

System Identification

6.435

SET 5

- Least Squares
- Statistical Properties

Munther A. Dahleh

Least Squares

- Linear regressions
- LS Estimates: Statistical properties
- Bias, variance, covariance
- Noise-variance estimation
- Introduction to model structure determination:

Statistical analysis and hypothesis testing

- $y(t) = \Phi^T(t)\theta$
 - ← unknown parameters
 - ← known ($1 \times n$)
 - ↳ measured

- Multivariable $y : p \times 1, \quad \Phi^T : p \times n, \quad \theta : n \times 1$

- Examples:

- $u(t) = a_0 + a_1 t + \dots + a_r t^r$
- $\Phi^T(t) = \begin{pmatrix} 1 & t & \dots & t^r \end{pmatrix} \quad \theta = \begin{pmatrix} a_0 \\ \vdots \\ a_r \end{pmatrix}$

- $y(t) = h_0(t) + h_1 u(t-1) + \dots + h_{\mu-1} u(t-\mu+1)$
- $\Phi^T(t) = \begin{pmatrix} u(t) & u(t-1) & \dots & u(t-\mu+1) \end{pmatrix}$
- $\theta = \begin{pmatrix} h_0 & h_1 & \dots & h_{\mu-1} \end{pmatrix}$

- In a matrix form

$$Y = \begin{pmatrix} y(1) \\ \vdots \\ y(N) \end{pmatrix} \quad \Phi = \begin{pmatrix} \Phi^T(1) \\ \vdots \\ \Phi^T(N) \end{pmatrix} \quad Y = \Phi\theta$$

- In general, this system is perturbed by noise

$$Y = \Phi\theta + e$$

e – noise, stochastic

- For $N > n$, the system is overdetermined.

- Define $\varepsilon = \begin{pmatrix} \varepsilon(1) \\ \vdots \\ \varepsilon(N) \end{pmatrix}$, then $\varepsilon = Y - \Phi\theta$

- LS: $\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \frac{1}{2} \|\varepsilon\|_2^2 = \underset{\theta}{\operatorname{argmin}} V(\theta)$

Solution to LS

- $\hat{\theta} = (\Phi^T \Phi)^{-1} (\Phi^T Y)$
- $V(\hat{\theta}) = \frac{1}{2} \left[y^T y - y^T \Phi (\Phi^T \Phi)^{-1} \Phi^T y \right]$
- Proof: Standard pseudo-inverse formula

$V(\hat{\theta})$: by substitution

- Equivalently

$$\hat{\theta} = \left(\sum_{t=1}^N \Phi(t) \Phi^T(t) \right)^{-1} \left(\sum_{t=1}^N \Phi(t) y(t) \right)$$

- Example

$$y(t) = b \quad \Rightarrow \quad \Phi^T(t) = 1 \quad , \quad \theta = b$$

$$\hat{\theta} = \frac{1}{N} \sum_{t=1}^N y(t)$$

Analysis of LS

- $Y = \Phi\theta + e$ $e = \begin{pmatrix} e(1) \\ \vdots \\ e(N) \end{pmatrix}$

- Assume that $E(ee^T) = \lambda^2 I$, $E(e) = 0$

\therefore white noise, zero mean, λ^2 .
variance

- Theorem:

- 1) $\hat{\theta}$: is unbiased estimate of θ

- 2) $\text{Cov}(\hat{\theta}) = E(\hat{\theta} - \theta)(\hat{\theta} - \theta)^T = \lambda^2 (\Phi^T \Phi)^{-1}$

- 3) an unbiased estimate of λ^2 is given by

$$s^2 = \frac{2V(\hat{\theta})}{N - n}$$

- Basic assumption: Φ is fixed.

- Proof:

- 1)
$$\begin{aligned}\hat{\theta} &= (\Phi^T \Phi)^{-1} (\Phi^T \Phi \theta + \Phi^T e) \\ &= \theta + (\Phi^T \Phi)^{-1} \Phi^T e.\end{aligned}$$

$$E\hat{\theta} = \theta + (\Phi^T \Phi)^{-1} \Phi^T Ee = \theta$$

$$\begin{aligned}
2) \quad E (\hat{\theta} - \theta) (\hat{\theta} - \theta)^T &= E (\Phi^T \Phi)^{-1} \Phi^T e e^T \Phi (\Phi^T \Phi)^{-1} \\
&= (\Phi^T \Phi)^{-1} \Phi^T E (e e^T) \Phi (\Phi^T \Phi)^{-1} \\
&= \lambda^2 (\Phi^T \Phi)^{-1}
\end{aligned}$$

