# Advanced Control Systems
# Detection, Estimation, and Filtering

## Graduate Course on the
## MEng PhD Program
## Spring 2012/2013

**Instructor:**
**Prof. Paulo Jorge Oliveira**
**p.oliveira@dem.ist.utl.pt or pjcro @ isr.ist.utl.pt**
**Phone: 21 8419511 or 21 8418053 (3511 or 2053 inside IST)**

# *Objectives:*

- Motivation for estimation, detection, filtering, and identification in stochastic signal processing

- Methodologies on design and synthesis of optimal estimation algorithms

- Characterization of estimators and tools to study their performance

- To provide an overview in all principal estimation approaches and the rationale for choosing a particular technique

**Both for parameter and state estimation,**

**always on the presence of stochastic disturbances**

In RADAR (Radio Detection and Ranging), SONAR (sound navigation and ranging), speech, image, sensor networks, geo-physical sciences,…

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Pre-requisites:*

• Random Variables and Stochastic Processes

Joint, marginal, and conditional probability density functions: Gaussian / normal distributions; Moments of random variables (mean and variance); Wide-sense stationary processes; Correlation and covariance; Power spectral density;

• Linear Algebra

Vectors: orthogonality, linear independence, inner product; norms;

Matrices: eigenvectors, rank, inverse, and pseudo-inverse;

• Linear Systems

LTIS and LTVs; ODEs and solutions; Response of linear systems; Transition matrix; Observability and controlability; Lyapunov stability.

The implementation of solutions for problems require the use of *MATLAB* and *Simulink*.

# *Syllabus:*

## Classical Estimation Theory

Chap. 1 - ***Motivation for Estimation in Stochastic Signal Processing*** [1/2 week]

Motivating examples of signals and systems in detection and estimation problems;

Chap. 2 - ***Minimum Variance Unbiased Estimation*** [1 week]

Unbiased estimators; Minimum Variance Criterion; Extension to vector parameters; Efficiency of estimators;

Chap. 3 - ***Cramer-Rao Lower Bound*** [1 week]
Estimator accuracy; Cramer-Rao lower bound (CRLB); CRLB for signals in white Gaussian noise;  Examples;

continues…

# Syllabus (cont.):

Chap. 4 - *Linear Models in the Presence of Stochastic Signals* [1/2 week]

Stationary and transient analysis; White Gaussian noise and linear systems;  Examples; Sufficient Statistics; Relation with MVU Estimators;

Chap. 5 - *Best Linear Unbiased Estimators* [1 week]

Definition of BLUE estimators;   White Gaussian noise and bandlimited systems; Examples; Generalized minimum variance unbiased estimation;

Chap. 6 - *Maximum Likelihood Estimation* [1 week]

The maximum likelihood estimator; Properties of the ML estimators; Solution for ML estimation; Examples; Monte-Carlo methods;

Chap. 7 - *Least Squares* [1 week]

The least squares approach; Linear and nonlinear least squares; Geometric interpretation; Constrained least squares;  Examples;

continues…

PO 1213

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Syllabus (cont.):*

Bayesian Estimation Theory

Chap. 8 – *Bayesian Estimation* [1/2 week]

Philosophy and estimator design; Prior knowledge; Bayesian linear model; Bayesian estimation on the presence of Gaussian pdfs; Minimum Mean Square Estimators;

Chap. 9 – *Wiener Filtering* [1/2 week]

The Wiener filter problem; Causal and non-causal solutions; Complementary filters;

Chap. 10 – *Kalman Filtering* [2 weeks]

Optimal estimator  in the presence of white Gaussian noise – the Kalman filter; Stability, convergence and robustness for LTV and LTI systems; Kalman and Wiener filters; Optimal smoothers; Examples; Extended Kalman Filters;

continues…

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Syllabus (cont.):*

Advanced Estimation Topics

Chap. 11 – *Multiple Model Adaptive Estimation* [1 week]

Joint system identification and parameter/state estimation using multiple models.

Chap. 12 – *Optimal Smoothing* [1 week]

Fixed point, fixed interval, and fixed lag smoothers.

Chap. 13 – *Advanced Topics* [2 weeks]

To be detailed later, e.g. Positioning and navigation systems; Failure detection and isolation; Multiple model adaptive estimation; Discretization; Missing data estimation; Outlier detection and removal; Feature based estimation; Principal component analysis; Nonlinear signal processing; Compressive sensing;…

End.

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# Grading:

- **Five problem sets** (50%)

  Due dates: weeks of 25-03, 08-04, 22-04, 06-05, 06-20 and 03-06 (tentative).

and

- **Term paper** (50%)

  Topic selected randomly by the student. Worked jointly by the faculty/student.
  To be completed in the final 3-4 weeks, i.e.  week of 5-07.

or

- Final exam (50%)

  Week of 15-07.

## Classes:

| | |
|---|---|
| Tuesdays: | 16h00 – 17h30, room C11 |
| Thurdays: | 16h00 – 17h30, room P9 |

**To discuss issues: i) e-mail, ii) phone, or iii) schedule an interview.**

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Bibliography:*

**Main references**

- Steven M. Kay, **Fundamentals of Statistical Signal Processing: Estimation Theory, Vol. I**, Prentice Hall Signal Processing Series, 1993.

- A. Gelb, **Applied Optimal Estimation,** MIT Press, 1974.

**Complementary reading**

- Steven M. Kay, **Fundamentals of Statistical Signal Processing: Detection Theory, Vol. II**, Prentice Hall Signal Processing Series, 1998.

- Harry L. Van Trees, **Detection, Estimation, and Modulation Theory, Parts I to IV,** John Wiley, 2001.

- Athanasios Papoulis and S. Unnikrishna Pillai, **Probability, Random Variables and Stochastic Processes,** McGraw Hill, 2001.

- Robert Brown and Patrick Hwang, **Introduction to Random Signals and Applied Kalman Filtering,** John Wiley, 1997.

- Gonzalo Arce, **Nonlinear Signal Processing: A Statistical Approach,** John Wiley, 2005.

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Advanced Control Systems Detection, Estimation, and Filtering*

*Graduate Course on the*
*MEng PhD Program*
*Spring 2012/2013*

## Chapter 1
## Motivation

*Instructor:*
*Prof. Paulo Jorge Oliveira*
*p.oliveira@dem.ist.utl.pt or pjcro @ isr.ist.utl.pt*
*Phone: 21 8419511 or 21 8418053 (3511 or 2053 inside IST)*

**DEM**
DEPARTAMENTO
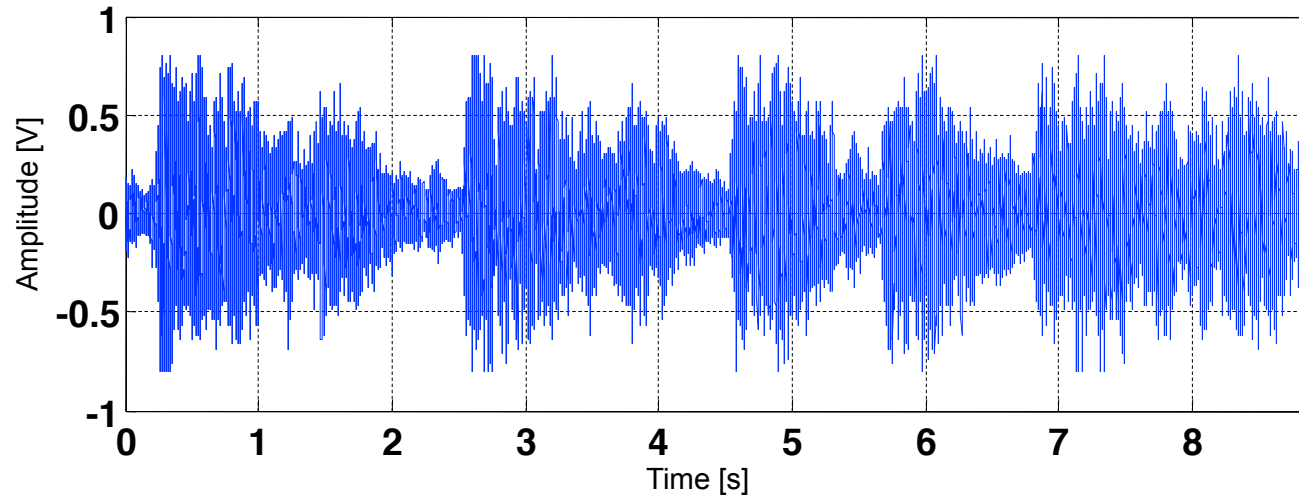DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Summary:*

- Motivation for estimation, detection, filtering, and identification in stochastic signal processing

- Methodologies on how to design optimal estimation algorithms

- Characterization of estimators and tools to study their performance

- To provide an overview in all principal estimation approaches and the rationale for choosing a particular technique

**Both for parameter and state estimation,**

**always on the presence of stochastic disturbances**

In RADAR (Radio Detection and Ranging), SONAR(sound navigation and ranging), speech, image, sensor networks, geo-physical sciences,…

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Speech*



Signals can be represented by functions (continuous time) or by vectors (where a sampling operation takes place)

Examples of speech/sound processing:

Automatic systems commanded by voice;

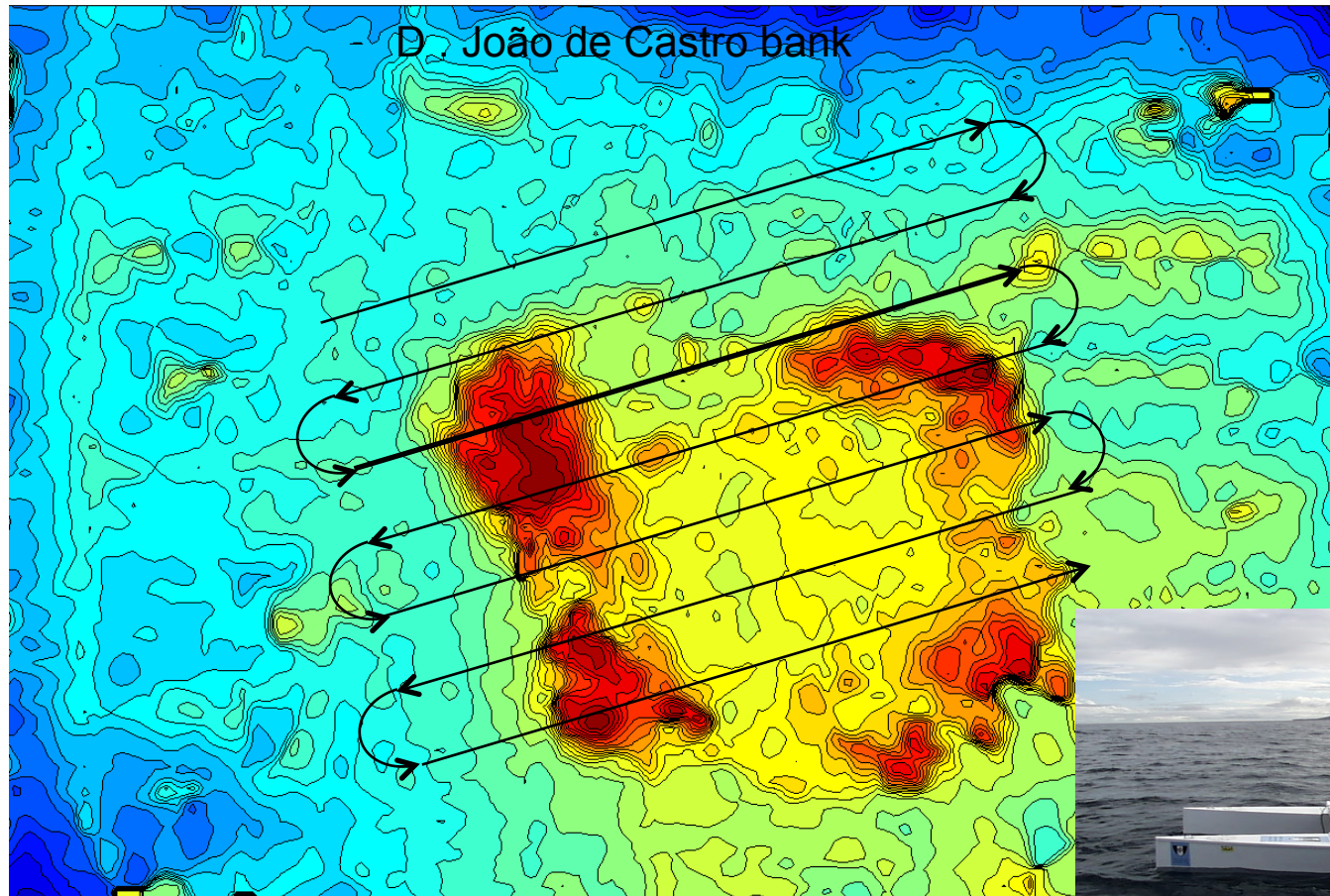Automatic translation; Voice recognition

Synthesis of voice

PO 1213

# *Echograms*

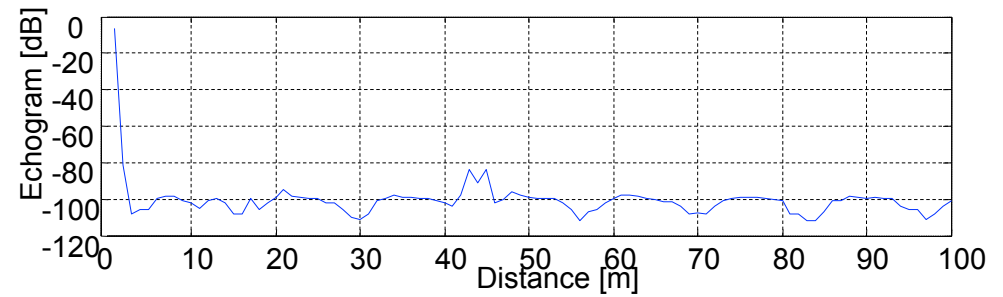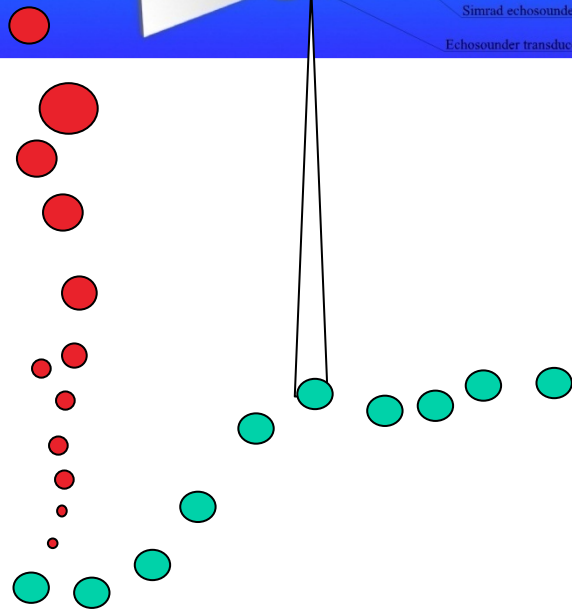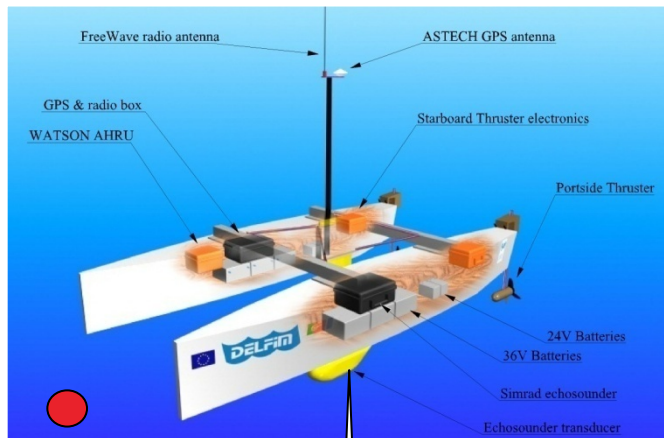**The quest for hydrothermal vents**

D. João de Castro bank

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Echograms*

**The quest for hydrothermal vents (cont.)**



D . João de Castro bank
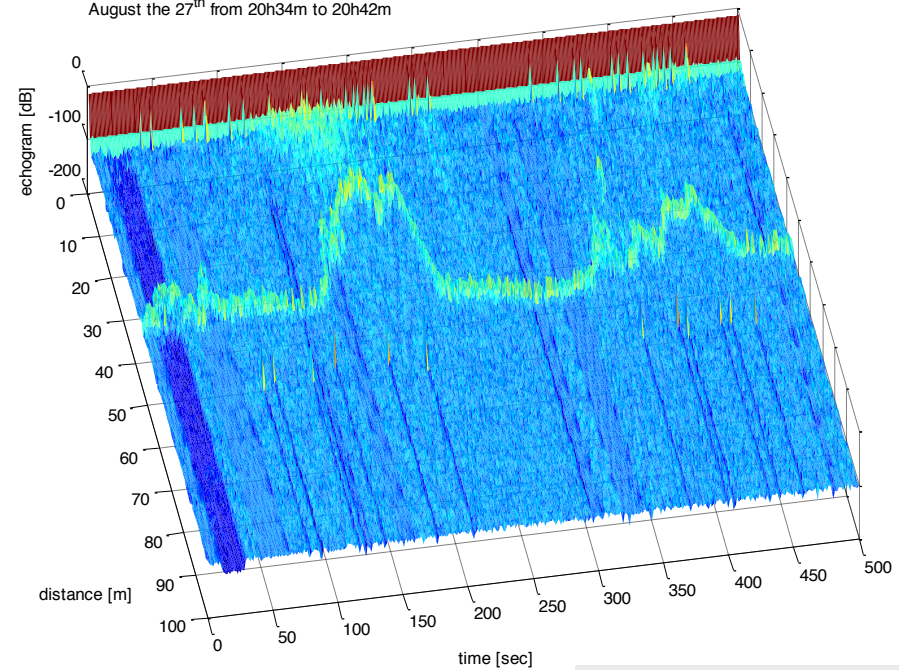
DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Echograms*

## *The quest for hydrothermal vents (cont.)*





Delfim at D. João de Castro Bank

August the 27th from 20h34m to 20h42m

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# Sidescan Sonar Imaging

# GPS – Global Positioning System

Latitude [m]                    σ = 0.00501



PO 1213

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Image with missing data*

## *Reconstruction*



η=0.1 η=0.2 η=0.3 η=0.4 η=0.5 η=0.6 η=0.7

PO 1213

# *Bathymetric survey*



Geo-referenced data

Reconstructed

Uncertainty on data

Final uncertainty

PO 1213

# *Deblurring an image*

Original Image

Blurred Image

Restored image

Causes:

• Out of focus acquisition

• Camera-object movement

• Shaking

• Shallow field of view …

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Stock Exchange*

***Models that explain evolution of phenomena***

- Causality

- Number of parameters

- Type of model

- Uncertainty

Is it possible to predict the

market price tomorrow,

next week, next month,

next year,?…

# GPS Intelligent Buoys(GIB)-ACSA/ORCA

**Tracking with a Sensor Network**



**Surface buoys with**

- *DGPS receivers*
- *Hydrophones*
- *Radio link*

**Control Station**

- *DGPS receiver*
- *Radio link*
- *PC with tracking software*

**Acoustic pinger**

PO 1213

**DEM**
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

**Tracking with a Sensor Network (cont.)**



PO 1213

# *The mathematical estimation problem*

**To be possible to design estimators, first the data must be modeled.**

*Example I:*

Assume that one sample *x* is available

(scalar example, i.e. N=1) with constant

**unknown** mean *θ*.

The probability density function (PDF) is

$$p(x;\theta) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\theta)^2}{2\sigma^2}}$$



For instance if x[0]<0 it is doubtful that the unknown parameter is >>0.

