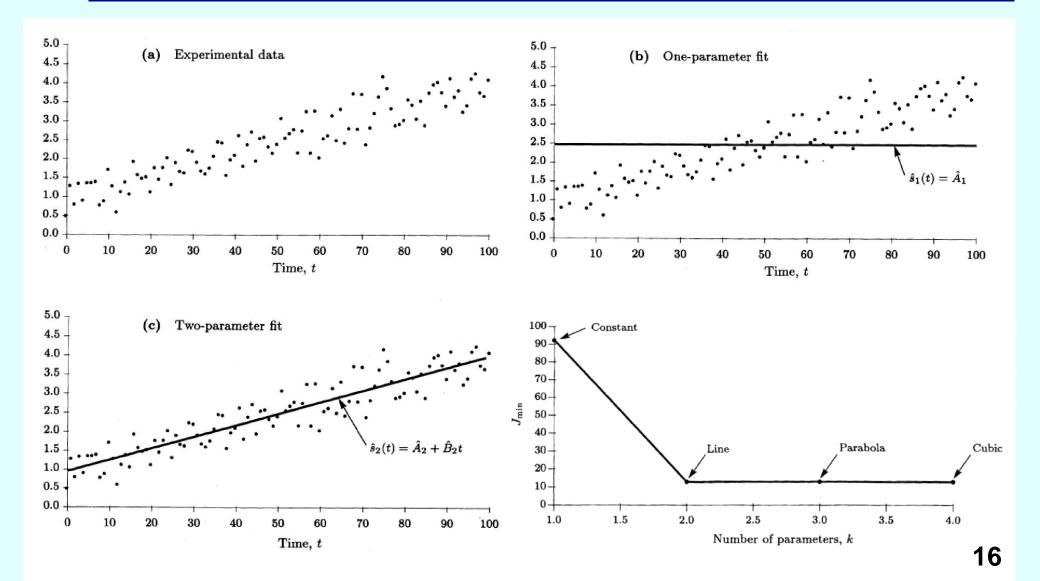
1. If
$$k = 1$$
,

$$\hat{\theta}(1) = \{ \mathbf{H}^{T}(1)\mathbf{H}(1) \}^{-1}\mathbf{H}^{T}(1)\mathbf{s}$$

$$J_{min}(1) = \{ \mathbf{s} - \mathbf{H}(1)\hat{\theta}(1) \}^{T} \{ s - H(1)\hat{\theta}(1) \}$$

2. As the order is increased $(k \to k+1)$, $\hat{\theta}(k)$, $J_{min}(k)$, D_k are updated accordingly.

Example: Curve Fitting



6. Sequential Least Squares

For observation data s(n-1) ($\in R^n$) and observation matrix H(n-1) ($\in R^{n\times p}$), consider the linear estimation:

$$\hat{\mathbf{s}} = \mathbf{H}(n-1)\theta.$$

We denote the least square solution as $\hat{\theta}(n-1)$. Suppose that additional data is observed as

$$s(n) = [s[0] \ s[1] \cdots s[n-1] \ s[n]]^T$$

$$H(n) = \begin{bmatrix} H(n-1) \\ h^T(n) \end{bmatrix} = \begin{bmatrix} n \times p \\ 1 \times p \end{bmatrix}$$

Then, the least square solution $\hat{\theta}(n)$ can be updated sequentially as follows.

Estimator Update:

$$\hat{\theta}(n) = \hat{\theta}(n-1) + K(n) \{s[n] - h^T(n)\hat{\theta}(n-1)\}$$

$$K(n) = \frac{\Sigma(n-1)h(n)}{\sigma_n^2 + h^T(n)\Sigma(n-1)h(n)}$$

Variance Update:

$$\Sigma(n) = \{ I - K(n)h^{T}(n) \} \Sigma(n-1)$$

Error Update:

$$J_{min}(n) = J_{min}(n-1) + \frac{\left\{s[n] - \boldsymbol{h}^{T}(n)\hat{\theta}(n-1)\right\}^{2}}{\sigma_{n}^{2} + \boldsymbol{h}^{T}(n)\boldsymbol{\Sigma}(n-1)\boldsymbol{h}(n)}$$

- Noise $s \hat{s}$ are independent; their covariance matrix has $\operatorname{diag}(\sigma_0^2, \sigma_1^2, \dots, \sigma_n^2)$.
- **\Sigma** stands for covariance matrix of $\hat{\theta}(\Sigma = C_{\theta})$.
- No need to compute inverse matrix.
- On-line process is possible.

Sequential least squares for DC level in white Gaussian noise

