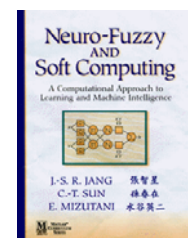


Chapter 5: Least-Square Methods for System Identification

- ◆ System Identification: an Introduction (5.1)
- ◆ Least-Squares Estimators (5.3)
- ◆ Statistical Properties & the Maximum Likelihood Estimator (5.7)
- ◆ LSE for Nonlinear Models (5.8)



Jyh-Shing Roger Jang et al., *Neuro-Fuzzy and Soft Computing: A Computational Approach to Learning and Machine Intelligence*, First Edition, Prentice Hall, 1997

System Identification: Introduction (5.1)

- ◆ Goal
 - Determine a mathematical model for an unknown system (or target system) by observing its input-output data pairs
- ◆ Purposes
 - To predict a system's behavior, as in time series prediction & weather forecasting
 - To explain the interactions & relationships between inputs & outputs of a system

System Identification: Introduction (5.1) (cont.)

◆ Purposes (cont.)

- To design a controller based on the model of a system, as an aircraft or ship control
- Simulate the system under control once the model is known

System Identification: Introduction (5.1) (cont.)

◆ There are 2 main steps that are involved

- Structure identification
- Parameter identification

System Identification: Introduction (5.1) (cont.)

– Structure identification

Apply a-priori knowledge about the target system to determine a class of models within which the search for the most suitable model is to be conducted; this class of model is denoted by a function $y = f(u, \theta)$ where:

- y is the model output
- u is the input vector
- θ is the parameter vector

f depends on the problem at hand & on the designer's experience & the laws of nature governing the target system

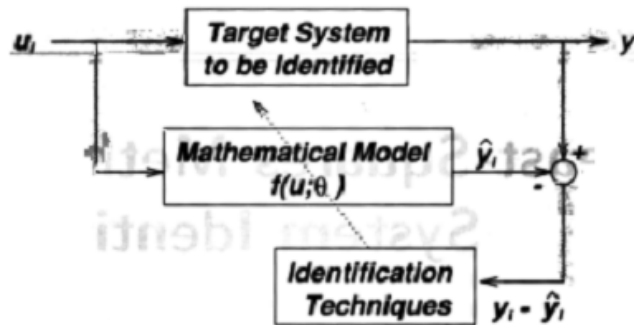
System Identification: Introduction (5.1) (cont.)

– Parameter identification

- The structure of the model is known, however we need to apply optimization techniques in order to determine the parameter vector $\theta = \hat{\theta}$ such that the resulting model $\hat{y} = f(u, \hat{\theta})$ describes the system appropriately:

$$|y_i - \hat{y}| \rightarrow 0 \text{ with } y_i \text{ assigned to } u_i$$

System Identification: Introduction (5.1) (cont.)



Block diagram for parameter identification

System Identification: Introduction (5.1) (cont.)

- ◆ The data set composed of m desired input-output pairs $(u_i; y_i)$ ($i = 1, \dots, m$) is called the training data
- ◆ System identification needs to do both structure & parameter identification repeatedly until satisfactory model is found: it does this as follows:
 - Specify & parameterize a class of mathematical models representing the system to be identified
 - Perform parameter identification to choose the parameters that best fit the training data set
 - Conduct validation set to see if the model identified responds correctly to an unseen data set
 - Terminate the procedure once the results of the validation test are satisfactory. Otherwise, another class of model is selected & repeat step 2 to 4

Least-Squares Estimators (5.3)

◆ General form:

$$y = \theta_1 f_1(u) + \theta_2 f_2(u) + \dots + \theta_n f_n(u) \quad (*)$$

where:

- $u = (u_1, \dots, u_p)^T$ is the model input vector
- f_1, \dots, f_n are known functions of u
- $\theta_1, \dots, \theta_n$ are unknown parameters to be estimated

Least-Squares Estimators (5.3) (cont.)