*In a actual problem, we are not given a PDF, but must be chosen to be*

*consistent with the data and with the prior knowledge.*

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

*Example II:*

Now the following sequence **x** is given.

Note that the value along time appears

to be decreasing. Lets consider that

the phenomena is described by



$$x[n] = A + Bn + w[n] \qquad n = 0, 1, ..., N-1$$

where A and B are constant unknown parameters and w[n] is assumed to be white Gaussian noise, with PDF $N(0, \sigma^2)$. For $\theta = [A\ B]$ and **x**=[x[0] x[1] …x[n]] *the data PDF is*

$$p(\mathbf{x} \mid \theta) = \frac{1}{(2\pi\sigma^2)^{\frac{N}{2}}} e^{-\frac{1}{2\sigma^2} \sum_{n=0}^{N-1} (x[n] - A - Bn)^2}$$

Where the uncertainty in the samples

is assumed to be uncorrelated.

**The performance of the estimators is dependent on the models used, so they must be mathematically treatable.**.

PO 1213

# *The mathematical estimation problem*

**Classical estimation techniques**

      **Parameters are assumed deterministic but unknown**

**Bayesian techniques**

      **Parameters are used to be unknown but are stochastic also described by a PDF.**

The joint PDF would then be         $p(\mathbf{x}, \theta) = p(\mathbf{x} \mid \theta) \, p(\theta)$

Dependence of data
on the parameters

Prior knowledge

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Exploiting simple estimators*

*Example III (Quiz):*

*Given a data sequence from a signal with*

*PDF as described by one of three models*

**Which one is the correct model?**



*First scenario:*

$$N = 100$$

$$\theta \in \{-40, 0, 40\}$$

$$\sigma^2 = 10^2$$

*For the signal*

$$x[n] = \theta + w[n] \qquad n = 0, 1, ..., N - 1$$

*The answer is obvious:*

*θ = 40!*

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Exploiting simple estimators*

Second scenario

(lousy sensor quality or lousy data):

Lets *repeat the problem with*

$$\sigma^2 = 100^2$$

$$N = 100$$
$$\theta \in \{-40, 0, 40\}$$

*The answer is not obvious anymore!*

*Lets propose a couple of estimators*

*and to study them…*

$$\hat{\theta}_1 = x[0]$$

$$\hat{\theta}_2 = \frac{1}{N}\sum_{n=0}^{N-1} x[n]$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Assessing estimator performance*

Estimators depend only on observed data thus can be viewed as a function

$$\theta = g\left(x[0], \quad \dots \quad x[N]\right) = g\left(x\right)$$

The study of estimator properties must be done resorting to statistic tools.

Is it exact?, i.e. Does it return the true value of the unknown parameters?

Is this a good estimator?  If many experiments can be performed, is it expected that the unknown parameter is achievable? Or are the results expected to be biased?

$$E\left[\hat{\theta}_1\right] = E\left[x[0]\right] = \theta$$

$$E\left[\hat{\theta}_2\right] = E\left[\frac{1}{N}\sum_{n=0}^{N-1}x[n]\right] =$$

$$= \frac{1}{N}\sum_{n=0}^{N-1}E\left[x[n]\right] =$$

$$= \frac{1}{N}N\theta = \theta$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Assessing estimator performance*

How good is an estimator? How much uncertainty corresponds to the computed value?



The use of computational tools is a good idea? No!

Formal methods are required



For our quiz:

$$\mathrm{var}\left(\hat{\theta}_1\right) = \mathrm{var}\left(x[0]\right) = \sigma^2$$

$$\mathrm{var}\left(\hat{\theta}_2\right) = \mathrm{var}\left(\frac{1}{N}\sum_{n=0}^{N-1}x[n]\right) = \frac{1}{N^2}\sum_{n=0}^{N-1}\mathrm{var}\left(x[n]\right) == \frac{1}{N^2}N\sigma^2 = \frac{\sigma^2}{N}!$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Assessing estimator performance*

Questions triggered from this simple example but valid to all our problems:

• **The second estimator is much better than the first estimator.**

• **The quality of the estimate increases with the number of points. Is it reasonable? Is it plausible?**

• **Do we have always data available? How to get data?**

• **Is this the best one can do with N samples?**

• **Are there better estimators that we can exploit?**

*Answers to this questions will be provided along the course…*

# *Bibliography:*

**Further reading**

• Paul Etter, *Underwater Acoustic Modelling and Simulation,* Taylor & Francis, 2003.

• François Le Chevalier, *Principles of Radar and Sonar Signal Processing*, Artech House, 2002.

• *Peter Wille, Sound Images of the Ocean in Research and Monitoring,* Springer, 2005.

• Lawrence Rabiner, Biing Juang, *Fundamentals of Speech Recognition,* Prentice Hall, 1993.

• Gilbert Strang, Kai Borre, *Linear Algebra, Geodesy, and GPS,* SIAM, 1997.

• Rafael Gonzales, Richard Woods, *Digital Image Processing,* Prentice Hall, 2001.

• Joseph Boccuzzi, *Signal Processing for Wireless Communications,* McGraw Hill, 2008.

• Venkatesh Saligrama, *Networked Sensing Information and Control,* Springer, 2008.

• Ching-Fang Lin, *Modern Navigation, Guidance, and Control Processing*, Prentice Hall, 1991.

See for instance http://www.**ieee**.org/portal/site

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Advanced Control Systems*
# *Detection, Estimation, and Filtering*

# *Chapter 2*
# *Minimum Variance Unbiased Estimation*

*Instructor:*
*Prof. Paulo Jorge Oliveira*
*p.oliveira@dem.ist.utl.pt or pjcro @ isr.ist.utl.pt*
*Phone: 21 8419511 or 21 8418053 (3511 or 2053 inside IST)*

**DEM**
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Syllabus:*

## Classical Estimation Theory

Chap. 1 - **Motivation for Estimation in Stochastic Signal Processing** [1/2 week] Motivating examples of signals and systems in detection and estimation problems;

## Chap. 2 - *Minimum Variance Unbiased Estimation* [1/2 week]

Unbiased estimators; Minimum Variance Criterion; Extension to vector parameters; Efficiency of estimators;

Chap. 3 - **Cramer-Rao Lower Bound** [1 week]
Estimator accuracy; Cramer-Rao lower bound (CRLB); CRLB for signals in white Gaussian noise; Examples;

continues…

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Unbiased estimators:*

The search for good estimators for unknown deterministic parameters begins

Example (revisited):

Unbiased estimator for DC level in white Gaussian noise. Signal model is

$$x[n] = A + w[n] \qquad n = 0,1,...,N-1$$

A reasonable estimator is $\quad \hat{A} = k\frac{1}{N}\sum_{n=0}^{N-1} x[n]$

Due to the linearity properties of the expectation operator *E[.]*:

$$E\left[\hat{A}\right] = E\left[k\frac{1}{N}\sum_{n=0}^{N-1} x[n]\right] = k\frac{1}{N}\sum_{n=0}^{N-1} E[x[n]] = k\frac{1}{N}NA = kA$$

### *Unbiased estimator iff k=1!*

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Unbiasness:*

Let the vector of deterministic unknown parameters, with p components, be described as

$$\theta = \begin{bmatrix} \theta_1 \\ \theta_2 \\ M \\ \theta_p \end{bmatrix} = \begin{bmatrix} \theta_1 & \theta_2 & ... & \theta_p \end{bmatrix}^T$$

An estimator must have the same dimensions, i.e.

$$\hat{\theta} = \begin{bmatrix} \hat{\theta}_1 \\ \hat{\theta}_2 \\ M \\ \hat{\theta}_p \end{bmatrix} = \begin{bmatrix} \hat{\theta}_1 & \hat{\theta}_2 & ... & \hat{\theta}_p \end{bmatrix}^T$$

The estimator is **unbiased** iff

$$E\left[\hat{\theta}\right] = \theta$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Minimum variance criterion:*

In searching for estimators some optimality criterion must be adopted.

A natural one is the **mean square error (MSE),** defined as

$$mse\left(\hat{\theta}\right) = E\left[\left(\hat{\theta} - \theta\right)^2\right]$$

Unfortunately the choice of the criterion leads to unrealizable estimators, i.e. not only function of the data

$$mse\left(\hat{\theta}\right) = E\left\{\left[\left(\hat{\theta} - E\left[\hat{\theta}\right]\right) + \left(E\left[\hat{\theta}\right] - \theta\right)\right]^2\right\} = \mathrm{var}\left(\hat{\theta}\right) + \left[\left(E\left[\hat{\theta}\right] - \theta\right)\right]^2 =$$

$$= \mathrm{var}\left(\hat{\theta}\right) + b^2\left(\theta\right), \qquad where \quad b^2\left(\theta\right) = E\left[\hat{\theta}\right] - \theta.$$

**Fortunately, after differentiation an estimator depending on θ will result.**

# *Minimum variance criterion:*

Example:

Find the value of k such that a realizable mse estimator results

$$x[n] = A + w[n] \quad n = 0,1,...,N-1, \qquad \hat{A} = k\frac{1}{N}\sum_{n=0}^{N-1} x[n]$$

From the previous page

$$\operatorname{var}\left(\hat{A}\right) = \frac{k^2\sigma^2}{N}, \quad and \quad b^2(A) = E\left[\hat{A}\right] - A = (kA - A)^2$$

$$mse\left(\hat{A}\right) = \operatorname{var}\left(\hat{A}\right) + b^2(A) = \frac{k^2\sigma^2}{N} + (k-1)^2 A^2$$

Lets find the minimum

$$\frac{d}{dk}mse\left(\hat{A}\right) = \frac{2k\sigma^2}{N} + 2(k-1)A^2 = 0, \quad results \quad in \quad k_{opt=}\frac{A^2}{A^2 + \sigma^2/N}$$

Unfortunately depends on the unknown parameter A.

# Minimum variance unbiased estimator:

The Minimum Variance Unbiased (MVU) Estimator **must have** smallest variance for all values of $\theta$.

In general, the MVU estimator does not always exist.

There is no "turn-the-crank" method.

Future approaches:

Chp. 3 – Cramer-Rao lower bound

Chp. 5 – Sufficient statistics

Chp. 6 – Restrict to linear estimators: BLUE

# *Bibliography:*

**Further reading**

• Harry L. Van Trees, *Detection, Estimation, and Modulation Theory, Parts I to IV,* John Wiley, 2001.

• J. Bibby, H. Toutenburg, *Prediction and Improved Estimation in Linear Models,* John Wiley, 1977.

• C.Rao, *Linear Statistical Inference and Its Applications,* John Wiley, 1973.

• P. Stoica, R. Moses, *"On Biased Estimators and the Unbiased Cramer-Rao Lower Bound,"* *Signal Process,* vol.21, pp. 349-350, 1990.

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Advanced Control Systems*
# *Detection, Estimation, and Filtering*

*Graduate Course on the*
*MEng PhD Program*
*Spring 2012/2013*

## *Chapter 3*
## *Cramer-Rao Lower Bounds*

*Instructor:*
*Prof. Paulo Jorge Oliveira*
*p.oliveira@dem.ist.utl.pt or pjcro @ isr.ist.utl.pt*
*Phone: 21 8419511 or 21 8418053 (3511 or 2053 inside IST)*

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# Syllabus:

Classical Estimation Theory

...

Chap. 2 - *Minimum Variance Unbiased Estimation* [1/2 week]

Unbiased estimators; Minimum Variance Criterion; Extension to vector parameters; Efficiency of estimators;

Chap. 3 - *Cramer-Rao Lower Bound* [1 week]
Estimator accuracy; Cramer-Rao lower bound (CRLB); CRLB for signals in white Gaussian noise;  Examples;

Chap. 4 - *Linear Models in the Presence of Stochastic Signals* [1 week]

Stationary and transient analysis; White Gaussian noise and linear systems;  Examples; Sufficient Statistics; Relation with MVU Estimators;

continues…

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Estimator accuracy:*

The accuracy on the estimates dependents very much on the PDFs

Example (revisited):

Model of signal $x[0] = A + w[0]$,

Observation PDF $p(x[0]; A) = \dfrac{1}{\sqrt{2\pi\sigma^2}} e^{-\dfrac{(x[0]-A)^2}{2\sigma^2}}$

for a disturbance $N(0, \sigma^2)$



*Remarks:*

If $\sigma^2$ is **Large** then the performance of the estimator is **Poor**;

If $\sigma^2$ is **Small** then the performance of the estimator is **Good**; or

If PDF concentration is **High** then the parameter accuracy is **High.**

**How to measure sharpness of PDF (or concentration)?**

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# Estimator accuracy:

When PDFs are seen as function of the unknown parameters, for x fixed, they are called as **Likelihood function.** To measure the sharpness note that (and ln is monotone…)

$$\ln p\left(x[0]; A\right) = -\ln \sqrt{2\pi\sigma^2} - \frac{1}{2\sigma^2}\left(x[0] - A\right)^2$$

Its first and second derivatives are respectively:

$$\frac{\partial}{\partial A}\ln p\left(x[0]; A\right) = \frac{1}{\sigma^2}\left(x[0] - A\right) \qquad \text{and} \qquad -\frac{\partial^2}{\partial A^2}\ln p\left(x[0]; A\right) = \frac{1}{\sigma^2}.$$

As we know that the estimator $\hat{A}$ has variance $\sigma^2$ (at least for this example)

$$\mathrm{var}\left(\hat{A}\right) = \frac{1}{-\dfrac{\partial^2}{\partial A^2}\ln p\left(x; A\right)} = \frac{1}{curvature}$$

We are now ready to present an important theorem…

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Cramer-Rao lower bound:*

**Theorem 3.1 (Cramer-Rao lower bound, scalar parameter)** – *It is assumed that the PDF p(**x**; θ) satisfies the "regularity" condition*

$$E\left[\frac{\partial}{\partial\theta}\ln p\left(\mathbf{x};\theta\right)\right] = 0 \qquad \text{for all} \quad \theta \tag{1}$$

*where the expectation is taken with respect to p(**x**; θ). Then, the variance of any unbiased estimator $\hat{\theta}$ must satisfy*

$$\operatorname{var}\left(\hat{\theta}\right) \geq \frac{1}{-E\left[\dfrac{\partial^2}{\partial\theta^2}\ln p\left(\mathbf{x};\theta\right)\right]} \tag{2}$$

*where the derivative is evaluated at the true value of θ and the expectation is taken with respect to p(**x**,  θ). Furthermore, an unbiased estimator can be found that attains the bound for all θ if and only if*

$$\frac{\partial}{\partial\theta}\ln p\left(\mathbf{x};\theta\right) = I\left(\theta\right)\left(g\left(\mathbf{x}\right) - \theta\right) \tag{3}$$

*for some functions g(.) and $I$(.). The estimator, which is the MVU estimator, is $\hat{\theta} = g\left(\mathbf{x}\right)$, and the minimum variance 1/ $I$(θ).*

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Cramer-Rao lower bound:*

**Proof outline:**

Lets derive the CRLB for a scalar parameter $\alpha=g(\theta)$. We consider all unbiased estimators

$$E\left[\hat{\alpha}\right] = \alpha = g\left(\theta\right) \qquad \text{or} \qquad \int \hat{\alpha} p\left(\boldsymbol{x};\theta\right) \, d\boldsymbol{x} = g\left(\theta\right).$$

(p.1)

Lets examine the regularity condition (1)

$$E\left[\frac{\partial}{\partial\theta}\ln p\left(\boldsymbol{x};\theta\right)\right] = \int \frac{\partial \ln p\left(\boldsymbol{x};\theta\right)}{\partial\theta} p\left(\boldsymbol{x};\theta\right) d\boldsymbol{x} = \int \frac{\partial p\left(\boldsymbol{x};\theta\right)}{\partial\theta} d\boldsymbol{x}$$

$$= \frac{\partial}{\partial\theta} \int p\left(\boldsymbol{x};\theta\right) d\boldsymbol{x} = \frac{\partial 1}{\partial\theta} = 0.$$

Remark: differentiation and integration are required to be interchangeable (Leibniz Rule)!

Lets differentiate (p.1) with respect to $\theta$ and use the previous results

$$\int \hat{\alpha} \frac{\partial p\left(\boldsymbol{x};\theta\right)}{\partial\theta} d\boldsymbol{x} = \frac{\partial g\left(\theta\right)}{\partial\theta} \qquad or \qquad \int \hat{\alpha} \frac{\partial \ln p\left(\boldsymbol{x};\theta\right)}{\partial\theta} p\left(\boldsymbol{x};\theta\right) d\boldsymbol{x} = \frac{\partial g\left(\theta\right)}{\partial\theta} \quad .$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Cramer-Rao lower bound:*

**Proof outline (cont.):**

This can be modified to

$$\int (\alpha - \hat{\alpha}) \frac{\partial \ln p(\boldsymbol{x};\theta)}{\partial \theta} p(\boldsymbol{x};\theta) d\boldsymbol{x} = \frac{\partial g(\theta)}{\partial \theta},$$

as

$$\int \alpha \frac{\partial \ln p(\boldsymbol{x};\theta)}{\partial \theta} p(\boldsymbol{x};\theta) d\boldsymbol{x} = \alpha E\left[\frac{\partial \ln p(\boldsymbol{x};\theta)}{\partial \theta}\right] = 0.$$

Now applying the Cauchy-Schwarz inequality, i.e.

$$\left[\int w(\boldsymbol{x}) g(\boldsymbol{x}) h(\boldsymbol{x}) d\boldsymbol{x}\right]^2 \leq \int w(\boldsymbol{x}) g^2(\boldsymbol{x}) d\boldsymbol{x} \int w(\boldsymbol{x}) h^2(\boldsymbol{x}) d\boldsymbol{x}$$

considering $\quad w(\boldsymbol{x}) = p(\boldsymbol{x};\theta), \quad g(\boldsymbol{x}) = \hat{\alpha} - \alpha, \quad$ and $\quad h(\boldsymbol{x}) = \frac{\partial \ln p(\boldsymbol{x};\theta)}{\partial \theta}$

results

$$\left(\frac{\partial g(\theta)}{\partial \theta}\right)^2 \leq \int (\alpha - \hat{\alpha})^2 p(\boldsymbol{x};\theta) d\mathbf{x} \int \left(\frac{\partial \ln p(\boldsymbol{x};\theta)}{\partial \theta}\right)^2 p(\boldsymbol{x};\theta) d\mathbf{x}$$

# *Cramer-Rao lower bound:*

**Proof outline (cont.):**

It remains to relate this expression with the one in the Theorem $\int \left( \dfrac{\partial \ln p(x;\theta)}{\partial \theta} \right)^2 p(x;\theta)dx = ?$

Starting with the previous result

$$E\left[ \frac{\partial}{\partial \theta} \ln p(\boldsymbol{X};\theta) \right] = \int \frac{\partial}{\partial \theta} \ln p(\boldsymbol{X};\theta) p(\boldsymbol{X};\theta) d\boldsymbol{x} = 0$$

thus, this function identically null verifies

$$\frac{\partial}{\partial \theta} \int \frac{\partial}{\partial \theta} \ln p(\boldsymbol{X};\theta) p(\boldsymbol{X};\theta) d\boldsymbol{x} = \int \left[ \frac{\partial^2 \ln p(\boldsymbol{X};\theta)}{\partial \theta^2} p(\boldsymbol{X};\theta) + \frac{\partial \ln p(\boldsymbol{X};\theta)}{\partial \theta} \frac{\partial p(\boldsymbol{X};\theta)}{\partial \theta} \right] d\boldsymbol{x} =$$

$$\int \left[ \frac{\partial^2 \ln p(\boldsymbol{X};\theta)}{\partial \theta^2} p(\boldsymbol{X};\theta) + \frac{\partial \ln p(\boldsymbol{X};\theta)}{\partial \theta} \frac{\partial \ln p(\boldsymbol{X};\theta)}{\partial \theta} p(\boldsymbol{X};\theta) \right] d\boldsymbol{x} = 0$$

And finally

$$E\left[ \frac{\partial^2 \ln p(\boldsymbol{X};\theta)}{\partial \theta^2} \right] = -E\left[ \left( \frac{\partial \ln p(\boldsymbol{X};\theta)}{\partial \theta} \right)^2 \right]$$

PO 1213

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Cramer-Rao lower bound:*

**Proof outline (cont.):**

Taking this into consideration, i.e.

$$E\left[\left(\frac{\partial \ln p(\boldsymbol{x};\theta)}{\partial \theta}\right)^2\right] = -E\left[\frac{\partial^2 \ln p(\boldsymbol{x};\theta)}{\partial \theta^2}\right]$$

expression (2) results, in the case where *g(θ)=θ*.

