



Detection & Estimation Theory: Lectures 1 and 2

Prof. Aggelos Bletsas

Technical University of Crete
School of Electrical and Computer Engineering



European Union
European Social Fund

Operational Programme
**Human Resources Development,
Education and Lifelong Learning**

Co-financed by Greece and the European Union





“Support for the Internationalization of Higher Education, School of Electrical and Computer Engineering, Technical University of Crete” (MIS 5150766), under the call for proposals “Support of the Internationalization of Higher Education - Technical University of Crete” (EDULLL 153).

The project is co-financed by Greece and the European Union (European Social Fund- ESF) through the Operational Programme “Human Resources Development, Education and Lifelong Learning 2014-2020.



Operational Programme
**Human Resources Development,
Education and Lifelong Learning**
Co-financed by Greece and the European Union



- Problem Definition & Basic Assumptions
- Formulation of (Bayesian) Risk
- Minimization of Bayesian Risk: Likelihood Ratio Test (LRT)
- LRT Special Cases: Maximum A posteriori Probability (MAP) Rule
 - Maximum Likelihood (ML)
- Example

Binary Hypothesis Testing: Problem Definition

- ▶ Given a collection of measurements $\mathbf{y} \in \mathcal{Y}$, find *optimal* decision function $\delta(\mathbf{y})$ that splits observation space \mathcal{Y} in two disjoint regions $\mathcal{Y}_0, \mathcal{Y}_1$:

$$\delta(\mathbf{y}) = \begin{cases} 0, & \text{under hypothesis } H_0 \text{ (i.e., } \mathbf{y} \in \mathcal{Y}_0) \\ 1, & \text{under hypothesis } H_1 \text{ (i.e., } \mathbf{y} \in \mathcal{Y}_1) \end{cases} \quad (1)$$

where hypothesis 0 and hypothesis 1 are denoted by H_0, H_1 , respectively.

- ▶ Vector \mathbf{y} denotes a collection of measurements.
- ▶ **Continuous case:** $f_{\mathbf{y}|H_j}(\mathbf{y}|H_j)$, $j \in \{0, 1\}$, i.e., conditional probability density function (pdf) is known and $f_{\mathbf{y}|H_0}(\mathbf{y}|H_0) \neq f_{\mathbf{y}|H_1}(\mathbf{y}|H_1)$.
- ▶ **Discrete case:** $\Pr(\mathbf{y}|H_j)$, $j \in \{0, 1\}$, i.e., conditional probability mass function (pmf) is known and $\Pr(\mathbf{y}|H_0) \neq \Pr(\mathbf{y}|H_1)$.
- ▶ Priors $\pi_0 \triangleq \Pr(H_0) = 1 - \pi_1$, $\pi_1 \triangleq \Pr(H_1)$ are known.

Binary Hypothesis Testing: Problem Definition

- ▶ Given a collection of measurements $\mathbf{y} \in \mathcal{Y}$, find *optimal* decision function $\delta(\mathbf{y})$ that splits observation space \mathcal{Y} in two disjoint regions $\mathcal{Y}_0, \mathcal{Y}_1$:

$$\delta(\mathbf{y}) = \begin{cases} 0, & \text{under hypothesis } H_0 \text{ (i.e., } \mathbf{y} \in \mathcal{Y}_0) \\ 1, & \text{under hypothesis } H_1 \text{ (i.e., } \mathbf{y} \in \mathcal{Y}_1) \end{cases} \quad (2)$$

where hypothesis 0 and hypothesis 1 are denoted by H_0, H_1 , respectively.

- ▶ We need an *optimality* criterion!
- ▶ Bayes comes to help: all uncertainties are quantifiable, all costs and benefits of decision can be measured!

Formulation of (Bayesian) Risk

- ▶ Define cost C_{ij} of deciding that H_i holds, when hypothesis H_j is true, $i, j \in \{0, 1\}$.
- ▶ Define

$$\begin{aligned} \Pr(\delta(\mathbf{y}) = i | H_j) &= \\ \triangleq \Pr(\mathbf{y} \in \mathcal{Y}_i | H_j) &= \begin{cases} \int_{\mathbf{y} \in \mathcal{Y}_i} f_{\mathbf{y} | H_j}(\mathbf{y} | H_j) d\mathbf{y}, & \text{(continuous case)} \\ \sum_{\mathbf{y} \in \mathcal{Y}_i} \Pr(\mathbf{y} | H_j). & \text{(discrete case)} \end{cases} \end{aligned} \quad (3)$$