$$\begin{aligned}
3) \quad E s^2 &= \frac{2V(\hat{\theta})}{N-n} = \frac{1}{N-n} E \left(Y^T \left(I - \Phi (\Phi^T \Phi)^{-1} \Phi^T \right) Y \right) \\
&= \frac{1}{N-n} E \left(e^T \left(I - \Phi (\Phi^T \Phi)^{-1} \Phi^T \right) e \right) \\
&= \frac{1}{N-n} E \operatorname{tr} \left(e^T \left(I_N - \Phi (\Phi^T \Phi)^{-1} \Phi^T \right) e \right) \\
&= \frac{1}{N-n} E \operatorname{tr} \left(I_N - \Phi (\Phi^T \Phi)^{-1} \Phi^T \right) \lambda^2 \\
&= \frac{1}{N-n} E \operatorname{tr} \left(I_N - \underbrace{(\Phi^T \Phi)^{-1} \Phi^T \Phi}_{I_N} \right) \lambda^2 \\
&= \frac{1}{N-n} (N-n) \lambda^2 = \lambda^2
\end{aligned}$$

Best Linear Unbiased Estimate

- $Y = \Phi\theta + e$

$$E(ee^T) = R \quad \text{correlated noise}$$

- Analysis of LS estimate:

- 1) $E(\hat{\theta}) = \theta$

- 2) $\text{Cov}(\hat{\theta}) = (\Phi^T \Phi)^{-1} \Phi^T R \Phi (\Phi^T \Phi)^{-1}$

- Consider general linear estimators:

$$\hat{\theta} = Z^T Y$$

- Want to find Z such that the estimate is unbiased and $\text{Cov}(\hat{\theta})$ is minimized.

Solution of BLUE

- Solution: $Z^* = R^{-1}\Phi (\Phi^T R^{-1}\Phi)^{-1}$
- $\text{Cov}_{Z^*}(\hat{\theta}) = (\Phi^T R^{-1}\Phi)^{-1} \leq \text{Cov}_Z(\hat{\theta})$ for any unbiased estimate.
- Proof:

$$\hat{\theta} = Z^T Y = Z^T \Phi \theta_o + Z^T e$$

$$E(\hat{\theta}) = \theta_o \quad \Rightarrow \quad Z^T \Phi = I$$

$$\text{Cov}_Z(\hat{\theta}) = E(Z^T Y - \theta_o)(Z^T Y - \theta_o)^T = Z^T R Z$$

$$\begin{aligned}\text{Cov}_{Z^*}(\hat{\theta}) &= (\Phi^T R^{-1}\Phi)^{-1} \Phi^T R^{-1} R R^{-1} \Phi (\Phi^T R^{-1}\Phi)^{-1} \\ &= (\Phi^T R^{-1}\Phi)^{-1} \Phi^T R^{-1} \Phi (\Phi^T R^{-1}\Phi)^{-1} \\ &= (\Phi^T R^{-1}\Phi)^{-1}\end{aligned}$$

- To show that $\text{Cov}_{Z^*}(\hat{\theta}) \leq \text{Cov}_Z(\hat{\theta})$

$$\begin{aligned}
\text{Cov}_Z(\hat{\theta}) - \text{Cov}_{Z^*}(\hat{\theta}) &= Z^T R Z - (\Phi^T R^{-1} \Phi)^{-1} \\
&= Z^T R Z - Z^T \Phi (\Phi^T R^{-1} \Phi)^{-1} \Phi^T Z \\
&= Z^T \left[R - \Phi (\Phi^T R^{-1} \Phi)^{-1} \Phi^T \right] Z
\end{aligned}$$

However,

$$\begin{aligned}
&R - \Phi (\Phi^T R^{-1} \Phi)^{-1} \Phi^T \\
&= \left(R - \Phi (\Phi^T R^{-1} \Phi)^{-1} \Phi^T \right) R^{-1} \left(R - \Phi (\Phi^T R^{-1} \Phi)^{-1} \Phi^T \right) \\
&\geq 0
\end{aligned}$$

result follows.

- If $R = \lambda^2 \Rightarrow Z^* = \Phi (\Phi^T \Phi)^{-1}$ and is equal to the least squares estimate. Hence LS is the BLUE when e is white.
- Can nonlinear estimates help? Not if the distribution of e is Gaussian
- BLUE of $A\theta_o$ for any constant matrix is $A\hat{\theta}$.