The result (3) will be obtained next…

$\square$

See also appendix 3.B for the derivation in the vector case.

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Cramer-Rao lower bound:*

*Summary:*

• *Being able to place a lower bound on the variance of any unbiased estimator is very useful.*

• *It allow us to assert that an estimator is the MVU estimator (if it attains the bound for all values of the unknown parameter).*

• *It provides in all cases a benchmark for the unbiased estimators that we can design.*

• *It alerts to impossibility of finding unbiased estimators with variance lower than the bound.*

•*Provides a systematic way of finding the MVU estimator, if it exists and if an extra condition is verified.*

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Example:*

Example (DC level in white Gaussian noise):

**Problem:** Find MVU estimator.          **Approach:** Compute CRLB, if right form we have it.

Signal model:   $x[n] = A + w[n],$          $n = 0, ..., N-1,$          $w[n]: N(0, \sigma^2)$

Likelihood function:          $p(\mathbf{x}; A) = \dfrac{1}{(2\pi\sigma^2)^{\frac{N}{2}}} e^{-\frac{1}{2\sigma^2}\sum_{n=0}^{N-1}(x[n]-A)^2}$

$$\frac{\partial}{\partial A}\ln p(\mathbf{x}; A) = \frac{\partial}{\partial A}\left(-\frac{1}{2\sigma^2}\sum_{n=0}^{N-1}(x[n]-A)^2\right) = \frac{1}{\sigma^2}\sum_{n=0}^{N-1}(x[n]-A) = \frac{N}{\sigma^2}(\bar{x}-A)$$

$$\frac{\partial^2}{\partial A^2}\ln p(\mathbf{x}; A) = -\frac{N}{\sigma^2}$$          **CRLB:**          $\text{var}(\hat{A}) \geq \dfrac{1}{N/\sigma^2} = \dfrac{\sigma^2}{N}$

The estimator is unbiased and has the same variance, **thus it is a MVU** estimator! And it has the form:

$$\frac{\partial}{\partial A}\ln p(x; A) = I(\theta)(g(x)-\theta), \qquad \text{for} \quad I(\theta) = \frac{N}{\sigma^2} \qquad g(x) = \bar{\mathbf{x}}.$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Cramer-Rao lower bound:*

**Proof outline (second part of the theorem):**

Still remains to prove that the CRLB is attained for the estimator $\hat{\theta} = g(\mathbf{x})$

$$\mathrm{var}\left(\hat{\theta}\right) = \frac{1}{I(\theta)}, \qquad \text{for} \quad I(\theta) = -E\left[\frac{\partial^2}{\partial \theta^2} \ln p(x;\theta)\right]$$

If

$$\frac{\partial}{\partial \theta} \ln p(x;\theta) = I(\theta)\left(g(\boldsymbol{x}) - \theta\right)$$

differentiation relative to the parameter gives

$$\frac{\partial^2}{\partial \theta^2} \ln p(\boldsymbol{x};\theta) = \frac{\partial I(\theta)}{\partial \theta}\left(g(\mathbf{x}) - \theta\right) - I(\theta)$$

and then

$$-E\left[\frac{\partial^2}{\partial \theta^2} \ln p(\boldsymbol{x};\theta)\right] = -\frac{\partial I(\theta)}{\partial \theta}\left(E\left[g(\mathbf{x})\right] - \theta\right) + I(\theta) = I(\theta)$$

i.e. the bound is attained.

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Example:*

Example (phase estimation):

Signal model:  $x[n] = A\cos(2\pi f_0 n + \phi) + w[n], \qquad n = 0, ..., N-1$

$A, f_0$ known

Likelihood function:  $p(\mathbf{x}; \phi) = \dfrac{1}{(2\pi\sigma^2)^{\frac{N}{2}}} e^{-\frac{1}{2\sigma^2}\sum_{n=0}^{N-1}(x[n] - A\cos(2\pi f_0 n + \phi))^2}$

$$\frac{\partial}{\partial \phi} \ln p(\mathbf{x}; \phi) = -\frac{A}{\sigma^2} \sum_{n=0}^{N-1} \left( x[n]\sin(2\pi f_0 n + \phi) - \frac{A}{2}\sin(4\pi f_0 n + 2\phi) \right) A$$

$$\frac{\partial^2}{\partial \phi^2} \ln p(\mathbf{x}; \phi) = -\frac{A}{\sigma^2} \sum_{n=0}^{N-1} \left( x[n]\cos(2\pi f_0 n + \phi) - A\cos(4\pi f_0 n + 2\phi) \right) A$$

$$-E\left[ \frac{\partial^2}{\partial \phi^2} \ln p(\mathbf{x}; \phi) \right] = -\frac{A^2}{\sigma^2} \sum_{n=0}^{N-1} \left( \frac{1}{2} - \frac{1}{2}\cos(4\pi f_0 n + 2\phi) \right) \approx \frac{NA^2}{2\sigma^2}$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Example:*

Example (phase estimation cont.):

$$-E\left[\frac{\partial^2}{\partial \phi^2}\ln p\left(\mathbf{x}\mid\phi\right)\right] = -\frac{A^2}{\sigma^2}\sum_{n=0}^{N-1}\left(\frac{1}{2}-\frac{1}{2}\cos\left(4\pi f_0 n + 2\phi\right)\right) \approx \frac{NA^2}{2\sigma^2}$$

as $\qquad \sum_{n=0}^{N-1}\cos\left(4\pi f_0 n + 2\phi\right) \approx 0$ for $f_0$ not near 0 or 1/2.

$$\frac{1}{N}\sum_{n=0}^{N-1}\cos\left(\alpha n + \beta\right) = \frac{1}{N}\operatorname{Re}\left\{\sum_{n=0}^{N-1}e^{j(\alpha n+\beta)}\right\} = \ldots = \frac{\sin\left(N\alpha/2\right)}{N\sin\left(\alpha/2\right)}\cos\left(\alpha\frac{N-1}{2}+\beta\right)$$

for large *N.*

$$\operatorname{var}\left(\hat{\phi}\right) \geq \frac{2\sigma^2}{NA^2}$$

- Bound decreases as $SNR = A^2/2\sigma^2$ increases
- Bound decreases as *N* increases

Does an efficient estimator exists? Does a MVUE estimator exists?

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Fisher information:*

*We define the Fisher Information (Matrix) as*

$$I\left(\hat{\theta}\right) = -E\left[\frac{\partial^2}{\partial\theta^2}\ln p\left(x;\theta\right)\right]$$

*Note:*

- *$I(\theta) \geq 0$*

- *It is additive for independent observations*

$$\ln p\left(\mathbf{x};\theta\right) = \ln \pi_{n=0}^{N-1} p\left(x[n];\theta\right) = \sum_{n=0}^{N-1}\ln p\left(x[n];\theta\right)$$

$$I\left(\theta\right) = -E\left[\frac{\partial^2}{\partial\theta^2}\ln p\left(\mathbf{x};\theta\right)\right] = -\sum_{n=0}^{N-1}E\left[\frac{\partial^2}{\partial\theta^2}\ln p\left(x[n];\theta\right)\right]$$

- *If identically distributed (same PDF for each x[n])*

$$I\left(\theta\right) = Ni\left(\theta\right) = -N\left[\frac{\partial^2}{\partial\theta^2}\ln p\left(x[.];\theta\right)\right]$$

*As N->∞, for iid => CRLB-> 0*

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Other estimator characteristic:*

Efficiency:

An estimator that is unbiased and attains the
CRLB is said to be **efficient.**

$\mathrm{var}\left(\hat{\theta}\right)$

$\hat{\theta}_1$

$\hat{\theta}_2$

CRLB

$\hat{\theta}_3 = MVU$

$\hat{\theta}_3$ MVU and efficient $\theta$

$\mathrm{var}\left(\hat{\theta}\right)$

$\hat{\theta}_1$

$\hat{\theta}_2$

CRLB

$\hat{\theta}_3$

$\hat{\theta}_3$ MVU but not efficient $\theta$

$\mathrm{var}\left(\hat{\theta}\right)$

$\hat{\theta}_1$

$\hat{\theta}_2$

$\hat{\theta}_3$

CRLB

No MVU available. $\theta$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Transformation of parameters:*

Imagine that the CRLB is known for the parameter $\theta$. Can we compute easily the CRLB

for a linear transformation of the form $\alpha = g(\theta) = a\theta + \beta$ ?

$$\hat{\alpha} = a\hat{\theta} + b, \qquad\qquad E\left[a\hat{\theta} + b\right] = aE\left[\hat{\theta}\right] + b = \alpha$$

$$\operatorname{var}\left[a\hat{\theta} + b\right] = a^2 \operatorname{var}\left[\hat{\theta}\right] = \frac{\left(\dfrac{\partial g(\theta)}{\partial \theta}\right)^2}{-E\left[\dfrac{\partial^2}{\partial \theta^2} \ln p(\boldsymbol{x}; \theta)\right]}$$

**Linear transformations preserve biasness and efficiency.**

And for a nonlinear transformation of the form $\alpha = g(\theta)$?

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Transformation of parameters:*

**Remark: after a nonlinear transformation, the good properties can be lost.**

$$\text{var}\left(\hat{\theta}\right) \geq \frac{\left(\dfrac{\partial g\left(\theta\right)}{\partial \theta}\right)^2}{-E\left[\dfrac{\partial^2}{\partial \theta^2}\ln p\left(\boldsymbol{x};\theta\right)\right]}$$

Example: Suppose that given a stochastic variable $\bar{x} : N\left(A, \dfrac{\sigma^2}{N}\right)$ we desire to have an estimator for $\alpha = g(\mathrm{A}) = A^2$ (power estimator). Note that

$$\text{var}\left(\bar{\mathbf{x}}\right) = E\left[\left(\bar{\mathbf{x}} - E\left[\left(\bar{\mathbf{x}}\right)\right]\right)^2\right] = E\left[\bar{\mathbf{x}}^2 - 2\bar{\mathbf{x}}E\left[\bar{\mathbf{x}}\right] + E^2\left[\bar{\mathbf{x}}\right]\right] = E\left[\bar{\mathbf{x}}^2\right] - 2E^2\left[\bar{\mathbf{x}}\right] + E^2\left[\bar{\mathbf{x}}\right] = E\left[\bar{\mathbf{x}}^2\right] - E^2\left[\bar{\mathbf{x}}\right]$$

$$E\left[\bar{\mathbf{x}}^2\right] = \text{var}\left(\bar{\mathbf{x}}\right) + E^2\left[\bar{\mathbf{x}}\right]$$

**A bias estimate results. Efficiency is lost.**

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# Cramer-Rao lower bound:

**Theorem 3.1 (Cramer-Rao lower bound, Vector parameter)** – *It is assumed that the PDF p(**x**;**θ**) satisfies the "regularity" condition*

$$E\left[\frac{\partial}{\partial \theta}\ln p(\boldsymbol{x};\boldsymbol{\theta})\right] = \mathbf{0} \qquad \text{for all } \theta$$

*where the expectation is taken with respect to p(**x**, θ). Then, the variance of any unbiased estimator $\hat{\theta}$ must satisfy*

$$\mathbf{C}_{\hat{\theta}} - I^{-1}(\boldsymbol{\theta}) \geq 0,$$

*where ≥ is interpreted as meaning the matrix is positive semi-definite. The Fisher information matrix I(**θ**) is given as*

$$\left[I(\boldsymbol{\theta})\right]_{ij} = -E\left[\frac{\partial^2}{\partial \theta_i \partial \theta_j}\ln p(\boldsymbol{x};\boldsymbol{\theta})\right],$$

*where the derivatives are evaluated at the true value of θ and the expectation is taken with respect to p(**x**;θ). Furthermore, an unbiased estimator may be found that attains the bound for all θ if and only if*

$$\frac{\partial}{\partial \theta}\ln p(\mathbf{x};\boldsymbol{\theta}) = I(\boldsymbol{\theta})(\mathbf{g}(\mathbf{x}) - \boldsymbol{\theta})$$

(3)

*for some functions p dimensional function **g(.)** and some p x p matrix **I (.)**. The estimator, which is the MVU estimator, is $\hat{\boldsymbol{\theta}} = \mathbf{g}(x)$, and its covariance matrix is **I**-1(**θ**).*

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Vector Transformation of parameters:*

The vector transformation of parameters $\boldsymbol{\alpha} = \mathbf{g}(\boldsymbol{\theta})$ impacts on the CRLB computation as

$$\mathbf{C}_{\hat{\alpha}} - \frac{\partial \mathbf{g}(\theta)}{\partial \theta} I^{-1}(\theta) \frac{\partial \mathbf{g}(\theta)^T}{\partial \theta} \geq 0$$

where the Jacobian is

$$\frac{\partial \mathbf{g}(\theta)}{\partial \theta} = \begin{bmatrix} \dfrac{\partial g_1(\theta)}{\partial \theta_1} & \cdots & \dfrac{\partial g_1(\theta)}{\partial \theta_p} \\ \cdots & \cdots & \cdots \\ \dfrac{\partial g_r(\theta)}{\partial \theta_1} & & \dfrac{\partial g_r(\theta)}{\partial \theta_p} \end{bmatrix}$$

In the Gaussian general case for *x[n]=s[n]+w[n]*, where $\mathbf{w} \sim N\left(\mu(\theta), \mathbf{C}_\theta\right)$

the Fisher information matrix is

$$[I(\boldsymbol{\theta})]_{ij} = \left[\frac{\partial \mu(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_i}\right]^T \mathbf{C}^{-1}(\boldsymbol{\theta}) \left[\frac{\partial \mu(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_i}\right] \mathbf{C}_{\hat{\alpha}} + \frac{1}{2} tr\left[\mathbf{C}^{-1}(\boldsymbol{\theta}) \frac{\partial C(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_i} \mathbf{C}^{-1}(\boldsymbol{\theta}) \frac{\partial C(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}_j}\right].$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Example:*

Example (line fitting):

Signal model:
$$x[n] = A + Bn + w[n], \qquad n = 0, ..., N-1$$

$A, B$ deterministic unknown quantities

Likelihood function:
$$p(\mathbf{x};\theta) = \frac{1}{(2\pi\sigma^2)^{\frac{N}{2}}} e^{-\frac{1}{2\sigma^2}\sum_{n=0}^{N-1}(x[n]-A-Bn)^2}, \qquad where \quad \theta = \begin{bmatrix} A & B \end{bmatrix}^T$$

The Fisher Information Matrix is

$$\mathbf{I}(\theta) = \begin{bmatrix} -E\left[\dfrac{\partial^2 \ln p(\mathbf{x};\theta)}{\partial A^2}\right] & -E\left[\dfrac{\partial^2 \ln p(\mathbf{x};\theta)}{\partial A \partial B}\right] \\ -E\left[\dfrac{\partial^2 \ln p(\mathbf{x};\theta)}{\partial B \partial A}\right] & -E\left[\dfrac{\partial^2 \ln p(\mathbf{x};\theta)}{\partial B^2}\right] \end{bmatrix}$$

where

$$\frac{\partial \ln p(\mathbf{x};\theta)}{\partial A} = \frac{1}{\sigma^2}\sum_{n=0}^{N-1}(x[n]-A-Bn), \qquad and \quad \frac{\partial \ln p(\mathbf{x};\theta)}{\partial B} = \frac{1}{\sigma^2}\sum_{n=0}^{N-1}(x[n]-A-Bn)n.$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Example:*

Example (cont.):

Moreover

$$\frac{\partial^2 \ln p(\mathbf{x};\theta)}{\partial A^2} = -\frac{N}{\sigma^2}, \frac{\partial^2 \ln p(\mathbf{x};\theta)}{\partial A \partial B} = -\frac{1}{\sigma^2}\sum_{n=0}^{N-1} n, \quad \text{and} \quad \frac{\partial^2 \ln p(\mathbf{x};\theta)}{\partial B^2} = -\frac{1}{\sigma^2}\sum_{n=0}^{N-1} n^2.$$

Since the second order derivatives do not depend on x, we have immediately that

$$\mathbf{I}(\theta) = \frac{1}{\sigma^2}\begin{bmatrix} N & \dfrac{N(N-1)}{2} \\ \dfrac{N(N-1)}{2} & \dfrac{N(N-1)(2N-1)}{6} \end{bmatrix}$$

And also,

$$\mathbf{I}^{-1}(\theta) = \sigma^2\begin{bmatrix} \dfrac{2(2N-1)}{N(N+1)} & -\dfrac{6}{N(N+1)} \\ -\dfrac{6}{N(N+1)} & \dfrac{12}{N(N^2-1)} \end{bmatrix}, \qquad \begin{array}{l} \text{var}(\hat{A}) \geq \dfrac{2(2N-1)}{N(N+1)}\sigma^2 \\[2em] \text{var}(\hat{B}) \geq \dfrac{12}{N(N^2-1)}\sigma^2 \end{array}$$

# *Example:*

Example (cont.):

Remarks:

For only one parameter to be determined $\mathrm{var}\left(\hat{A}\right) \geq \dfrac{\sigma^2}{N}$. Thus a general results was obtained: ***when more parameters are to be estimated the CRLB always degrades.***

*Moreover*

$$\frac{CRLB(\hat{A})}{CRLB(\hat{B})} = \frac{(2N-1)(N-1)}{6} > 1, \quad \text{for} \quad N \geq 3.$$

The parameter *B* is easier to be determined, as its CRLB decreases with $1/N^3$. This means that *x[n]* is more sensitive to changes in *B* than changes in *A*.

$$\Delta x\left[n\right] \approx \frac{\partial x\left[n\right]}{\partial A} \Delta A = \Delta A$$

$$\Delta x\left[n\right] \approx \frac{\partial x\left[n\right]}{\partial B} \Delta B = n\Delta B.$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Bibliography:*

**Further reading**

• Harry L. Van Trees, *Detection, Estimation, and Modulation Theory, Parts I to IV,* John Wiley, 2001.