- ▶ We are ready to define the *conditional* Bayesian Risk $R(\cdot | \cdot)$ for decision rule $\delta(\mathbf{y})$ under hypothesis H_j :

$$\begin{aligned} R(\delta(\mathbf{y}) | H_j) &= C_{1j} \Pr(\delta(\mathbf{y}) = 1 | H_j) + C_{0j} \Pr(\delta(\mathbf{y}) = 0 | H_j) \\ &= \sum_{i=0}^1 C_{ij} \Pr(\delta(\mathbf{y}) = i | H_j). \end{aligned} \quad (4)$$

Formulation of (Bayesian) Risk

- ▶ *Conditional* Bayesian Risk $R(\cdot|\cdot)$ for decision rule $\delta(\mathbf{y})$ under hypothesis H_j :

$$R(\delta(\mathbf{y})|H_j) = \sum_{i=0}^1 C_{ij} \Pr(\delta(\mathbf{y}) = i|H_j). \quad (5)$$

- ▶ Thus, the average unconditional Bayesian cost of decision rule $\delta(\mathbf{y})$ follows:

$$R(\delta(\mathbf{y})) = R(\delta(\mathbf{y})|H_0) \Pr(H_0) + R(\delta(\mathbf{y})|H_1) \Pr(H_1) \quad (6)$$

$$= R(\delta(\mathbf{y})|H_0) \pi_0 + R(\delta(\mathbf{y})|H_1) \pi_1 \quad (7)$$

$$= \sum_{j=0}^1 R(\delta(\mathbf{y})|H_j) \pi_j \quad (8)$$

$$\stackrel{\text{Eq. (5)}}{=} \sum_{i=0}^1 \sum_{j=0}^1 \pi_j C_{ij} \Pr(\delta(\mathbf{y}) = i|H_j) \quad (9)$$

Formulation of (Bayesian) Risk

- ▶ Exploiting the fact that $\mathcal{Y}_0 \cup \mathcal{Y}_1 = \mathcal{Y}$ and $\mathcal{Y}_0 \cap \mathcal{Y}_1 = \emptyset$:

$$\Pr(\delta(\mathbf{y}) = 0|H_j) + \Pr(\delta(\mathbf{y}) = 1|H_j) = 1 \Leftrightarrow \quad (10)$$

$$\Pr(\mathbf{y} \in \mathcal{Y}_0|H_j) = 1 - \Pr(\mathbf{y} \in \mathcal{Y}_1|H_j) \quad (11)$$

- ▶ Average Bayesian cost of decision rule $\delta(\mathbf{y})$:

$$R(\delta(\mathbf{y})) = \sum_{i=0}^1 \sum_{j=0}^1 \pi_j C_{ij} \Pr(\delta(\mathbf{y}) = i|H_j) \quad (12)$$

$$= \sum_{j=0}^1 \sum_{i=0}^1 \pi_j C_{ij} \Pr(\mathbf{y} \in \mathcal{Y}_i|H_j) \quad (13)$$

$$= \sum_{j=0}^1 \pi_j C_{0j} \Pr(\mathbf{y} \in \mathcal{Y}_0|H_j) + \pi_j C_{1j} \Pr(\mathbf{y} \in \mathcal{Y}_1|H_j) \quad (14)$$

$$\stackrel{(11)}{=} \sum_{j=0}^1 \pi_j C_{0j} + \sum_{j=0}^1 \pi_j (C_{1j} - C_{0j}) \Pr(\mathbf{y} \in \mathcal{Y}_1|H_j) \quad (15)$$

Minimization of Bayesian Risk

- ▶ Average Bayesian cost (Bayesian Risk) of decision rule $\delta(\mathbf{y})$:

$$R(\delta(\mathbf{y})) = \sum_{j=0}^1 \pi_j C_{0j} + \sum_{j=0}^1 \pi_j (C_{1j} - C_{0j}) \Pr(\mathbf{y} \in \mathcal{Y}_1 | H_j) \quad (16)$$

- ▶ Notice that the first sum is independent of the measurement data \mathbf{y} . The *optimal* decision rule should perform the following minimization:

$$\min_{\delta(\mathbf{y})} R(\delta(\mathbf{y})) \quad (17)$$

- ▶ Two cases: \mathbf{y} continuous or discrete (solution will be found the same!).