◆ The task of fitting data using a linear model is referred to as linear regression

◆ We collect a training data set

$\{(u_i; y_i), i = 1, \dots, m\}$

Equation (*) becomes:

$$\begin{cases} f_1(u_1)\theta_1 + f_2(u_1)\theta_2 + \dots + f_n(u_1)\theta_n = y_1 \\ f_1(u_2)\theta_1 + f_2(u_2)\theta_2 + \dots + f_n(u_2)\theta_n = y_2 \\ \vdots \\ f_1(u_m)\theta_1 + f_2(u_m)\theta_2 + \dots + f_n(u_m)\theta_n = y_m \end{cases}$$

Which is equivalent to: $A \theta = y$

Least-Squares Estimators (5.3) (cont.)

Where: A is an $m \times n$ matrix which is:
$$\mathbf{A} = \begin{bmatrix} \mathbf{f}_1(\mathbf{u}_1) \cdots \mathbf{f}_n(\mathbf{u}_1) \\ \vdots \\ \mathbf{f}_1(\mathbf{u}_m) \cdots \mathbf{f}_n(\mathbf{u}_m) \end{bmatrix}$$

θ is $n \times 1$ unknown parameter vector:
$$\theta = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_n \end{bmatrix}$$

and y is an $m \times 1$ output vector:

$$y = \begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} \text{ and } \mathbf{a}_i^T = [\mathbf{f}_1(\mathbf{u}_i), \dots, \mathbf{f}_n(\mathbf{u}_i)]$$

$$\mathbf{A} \theta = y \Leftrightarrow \theta = \mathbf{A}^{-1}y \text{ (solution)}$$

Least-Squares Estimators (5.3) (cont.)

- ◆ We have m outputs & n fitting parameters to find (or m equations & n unknown variables)
- ◆ Usually m is greater than n , since the model is just an approximation of the target system & the data observed might be corrupted, therefore an exact solution is not always possible!
- ◆ To overcome this inherent conceptual problem, an error vector e is added to compensate

$$\mathbf{A} \theta + e = y$$

Least-Squares Estimators (5.3) (cont.)

- ◆ Our goal consists now of finding $\hat{\theta}$ that reduces the errors between y_i & $\hat{y}_i = \mathbf{a}_i^T \theta$

$$\underset{\theta}{\text{minimize}} \sum_{i=1}^m |y_i - \mathbf{a}_i^T \theta| \quad (\text{not derivable!})$$

rather

$$\underset{\theta}{\text{minimize}} \sum_{i=1}^m (y_i - \mathbf{a}_i^T \theta)^2$$

Least-Squares Estimators (5.3) (cont.)

- ◆ If $\mathbf{e} = \mathbf{y} - \mathbf{A}\theta$ then:

$$\mathbf{E}(\theta) = \sum_{i=1}^m (y_i - \mathbf{a}_i^T \theta)^2 = \mathbf{e}^T \mathbf{e} = (\mathbf{y} - \mathbf{A}\theta)^T (\mathbf{y} - \mathbf{A}\theta)$$

We need to compute:

$$\underset{\theta}{\text{min}} (\mathbf{y} - \mathbf{A}\theta)^T (\mathbf{y} - \mathbf{A}\theta)$$

Least-Squares Estimators (5.3) (cont.)

◆ Theorem [least-squares estimator]

The squared error is minimized when $\theta = \hat{\theta}$ (called the least-squares estimators, LSE) satisfies the normal equation

$$A^T A \hat{\theta} = A^T y,$$

if $A^T A$ is nonsingular, $\hat{\theta}$ is unique & is given by

$$\hat{\theta} = (A^T A)^{-1} A^T y$$

Least-Squares Estimators (5.3) (cont.)

◆ Example

- The relationship is between the spring length & the force applied $L = k_1 f + k_0$ (linear model)
- Goal: find $\hat{\theta} = (k_0, k_1)^T$ that best fits the data for a given force f_0 , we need to determine the corresponding spring length L_0
 - Solution 1: provide 2 pairs (L_0, f_0) and (L_1, f_1) and solve a linear system of 2 equations and 2 variables k_1 and k_0
 $\Rightarrow \hat{k}_1$ & \hat{k}_0 . However, because of noisy data, this solution is not reliable!
 - Solution 2: use a larger training set (L_i, f_i)

Least-Squares Estimators (5.3) (cont.)