• J. Bibby, H. Toutenburg, *Prediction and Improved Estimation in Linear Models,* John Wiley, 1977.

• C.Rao, *Linear Statistical Inference and Its Applications,* John Wiley, 1973.

• P. Stoica, R. Moses, *"On Biased Estimators and the Unbiased Cramer-Rao Lower Bound,"* *Signal Process,* vol.21, pp. 349-350, 1990.

**DEM**
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Advanced Control Systems*
# *Detection, Estimation, and Filtering*

*Graduate Course on the*
*MEng PhD Program*
*Spring 2012/2013*

## *Chapter 4*
## *Linear Models in the Presence of Stochastic Signals*

*Instructor:*
*Prof. Paulo Jorge Oliveira*
*p.oliveira@dem.ist.utl.pt or pjcro @ isr.ist.utl.pt*
*Phone: 21 8419511 or 21 8418053 (3511 or 2053 inside IST)*

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Syllabus:*

Classical Estimation Theory

...

Chap. 3 - **Cramer-Rao Lower Bound** [1 week]
Estimator accuracy; Cramer-Rao lower bound (CRLB); CRLB for signals in white Gaussian noise; Examples;

Chap. 4 - **Linear Models in the Presence of Stochastic Signals** [1 week] Stationary and transient analysis; White Gaussian noise and linear systems; Examples; Sufficient Statistics; Relation with MVU Estimators;

Chap. 5 - **Best Linear Unbiased Estimators** [1 week]
Definition of BLUE estimators; White Gaussian noise and bandlimited systems; Examples; Generalized minimum variance unbiased estimation;

continues…

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *A very special class of systems:*

FACT:

**The determination of the MVU Estimator is in general a difficult task.**

**A class of systems that allows the determination of this estimator easily…**

**LINEAR SYSTEMS**

**The statistical performance is also easy to compute**

**and an efficient solution is obtained.**

**The key point is on the formulation of a problem as a linear one.**

# MVU Estimator for the Linear Model:

**Theorem 4.1** – If the data observed can be modeled as

$$\mathbf{x} = \mathbf{H}\theta + \mathbf{w}$$

where $\mathbf{x}$ is a N x 1 vector of observations, $\mathbf{H}$ is a known N x p observation matrix (with N>p) and rank p, $\theta$ is a p x 1 vector of parameters to be estimated, and $\mathbf{w}$ is an N x 1 noise vector with PDF $N(0, \sigma^2\mathbf{I})$, then the MVU is

$$\hat{\theta} = \left(\mathbf{H}^T\mathbf{H}\right)^{-1}\mathbf{H}^T\mathbf{x} \qquad (1)$$

and the covariance matrix of estimate is

$$C_{\hat{\theta}} = \sigma^2\left(\mathbf{H}^T\mathbf{H}\right)^{-1}. \qquad (2)$$

**Proof outline:**

As discussed in Chapter 3, it is possible to determine the MVU estimator if the equality constraints of the CRLB are satisfied.

From the signal model, it follows that the log-likelihood function is

$$\ln p\left(\mathbf{x};\boldsymbol{\theta}\right) = -\ln\left(2\pi\sigma^2\right)^{N/2} - \frac{\left(\mathbf{x} - \mathbf{H}\boldsymbol{\theta}\right)^T\left(\mathbf{x} - \mathbf{H}\boldsymbol{\theta}\right)}{2\sigma^2}$$

And

$$\frac{\partial \ln p\left(\mathbf{x};\boldsymbol{\theta}\right)}{\partial \boldsymbol{\theta}} = -\frac{1}{2\sigma^2}\frac{\partial}{\partial \boldsymbol{\theta}}\left[\mathbf{x}^T\mathbf{x} - 2\mathbf{x}^T\mathbf{H}\boldsymbol{\theta} + \boldsymbol{\theta}^T\mathbf{H}^T\mathbf{H}\boldsymbol{\theta}\right].$$

Using the relations (deduce them, good exercise…)

$$\frac{\partial \mathbf{b}^T\boldsymbol{\theta}}{\partial \boldsymbol{\theta}} = \mathbf{b} \qquad \frac{\partial \boldsymbol{\theta}^T\mathbf{A}\boldsymbol{\theta}}{\partial \boldsymbol{\theta}} = 2\mathbf{A}\boldsymbol{\theta}$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *MVU Estimator for the Linear Model:*

**Proof outline (cont):**

It follows

$$\frac{\partial \ln p\left(\mathbf{x};\theta\right)}{\partial \theta} = \frac{1}{\sigma^2}\left[\mathbf{H}^T\mathbf{x} - \mathbf{H}^T\mathbf{H}\theta\right].$$

Under the assumptions of the theorem, $\mathbf{H}^T\mathbf{H}$ is invertible

$$\frac{\partial \ln p\left(\mathbf{x};\theta\right)}{\partial \theta} = \frac{\mathbf{H}^T\mathbf{H}}{\sigma^2}\left[\left(\mathbf{H}^T\mathbf{H}\right)^{-1}\mathbf{H}^T\mathbf{x} - \theta\right]. \qquad \left(\frac{\partial \ln p\left(\mathbf{x};\theta\right)}{\partial \theta} = \mathbf{I}(\theta)\left[\mathbf{g}(\mathbf{x}) - \theta\right]\right)$$

Note that it is in the format introduced in the previous chapter, from where (1) and (2) follows immediately.

☐

Major constraints:

      what if $\mathbf{H}^T\mathbf{H}$ is not invertible?

      what if $\mathbf{H}^T\mathbf{H}$ is ill-conditioned?

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Example - Fourier Analysis:*

Cyclic components in white Gaussian noise

Signal model:

$$x[n] = \sum_{k=1}^{M} a_k \cos\left(\frac{2\pi kn}{N}\right) + \sum_{k=1}^{M} b_k \sin\left(\frac{2\pi kn}{N}\right) + w[n], \qquad n = 0,...,N-1, \quad w[n]: N\left(0,\sigma^2\right)$$

Defining

$$\theta = \begin{bmatrix} a_1 & ... & a_M & b_1 & ... & b_M \end{bmatrix}^T, \mathbf{w} = \begin{bmatrix} w_0 & & w_{N-1} \end{bmatrix}^T, \text{and}$$

$$\mathbf{H} = \begin{bmatrix} 1 & ... & 1 & 0 & ... & 0 \\ c\left(\dfrac{2\pi}{N}\right) & ... & c\left(\dfrac{2\pi M}{N}\right) & s\left(\dfrac{2\pi}{N}\right) & ... & s\left(\dfrac{2\pi M}{N}\right) \\ \vdots & \vdots & \ddots & \vdots & \vdots & \vdots \\ c\left(\dfrac{2\pi(N-1)}{N}\right) & ... & c\left(\dfrac{2\pi M(N-1)}{N}\right) & s\left(\dfrac{2\pi(N-1)}{N}\right) & ... & s\left(\dfrac{2\pi M(N-1)}{N}\right) \end{bmatrix}$$

The model can be reformulated as a linear system, with solution if *M < N/2*

$$\mathbf{x} = \mathbf{H}\theta + \mathbf{w}$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Example - Fourier Analysis (cont.):*

Important fact: The columns of $\mathbf{H}$ are orthogonal.

Define

$$\mathbf{H} = \begin{bmatrix} \mathbf{h_1} & \mathbf{h_2} & ... & \mathbf{h_{2M}} \end{bmatrix}, \qquad \text{it follows} \quad \mathbf{h}_i^T \mathbf{h}_j = 0, \qquad i \neq j \quad .$$

Moreover, the discrete Fourier Transform (DFT) relations can be applied, i.e.

$$\sum_{n=0}^{N-1} \cos\left(\frac{2\pi in}{N}\right)\cos\left(\frac{2\pi jn}{N}\right) = \frac{N}{2}\delta_{ij}$$

$$\sum_{n=0}^{N-1} \sin\left(\frac{2\pi in}{N}\right)\sin\left(\frac{2\pi jn}{N}\right) = \frac{N}{2}\delta_{ij}$$

$$\sum_{n=0}^{N-1} \cos\left(\frac{2\pi in}{N}\right)\sin\left(\frac{2\pi jn}{N}\right) = 0, \qquad \text{for all} \quad i, j.$$

From where it follows

$$\mathbf{H}^T\mathbf{H} = \frac{N}{2}\mathbf{I}.$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Example - Fourier Analysis (cont.):*

The MVU estimator is

$$\hat{\theta} = \left(\mathbf{H}^T\mathbf{H}\right)^{-1}\mathbf{H}^T\mathbf{x} = \frac{2}{N}\mathbf{H}^T\mathbf{x} = \begin{bmatrix} \dfrac{2}{N}h_1^T\mathbf{x} \\ \vdots \\ \dfrac{2}{N}h_{2M}^T\mathbf{x} \end{bmatrix},$$

§

or finally

$$\hat{a}_k = \frac{2}{N}\sum_{n=0}^{N-1}x[n]\cos\left(\frac{2\pi kn}{N}\right),$$

$$\hat{b}_k = \frac{2}{N}\sum_{n=0}^{N-1}x[n]\sin\left(\frac{2\pi kn}{N}\right).$$

with covariance

$$C_{\hat{\theta}} = \frac{2\sigma^2}{N}\mathbf{I}.$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Example: System Identification*

Signal model, where the Finite Impulse Response (FIR) is to be estimated.

The user can apply the input signal $u$:

$w[n]$

$u[n] \longrightarrow \boxed{H(z)} \longrightarrow \overset{+}{\bigcirc} \longrightarrow x[n]$

$$x[n] = \sum_{k=0}^{p-1} h[k] u[n-k] + w[n], \qquad n = 0,...,N-1, \quad w[n] \sim N(0,\sigma^2) \ .$$

In matrix form, considering $\mathbf{x}=[x_0 \ ... \ x_{N-1}]^T$, the input/output relations of this linear system can be written as

$$\mathbf{X} = \begin{bmatrix} u[0] & 0 & ... & 0 \\ u[1] & u[0] & ... & 0 \\ \vdots & \vdots & \ddots & \vdots \\ u[N-1] & u[N-2] & ... & u[N-p] \end{bmatrix} \begin{bmatrix} h[0] \\ h[1] \\ \vdots \\ h[p-1] \end{bmatrix} + \mathbf{w} \qquad \mathbf{w} = \begin{bmatrix} w_0 & w_{N-1} \end{bmatrix}^T$$

Or in compact form, once again

$$\mathbf{x} = \mathbf{H}\theta + \mathbf{w}$$

PO 1213

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

The MVU estimator is once again

$$\hat{\boldsymbol{\theta}} = \left(\mathbf{H}^T \mathbf{H}\right)^{-1} \mathbf{H}^T \mathbf{x}, \quad \text{with covariance} \quad C_{\hat{\boldsymbol{\theta}}} = \sigma^2 \left(\mathbf{H}^T \mathbf{H}\right)^{-1}.$$

Note that accuracy depends on the input signal applied. How to choose it?

Problem: Choose u[n] to minimize $\operatorname{var}\left(\hat{\theta}_i\right) = \left[\mathbf{C}_{\hat{\boldsymbol{\theta}}}\right]_{ii}, i = 1,...,p,$ subject to the constraint that $\sum_{n=0}^{N-1} u[n]$ is fixed.

Introducing the crosscorrelation (autocorrelation)

$$r_{ux}[i] = \frac{1}{N}\sum_{n=0}^{N-1-i} u[n]x[n+i] \qquad \mathbf{H}^T\mathbf{H} = \begin{bmatrix} r_{uu}[0] & r_{uu}[1] & \cdots & r_{uu}[p-1] \\ r_{uu}[1] & r_{uu}[0] & \cdots & r_{uu}[p-2] \\ \vdots & \vdots & \ddots & \vdots \\ r_{uu}[p-1] & r_{uu}[p-2] & \cdots & r_{uu}[0] \end{bmatrix}$$

Choosing a Pseudorandom Noise (PRN) makes this last matrix diagonal

$$C_{\hat{\boldsymbol{\theta}}} = \sigma^2 \left(r_{uu}[0]\mathbf{I}\right)^{-1} = \frac{\sigma^2}{r_{uu}[0]}\mathbf{I}.$$

Input Signal Energy

# *Extension to non-white Gaussian noise:*

*Theorem (Generalization of Theorem 4.1) – If the data observed can be modeled as*

$$\mathbf{x} = \mathbf{H}\theta + \mathbf{w}$$

*where $\mathbf{x}$ is a N x 1 vector of observations, $\mathbf{H}$ is a known N x p observation matrix (with N>p) and rank p, $\theta$ is a p x 1 vector of parameters to be estimated, and $\mathbf{w}$ is an N x 1 colored noise vector with PDF $N(0, C)$ ($C \neq \sigma^2 I$), then the MVU is*

$$\hat{\theta} = \left(\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H}\right)^{-1} \mathbf{H}^T \mathbf{C}^{-1} \mathbf{x}$$

(1bis)

*and the covariance matrix of estimate is*

$$C_{\hat{\theta}} = \left(\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H}\right)^{-1}.$$

(2bis)

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Extension to non-white Gaussian noise:*

*Proof: The covariance matrix and its inverse are both positive semi-definite. Thus*

$$\mathbf{C}^{-1} = \mathbf{D}^T \mathbf{D}, \qquad \text{where} \qquad \mathbf{D} \in R^{NxN}$$

*A noise whitening operation can be performed. For that purpose lets compute the covariance of*

$$E\left[\left(\mathbf{D}w\right)\left(\mathbf{D}w\right)^T\right] = \mathbf{D}\mathbf{C}\mathbf{D}^T = \mathbf{D}\mathbf{D}^{-1}\mathbf{D}^{T^{-1}}\mathbf{D}^T = I.$$

*If we define the new variable x' as*

$$\mathbf{x}' = \mathbf{D}\mathbf{x} = \mathbf{D}\mathbf{H}\theta + \mathbf{D}\mathbf{w} = \mathbf{H}'\theta + \mathbf{w}'.$$

*Applying the usual solution to this linear model (transformed) results in*

$$\hat{\theta} = \left(\mathbf{H}'^T \mathbf{H}'\right)^{-1} \mathbf{H}'^T \mathbf{x}' = \left(\mathbf{H}^T \mathbf{D}^T \mathbf{D}\mathbf{H}\right)^{-1} \mathbf{H}^T \mathbf{D}^T \mathbf{D}\mathbf{x} = \left(\mathbf{H}^T \mathbf{C}^{-1}\mathbf{H}\right)^{-1} \mathbf{H}^T \mathbf{C}^{-1}\mathbf{x}$$

*For a covariance*

$$C_{\hat{\theta}} = \left(\mathbf{H}'^T \mathbf{H}'\right)^{-1} = \left(\mathbf{H}^T \mathbf{D}^T \mathbf{D}\mathbf{H}'\right)^{-1} = \left(\mathbf{H}^T \mathbf{C}^{-1}\mathbf{H}'\right)^{-1}.$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Sufficient Statistics:*

*General MVU Estimation:*

*Assume that the CRLB is not satisfied with equality!*

*There is no efficient estimator.*

*How do we find the MVU estimator (if it exists)?*

**Use the concept of Sufficient Statistics.**

*Example: To compute the value of a DC signal in noise, given n samples, i=0,…,N-1.*

*Consider*

$$S_1 = \left\{ x[0], x[1], ..., x[N-1] \right\}$$

$$S_2 = \left\{ x[0] + x[1], ..., x[N-1] \right\}$$

$$S_3 = \left\{ \sum_{n=0}^{N-1} x[n] \right\}$$

*All sets are sufficient since the unknown parameter can be found. $S_3$ is the minimal one.*

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# Sufficient Statistics:

*Theorem 5.1 (Neyman-Fisher Factorization) – If we can factor the PDF p ($x$;θ) as*

$$p(\mathbf{x};\theta) = g\big(T(\mathbf{x}),\theta\big)h(\mathbf{x})$$

(3)

*where g(.) is a function depending on $x$ only through $T(x)$ and $h$ (.) is a function depending only on $x$, then $T(x)$ is sufficient statistic for θ. Conversely, if $T(x)$ is a sufficient statistic for θ then the PDF can be factored as in (3).*

*Proof outline (=>):*

- *p($x$,$T(x)$;θ) must have a minimum at $x$=$x_0$, denoted as $T(x_0)$=$T_0$;*
- *If $y$=g($x$), for the vector random variable $x$,* $p(y) = \int p(\mathbf{x})\delta\big(y - g(\mathbf{x})\big)d\mathbf{x}.$

$$p\big(\mathbf{x}\,|\,T(\mathbf{x}) = T_0;\theta\big)$$

- *Knowledge of the value of a sufficient statistics makes the conditional PDF not to depend on the parameters*

DEM
DEPARTAMENTO
DE ENGENHARIA MECÃNICA
TÉCNICO LISBOA

# Sufficient Statistics:

*Proof outline (cont):*

*Using conditional*

*probability definition:*

$$p\left(\mathbf{x}\,|\,T\left(\mathbf{x}\right)=T_0;\theta\right) = \frac{p\left(\mathbf{x},T\left(\mathbf{x}\right)=T_0;\theta\right)}{p\left(T\left(\mathbf{x}\right)=T_0;\theta\right)} = \frac{p\left(\mathbf{x};\theta\right)\delta\left(T\left(\mathbf{x}\right)-T_0\right)}{p\left(T\left(\mathbf{x}\right)=T_0;\theta\right)}$$

$$= \frac{g\left(\mathbf{x},T\left(\mathbf{x}\right)=T_0,\theta\right)h\left(\mathbf{x}\right)\delta\left(T\left(\mathbf{x}\right)-T_0\right)}{p\left(T\left(\mathbf{x}\right)=T_0;\theta\right)}.$$

*Where the factorization was used in the last step. The denominator can be written as*

$$p\left(T\left(\mathbf{x}\right)=T_0;\theta\right) = \int p\left(\mathbf{x};\theta\right)\delta\left(T\left(\mathbf{x}\right)-T_0\right)d\mathbf{x} =$$

$$= \int g\left(T\left(\mathbf{x}\right)=T_0,\theta\right)h\left(\mathbf{x}\right)\delta\left(T\left(\mathbf{x}\right)-T_0\right)d\mathbf{x} = g\left(T\left(\mathbf{x}\right)=T_0,\theta\right)\int h\left(\mathbf{x}\right)\delta\left(T\left(\mathbf{x}\right)-T_0\right)d\mathbf{x}.$$