Minimization of Bayesian Risk

- ▶ Continuous case: $\Pr(\mathbf{y} \in \mathcal{Y}_1 | H_j) = \int_{\mathbf{y} \in \mathcal{Y}_1} f_{\mathbf{y}|H_j}(\mathbf{y}|H_j) d\mathbf{y}$,
- ▶ Bayesian Risk of decision rule $\delta(\mathbf{y})$:

$$\begin{aligned} R(\delta(\mathbf{y})) &= \sum_{j=0}^1 \pi_j C_{0j} + \sum_{j=0}^1 \pi_j (C_{1j} - C_{0j}) \int_{\mathbf{y} \in \mathcal{Y}_1} f_{\mathbf{y}|H_j}(\mathbf{y}|H_j) d\mathbf{y} \\ &= \sum_{j=0}^1 \pi_j C_{0j} + \int_{\mathbf{y} \in \mathcal{Y}_1} \sum_{j=0}^1 \pi_j (C_{1j} - C_{0j}) f_{\mathbf{y}|H_j}(\mathbf{y}|H_j) d\mathbf{y} \quad (18) \end{aligned}$$

- ▶ Remember that $\delta(\mathbf{y})$ controls what data \mathbf{y} is allocated to \mathcal{Y}_1 and what data to \mathcal{Y}_0 . From Eq. (18), $R(\delta(\mathbf{y}))$ is minimized when then integrand of (18) is minimized (i.e., negative or zero):

$$\begin{aligned} \delta_B(\mathbf{y}) &= \arg \min_{\delta(\mathbf{y})} R(\delta(\mathbf{y})) \Leftrightarrow \\ \text{select } \mathcal{Y}_1 &: \left\{ \mathbf{y} \in \mathcal{Y} : \sum_{j=0}^1 \pi_j (C_{1j} - C_{0j}) f_{\mathbf{y}|H_j}(\mathbf{y}|H_j) \leq 0 \right\} \quad (19) \end{aligned}$$

Minimization of Bayesian Risk

- ▶ Discrete case: $\Pr(\mathbf{y} \in \mathcal{Y}_1 | H_j) = \sum_{\mathbf{y} \in \mathcal{Y}_1} \Pr(\mathbf{y} | H_j)$,
- ▶ Bayesian Risk of decision rule $\delta(\mathbf{y})$:

$$\begin{aligned} R(\delta(\mathbf{y})) &= \sum_{j=0}^1 \pi_j C_{0j} + \sum_{j=0}^1 \pi_j (C_{1j} - C_{0j}) \sum_{\mathbf{y} \in \mathcal{Y}_1} \Pr(\mathbf{y} | H_j) \\ &= \sum_{j=0}^1 \pi_j C_{0j} + \sum_{\mathbf{y} \in \mathcal{Y}_1} \sum_{j=0}^1 \pi_j (C_{1j} - C_{0j}) \Pr(\mathbf{y} | H_j) d\mathbf{y} \end{aligned} \quad (20)$$

- ▶ Similarly to the continuous case, $R(\delta(\mathbf{y}))$ is minimized when then integrand of (20) is minimized (i.e., negative or zero):

$$\begin{aligned} \delta_B(\mathbf{y}) &= \arg \min_{\delta(\mathbf{y})} R(\delta(\mathbf{y})) \Leftrightarrow \\ \text{select } \mathcal{Y}_1 &: \left\{ \mathbf{y} \in \mathcal{Y} : \sum_{j=0}^1 \pi_j (C_{1j} - C_{0j}) \Pr(\mathbf{y} | H_j) \leq 0 \right\} \end{aligned} \quad (21)$$

Minimization of Bayesian Risk

- ▶ Define likelihood ratio (LR) for continuous or discrete case:

$$L(\mathbf{y}) = \frac{f_{\mathbf{y}|H_1}(\mathbf{y}|H_1)}{f_{\mathbf{y}|H_0}(\mathbf{y}|H_0)} \text{ (continuous case)}, L(\mathbf{y}) = \frac{\Pr(\mathbf{y}|H_1)}{\Pr(\mathbf{y}|H_0)} \text{ (discrete case)}.$$

- ▶ We can safely assume that the LR is finite positive for all cases of interest (see below):
 - ▶ LR numerator and denominator both positive: LRT finite positive.
 - ▶ LR numerator and denominator both zero: this is impossible, since in that case $\mathbf{y} \notin \mathcal{Y}$ (and we have also assumed that $\Pr(\mathbf{y}|H_1) \neq \Pr(\mathbf{y}|H_0)$).
 - ▶ either numerator or denominator (only one of the two) is zero; if numerator is zero then that particular \mathbf{y} cannot occur under H_1 ; similarly, if denominator is zero then that \mathbf{y} cannot occur under H_0 .