Experiment	Force (newtons)	Length of Spring (inches)
1	1.1	1.5
2	1.9	2.1
3	3.2	2.5
4	4.4	3.3
5	5.9	4.1
6	7.4	4.6
7	9.2	5.0

Training data for the spring example

Least-Squares Estimators (5.3) (cont.)

since $y = e + A\theta$, we can write:

$$\underbrace{\begin{bmatrix} 1 & 1.1 \\ 1 & 1.9 \\ 1 & 3.2 \\ 1 & 4.4 \\ 1 & 5.9 \\ 1 & 7.4 \\ 1 & 9.2 \end{bmatrix}}_A + \underbrace{\begin{bmatrix} k_0 \\ k_1 \\ \theta \end{bmatrix}}_{\theta} + \underbrace{\begin{bmatrix} e_1 \\ e_2 \\ e_3 \\ e_4 \\ \vdots \\ e_7 \\ e \end{bmatrix}}_e = \underbrace{\begin{bmatrix} 1.5 \\ 2.1 \\ 2.5 \\ 3.3 \\ 4.1 \\ 4.6 \\ 5.0 \end{bmatrix}}_y$$

therefore the LSE of $[k_0, k_1]^T$ which minimizes

$$e^T e = \sum_{i=1}^7 e_i^2 \text{ is equal to : } \begin{bmatrix} \hat{k}_0 \\ \hat{k}_1 \end{bmatrix} = (A^T A)^{-1} A^T y = \begin{bmatrix} 1.20 \\ 0.44 \end{bmatrix}$$

Least-Squares Estimators (5.3) (cont.)

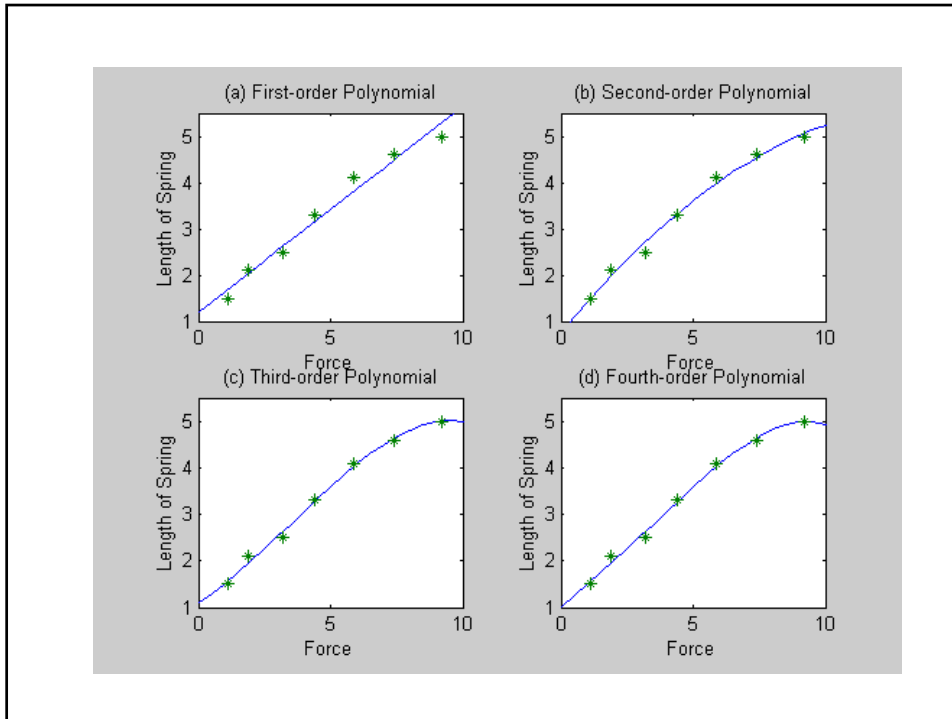
- ◆ We rely on this estimation because we have more data
- ◆ If we are not happy with the LSE estimators then we can increase the model's degree of freedom such that:

$$L = k_0 + k_1 f + k_2 f^2 + \dots + k_n f^n$$

(least square polynomial!)

Least-Squares Estimators (5.3) (cont.)