*The integral is zero in $R^n$ except over the surface where T($\mathbf{x}$)=$T_0$. where it is constant.*

$$p\left(\mathbf{x}\,|\,T\left(\mathbf{x}\right)=T_0;\theta\right) = \frac{h\left(\mathbf{x}\right)\delta\left(T\left(\mathbf{x}\right)-T_0\right)}{\int h\left(\mathbf{x}\right)\delta\left(T\left(\mathbf{x}\right)-T_0\right)d\mathbf{x}},$$

*Which does not depend on θ. Hence, we conclude that T($\mathbf{x}$) is a sufficient statistic.* ∎

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Sufficient Statistics:*

*Proof outline (<=):*

*Consider the joint PDF*

$$p\big(\mathbf{x}, T(\mathbf{x}) = T_0; \theta\big) = p\big(\mathbf{x} \mid T(\mathbf{x}) = T_0; \theta\big) p\big(T(\mathbf{x}) = T_0; \theta\big) = p(\mathbf{x}; \theta)\delta\big(T(\mathbf{x}) - T_0\big).$$

*Because T(x) is a sufficient statistic, the conditional PDF does not depend on θ. We can let*

$$p\big(\mathbf{x} \mid T(\mathbf{x}) = T_0\big) = w(\mathbf{x})\delta\big(T(\mathbf{x}) - T_0\big)$$

*Substituting in the previous expression*

$$p(\mathbf{x}; \theta)\delta\big(T(\mathbf{x}) - T_0\big) = w(\mathbf{x})\delta\big(T(\mathbf{x}) - T_0\big) p\big(T(\mathbf{x}) = T_0; \theta\big)$$

*Setting w(x) to*

$$w(\mathbf{x}) = \frac{h(\mathbf{x})}{\int h(\mathbf{x})\delta\big(T(\mathbf{x}) - T_0\big) d\mathbf{x}},$$

*Allows one to write*

$$p(\mathbf{x}; \theta)\delta\big(T(\mathbf{x}) - T_0\big) = \frac{h(\mathbf{x})\delta\big(T(\mathbf{x}) - T_0\big)}{\int h(\mathbf{x})\delta\big(T(\mathbf{x}) - T_0\big) d\mathbf{x}} \, p\big(T(\mathbf{x}) = T_0; \theta\big)$$

*Thus based on the factorization a sufficient statistic can be found*

$$p(\mathbf{x}; \theta) = g\big(T(\mathbf{x}) = T_0; \theta\big) h(\mathbf{x})$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

DC Level in WGN: Assuming that we have found the sufficient statistics, the **Rao - Blackwell-Lehmann-Scheffe Theorem** can be used to find the MVU estimator in two different ways:

1) Find any unbiased estimator of A, say $\breve{A} = x[0]$, and determine $\hat{A} = E\left[\breve{A} \,|\, T\right]$.
The expectation is taken with respect to $p\left(\breve{A} \,|\, T\right)$.

2) Find some function *g(.)* so that $\hat{A} = g(T)$ is an unbiased estimator of *A*.

First approach:

Let $\breve{A} = x[0]$ and determine $\hat{A} = E\left[x[0] \,|\, \sum_{n=0}^{N-1} x[n]\right]$

We need auxiliary results for [x y]$^\mathsf{T}$ a Gaussian random vector with mean $\boldsymbol{\mu}$=[E[x] E[x]]$^\mathsf{T}$

$$E[x \,|\, y] = E[x] + \frac{\operatorname{cov}(x, y)}{\operatorname{var}(y)}\left(y - E[y]\right)$$

(see Appendix 10A for details.)

# *Motivating Example:*

DC Level in WGN (cont.):

1) Find any unbiased estimator of A, say $\breve{A} = x[0]$ , and determine $\hat{A} = E[\breve{A}|T]$ .

The expectation is taken with respect to $p(\breve{A}|T)$.

Applying the previous results to *x=x[0]* and $y = \sum_{n=0}^{N-1} x[n]$

$$\begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} x[0] \\ \sum_{n=0}^{N-1} x[n] \end{bmatrix} = \underbrace{\begin{bmatrix} 1 & 0 & \cdots & 0 \\ 1 & 1 & \cdots & 1 \end{bmatrix}}_{\mathbf{L}} \begin{bmatrix} x[0] \\ x[1] \\ \vdots \\ x[N-1] \end{bmatrix}$$

Hence the PDF of [x y]ᵀ is *N(μ, C),* where

$$\boldsymbol{\mu} = \mathbf{L}E[\mathbf{x}] = \mathbf{L}A\mathbf{1} = \begin{bmatrix} A \\ NA \end{bmatrix},$$

$$C = \sigma^2 \mathbf{L}\mathbf{L}^T = \sigma^2 \begin{bmatrix} 1 & 1 \\ 1 & N \end{bmatrix}.$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Motivating Example:*

DC Level in WGN (cont.):

1) Find any unbiased estimator of A, say $\breve{A} = x\left[0\right]$ , and determine $\hat{A} = E\left[\breve{A} \mid T\right]$ .

The expectation is taken with respect to $p\left(\breve{A} \mid T\right)$.

Hence we have finally

$$\hat{A} = E\left[x \mid y\right] = A + \frac{\sigma^2}{N\sigma^2}\left(\sum\nolimits_{n=0}^{N-1} x\left[n\right] - NA\right) = \frac{1}{N}\sum\nolimits_{n=0}^{N-1} x\left[n\right].$$

Which is the MVU estimator. Usually this option is mathematically intractable.

2) Find some function *g(.)* so that $\hat{A} = g\left(T\right)$ is an unbiased estimator of *A.*

*We need to find some function* $\hat{A} = g\left(\sum\nolimits_{n=0}^{N-1} x\left[n\right]\right)$ *so that it is an unbiased estimator.*

*That is the case of*

$$\hat{A} = \frac{1}{N}\sum\nolimits_{n=0}^{N-1} x\left[n\right].$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# RBLS Theorem:

**Definition:** a statistic is complete if there is only one function of the statistic that is unbiased.

**Theorem 5.1 (Rao-Blackwell-Lehmann-Scheffe)** – If $\breve{\theta}$ is an unbiased estimator of $\theta$ and T(**x**) is a sufficient statistic for $\theta$, then $\hat{\theta} = E\left[\breve{\theta} \mid T(\mathbf{x})\right]$ is

1. A valid estimator for $\theta$

2. Unbiased

3. Of lesser or equal variance than that of $\breve{\theta}$ , for all $\theta$.

Additionally, if the sufficient statistic is complete, then $\hat{\theta}$ is the MVU estimator.

To validate that a statistic is complete is in general very difficult, (see examples 5.6 and 5.7). It must verify

$$\int_{-\infty}^{+\infty} v(T) p(T;\theta) dT = 0, \quad \text{for all} \quad \theta. \quad (5.8)$$

Only for the zero function and for *v(T)*.

*Note: -* For an example of an incomplete statistic check Example 5.7

DEM
DEPARTAMENTO
DE ENGENHARIA MECÃNICA
TÉCNICO LISBOA

# *Methodology:*

Use Neyman-Fisher factorization theorem (5.1) to find sufficient statistic

$\downarrow T(\boldsymbol{x})$

Determine if $T(\boldsymbol{x})$ is complete see (5.8)

Find function of $T(\boldsymbol{x})$ that is unbiased

$$\hat{\theta} = g\big(T(\mathbf{x})\big) = \text{MVU Estimator}$$

# *Example:*

Mean of Uniform Noise:

Data model:        *x[n]=w[n], n=0,1,…,N-1*

Where *w[n]* is IID noise with *PDF  U[0,β], for β>0.*

We wish to find the MVU estimator for the mean *θ=β/2.*

The approach to find the CRLB can not be followed as the PDF does not satisfy the regularity conditions. A natural estimator is

$$\hat{\theta} = \frac{1}{N} \sum_{n=0}^{N-1} x[n], \qquad with \quad \mathrm{var}(\hat{\theta}) = \frac{1}{N} \mathrm{var}(x[n]) = \frac{\beta^2}{12N}.$$

To determine if the sample mean is the MVU we will follow the methodology previously presented.

# *Example:*

Lets define the unit step function:

$$u(x) = \begin{cases} 1 & for \quad x > 0 \\ 0 & for \quad x < 0 \end{cases}.$$

Then,

$$p(x[n];\beta) = \frac{1}{\beta}\left[u(x[n]) - u(x[n] - \beta)\right], \quad where \quad \beta = 2\theta.$$

and the PDF is

$$p(x[n];\beta) = \frac{1}{\beta^N}\prod_{n=0}^{N-1}\left[u(x[n]) - u(x[n] - \beta)\right] = \begin{cases} \dfrac{1}{\beta^N} & 0 < x[n] < \beta \quad n = 0,1,...,N-1 \\ 0 & otherwise \end{cases}.$$

Alternative, we can write

$$p(x[n];\beta) = \begin{cases} \dfrac{1}{\beta^N} & \max(x[n]) < \beta, \min(x[n]) > 0 \\ 0 & otherwise \end{cases},$$

So that

$$p(x[n];\beta) = \frac{1}{\beta^N}u\big(\beta - \max(x[n])\big)u\big(\min(x[n])\big)$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Example:*

Note that can be identified

$$p\big(x[n];\beta\big) = \underbrace{\frac{1}{\beta^N} u\big(\beta - \max\big(x[n]\big)\big)}_{g\big(T(x),\beta\big)} \quad \underbrace{u\big(\min\big(x[n]\big)\big)}_{h(x)}$$

By the Neyman-Fisher factorization theorem, *T(x)=max(x[n])* is a sufficient statistic for θ. Furthermore, it can be shown that the sufficient statistic is complete. We need next to find a function of *T(x)* that is not biased (denominated as order statistics). Lets write the cumulative distribution function

$$\Pr\big\{T \le \xi\big\} = \Pr\big\{x[0] \le \xi, x[1] \le \xi, ..., x[N-1] \le \xi, \big\} = \prod_{n=0}^{N-1} \Pr\big\{x[n] \le \xi\big\} = \Pr\big\{x[n] \le \xi\big\}^N.$$

The PDF follows as

$$p_T(\xi) = \frac{d\,\Pr\big\{T \le \xi\big\}}{d\xi} = N \Pr\big\{x[n] \le \xi\big\}^{N-1} \frac{d\,\Pr\big\{x[n] \le \xi\big\}}{d\xi}.$$

# *Example:*

But
$$\frac{d\,\mathrm{Pr}\{x[n] \le \xi\}}{d\xi} = p_{x[n]}(\xi) = \begin{cases} \dfrac{1}{\beta} & 0 < \xi < \beta \\ 0 & otherwise \end{cases},$$

Integrating we obtain

$$p_T(\xi) = \begin{cases} 0 & \xi < 0 \\ N\left(\dfrac{\xi}{\beta}\right)^{N-1}\dfrac{1}{\beta} & 0 < \xi < \beta \\ 0 & \xi > \beta \end{cases}, \quad and \quad E[T] = \int_0^{\beta} \xi N\left(\frac{\xi}{\beta}\right)^{N-1}\frac{1}{\beta}d\xi$$

From where it results

$$E[T] = \frac{N}{N+1}\beta = \frac{2N}{N+1}\theta, \qquad thus \quad \hat{\theta} = \frac{N+1}{2N}T \text{ makes the expected value unbiased.}$$

The MVU estimator is
$$\hat{\theta} = \frac{N+1}{2N}\max(x[n])$$

with a variance... $\mathrm{var}(\hat{\theta}) = \dfrac{\beta^2}{4N(N+2)} << \dfrac{\beta^2}{12N}$ (sample mean var)   for large N!

# *Bibliography:*

**Further reading**

• Thomas Kailath, *Linear Systems,* Prentice Hall, 1980.

• Thomas Kailath, Ali Sayed, and Babak Hassibi, *Linear Estimation,* Prentice Hall, 2000.

• Harry L. Van Trees, *Detection, Estimation, and Modulation Theory, Parts I to IV,* John Wiley, 2001.

• J. Bibby, H. Toutenburg, *Prediction and Improved Estimation in Linear Models,* John Wiley, 1977.

• C.Rao, *Linear Statistical Inference and Its Applications,* John Wiley, 1973.

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Advanced Control Systems*
# *Detection, Estimation, and Filtering*

## *Graduate Course on the*
## *MEng PhD Program*
## *Spring 2012/2013*

## *Chapter 5*
## *Best Linear Unbiased Estimators*

*Instructor:*
*Prof. Paulo Jorge Oliveira*
*p.oliveira@dem.ist.utl.pt or pjcro @ isr.ist.utl.pt*
*Phone: 21 8419511 or 21 8418053 (3511 or 2053 inside IST)*

# *Syllabus:*

Classical Estimation Theory

...

Chap. 4 - *Linear Models in the Presence of Stochastic Signals* [1 week]

Stationary and transient analysis; White Gaussian noise and linear systems;  Examples;

Chap. 5 - *Best Linear Unbiased Estimators* [1 week]

Definition of BLUE estimators;  White Gaussian noise and bandlimited

systems; Examples; Generalized MVU estimation;

Chap. 6 - *Maximum Likelihood Estimation* [1 week]

The maximum likelihood estimator; Properties of the ML estimators; Solution for ML

estimation; Examples; Monte-Carlo methods;

continues…

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *An alternative strategy:*

FACT:

**It occurs that the MVU estimator, if it exists, can not be found.**

e.g. the PDF for the data is not known, the user would not like to assume a model for the PDF, or the problem can be mathematically untreatable.

**An alternative strategy can be pursued is to study the class of**

**Best Linear Unbiased Estimators**

**Only suboptimal performance can be achieved.**

The performance degradation, relative to the MVU estimator, is unknown but the resulting performance can be enough for the problem at hand.

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# BLUE structure:

The Best Linear Unbiased Estimator consists of restrict the estimator to be a linear function of the data, i.e.

$$\hat{\theta} = \sum_{n=0}^{N-1} a_n x[n]$$
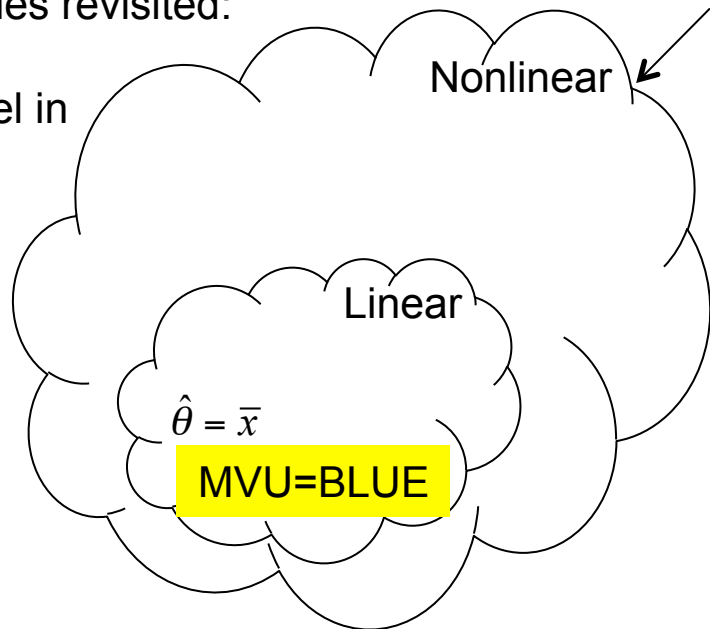
where the $a_n$'s are constants to be determined.

Optimality in general is lost.

Examples revisited:

DC level in WGN

All unbiased estimators

Nonlinear

Nonlinear

$$\hat{\theta} = \frac{N+1}{2N} \max x(n)$$

MVU

Linear

Linear

$$\hat{\theta} = \bar{x}$$

MVU=BLUE

$$\hat{\theta} = \bar{x}$$

BLUE

Mean of Uniform noise

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Finding the BLUE:*

To find the BLUE we constrain the estimator

- to be linear

$$\hat{\theta} = \sum_{n=0}^{N-1} a_n x[n]$$

- to be unbiased

$$E\left[\hat{\theta}\right] = \sum_{n=0}^{N-1} a_n E\left[x[n]\right] = \theta \qquad (6.2)$$

- to minimize its variance

$$\mathrm{var}\left(\hat{\theta}\right) = E\left[\left(\sum_{n=0}^{N-1} a_n x[n] - E\left[\sum_{n=0}^{N-1} a_n x[n]\right]\right)^2\right]$$

Defining $\boldsymbol{a}=[a_0\ a_1\ \dots\ a_{N-1}]^T$ and $\boldsymbol{x}=[x_0\ x_1\ \dots\ x_{N-1}]^T$ this last expression can be simplified:

$$\mathrm{var}\left(\hat{\theta}\right) = E\left[\left(\mathbf{a}^T\mathbf{x} - \mathbf{a}^T E\left[\mathbf{x}\right]\right)^2\right] = E\left[\left(\mathbf{a}^T\left(\mathbf{x} - E\left[\mathbf{x}\right]\right)\right)^2\right] =$$

$$= E\left[\mathbf{a}^T\left(\mathbf{x} - E\left[\mathbf{x}\right]\right)\left(\mathbf{x} - E\left[\mathbf{x}\right]\right)^T \mathbf{a}\right] = \mathbf{a}^T C \mathbf{a}. \qquad (6.3)$$

The problem of finding the BLUE can be stated as, for $\mathbf{a} \in R^N$

$$\min \qquad \mathbf{a}^T C \mathbf{a}$$
$$\text{subject to} \quad \mathbf{a}^T E\left[\mathbf{x}\right] = \theta$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Finding the BLUE:*

Given the scalar parameter *θ,* the expected value of the samples can be assumed as

$$E[x[n]]=s[n]\theta,$$

where *s[n]* is known.

$$\sum_{n=0}^{N-1} a_n E\left[x\left[n\right]\right] = \sum_{n=0}^{N-1} a_n s\left[n\right]\theta = \theta, \qquad \text{(from 6.2)}$$

Thus the previous problem can be stated as

$$\min_{\mathbf{a} \in R^N} \quad \mathbf{a}^T \mathbf{C} \mathbf{a}$$
$$\text{s.t. } \mathbf{a}^T \mathbf{s} = 1$$

The method of Lagrangian multipliers can be used to solve this problem. Define the Lagrangian function as

$$J\left(\mathbf{a}, \lambda\right) = \mathbf{a}^T \mathbf{C} \mathbf{a} + \lambda\left(\mathbf{a}^T \mathbf{s} - 1\right), \qquad \lambda \in R$$

The gradient of *J* relative to **a** is

$$\frac{\partial J\left(\mathbf{a}, \lambda\right)}{\partial \mathbf{a}} = 2\mathbf{C}\mathbf{a} + \lambda\mathbf{s} = 0$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Finding the BLUE:*

Solving for **a** produces

$$\mathbf{a} = -\frac{\lambda}{2}\mathbf{C}^{-1}\mathbf{s}$$

Using the constrain as

$$\mathbf{a}^T\mathbf{s} = -\frac{\lambda}{2}\mathbf{s}^T\mathbf{C}^{-1}\mathbf{s} = 1, \quad \text{or} \quad \lambda = -\frac{2}{\mathbf{s}^T\mathbf{C}^{-1}\mathbf{s}} \quad.$$

Finally, the solution is

$$\mathbf{a}_{opt} = -\frac{\mathbf{C}^{-1}\mathbf{s}}{\mathbf{s}^T\mathbf{C}^{-1}\mathbf{s}}, \qquad \text{with a variance} \quad \text{var}\left(\hat{\theta}\right) = \mathbf{a}_{opt}^T\mathbf{C}\mathbf{a}_{opt} = \frac{\mathbf{s}^T\mathbf{C}^{-1}\mathbf{C}\mathbf{C}^{-1}\mathbf{s}}{\left(\mathbf{s}^T\mathbf{C}^{-1}\mathbf{s}\right)^2} = \frac{1}{\mathbf{s}^T\mathbf{C}^{-1}\mathbf{s}}.$$

Taking into account that *E[x]=sθ,* finally the estimator

$$\hat{\theta} = \mathbf{a}_{opt}^T\mathbf{x} = \frac{\mathbf{s}^T\mathbf{C}^{-1}\mathbf{x}}{\mathbf{s}^T\mathbf{C}^{-1}\mathbf{s}}, \qquad \text{Its expected value is} \quad E\left(\hat{\theta}\right) = \frac{\mathbf{s}^T\mathbf{C}^{-1}E\left[\mathbf{x}\right]}{\mathbf{s}^T\mathbf{C}^{-1}\mathbf{s}} = \frac{\mathbf{s}^T\mathbf{C}^{-1}\mathbf{s}\theta}{\mathbf{s}^T\mathbf{C}^{-1}\mathbf{s}} = \theta!$$

Thus it is unbiased, as required.