Minimization of Bayesian Risk: Likelihood Ratio Test

- ▶ We further assume that $C_{01} > C_{11}$, i.e., the cost of wrong decision is strictly higher than the cost of correct decision.
- ▶ Continuous case:

$$\begin{aligned} \text{select } \mathcal{Y}_1 : & \left\{ \mathbf{y} \in \mathcal{Y} : \sum_{j=0}^1 \pi_j (C_{1j} - C_{0j}) f_{\mathbf{y}|H_j}(\mathbf{y}|H_j) \leq 0 \right\} \Leftrightarrow \\ & \Leftrightarrow \sum_{j=0}^1 \pi_j (C_{1j} - C_{0j}) f_{\mathbf{y}|H_j}(\mathbf{y}|H_j) \stackrel{H_1}{\leq} 0 \end{aligned} \quad (22)$$

$$\begin{aligned} & \Leftrightarrow \pi_0 (C_{10} - C_{00}) f_{\mathbf{y}|H_0}(\mathbf{y}|H_0) \stackrel{H_1}{\leq} -\pi_1 (C_{11} - C_{01}) f_{\mathbf{y}|H_1}(\mathbf{y}|H_1) \\ & \stackrel{(C_{01}-C_{11})>0}{\Leftrightarrow} \frac{f_{\mathbf{y}|H_1}(\mathbf{y}|H_1)}{f_{\mathbf{y}|H_0}(\mathbf{y}|H_0)} \stackrel{H_1}{\geq} \frac{C_{10} - C_{00}}{C_{01} - C_{11}} \frac{\pi_0}{\pi_1} \triangleq \tau \end{aligned} \quad (23)$$

$$\Leftrightarrow L(\mathbf{y}) \stackrel{H_1}{\geq} \tau \quad (24)$$

Minimization of Bayesian Risk: Likelihood Ratio Test

- ▶ Continuous case:

$$L(\mathbf{y}) \triangleq \frac{f_{\mathbf{y}|H_1}(\mathbf{y}|H_1)}{f_{\mathbf{y}|H_0}(\mathbf{y}|H_0)} \stackrel{H_1}{\geq} \frac{C_{10} - C_{00}}{C_{01} - C_{11}} \frac{\pi_0}{\pi_1} \triangleq \tau \quad (25)$$

$$\Leftrightarrow L(\mathbf{y}) \stackrel{H_1}{\geq} \tau \quad (26)$$

Notice that the values of \mathbf{y} where the integrand goes to zero (or equivalently the (LR) ratio is equal to τ) do not matter; some can be allocated to \mathcal{Y}_1 and some (or none) to \mathcal{Y}_0 .

- ▶ Discrete case - with similar reasoning, minimization in Eq. (21) offers the following LR test:

$$L(\mathbf{y}) \triangleq \frac{\Pr(\mathbf{y}|H_1)}{\Pr(\mathbf{y}|H_0)} \stackrel{H_1}{\geq} \frac{C_{10} - C_{00}}{C_{01} - C_{11}} \frac{\pi_0}{\pi_1} \triangleq \tau \quad (27)$$

$$\Leftrightarrow L(\mathbf{y}) \stackrel{H_1}{\geq} \tau \quad (28)$$

Minimization of Bayesian Risk: Likelihood Ratio Test

- ▶ Thus, optimum Bayesian decision rule $\delta_B(\mathbf{y})$, i.e., rule that minimizes Bayes risk, can be written as follows:

$$\delta_B(\mathbf{y}) = \begin{cases} 1, & \text{if } L(\mathbf{y}) \geq \tau, \\ 0, & \text{if } L(\mathbf{y}) < \tau, \end{cases} \quad (29)$$

or more compactly,

$$L(\mathbf{y}) \stackrel{H_1}{\geq} \tau. \quad (30)$$

Likelihood Ratio Test (LRT) & Symmetric Costs

- ▶ Set symmetric costs, i.e., 1 for (any) erroneous detection and 0 for (any) correct decision:

$$C_{ij} = 1 - \delta_{ij} = \begin{cases} 0, & i = j, \\ 1, & i \neq j, \end{cases} \quad (31)$$

where δ_{ij} denotes the Kronecker delta. For such costs, Bayesian Risk is equivalent to probability of error! From Eq. (4):

$$\begin{aligned} R(\delta(\mathbf{y})|H_j) &= C_{1j} \Pr(\delta(\mathbf{y}) = 1|H_j) + C_{0j} \Pr(\delta(\mathbf{y}) = 0|H_j) \Rightarrow \\ R(\delta(\mathbf{y})|H_0) &= C_{10} \Pr(\delta(\mathbf{y}) = 1|H_0) + C_{00} \Pr(\delta(\mathbf{y}) = 0|H_0) \\ &= \Pr(\delta(\mathbf{y}) = 1|H_0) \equiv \Pr(\text{error}|H_0). \end{aligned} \quad (32)$$