- ◆ Higher order models fit better the data but they do not always reflect the inner law that governs the system
- ◆ For example, when f is increasing toward 10N, the length is decreasing!



Statistical Properties & the Maximum Likelihood Estimator (5.7)

◆ Statistical qualities of LSE

– Definition [unbiased estimator]

An estimator $\hat{\theta}$ of the parameter θ is unbiased if $E(\hat{\theta}) = \theta$, where $E[\cdot]$ is the statistical expectation

– Definition [minimal variance]

An estimator $\hat{\theta}$ is a minimum variance estimator if for any other estimator θ^* : $\text{cov}(\hat{\theta}) \leq \text{cov}(\theta^*)$, where $\text{cov}(\theta)$ is the covariance matrix of the random vector θ .

Statistical Properties & the Maximum Likelihood Estimator (5.7) (cont.)

◆ Statistical qualities of LSE (cont.)

– Gauss-Markov conditions:

- The error vector e is a vector of m uncorrelated random variables, each with zero mean & the same variance σ^2 . This means that:

$$E[e] = 0 \text{ \& } E[ee^T] = \sigma^2 I$$

$$\left[E \begin{pmatrix} e_1 \\ e_2 \end{pmatrix} (e_1, e_2) = \begin{pmatrix} E(e_1^2) & E(e_1 e_2) \\ E(e_2 e_1) & E(e_2^2) \end{pmatrix} = \sigma^2 I \right]$$

Statistical Properties & the Maximum Likelihood Estimator (5.7) (cont.)

◆ *Theorem [Gauss-Markov]:*

LSE is unbiased & has minimum variance.

$$\begin{aligned} \text{Proof: } E[\hat{\theta}] &= E[(A^T A)^{-1} A^T y] = (A^T A)^{-1} A^T E[y] \\ &= (A^T A)^{-1} A^T A \theta = \theta \end{aligned}$$

since $E[y] = E[A\theta] + E(e) = A\theta$ and

$$(A^T A)^{-1} = A^{-1} (A^T)^{-1}$$

Statistical Properties & the Maximum Likelihood Estimator (5.7) (cont.)

◆ Maximum likelihood (ML) estimator

- ML is one of the most widely used technique for parameter estimation of a statistical distribution

- *ML definition:*

For a sample of n observations (of a probability density function) x_1, x_2, \dots, x_n , the likelihood function L is defined by:

$$L = f(x_1; \theta) f(x_2; \theta) \dots f(x_n; \theta)$$

Statistical Properties & the Maximum Likelihood Estimator (5.7) (cont.)

◆ Maximum likelihood (ML) estimator (cont.)

- The criterion for choosing θ is:
“pick a value of θ that provides a high probability of obtaining the actual observed data x_1, x_2, \dots, x_n ”

- Therefore, ML estimator $\hat{\theta}$ is defined as the value of θ which maximizes L : $\frac{\partial L}{\partial \theta}(\hat{\theta}) = 0$

or equivalently: $\frac{\partial \ln L}{\partial \theta}(\hat{\theta}) = 0$

Statistical Properties & the Maximum Likelihood Estimator (5.7) (cont.)

◆ Maximum likelihood (ML) estimator (cont.)

– Example 1: ML estimation for exponential distribution

$$f(x; \theta) = \theta^{-1} e^{-x/\theta} \quad L = (\theta^{-1} e^{-x_1/\theta}) (\theta^{-1} e^{-x_2/\theta}) \dots (\theta^{-1} e^{-x_m/\theta}) \\ = \theta^{-m} \exp(-\theta^{-1} \sum x_i)$$

is the likelihood function for m observations x_1, x_2, \dots, x_m .

$$\ln L = -m \ln \theta - \frac{1}{\theta} \sum_{i=1}^m x_i \quad \frac{\partial \ln L}{\partial \theta} = -\frac{m}{\theta} + \frac{\sum x_i}{\theta^2} = 0$$

$$\text{Therefore: } \hat{\theta} = \frac{\sum x_i}{m}$$

Statistical Properties & the Maximum Likelihood Estimator (5.7) (cont.)

◆ Maximum likelihood (ML) estimator (cont.)