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Example:*

Example (DC level in white ~~Gaussian~~ noise revisited):

Model of signal: $x[n] = A + w[n], \qquad n = 0,...,N-1$

Where $w[n]$ is zero mean white noise with variance $\sigma^2$ (and an unspecified PDF), the problem is to estimate $A$.

Because $E[x[n]]=A$, we have $\mathbf{s=1}$, where $\mathbf{1}=[1\ 1\ ...\ 1]^T$.

The BLUE is $\hat{A} = \dfrac{\mathbf{1}^T \dfrac{1}{\sigma^2} \mathbf{Ix}}{\mathbf{1}^T \dfrac{1}{\sigma^2} \mathbf{I1}} = \dfrac{1}{N}\sum_{N=0}^{N-1} x[n] = \bar{x}$, and the variance is $\mathrm{var}(\hat{A}) = \dfrac{1}{\mathbf{1}^T \dfrac{1}{\sigma^2} \mathbf{I1}} = \dfrac{\sigma^2}{N}$.

Hence the sample mean is the BLUE independent of the PDF of data. It is the MVU estimator for the Gaussian case.

And in general: is it optimal?...

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Example:*

Example (DC level UNCORRELATED noise):

Model of signal: $x[n] = A + w[n], \qquad n = 0,...,N-1$

Where $w[n]$ is zero mean uncorrelated noise with $var(w[n])=\sigma^2$. Once again, the problem is to estimate $A$.

We have again $s=1$, and $C=diag(\sigma_0^2 \; \sigma_1^2 ... \sigma_{N-1}^2)$, and $C^{-1}=diag(1/\sigma_0^2 \; 1/\sigma_1^2 ... 1/\sigma_{N-1}^2)$..

The BLUE is
$$\hat{A} = \frac{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{x}}{\mathbf{1}^T \mathbf{C}^{-1} \mathbf{1}} = \frac{\sum_{N=0}^{N-1} \dfrac{x[n]}{\sigma_n^2}}{\sum_{N=0}^{N-1} \dfrac{1}{\sigma_n^2}}$$
and the variance is
$$\mathrm{var}\left(\hat{A}\right) = \frac{1}{\sum_{N=0}^{N-1} \dfrac{1}{\sigma_n^2}}$$

The BLUE weights those samples more heavily with smallest variances, in an attempt to equalize the noise contribution from each sample...

Is it optimal? In what cases?...

# *Extending BLUE to a Vector Parameter:*

To find the BLUE for a *p x 1* vector parameter, we constrain the estimator

- to be linear $\hat{\theta}_i = \sum_{n=0}^{N-1} a_{in} x[n]$     $i = 1, 2, ..., p$   or   $\hat{\theta} = \mathbf{A}\mathbf{x}$,     $\mathbf{A} \in R^{p \text{ x } N}$

- to be unbiased     $E\left[\hat{\boldsymbol{\theta}}\right] = \mathbf{A}E\left[\hat{\boldsymbol{\theta}}\right] = \boldsymbol{\theta}$

- to minimize its variance     $\mathrm{var}\left(\hat{\theta}_i\right) = E\left[\left(\sum_{n=0}^{N-1} a_{in} x[n] - E\left[\sum_{n=0}^{N-1} ai_n x[n]\right]\right)^2\right]$

The problem of finding the BLUE can be stated as, for

$$\min \quad \mathrm{var}\left(\hat{\theta}_i\right) = \mathbf{a}_i^T C \mathbf{a}_i$$

$$\text{subject to} \quad \mathbf{a}_i^T E\left[\mathbf{x}\right] = \theta$$

where

$$E\left[\mathbf{x}\right] = \mathbf{H}\boldsymbol{\theta}.$$

# *Extension to non-white Gaussian noise:*

*Theorem 6.1(Gauss-Markov Theorem)* – *If the data observed are of the general linear model form*

$$\mathbf{x} = \mathbf{H}\theta + \mathbf{w}$$

*where* **H** *is a known N x p observation matrix (with N>p) and rank p,* **x** *is a N x 1 vector of observations,* **θ** *is a p x 1 vector of parameters to be estimated, and* **w** *is an N x 1 noise vector with zero mean and covariance* **C** *(for an arbitrary PDF), then the BLUE is*

$$\hat{\boldsymbol{\theta}} = \left(\mathbf{H}^T\mathbf{C}^{-1}\mathbf{H}\right)^{-1}\mathbf{H}^T\mathbf{C}^{-1}\mathbf{x} \qquad (1bis)$$

*and the covariance matrix of estimate is*

$$C_{\hat{\theta}} = \left(\mathbf{H}^T\mathbf{C}^{-1}\mathbf{H}\right)^{-1}. \qquad (2bis)$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Bibliography:*

**Further reading**

• Thomas Kailath, *Linear Systems,* Prentice Hall, 1980.

• Thomas Kailath, Ali Sayed, and Babak Hassibi, *Linear Estimation,* Prentice Hall, 2000.

• Harry L. Van Trees, *Detection, Estimation, and Modulation Theory, Parts I to IV,* John Wiley, 2001.

• J. Bibby, H. Toutenburg, *Prediction and Improved Estimation in Linear Models,* John Wiley, 1977.

• C.Rao, *Linear Statistical Inference and Its Applications,* John Wiley, 1973.

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Advanced Control Systems*
# *Detection, Estimation, and Filtering*

### *Graduate Course on the*
### *MEng PhD Program*
### *Spring 2012/2013*

### *Chapter 6*
### *Maximum Likelihood Estimation*

### *Instructor:*
### *Prof. Paulo Jorge Oliveira*
### *p.oliveira@dem.ist.utl.pt or pjcro @ isr.ist.utl.pt*
### *Phone: 21 8419511 or 21 8418053 (3511 or 2053 inside IST)*

**DEM**
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Syllabus:*

Classical Estimation Theory

…

Chap. 5 - **Best Linear Unbiased Estimators** [1 week]

Definition of BLUE estimators;    White Gaussian noise and bandlimited systems;
Examples; Generalized MVU estimation;

## Chap. 6 - *Maximum Likelihood Estimation* [1 week]

The maximum likelihood estimator; Properties of the ML estimators;
Solution for ML estimation; Examples; Monte-Carlo methods;

Chap. 7 - **Least Squares** [1 week]

The least squares approach; Linear and nonlinear least squares; Geometric
interpretation; Constrained least squares;  Examples;

continues…

# *Motivating example:*

Example (DC level in white Gaussian noise modified):

For this example the methods previously introduced will not work…

Signal model: $\qquad x\left[n\right] = A + w\left[n\right], \qquad n = 0,...,N-1$

Where A is the unknown level to be estimated and *w[n]* is zero mean white Gaussian with unknown variance A.

First, lets try to find the CRLB. The PDF is:

$$p\left(\mathbf{x};A\right) = \frac{1}{\left(2\pi A\right)^{N/2}} \exp\left( -\frac{1}{2A} \sum_{n=0}^{N-1}\left(x[n] - A\right)^2 \right) \qquad (1)$$

The derivative of the log-likelihood function is

$$\frac{\partial}{\partial A} \ln p\left(\mathbf{x};A\right) = -\frac{N}{2A} + \frac{1}{A}\sum_{n=0}^{N-1}\left(x[n]-A\right) + \frac{1}{2A^2}\sum_{n=0}^{N-1}\left(x[n]-A\right)^2$$

$$\overset{?}{=} I\left(A\right)\left(g\left(\mathbf{x}\right) - A\right)$$

It appears that it is not possible…
So, an efficient estimator does not exist.

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *An example:*

Example (DC level in white Gaussian noise modified) (cont):

However, from the second derivative, it is possible to compute the CRLB to be

$$\mathrm{var}\left(\hat{A}\right) \geq \frac{A^2}{N\left(A + 1/2\right)}.$$

Secondly, to find the MVU estimator based on the theory of sufficient statistics, one must factorize (1) in the form

$$p\left(\mathbf{x};\theta\right) = g\left(T\left(\mathbf{x}\right),\theta\right)h\left(\mathbf{x}\right)$$

It is possible, if one considers

$$p\left(\mathbf{x};A\right) = \underbrace{\frac{1}{\left(2\pi A\right)^{N/2}}\exp\left(-\frac{1}{2}\left(\frac{1}{A}\sum_{n=0}^{N-1}x^2[n] + NA\right)\right)}_{g\left(\sum_{n=0}^{N-1}x^2[n], A\right)}\underbrace{\exp\left(\sum_{n=0}^{N-1}x[n]\right)}_{h\left(\mathbf{x}\right)}$$

# *An example:*

Example (DC level in white Gaussian noise modified) (cont):

So a sufficient statistics is

$$T\left(\mathbf{x}\right) = \sum_{n=0}^{N-1} x^2[n]$$

It is required to find a function of the sufficient statistics that produces an unbiased estimator, i.e.

$$E\left[g\left(\sum_{n=0}^{N-1} x^2[n]\right)\right] = A$$

Taking into account the auxiliary result

$$\mathrm{var}\left(x\left[n\right]\right) = E\left[\left(x\left[n\right] - E\left(x\left[n\right]\right)\right)^2\right] = E\left[x^2\left[n\right]\right] - 2E\left(x\left[n\right]\right)E\left(x\left[n\right]\right) + E^2\left(x\left[n\right]\right)$$

We have that

$$E\left[x^2\left[n\right]\right] = \mathrm{var}\left(x\left[n\right]\right) + E^2\left(x\left[n\right]\right) \qquad \text{(in our case } E\left[x^2\left[n\right]\right] = A + A^2\text{)}$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *An example:*

Example (DC level in white Gaussian noise modified) (cont):

Since

$$E\left[\sum_{n=0}^{N-1} x^2[n]\right] = NE\left[\sum_{n=0}^{N-1} x^2[n]\right] = N\left[\text{var}(x[n]) + E^2(x[n])\right] = N\left[A + A^2\right]!$$

It is impossible to find a solution for a generic unknown parameter $A$, i.e.

$$N\left[A + A^2\right] \neq A!$$

A final alternative is to find the optimal estimator would be to determine

$$E\left[\hat{A} \mid \sum_{n=0}^{N-1} x^2[n]\right] = \text{???}$$

That appears to be a formidable task!

We exhausted the optimal approaches studied… We can propose other estimators, but without any guarantee of optimality.

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *An example:*

Example (DC level in white Gaussian noise modified) (cont):

Those estimators should be at least *approximately optimal,* i.e.

$$E\left[\hat{A}\right] \to A$$
$$\mathrm{var}\left(\hat{A}\right) \to \mathrm{CRLB}$$

For instance, lets consider the estimator (why? explanation will be provided next…)

$$\hat{A} = -\frac{1}{2} + \sqrt{\frac{1}{N}\sum_{n=0}^{N-1} x^2[n] + \frac{1}{4}}$$

This estimator is biased, since

$$E\left[\hat{A}\right] = E\left[-\frac{1}{2} + \sqrt{\frac{1}{N}\sum_{n=0}^{N-1} x^2[n] + \frac{1}{4}}\right] \neq -\frac{1}{2} + \sqrt{E\left[\frac{1}{N}\sum_{n=0}^{N-1} x^2[n]\right] + \frac{1}{4}} =$$

$$= -\frac{1}{2} + \sqrt{A + A^2 + \frac{1}{4}} = A!$$

But it can be verified that is is consistent, i.e.

$$\frac{1}{N}\sum_{n=0}^{N-1} x^2[n] \to E\left[x^2[n]\right] = A + A^2 \qquad \text{and therefore} \quad \hat{A} \to A$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *An example:*

Example (DC level in white Gaussian noise modified) (cont):

Consider that $\hat{A} = g(u)$, where $g(u) = -\dfrac{1}{2} + \sqrt{u + \dfrac{1}{4}}$ and lets linearise this function, near $u_0 = E[u] = A + A^2$.

$$g(u) \approx g(u_0) + \left.\frac{dg(u)}{du}\right|_{u=u_0} (u - u_0) \qquad \text{(using Taylor's series expansion)}$$

$$\hat{A} \approx A + \frac{\dfrac{1}{2}}{A + \dfrac{1}{2}} \left[ \frac{1}{N} \sum_{n=0}^{N-1} x^2[n] - \left( A + A^2 \right) \right]$$

$$E\left[ \hat{A} \right] \approx A.$$

Thus this estimator is asymptotically unbiased.

And what about its variance?…

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *An example:*

Example (DC level in white Gaussian noise modified) (cont):

It is given by

$$\operatorname{var}\left(\hat{A}\right) \approx \left(\frac{\frac{1}{2}}{A + \frac{1}{2}}\right)^2 \operatorname{var}\left[\frac{1}{N}\sum_{n=0}^{N-1} x^2[n] - \left(A + A^2\right)\right] \approx \frac{\frac{1}{4}}{N\left(A + \frac{1}{2}\right)^2}\operatorname{var}\left(x^2[n]\right)$$

But $var(x^2[n]) = 4A^3 + 2A^2$, so that

$$\operatorname{var}\left(\hat{A}\right) \approx \frac{\frac{1}{4}}{N\left(A + \frac{1}{2}\right)^2} 4A^2\left(A + \frac{1}{2}\right) \approx \frac{A^2}{N\left(A + \frac{1}{2}\right)}$$

Thus this estimator asymptotically equals the CRLB!!!

Discuss the impact of one such methodology that provides asymptotic results.
The value for science and for engineering

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *An asymptotically optimal solution:*

**What to do, if the MVU estimator does not exist or can not be found?**

**An alternative consists of exploiting the…**

**Maximum Likelihood Principle.**

**It can be understood as a "turn the crank" method.**

**Only suboptimal performance can be achieved.**

It is the most popular approach to obtaining practical estimators.

Its optimality is verified  for large enough data sets.

**DEM**
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Motivating example revisited:*

Example (DC level in white Gaussian noise modified):

The method consists only on the computation of the maximum of the (log) likelihood function. In our case, it is required to solve:

$$\frac{\partial}{\partial A}\ln p\left(\mathbf{x};A\right) = -\frac{N}{2A} + \frac{1}{A}\sum\nolimits_{n=0}^{N-1}\left(x[n]-A\right) + \frac{1}{2A^2}\sum\nolimits_{n=0}^{N-1}\left(x[n]-A\right)^2 = 0$$

$$= -\frac{N}{2A} + \frac{1}{A}\sum\nolimits_{n=0}^{N-1}x[n] - \frac{1}{A}NA + \frac{1}{2A^2}\sum\nolimits_{n=0}^{N-1}\left(x^2[n]-2Ax[n]+A^2\right) =$$

$$= -\frac{N}{2A} + \frac{1}{A}\sum\nolimits_{n=0}^{N-1}x[n] - N + \frac{1}{2A^2}\sum\nolimits_{n=0}^{N-1}x^2[n] - \frac{1}{2A^2}2A\sum\nolimits_{n=0}^{N-1}x[n] + \frac{1}{2A^2}NA^2 =$$

$$= -\frac{N}{2A} - \frac{N}{2} + \frac{1}{2A^2}\sum\nolimits_{n=0}^{N-1}x^2[n] = -\frac{A^2 + A - \dfrac{1}{N}\sum\nolimits_{n=0}^{N-1}x^2[n]}{2A^2N} = 0$$

From where our previous unexplained estimator results

$$\hat{A} = -\frac{1}{2} + \sqrt{\frac{1}{N}\sum\nolimits_{n=0}^{N-1}x^2[n] + \frac{1}{4}}$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Maximum Likelihood Principle:*

***Theorem 7.1 (Asymptotic Properties of the MLE)*** *– If the PDF $p(x;\theta)$ of the data $x$ satisfies some regularity conditions, then the MLE of the unknown parameter is asymptotically distributed (for large data records) according to*

$$\hat{\theta} \overset{a}{\sim} N\left(\theta, I^{-1}(\theta)\right)$$

*where $I(\theta)$ is the Fisher information evaluated at the true value of the unknown parameter*

**In practice it is seldom known in advance how large N must be.**

**Analytical expression for the PDF of the MLE is usually impossible to derive.**

**Thus, to assess the MLE performance, computer simulations are usual.**

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Properties of MLE:*

*Proof outline:*

*The following regularity conditions are assumed:*

*1) The first and second-order derivative of the log-likelihood are well defined.*

*2)*
$$E\left[\frac{\partial \ln p\big(x[n];\theta\big)}{\partial \theta}\right] = 0$$

*First, it is required to show that the MLE is consistent. Related with the Kullbak_Leibner information (and also with measure of the difference between two probability distributions)*

$$\int \ln\left[\frac{p\big(x[n];\theta_1\big)}{p\big(x[n];\theta_2\big)}\right] p\big(x[n];\theta_1\big)dx[n] \geq 0 \qquad (1)$$

*Where equality occurs for* **θ₁=θ₂.**

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Properties of MLE:*

*Proof outline:*

*Now, maximizing the log-likelihood*

$$\frac{1}{N}\ln p(\mathbf{x};\theta) = \frac{1}{N}\sum_{n=0}^{N-1}\ln p(x[n];\theta) \rightarrow \int \ln p(x[n];\theta) p(x[n];\theta_0) dx[n]$$

*Where the last relation is due to the fact that, by the law of large numbers, it converges to the expected value. The MLE is <span style="color:red">consistent</span> and is maximized for $\hat{\theta} = \theta_0$, i.e.*

$$\int \ln p(x[n];\theta_0) p(x[n];\theta_0) dx[n] \geq \int \ln p(x[n];\theta_1) p(x[n];\theta_0) dx[n]$$

*Moreover is the maximum, due to suitable continuity argument and the relation (1). Using the Taylor series expansion, one obtains*

$$\left.\frac{\partial \ln p(\mathbf{x};\theta)}{\partial \theta}\right|_{\theta=\hat{\theta}} \approx \left.\frac{\partial \ln p(\mathbf{x};\theta)}{\partial \theta}\right|_{\theta=\theta_0} + \left.\frac{\partial^2 \ln p(\mathbf{x};\theta)}{\partial \theta^2}\right|_{\theta=\theta_0} (\hat{\theta} - \theta_0) \approx 0$$

*Where the last quantity is approx. 0 if near an maximum.*

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Properties of MLE:*

*Proof outline:*

*This relation can therefore be approximately written as*

$$\sqrt{N}\left(\hat{\theta}-\theta_0\right) = \frac{\left.\frac{1}{\sqrt{N}}\frac{\partial \ln p\left(\mathbf{x};\theta\right)}{\partial \theta}\right|_{\theta=\theta_0}}{\left.-\frac{1}{N}\frac{\partial^2 \ln p\left(\mathbf{x};\theta\right)}{\partial \theta^2}\right|_{\theta=\hat{\theta}}} \rightarrow \frac{\left.\frac{1}{\sqrt{N}}\sum_{n=0}^{N-1}\frac{\partial \ln p\left(x\left[n\right];\theta\right)}{\partial \theta}\right|_{\theta=\theta_0}}{\left.-\frac{1}{N}\sum_{n=0}^{N-1}\frac{\partial^2 \ln p\left(x\left[n\right];\theta\right)}{\partial \theta^2}\right|_{\theta=\hat{\theta}}} \sim N\left(0,i^{-1}\left(\theta_0\right)\right)$$

*From where it can be concluded, using the law of large numbers and the IID of the samples, that*

$$\hat{\boldsymbol{\theta}} \overset{a}{\sim} N\left(\boldsymbol{\theta}_0,\boldsymbol{I}^{-1}\left(\boldsymbol{\theta}_0\right)\right)$$

# *MLE PDF:*

In general is very difficult (or impossible) to obtain the PDF of the MLE.