$$\begin{aligned} R(\delta(\mathbf{y})|H_1) &= C_{11} \Pr(\delta(\mathbf{y}) = 1|H_1) + C_{01} \Pr(\delta(\mathbf{y}) = 0|H_1) \\ &= \Pr(\delta(\mathbf{y}) = 0|H_1) \equiv \Pr(\text{error}|H_1) \Rightarrow \end{aligned} \quad (33)$$

$$\begin{aligned} R(\delta(\mathbf{y})) &= R(\delta(\mathbf{y})|H_0) \pi_0 + R(\delta(\mathbf{y})|H_1) \pi_1 \\ &= \Pr(\text{error}|H_0) \pi_0 + \Pr(\text{error}|H_1) \pi_1 \equiv \Pr(\text{error}). \end{aligned} \quad (34)$$

- ▶ Set symmetric costs $C_{ij} = 1 - \delta_{ij}$, as before. For such costs, Bayesian Risk is equivalent to probability of error!
- ▶ Continuous case:¹

$$L(\mathbf{y}) \triangleq \frac{f_{\mathbf{y}|H_1}(\mathbf{y}|H_1)}{f_{\mathbf{y}|H_0}(\mathbf{y}|H_0)} \stackrel{H_1}{\geq} \tau \triangleq \frac{C_{10} - C_{00}}{C_{01} - C_{11}} \frac{\pi_0}{\pi_1} = \frac{(1-0)}{(1-0)} \frac{\pi_0}{\pi_1} = \frac{\pi_0}{\pi_1}$$

$$\Leftrightarrow f_{\mathbf{y}|H_1}(\mathbf{y}|H_1) \pi_1 \stackrel{H_1}{\geq} f_{\mathbf{y}|H_0}(\mathbf{y}|H_0) \pi_0 \quad (35)$$

$$\Leftrightarrow \frac{f_{\mathbf{y}|H_1}(\mathbf{y}|H_1) \pi_1}{f_{\mathbf{y}}(\mathbf{y})} \stackrel{H_1}{\geq} \frac{f_{\mathbf{y}|H_0}(\mathbf{y}|H_0) \pi_0}{f_{\mathbf{y}}(\mathbf{y})} \quad (36)$$

$$\stackrel{\text{Bayes}^{(*)}}{\Leftrightarrow} \Pr(H_1|\mathbf{y}) \stackrel{H_1}{\geq} \Pr(H_0|\mathbf{y}) \quad (\text{MAP Rule}) \quad (37)$$

- ▶ Discrete case: same rule as above!

¹(*) holds because the Bayes property holds for continuous distributions as well.

- ▶ Set symmetric costs $C_{ij} = 1 - \delta_{ij}$, as before and equal priors $\pi_0 = \pi_1$ (special MAP case):
- ▶ Continuous case:

$$\frac{f_{\mathbf{y}|H_1}(\mathbf{y}|H_1) \pi_1}{f_{\mathbf{y}}(\mathbf{y})} \stackrel{H_1}{\geq} \frac{f_{\mathbf{y}|H_0}(\mathbf{y}|H_0) \pi_0}{f_{\mathbf{y}}(\mathbf{y})} \quad (38)$$

$$\Rightarrow f_{\mathbf{y}|H_1}(\mathbf{y}|H_1) \stackrel{H_1}{\geq} f_{\mathbf{y}|H_0}(\mathbf{y}|H_0) \quad (\text{ML Rule}) \quad (39)$$

- ▶ Discrete case - same derivation as above:

$$\Pr(\mathbf{y}|H_1) \stackrel{H_1}{\geq} \Pr(\mathbf{y}|H_0) \quad (\text{ML Rule}) \quad (40)$$

- MAP and ML minimize Bayesian risk and probability of error.

Simple example

- ▶ Assume $y_k = m_0 + v_k$ under hypothesis H_0 and $y_k = m_1 + v_k$ under hypothesis H_1 , where $m_1 > m_0$ and variables $\{v_k\}$, $k \in \{1, 2, \dots, M\}$ are derived from white Gaussian noise (WGN), i.e. $v_i \perp v_j, i \neq j$ (statistically independent) and $v_k \sim \mathcal{N}(0, \sigma^2)$. Find optimal decision rule that detects which hypothesis holds.
- ▶ Solution: Affine transformation of Gaussian is also Gaussian:

$$y_k \sim \begin{cases} \mathcal{N}(m_0, \sigma^2), & \text{under } H_0 \\ \mathcal{N}(m_1, \sigma^2), & \text{under } H_1 \end{cases} \quad (41)$$