– Example 2: ML estimation for normal distribution

$$f(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2\right]$$

where μ and σ^2 are respectively the mean & the variance unknown parameters.

For m observations x_1, x_2, \dots, x_m , we have:

Statistical Properties & the Maximum Likelihood Estimator (5.7) (cont.)

- Example 2 (cont.)

$$L = \left(\frac{1}{\sqrt{2\pi}\sigma} \right)^m \exp \left[-\frac{1}{2\sigma^2} (x_i - \mu)^2 \right]$$

$$\text{and } \ln L = -m \ln(\sqrt{2\pi}\sigma) - \frac{1}{2\sigma^2} \sum (x_i - \mu)^2$$

$$\text{Therefore : } \frac{\partial \ln L}{\partial \mu} = 0 \Rightarrow \frac{1}{\sigma^2} \sum (x_i - \mu) = 0 \Rightarrow \hat{\mu} = \frac{\sum x_i}{m}$$

$$\frac{\partial \ln L}{\partial \sigma} = 0 \Rightarrow -\frac{m}{\sigma} - \frac{1}{2} (-2\sigma^{-3}) \sum (x_i - \mu)^2 = 0$$

$$\Rightarrow \hat{\sigma}^2 = \frac{\sum (x_i - \hat{\mu})^2}{m}$$

Statistical Properties & the Maximum Likelihood Estimator (5.7) (cont.)

◆ Maximum likelihood (ML) estimator (cont.)

- Equivalence between LSE & MLE
- Theorem: Under the Gauss conditions & if each component of the vector e follows a normal distribution then the LSE of θ = MLE of θ

LSE for Nonlinear Models (5.8)

◆ Nonlinear models are divided into 2 families

- Intrinsically linear
- Intrinsically nonlinear
 - Through appropriate transformations of the input-output variables & fitting parameters, an intrinsically linear model can become a linear model
 - By this transformation into linear models, LSE can be used to optimize the unknown parameters

LSE for Nonlinear Models (5.8) (cont.)

◆ Example: Decay of the radioactivity of chemicals

$y = ae^{bt}$ ($a, b < 0$ are to be determined)

(t, y) are the input-output data

given a training data $\{(t_i, y_i), i = 1, \dots, m\}$

$$E(\mathbf{a}, \mathbf{b}) = \sum_{i=1}^m (y_i - ae^{bt_i})^2$$

minimizing the error yields:

$$\begin{cases} \frac{\partial E}{\partial a} = 2 \sum_{i=1}^m (y_i - ae^{-bt_i}) (-e^{bt_i}) = 0 \\ \frac{\partial E}{\partial b} = 2 \sum_{i=1}^m (y_i - ae^{-bt_i}) (-at_i e^{-bt_i}) = 0 \end{cases}$$

we obtain 2 nonlinear equations with respect to a & b !
Therefore, we are better off to linearize the model before.

LSE for Nonlinear Models (5.8) (cont.)

$$\ln y = \ln a + bt$$

with $\ln a = a'$, $\ln y = y'$ and $b = b'$

$$\Rightarrow y' = a' + b't$$

which is an intrinsically linear model. The linear & nonlinear models are close but can sometimes differ significantly!

The following table shows a family of intrinsically linear models

LSE for Nonlinear Models (5.8) (cont.)

Nonlinear models	Transformation	Linear forms
$y = ae^{bx}$	Natural logarithm	$\ln y = \ln a + bx$
$y = ax^b$	Natural logarithm	$\ln y = \ln a + b \ln x$
$y = \frac{ax}{b+x}$	Reciprocal	$\frac{1}{y} = \frac{1}{a} + \frac{b}{a} \frac{1}{x}$
$y = \frac{a}{b+x}$	Reciprocal	$\frac{1}{y} = \frac{1}{a} + \frac{b}{a} x$

LSE for Nonlinear Models (5.8) (cont.)

Suppose that the underlying model is: $y = ae^{bt}$

t_i	0	0.80	1.84	2.90	4.06	4.81	6.07	7.06	8.15	8.87	9.98
y_i	0.98	0.69	0.47	0.46	0.29	0.16	0.23	0.10	0.03	0.12	0.01

LSE for Nonlinear Models (5.8) (cont.)