How to study its performance?

**Use Monte Carlo Method**

1 . Simulate the noise characteristics, the signal model, and compute the estimates.

2. Repeat M times these realizations. (How to select M?)

3. Compute the experiments ensemble mean and covariance, using

$$\widehat{E\left[\hat{A}\right]} = \frac{1}{M}\sum\nolimits_{i=1}^{M}\hat{A}_{i}$$

$$\widehat{\mathrm{var}\left(\hat{A}\right)} = \frac{1}{M}\sum\nolimits_{i=1}^{M}\left(\hat{A}_{i} - \widehat{E\left[\hat{A}\right]}\right)^{2}$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Invariance Property:*

**Theorem 7.2 (Invariance Property of the MLE)** – *The MLE of the parameter α=g(θ),*

*where the PDF* $p(\boldsymbol{x};\theta)$ *is parameterized by θ, is given by*

$$\hat{\alpha} = g\left(\hat{\theta}\right)$$

*Where $\hat{\theta}$ is the MLE of θ. The MLE is obtained by maximization of* $p(\boldsymbol{x};\theta)$, *If g is not a one-*

*to-one function, then $\hat{\alpha}$ maximized the modified likelihood fuction*

$$\overline{p}_T\left(\mathbf{x};\alpha\right) = \max_{\left\{\theta \,:\, \alpha \,=\, g(\theta)\right\}} p\left(\mathbf{x};\theta\right).$$

*Proof outline (simple case: g() one to one WGN, IID, expected value):*

*The MLE for the transformed parameter can be found minimizing the log-likelihood, i.e.*

$$\frac{\partial}{\partial\alpha}\sum_{n=0}^{N-1}\left(x[n] - g^{-1}\left(\alpha\right)\right)^2 = k + k'\sum_{n=0}^{N-1}\left(x[n] - g^{-1}\left(\alpha\right)\right)\frac{\partial}{\partial\alpha}g^{-1}\left(\alpha\right) = 0, \qquad k,k' > 0 \quad.$$

*Thus*

$$\sum_{n=0}^{N-1}x[n] - Ng^{-1}\left(\alpha\right) = 0, \qquad g^{-1}\left(\alpha\right) = \frac{1}{N}\sum_{n=0}^{N-1}x[n] = \overline{x} \qquad \alpha = g\left(\overline{x}\right).$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Numerical Determination of the MLE:*

The MLE in general can not be found in close form.

**But it can be found numerically. Grid search, gradient or Newton methods can be used.**

**Conditions for nonlinear optimization methods are central to that discussion.**

**For different data-sets, the target function changes and thus also the maximum changes.**

**In general there is not or maximum, but a number of local maxima.**

**How to avoid attraction to local maxima? Regions of attraction?...**

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Motivating example:*

Example (Exponential in white Gaussian noise):

Signal model: $x[n] = r^n + w[n], \qquad n = 0, ..., N-1$

Where $w[n]$ is zero mean white Gaussian noise with variance $\sigma^2$ and the exponential factor $r$ is to estimated.

For the likelihood function, the MLE is the value of $r$ that maximizes is :

$$p(\mathbf{x}; A) = \frac{1}{\left(2\pi\sigma^2\right)^{N/2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{n=0}^{N-1} \left(x[n] - r^n\right)^2\right) \qquad (1)$$

Or, equivalently, the value that minimizes

$$J(r) = \sum_{n=0}^{N-1} \left(x[n] - r^n\right)^2.$$

Differentiating $J(r)$ and setting to zero produces

$$\frac{\partial J(r)}{\partial r} = 2\sum_{n=0}^{N-1} \left(x[n] - r^n\right) nr^{n-1}.$$

It is a nonlinear equation in $r$ and cannot be solved directly.

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Numerical Solution (basics):*

*The use of iterative methods to maximize the log-likelihood function is an example of application of nonlinear optimization methods. See a good book (or class) on the field…*

$$\frac{\partial \ln p(\mathbf{x};\theta)}{\partial \theta} = g(\theta) = 0$$

*For instance, one of the most basic method, is the Newton-Raphson method. From an initial guess $\Theta_0$, and from a Taylor series expansion results*

$$g(\theta) \approx g(\theta_0) + \left.\frac{\partial g(\theta)}{\partial \theta}\right|_{\theta=\theta_0} (\theta - \theta_0) \approx 0$$

*The following recursion results* $\theta$

$$\theta_{k+1} = \theta_k - \left[\left.\frac{\partial^2 \ln p(\mathbf{x};\theta)}{\partial \theta^2}\right|_{\theta=\theta_k}\right]^{-1} \left.\frac{\partial \ln p(\mathbf{x};\theta)}{\partial \theta}\right|_{\theta=\theta_k}$$

# *Motivating example:*

Example (Exponential in white Gaussian noise):

Computer simulation

N=50, r=0.5, and $\sigma^2$=0.01

Maximum at r=0.493
(a specific realization)

| Iteration | Initial Guess, $r_0$ | | |
|---|---|---|---|
| | 0.8 | 0.2 | 1.2 |
| 1 | 0.723 | 0.799 | 1.187 |
| 2 | 0.638 | 0.722 | 1.174 |
| 3 | 0.561 | 0.637 | 1.161 |
| 4 | 0.510 | 0.560 | 1.148 |
| 5 | 0.494 | 0.510 | 1.136 |
| 6 | 0.493 | 0.494 | 1.123 |
| 7 | | 0.493 | 1.111 |
| 8 | | | 1.098 |
| 9 | | | 1.086 |
| 10 | | | 1.074 |
| $\vdots$ | | | $\vdots$ |
| 29 | | | 0.493 |

$$-\sum_{n=0}^{N-1}(x[n]-r^n)^2$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Numerical Solution (basics):*

*The importance of*

    **stability conditions,**

    **convergence rates, and**

    **domains of attraction**

*can hardly be overemphasized. Engineering/scientific content…*

*Other methods mentioned:*

    *Scoring*

    *Expectation / maximization (nice term paper subject)*

# *Invariance Property:*

**Theorem 7.5 (Optimality of the MLE for the Linear Model)** – *If the observed data $\mathbf{x}$ are*

*described by the general linear model*

$$\mathbf{x} = \mathbf{H}\theta + \mathbf{w}$$

*where H is a known N x p matrix with N>p and of rank p, θ is a p x 1 parameter vector to*

*be estimated, and w is the noise vector with PDF N(0,C), the the MLE of θ is*

$$\hat{\theta} = \left(\mathbf{H}^{\mathrm{T}}\mathbf{C}^{-1}\mathbf{H}\right)^{-1}\mathbf{H}^{\mathrm{T}}\mathbf{C}^{-1}\mathbf{x}.$$

*And is also an efficient estimator in that it attains the CRLB and hence is the MVU*

*estimator. The PDF of θ is*

$$\hat{\theta} \sim N\left(\theta, \left(\mathbf{H}^{\mathrm{T}}\mathbf{C}^{-1}\mathbf{H}\right)^{-1}\right).$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Method of Scoring:*

*The method of scoring is based on the approximation for one element found also in the Newton-Raphson method. Note that for IID samples we have*

$$\frac{\partial^2 \ln p(\mathbf{x};\theta)}{\partial \theta^2} = \sum_{n=0}^{N-1} \frac{\partial^2 \ln p(x[n];\theta)}{\partial \theta^2} = NE\left[\frac{\partial^2 \ln p(x[n];\theta)}{\partial \theta^2}\right] = -Ni(\theta) = -I(\theta).$$

*So the iterations on NR method can be transformed in*

$$\theta_{k+1} = \theta_k - I^{-1}(\theta) \left.\frac{\partial \ln p(\mathbf{x};\theta)}{\partial \theta}\right|_{\theta=\theta_k}$$

*Resulting in a method that is more stable. However it suffers from the same convergence problems as the NR method.*

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Maximum Likelihood Principle:*

*Theorem 7.1 (Asymptotic Properties of the MLE – Vector Parameter)* – *If the PDF $p(x;\theta)$ of the data **x** satisfies some "regularity" conditions, then the MLE of the unknown parameter **θ** is asymptotically distributed (for large data records) according to*

$$\hat{\boldsymbol{\theta}} \overset{a}{\sim} N\left(\boldsymbol{\theta}, \boldsymbol{I}^{-1}(\boldsymbol{\theta})\right)$$

*where $I(\boldsymbol{\theta})$ is the Fisher information evaluated at the true value of the unknown parameter*

**In practice it is seldom known in advance how large N must be.**

**In the cases where the number of parameters increases, relative to the number of samples available, the assumptions fails and the MLE estimator can provide very poor estimates.**

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Bibliography:*

**Further reading**

• Harry L. Van Trees, *Detection, Estimation, and Modulation Theory, Parts I to IV,* John Wiley, 2001.

• Anthony William Fairbank Edwards, **Likelihood,** Cambridge University Press, 1972.

• Jerry M. Mendel, **Lessons in Digital Estimation Theory,** Prentice Hall, 1987.

• Geoffrey J. McLachlan (Author), Thriyambakam Krishnan, **The EM Algorithm and Extensions,** Wiley, 1997.

**DEM**
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# Advanced Control Systems
# Detection, Estimation, and Filtering

## Graduate Course on the
## MEng PhD Program
## Spring 2012/2013

## Chapter 7
## Least Squares

*Instructor:*
*Prof. Paulo Jorge Oliveira*
*p.oliveira@dem.ist.utl.pt or pjcro @ isr.ist.utl.pt*
*Phone: 21 8419511 or 21 8418053 (3511 or 2053 inside IST)*

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Syllabus:*

Classical Estimation Theory

...

Chap. 6 - *Maximum Likelihood Estimation* [1 week]

The maximum likelihood estimator; Properties of the ML estimators; Solution for ML estimation; Examples; Monte-Carlo methods;

Chap. 7 - *Least Squares* [1 week]

The least squares approach; Linear and nonlinear least squares; Geometric interpretation; Constrained least squares; Examples;

Chap. 8 – *Bayesian Estimation* [1 week]

Philosophy and estimator design; Prior knowledge; Bayesian linear model; Bayesian estimation on the presence of Gaussian pdfs; Minimum Mean Square Estimators;

continues…

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Least Squares Approach:*



$$J(\theta) = \sum_{n=0}^{N-1} \left( x[n] - s[n] \right)^2 = \sum_{n=0}^{N-1} \varepsilon^2[n]$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Least Squares Approach:*

The least squares estimator (LSE) is obtained minimizing the LS error criterion

$$J(\theta) = \sum_{n=0}^{N-1} \left( x[n] - s[n] \right)^2 = \sum_{n=0}^{N-1} \varepsilon^2[n]$$

where the dependency on θ is via s[n].

Note:

No probabilistic assumptions have been made about the data x[n];

Method valid both for Gaussian and for non-Gaussian disturbances;

Performance optimality of the LSE can not be guaranteed;

Method applied when:

a precise statistical characterization of the data is unknown;

optimal estimator can not be found;

…

# *Linear Least Squares:*

The least squares approach for a **scalar parameter**, we must assume

$$s[n] = \theta h[n].$$

The criterion to minimize is

$$J(\theta) = \sum_{n=0}^{N-1} \left( x[n] - \theta h[n] \right)^2$$



It is immediate that

$$\frac{\partial J(\theta)}{\partial \theta} = -2 \sum_{n=0}^{N-1} \left( x[n] - \theta h[n] \right) h[n] = 0$$

With a solution given by

$$\hat{\theta} = \frac{\sum_{n=0}^{N-1} x[n] h[n]}{\sum_{n=0}^{N-1} h^2[n]}.$$

Thus the minimum cost of the criterion verifies

$$0 < J_{min}(\theta) = \sum_{n=0}^{N-1} x^2[n] - \frac{\left( \sum_{n=0}^{N-1} x[n] h[n] \right)^2}{\sum_{n=0}^{N-1} h^2[n]} < \sum_{n=0}^{N-1} x^2[n].$$

# *Linear Least Squares:*

The extension of the least squares approach for a **vector parameter** is immediate.

For the signal $s = [s[0] \; s[1] \; ... \; s[N-1]]$
The criterion to minimize is

$$J(\theta) = \sum_{n=0}^{N-1} \left( x[n] - s[n] \right)^2 = \left( \mathbf{x} - \mathbf{H}\theta \right)^T \left( \mathbf{x} - \mathbf{H}\theta \right) = \mathbf{x}^T \mathbf{x} - 2\mathbf{x}^T \mathbf{H}\theta + \theta^T \mathbf{H}^T \mathbf{H}\theta.$$

The gradient is

$$\frac{\partial J(\theta)}{\partial \theta} = -2\mathbf{H}^T \mathbf{x} + 2\mathbf{H}^T \mathbf{H}\theta.$$

With a solution given by

$$\hat{\theta} = \left( \mathbf{H}^T \mathbf{H} \right)^{-1} \mathbf{H}^T \mathbf{x}$$

The minimum cost of the criterion verifies

$$0 < J_{min}(\theta) = \mathbf{x}^T \mathbf{x} - \mathbf{x}^T \mathbf{H} \left( \mathbf{H}^T \mathbf{H} \right)^{-1} \mathbf{H}^T \mathbf{x} < \mathbf{x}^T \mathbf{x}.$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÃNICA
TÉCNICO LISBOA

# *Geometrical Interpretation:*

Note that the solution obtained

$$\hat{\theta} = \left(\mathbf{H}^T\mathbf{H}\right)^{-1}\mathbf{H}^T\mathbf{x}$$

can be rewritten as

$$\left(\mathbf{H}^T\mathbf{H}\right)\theta = \left(\mathbf{H}^T\mathbf{H}\right)\left(\mathbf{H}^T\mathbf{H}\right)^{-1}\mathbf{H}^T\mathbf{x}$$

$$\left(\mathbf{H}^T\mathbf{H}\right)\theta = \mathbf{H}^T\mathbf{x}$$

$$\mathbf{H}^T\left(\mathbf{H}\theta - \mathbf{x}\right) = 0$$

Denoting as the error vector $\varepsilon = \mathbf{H}\theta - \mathbf{x}$, the previous expression can be interpreted as that the error vector must be orthogonal to the columns of **H.**

# *Extensions to Least Squares:*

Other extensions of the least squares approach are also very popular

Weighted Least Squares:

criterion $\quad J_W(\boldsymbol{\theta}) = (\mathbf{x} - \mathbf{H}\boldsymbol{\theta})^T \mathbf{W}(\mathbf{x} - \mathbf{H}\boldsymbol{\theta})$

solution $\quad \hat{\boldsymbol{\theta}} = (\mathbf{H}^T \mathbf{W} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{W} \mathbf{x}$

minimum $\quad 0 < J_{min}(\boldsymbol{\theta}) = \mathbf{x}^T (\mathbf{W} - \mathbf{W}\mathbf{H}(\mathbf{H}^T \mathbf{W}\mathbf{H})^{-1} \mathbf{H}^T \mathbf{W})\mathbf{x} < \mathbf{x}^T \mathbf{W} \mathbf{x}.$

$\mathbf{W}$ can be set as the inverse covariance matrix, leading to an optimal solution in the case of correlated Gaussian noise.