Since $\{v_k\}$ are independent, observations $\{y_k\}$ are independent and the product of their conditional pdfs offers the conditional probability density of each hypothesis and their ratio:

$$\begin{aligned} f_{\mathbf{y}|H_j}(\mathbf{y} = [y_1 \ y_2 \ \dots \ y_M]|H_j) &= \prod_{k=1}^M f_{y_k|H_j}(y_k|H_j) = \prod_{k=1}^M \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left[-\frac{(y_k - m_j)^2}{2\sigma^2}\right] \\ &= \frac{1}{(2\pi\sigma^2)^{\frac{M}{2}}} \exp\left[-\frac{1}{2\sigma^2} \sum_{k=1}^M (y_k - m_j)^2\right]. \\ L(\mathbf{y}) &\triangleq \frac{f_{\mathbf{y}|H_1}(\mathbf{y}|H_1)}{f_{\mathbf{y}|H_0}(\mathbf{y}|H_0)} = \exp\left[-\frac{1}{2\sigma^2} \sum_{k=1}^M (y_k - m_1)^2 + \frac{1}{2\sigma^2} \sum_{k=1}^M (y_k - m_0)^2\right]. \end{aligned}$$

$$L(\mathbf{y}) = \exp \left[-\frac{1}{2\sigma^2} \sum_{k=1}^M (y_k - m_1)^2 + \frac{1}{2\sigma^2} \sum_{k=1}^M (y_k - m_0)^2 \right] \quad (42)$$

$$= \exp \left[-\frac{1}{2\sigma^2} \sum_{k=1}^M (y_k^2 - 2m_1 y_k + m_1^2) + \frac{1}{2\sigma^2} \sum_{k=1}^M (y_k^2 - 2m_0 y_k + m_0^2) \right] \quad (43)$$

$$= \exp \left[+\frac{m_1}{\sigma^2} \sum_{k=1}^M y_k - \frac{1}{2\sigma^2} \sum_{k=1}^M m_1^2 - \frac{m_0}{\sigma^2} \sum_{k=1}^M y_k + \frac{1}{2\sigma^2} \sum_{k=1}^M m_0^2 \right] \quad (44)$$

$$= \exp \left[+\frac{m_1 - m_0}{\sigma^2} \sum_{k=1}^M y_k - \frac{1}{2\sigma^2} \sum_{k=1}^M m_1^2 + \frac{1}{2\sigma^2} \sum_{k=1}^M m_0^2 \right] \quad (45)$$

$$= \exp \left[+\frac{m_1 - m_0}{\sigma^2} \sum_{k=1}^M y_k - \frac{1}{2\sigma^2} M m_1^2 + \frac{1}{2\sigma^2} M m_0^2 \right] \quad (46)$$

$$= \exp \left[-\frac{M(m_1^2 - m_0^2)}{2\sigma^2} + \frac{m_1 - m_0}{\sigma^2} \sum_{k=1}^M y_k \right] \quad (47)$$

$$L(\mathbf{y}) = \exp \left[-\frac{M(m_1^2 - m_0^2)}{2\sigma^2} + \frac{m_1 - m_0}{\sigma^2} \sum_{k=1}^M y_k \right] \quad (48)$$

$$L(\mathbf{y}) \stackrel{H_1}{\geq} \tau \Leftrightarrow \ln(L(\mathbf{y})) \stackrel{H_1}{\geq} \ln(\tau) \Leftrightarrow \quad (49)$$

$$-\frac{M(m_1^2 - m_0^2)}{2\sigma^2} + \frac{m_1 - m_0}{\sigma^2} \sum_{k=1}^M y_k \stackrel{H_1}{\geq} \ln(\tau) \Leftrightarrow \quad (50)$$

$$+ \frac{m_1 - m_0}{\sigma^2} \sum_{k=1}^M y_k \stackrel{H_1}{\geq} + \frac{M(m_1^2 - m_0^2)}{2\sigma^2} + \ln(\tau) \stackrel{m_1 > m_0}{\Leftrightarrow} \quad (51)$$

$$+ \frac{1}{M} \sum_{k=1}^M y_k \stackrel{H_1}{\geq} + \frac{(m_1 + m_0)}{2} + \frac{\sigma^2}{M(m_1 - m_0)} \ln(\tau), \quad (52)$$

where we used the fact that $(m_1 - m_0) > 0$.

- The left-hand side of the above inequality, i.e., the term $\frac{1}{M} \sum_{k=1}^M y_k$, is called the sufficient statistic.
- Observe that the sufficient statistic is the sample mean for $M \rightarrow +\infty$, under each hypothesis.

Bernard C. Levy, Principles of Signal Detection and Parameter Estimation, Springer 2008.

Thank you!