- Example:

$$y(t) = b_o + e(t) \quad Ee^2(t) = \lambda t^2$$

$$\Phi = \begin{pmatrix} \vdots \end{pmatrix} \quad R = \begin{pmatrix} \lambda_1^2 & & \\ & \cdots & \\ & & \lambda_N^2 \end{pmatrix}$$

BLUE of $b_o = \theta$ is

$$\hat{\theta} = \frac{1}{\sum_{j=1}^N \left(\frac{1}{\lambda_j^2} \right)} \sum_{i=1}^N \frac{1}{\lambda_i^2} y(i)$$

Maximum Likelihood Estimate

$$y = \Phi\theta_o + e \quad \theta_o \text{ is unknown}$$

$$P(y|\theta_o) = P(\Phi\theta_o + e|\theta_o)$$

Suppose: Maximum likelihood estimate

$$\hat{\theta}_{ML} = \arg \max_{\theta} P(y|\theta)$$

The parameter that makes y the “most likely event”.

Result: Suppose that y is a random variable with a distribution that depends on θ_o . Let $L(y, \theta)$ be the likelihood function, and $\hat{\theta}(y)$ is an unbiased estimate of θ . Then:

$$\text{Cov}(\hat{\theta}) \geq \left[E \left(\frac{\partial \log L}{\partial \theta} \right)^T \left(\frac{\partial \log L}{\partial \theta} \right) \right]^{-1} = - \left[E \frac{\partial^2 \log L}{\partial \theta^2} \right]^{-1}$$

Least Squares

$$y = \Phi\theta_o + e \quad e \sim N(0, \lambda^2 I)$$

$$L(Y, \theta, \lambda^2) = \frac{1}{(2\pi)^{\frac{N}{2}} (\det \lambda^2 I_N)^{\frac{1}{2}}} \exp\left(-\frac{1}{2}(Y - \Phi\theta)^T (\lambda^2 I_N)^{-1} (Y - \Phi\theta)\right)$$

$$\log L = -\frac{1}{2\lambda^2}(Y - \Phi\theta)^T (Y - \Phi\theta) - \frac{N}{2} \log 2\pi - \frac{N}{2} \log \lambda^2$$

Differentiate:

$$\frac{\partial \log L}{\partial \theta} = +\frac{1}{\lambda^2}(y - \Phi\theta)^T \Phi$$

$$\frac{\partial \log L}{\partial \lambda^2} = \frac{1}{2\lambda^4}(Y - \Phi\theta)^T (Y - \Phi\theta) - \frac{N}{2\lambda^2}$$

$$\frac{\partial^2 \log L}{\partial \theta^2} = -\frac{1}{\lambda^2} \Phi^T \Phi$$

$$\frac{\partial}{\partial \theta} \left(\frac{\partial \log L}{\partial \lambda^2} \right) = -\frac{1}{2\lambda^4} (Y - \Phi\theta)^T \Phi \quad \Rightarrow \quad E = 0$$

$$\frac{\partial^2 \log L}{\partial (\lambda^2)^2} = -\frac{1}{\lambda^6} (Y - \Phi\theta)^T (Y - \Phi\theta) + \frac{N}{2\lambda^4}$$

$$\begin{aligned} E \left(\frac{\partial^2 \log L}{\partial (\lambda^2)^2} \right) &= -\frac{1}{\lambda^6} N \lambda^2 + \frac{N}{2\lambda^4} \\ &= -\frac{N}{\lambda^4} + \frac{N}{2\lambda^4} = -\frac{1}{2} \frac{N}{\lambda^4} \end{aligned}$$

$$\text{Cov}(\hat{\theta}) \geq J^{-1} \quad \hat{\theta} = \begin{pmatrix} \hat{\theta} \\ \lambda^2 \end{pmatrix}$$

$$J = E \begin{pmatrix} \frac{1}{\lambda^2} \Phi^T \Phi & 0 \\ 0 & \frac{N}{2\lambda^4} \end{pmatrix}$$

$$\text{Cov}(\hat{\theta}) \geq \left(\frac{1}{\lambda^2} \Phi^T \Phi \right)^{-1} = \lambda^2 (\Phi^T \Phi)^{-1}$$

$$\text{Cov}(\lambda^2) \geq \frac{2\lambda^4}{N}$$

• Remarks:

- LS estimate for $\hat{\theta}$ is efficient, i.e. $\text{Cov}(\hat{\theta}) = J^{-1}$
= lower bound (termed efficient)
- LS estimate of λ^2 (guess) is asymptotically efficient.

To show that:

$$\text{Var}(s^2) = E(s^2 - \lambda^2)^2 = E(s^2)^2 - \lambda^4 \quad (\text{assuming Gaussian dist})$$

$$\text{and } s^2 = \frac{e^T P e}{N - n} \quad P = \left(\mathbf{1} - \Phi (\Phi^T \Phi)^{-1} \Phi^T \right)$$
$$P^2 = P$$

$$\text{Var}(s^2) = \frac{2\lambda^4}{N - n} > \frac{2\lambda^4}{N} = \text{Cramer-Rao bdd.}$$

- Result: BLUE for Gaussian distribution is best estimate over Nonlinear estimators.