Order-recursive Least Squares (see pp. 232)

same criterion but the observation and parameter matrices vary their length

$$\mathbf{H}_{k+1} = \begin{bmatrix} \mathbf{H}_k & h_{k+1} \end{bmatrix} = \begin{bmatrix} N \text{ x } k & N \text{ x } 1 \end{bmatrix}$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Extensions to Least Squares:*

Order-recursive Least Squares (cont.)

solution

$$\hat{\boldsymbol{\theta}}_{k+1} = \begin{bmatrix} \hat{\boldsymbol{\theta}}_k - \dfrac{\left(\mathbf{H}_k^T\mathbf{H}_k\right)^{-1}\mathbf{H}_k^T\mathbf{h}_{k+1}\mathbf{h}_{k+1}^T\mathbf{P}_k^{\perp}\mathbf{x}}{\mathbf{h}_{k+1}^T\mathbf{P}_k^{\perp}\mathbf{h}_{k+1}} \\[2em] \dfrac{\mathbf{h}_{k+1}^T\mathbf{P}_k^{\perp}\mathbf{x}}{\mathbf{h}_{k+1}^T\mathbf{P}_k^{\perp}\mathbf{h}_{k+1}} \end{bmatrix} = \begin{bmatrix} k \text{ x } 1 \\ 1 \text{ x } 1 \end{bmatrix}$$

where

$$\mathbf{P}_k^{\perp} = I - \mathbf{H}_k\left(\mathbf{H}_k^T\mathbf{H}_k\right)^{-1}\mathbf{H}_k^T$$

minimum

$$J_{\min}\left(\theta_{k+1}\right) = J_{\min}\left(\theta_k\right) - \dfrac{\left(\mathbf{h}_{k+1}^T\mathbf{P}_k^{\perp}\mathbf{x}\right)^2}{\mathbf{h}_{k+1}^T\mathbf{P}_k^{\perp}\mathbf{h}_{k+1}}$$
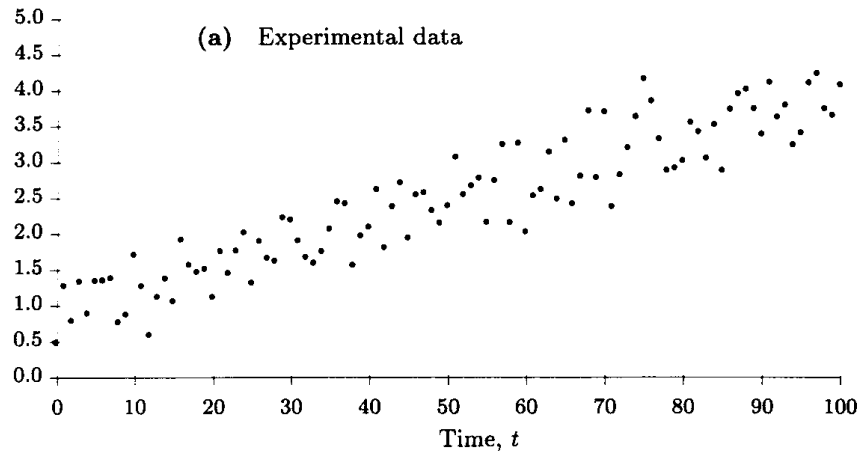
Example:
Line fitting

$$s_1[n] = A_1 \qquad s_2[n] = A_2 + B_c n \qquad H_1 = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \qquad H_2 = \begin{bmatrix} H_1 & \begin{matrix} 0 \\ 1 \\ \vdots \\ N-1 \end{matrix} \end{bmatrix}$$

# *Example:*



(a) Experimental data

(b) One-parameter fit

$\hat{s}_1(t) = \hat{A}_1$

(c) Two-parameter fit

$\hat{s}_2(t) = \hat{A}_2 + \hat{B}_2 t$

Constant

Line

Parabola

Cubic

$J_{min}$

Number of parameters, $k$

PO 1213

# *Sequential Least Squares:*

In many estimation, detection, or identification problems data are obtained as samples of the output of a process.

It would be advantageous that the least squares solution could be written as a recursive solution.

Lets revisit our old DC level in Gaussian noise example:

At time N-1, the data set available is $x$=[x[0] x[1] … x[N-1]] and the MVU estimator solution is given by

$$\hat{A}[N-1] = \frac{1}{N}\sum_{n=0}^{N-1} x[n]$$

If a new sample is obtained, i.e. x[n] is available, the estimator is given by

$$\hat{A}[N] = \frac{1}{N+1}\sum_{n=0}^{N} x[n] = \frac{1}{N+1}\left(\sum_{n=0}^{N-1} x[n] + x[N]\right) = \frac{N}{N+1}\hat{A}[N-1] + \frac{1}{N+1} x[N]$$

That can be rewritten as

$$\hat{A}[N] = \hat{A}[N-1] + \frac{1}{N+1}\left(x[N] - \hat{A}[N-1]\right).$$

Much remains to be said, see next chapters…

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Sequential Least Squares:*

$$\hat{A}[N] = \hat{A}[N-1] + \frac{1}{N+1}\left(x[N] - \hat{A}[N-1]\right)$$

Recursive solution

Correction term, reflecting that with more one sample more is known on the parameter.

The gain is decreasing thus preserving a memory on the past samples.

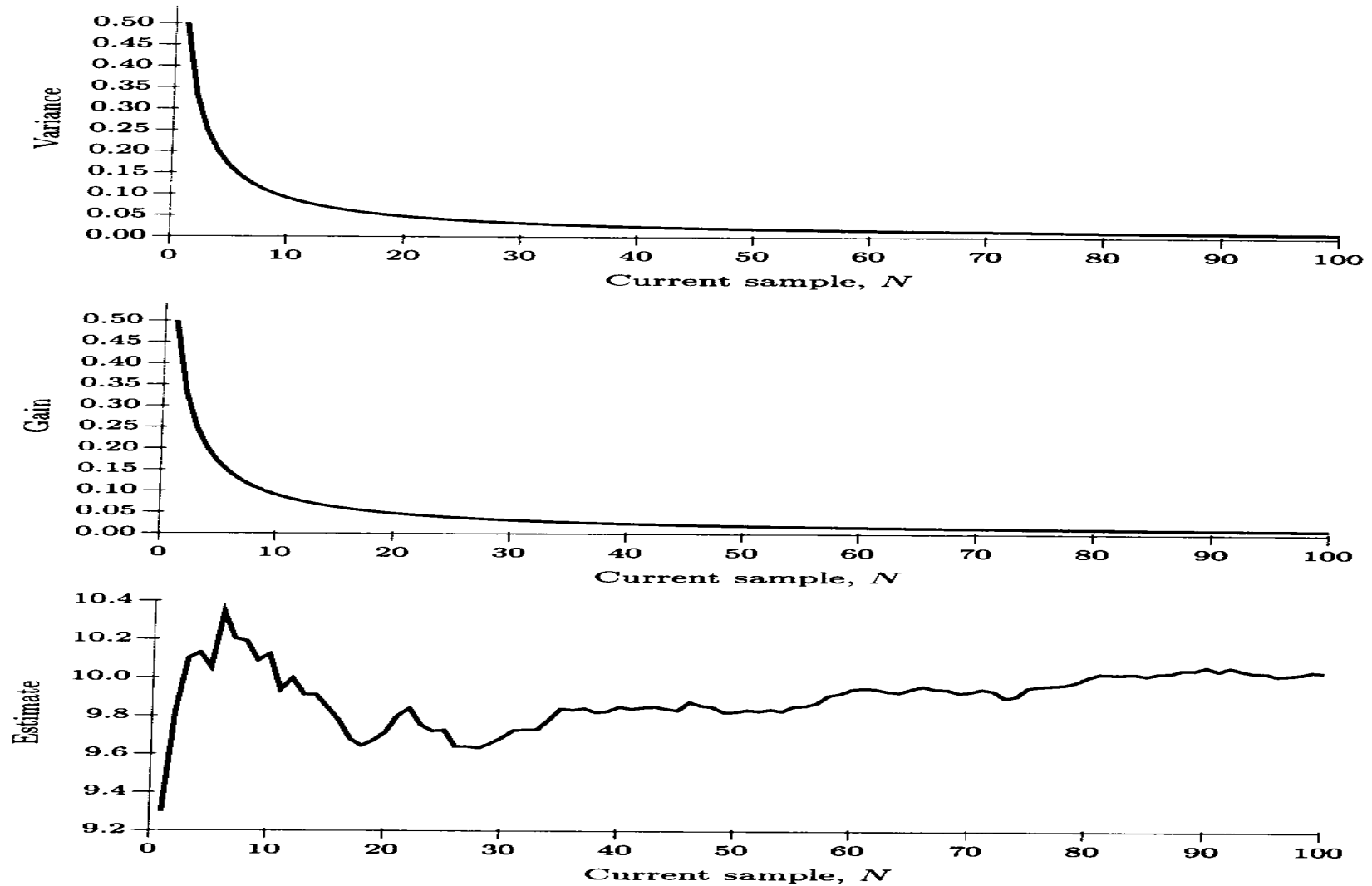The value of the criterion can also be written recursively, i.e.

$$J_{\min}(N) = J_{\min}(N-1) + \frac{N}{N+1}\left(x[N] - \hat{A}[N-1]\right)^2$$

Seems a paradox, but if our fitting is parfait does not increases…

More points to be fitted with the same number of parameters.

## It is an OPTIMAL solution!

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Sequential Least Squares:*

# *Sequential Least Squares:*

The optimal solution, in the case where a Gaussian noise occurs, with time varying variance

Signal Model $\qquad x[n]=\boldsymbol{h}[n]\boldsymbol{\theta}, \qquad n=0,...,N-1,...$

Estimator Update:

$$\hat{\boldsymbol{\theta}}\left[n\right]=\hat{\boldsymbol{\theta}}\left[n-1\right]+\mathbf{K}\left[n\right]\left(x\left[N\right]-h^{T}\left[n\right]\hat{\boldsymbol{\theta}}\left[n-1\right]\right)$$

Where

$$\mathbf{K}\left[n\right]=\frac{\Sigma\left[n-1\right]h\left[n\right]}{\sigma_{n}^{2}+h^{T}\left[n\right]\Sigma\left[n-1\right]h\left[n\right]}$$

Covariance Update:

$$\Sigma\left[n\right]=\left(\mathbf{I}-\mathbf{K}\left[n\right]h^{T}\left[n\right]\right)\Sigma\left[n-1\right]$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Sequential Least Squares:*

The signal model and the parameter estimation problem can be interpreted resorting to the dynamic model

$$\theta[n+1] = \theta[n]$$
$$x[n] = \mathbf{h}^T[n]\theta[n] + w[n]$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Constrained Least Squares:*

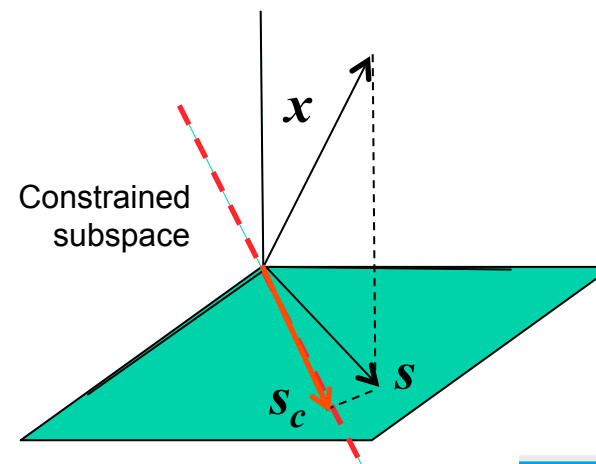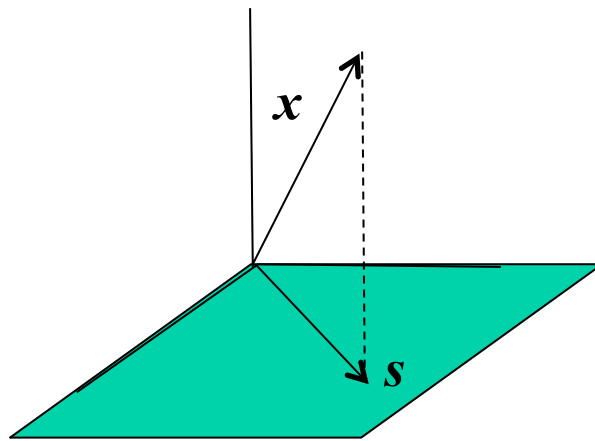This alternative method can be very useful if the problem at hand verifies some properties.

criterion
$$J_C(\theta) = (\mathbf{x} - \mathbf{H}\theta)^T (\mathbf{x} - \mathbf{H}\theta)$$

$$s.t. \qquad \mathbf{A}\theta = \mathbf{b}$$

solution $\hat{\theta}_C = \hat{\theta} - (\mathbf{H}^T\mathbf{H})^{-1} \mathbf{A}^T \left( \mathbf{A}(\mathbf{H}^T\mathbf{H})^{-1} \mathbf{A}^T \right)^{-1} (\mathbf{A}\hat{\theta} - \mathbf{b})$

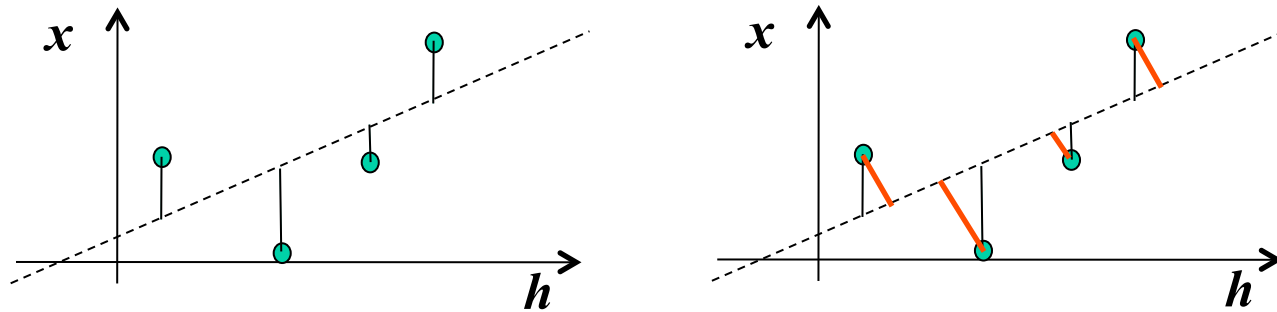The constrained LSE is a corrected version of the unconstrained LSE.

It can also be interpreted as the constrained signal estimate to be the projection of the unconstrained solution onto the constrained subspace.



PO 1213

# *Extensions to Least Squares:*

Other extensions:

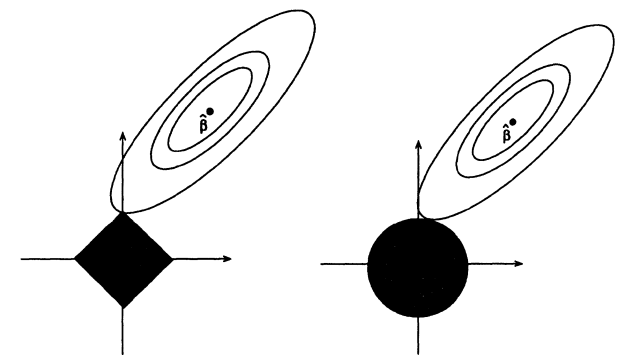Total Least Squares (errors in variables, or orthogonal regression)



When could also be errors in the independent variables.

Lasso – Least Absolute Shrinkage and Selection Operator

criterion $\quad J(\theta) = (\mathbf{x} - \mathbf{H}\theta)^T (\mathbf{x} - \mathbf{H}\theta)$

$$s.t. \quad \sum_j |\theta| \le t, \quad \text{with} \quad t > 0$$

solution $\quad \hat{\theta} = \left(\mathbf{H}^T\mathbf{H} + \lambda\mathbf{W}^-\right)^{-1} \mathbf{H}^T\mathbf{x}$

$\mathbf{W}$ diagonal matrix with elements $\left|\hat{\theta}_i\right|$, and $\mathbf{W}^-$ is the generalized inverse.

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Nonlinear Least Squares:*

In general the signal model is

model
$$\mathbf{x} = s(\boldsymbol{\theta})^T + \mathbf{w}$$

where s() is in general a nonlinear function of the unknown parameters. The criterion to be minimized can be written as (if a quadratic error is selected)

criterion
$$J(\boldsymbol{\theta}) = (\mathbf{x} - \mathbf{s}(\boldsymbol{\theta}))^T (\mathbf{x} - \mathbf{s}(\boldsymbol{\theta}))$$

termed also as nonlinear regression problem, in statistics.

Solution is general is not available, except if resorting to numerical methods.

Two methods than can reduce the complexity can be identified:

1 – Transformation of parameters;

2 – Separability of parameters;

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Nonlinear Least Squares:*

Transformation of parameters

We seek a one-to-one transformation that produces a linear signal model in the new space:

$$\alpha = \mathbf{g}\left(\theta\right)$$

Where $\mathbf{g}()$ is a p-dimensional function of the unknown parameters, with inverse:

$$\mathbf{s}\left(\theta\left(\alpha\right)\right) = \mathbf{s}\left(\mathbf{g}^{-1}\left(\alpha\right)\right) = \mathbf{H}\alpha.$$

Then the solution is

$$\hat{\theta} = \mathbf{g}^{-1}\left(\alpha\right) = \mathbf{g}^{-1}\left(\left(\mathbf{H}^{T}\mathbf{H}\right)^{-1}\mathbf{H}^{T}\mathbf{x}\right)$$

The transformation $\mathbf{g}()$, if it exists, is usually quite difficult.

Only a few nonlinear least squares problems may be solved in this manner.

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

# *Nonlinear Least Squares:*

Separability of parameters

Assume that the model is nonlinear but still is linear in some of the parameters. Thus

$$\mathbf{s} = \mathbf{H}(\alpha)\beta$$

Where

$$\theta = \begin{bmatrix} \alpha \\ \beta \end{bmatrix} = \begin{bmatrix} (p-q) \times 1 \\ q \times 1 \end{bmatrix}$$

The criterion

$$J(\alpha,\beta) = \left(\mathbf{x} - \mathbf{H}(\alpha)\beta\right)^T \left(\mathbf{x} - \mathbf{H}(\alpha)\beta\right)$$

is linear in **β** and nonlinear in **α**. For a given **α** can be minimized, with (partial) solution

$$\hat{\beta} = \left(\mathbf{H}^T(\alpha)\mathbf{H}(\alpha)\right)^{-1}\mathbf{H}^T(\alpha)\mathbf{x}$$

The problem now reduces to the maximization of

$$J(\alpha,\hat{\beta}) = \mathbf{x}^T\left(I - \mathbf{H}(\alpha)\left(\mathbf{H}^T(\alpha)\mathbf{H}(\alpha)\right)^{-1}\mathbf{H}^T(\alpha)\right)\mathbf{x}$$

over **α.**

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA

General case

When all the other methods fail, a Taylor series expansion can be used. The criterion is then approximated...

$$J(\theta) = \sum_{n=0}^{N-1} \left(x[n] - s[n;\theta]\right)^2 \approx \sum_{n=0}^{N-1} \left( x[n] - s[n;\theta_0] - \left.\frac{ds[n;\theta]}{d\theta}\right|_{\theta_0} (\theta - \theta_0) \right)^2$$

If we set up an iterative procedure (as in the Newton-Rawphson case)

$$\theta_{k+1} = \theta_k + \left(\mathbf{H}^T(\theta_k)\mathbf{H}(\theta_k)\right)^{-1}\mathbf{H}^T(\theta_k)\left(\mathbf{x} - \mathbf{s}(\theta_k)\right)$$

Where

$$\left[\mathbf{H}(\theta)\right]_{ij} = \frac{\partial s[i]}{\partial \theta_j}$$

The solution can be trivially generalized to the vector case:

$$\theta_{k+1} = \theta_k + \left(\mathbf{H}^T(\theta_k)\mathbf{H}(\theta_k)\right)^{-1}\mathbf{H}^T(\theta_k)\left(\mathbf{x} - \mathbf{s}(\theta_k)\right)$$

DEM
DEPARTAMENTO
DE ENGENHARIA MECÃNICA
TÉCNICO LISBOA

# *Bibliography:*

**Further reading**

• Thomas Kailath, *Linear Systems,* Prentice Hall, 1980.

• Thomas Kailath, Ali Sayed, and Babak Hassibi, *Linear Estimation,* Prentice Hall, 2000.

• Harry L. Van Trees, *Detection, Estimation, and Modulation Theory, Parts I to IV,* John Wiley, 2001.

• J. Bibby, H. Toutenburg, *Prediction and Improved Estimation in Linear Models,* John Wiley, 1977.

• C.Rao, *Linear Statistical Inference and Its Applications,* John Wiley, 1973.

DEM
DEPARTAMENTO
DE ENGENHARIA MECÂNICA
TÉCNICO LISBOA