Detection & Estimation Theory: Lecture 3

Prof. Aggelos Bletsas

Technical University of Crete
School of Electrical and Computer Engineering



European Union
European Social Fund

Operational Programme
**Human Resources Development,
Education and Lifelong Learning**

Co-financed by Greece and the European Union





“Support for the Internationalization of Higher Education, School of Electrical and Computer Engineering, Technical University of Crete” (MIS 5150766), under the call for proposals “Support of the Internationalization of Higher Education - Technical University of Crete” (EDULLL 153).

The project is co-financed by Greece and the European Union (European Social Fund- ESF) through the Operational Programme “Human Resources Development, Education and Lifelong Learning 2014-2020.



Operational Programme
**Human Resources Development,
Education and Lifelong Learning**
Co-financed by Greece and the European Union



- One more simple (binary hypothesis testing) example
- Common distributions for Sufficient Statistics
- Gaussian vectors (or jointly Gaussian random variables)

Another simple example

- ▶ Assume $y_k = v_k$, where $v_k \sim \mathcal{N}(0, \sigma_0^2)$ under hypothesis H_0 and $y_k = v_k$, where $v_k \sim \mathcal{N}(0, \sigma_1^2)$ under hypothesis H_1 , with $\sigma_1^2 > \sigma_0^2$ and variables $\{v_k\}$, $k \in \{1, 2, \dots, M\}$ are derived from white Gaussian noise (WGN), i.e. $v_i \perp v_j$, $i \neq j$ (statistically independent) and v_k is Gaussian. Find optimal decision rule that detects which hypothesis holds.
- ▶ Notice that in this example, the variance of measurements changes per hypothesis (and not the mean, as in the previous example).
- ▶ Solution: Affine transformation of Gaussian is also Gaussian:

$$y_k \sim \begin{cases} \mathcal{N}(0, \sigma_0^2), & \text{under } H_0 \\ \mathcal{N}(0, \sigma_1^2), & \text{under } H_1 \end{cases} \quad (1)$$

Since $\{v_k\}$ are independent, observations $\{y_k\}$ are independent and the product of their conditional pdfs offers the conditional probability density of each hypothesis and their ratio, as follows (with $j \in \{0, 1\}$):

$$\begin{aligned} f_{\mathbf{y}|H_j}(\mathbf{y} = [y_1 \ y_2 \ \dots \ y_M]|H_j) &= \prod_{k=1}^M f_{y_k|H_j}(y_k|H_j) = \prod_{k=1}^M \frac{1}{\sqrt{2\pi\sigma_j^2}} \exp\left[-\frac{y_k^2}{2\sigma_j^2}\right] \\ &= \frac{1}{(2\pi\sigma_j^2)^{\frac{M}{2}}} \exp\left[-\frac{1}{2\sigma_j^2} \sum_{k=1}^M y_k^2\right]. \\ L(\mathbf{y}) \triangleq \frac{f_{\mathbf{y}|H_1}(\mathbf{y}|H_1)}{f_{\mathbf{y}|H_0}(\mathbf{y}|H_0)} &= \frac{\sigma_0^M}{\sigma_1^M} \exp\left[\left(\frac{1}{2\sigma_0^2} - \frac{1}{2\sigma_1^2}\right) \sum_{k=1}^M y_k^2\right]. \end{aligned}$$

$$L(\mathbf{y}) = \frac{\sigma_0^M}{\sigma_1^M} \exp \left[\left(\frac{1}{2\sigma_0^2} - \frac{1}{2\sigma_1^2} \right) \sum_{k=1}^M y_k^2 \right]. \quad (2)$$

$$L(\mathbf{y}) \stackrel{H_1}{\geq} \tau \Leftrightarrow \ln(L(\mathbf{y})) \stackrel{H_1}{\geq} \ln(\tau) \Leftrightarrow \quad (3)$$

$$-M \ln \left(\frac{\sigma_1}{\sigma_0} \right) + \frac{\sigma_1^2 - \sigma_0^2}{2\sigma_0^2\sigma_1^2} \sum_{k=1}^M y_k^2 \stackrel{H_1}{\geq} \ln(\tau) \Leftrightarrow \quad (4)$$

$$\frac{1}{M} \frac{\sigma_1^2 - \sigma_0^2}{2\sigma_0^2\sigma_1^2} \sum_{k=1}^M y_k^2 \stackrel{H_1}{\geq} \frac{1}{M} \ln(\tau) + \ln \left(\frac{\sigma_1}{\sigma_0} \right) \stackrel{\sigma_1^2 > \sigma_0^2}{\Leftrightarrow} \quad (5)$$

$$\frac{1}{M} \sum_{k=1}^M y_k^2 \stackrel{H_1}{\geq} \frac{2\sigma_0^2\sigma_1^2}{\sigma_1^2 - \sigma_0^2} \left(\frac{1}{M} \ln(\tau) + \ln \left(\frac{\sigma_1}{\sigma_0} \right) \right), \quad (6)$$

where we used the fact in Eq. (6) that $(\sigma_1^2 - \sigma_0^2) > 0$.

- In this case, the sufficient statistic is the term $\frac{1}{M} \sum_{k=1}^M y_k^2$.
- Observe that the sufficient statistic is the sample variance (since $\mathbb{E}[y_k] = 0$) for $M \rightarrow +\infty$, under each hypothesis (and not the sample mean, as in the previous example).

$$S \triangleq \frac{1}{M} \sum_{k=1}^M y_k^2 \stackrel{H_1}{\geq} \frac{2\sigma_0^2\sigma_1^2}{\sigma_1^2 - \sigma_0^2} \left(\frac{1}{M} \ln(\tau) + \ln \left(\frac{\sigma_1}{\sigma_0} \right) \right). \quad (7)$$

Additional remarks:

1. Under H_j , for $\lim_{M \rightarrow +\infty} S = \sigma_j^2$, $j \in \{0, 1\}$.
2. Using $1 - \frac{1}{x} \leq \ln(x) \leq x - 1$, it can be easily shown that $\sigma_0^2 \leq \frac{\sigma_0^2\sigma_1^2}{\sigma_1^2 - \sigma_0^2} \ln \left(\frac{\sigma_1}{\sigma_0} \right)^2 \leq \sigma_1^2$.
3. For $M \rightarrow +\infty$, $\frac{2\sigma_0^2\sigma_1^2}{\sigma_1^2 - \sigma_0^2} \left(\frac{1}{M} \ln(\tau) + \ln \left(\frac{\sigma_1}{\sigma_0} \right) \right) \rightarrow \frac{\sigma_0^2\sigma_1^2}{\sigma_1^2 - \sigma_0^2} \ln \left(\frac{\sigma_1}{\sigma_0} \right)^2$.
4. Under hypothesis H_j , $\frac{MS}{\sigma_j^2} = \sum_{k=1}^M \left(\frac{y_k}{\sigma_j} \right)^2 =$ sum of independent squared zero-mean Gaussians of unit variance: Under hypothesis H_j , $\frac{MS}{\sigma_j^2}$ corresponds to Chi-squared distribution with M degrees of freedom [will explain it subsequently].

Useful distributions

- ▶ $z = \sum_{i=1}^M z_i^2$, $z_i \sim \mathcal{N}(0, 1)$ and $\{z_i\}$ independent, identically distributed (i.i.d.):
 - ▶ z distributed according to the Chi-squared distribution with M degrees of freedom and pdf as follows:

$$f_z(z) = \frac{1}{\Gamma(M/2) 2^{M/2}} z^{(M/2-1)} e^{-z/2} u(z), \quad (8)$$

with $u(z)$ the step function (i.e., $u(z) = 1$ for $z \geq 0$ and zero otherwise) and $\Gamma(z) = \int_0^{+\infty} t^{z-1} e^{-t} dt$ Euler's gamma function, defined everywhere apart from non-positive integers (and $\Gamma(n) = (n-1)!$ for any positive integer n).

- ▶ $\mathbb{E}[z] = M$, $\sigma_z^2 \triangleq$ variance of $z = \text{var}(z) = 2M$.
- ▶ Special case - $M = 2$: exponential distribution, with p.d.f as follows (since $\Gamma(M = 2/2) = \Gamma(1) = 0! = 1$):

$$f_z(z) = \frac{1}{2} e^{-z/2} u(z). \quad (9)$$

- ▶ In general, the pdf of a random variable according to the exponential distribution with parameter $\lambda > 0$ is given by:

$$f_z(z) = \lambda e^{-\lambda z} u(z), \quad (10)$$

with $\mathbb{E}[z] = 1/\lambda$ and $\sigma_z^2 \triangleq \text{var}(z) = 1/\lambda^2$.

- ▶ This is equivalent to $z = z_1^2 + z_2^2$, with z_1, z_2 independent and identically distributed according to $\mathcal{N}(0, \sigma^2)$ and $\mathbb{E}[z] = 1/\lambda = 2\sigma^2$.
 - ▶ For the special case of $y = \sqrt{z_1^2 + z_2^2}$, the Rayleigh pdf occurs:

$$f_y(y) = \frac{y}{\sigma^2} \exp\left(-\frac{y^2}{2\sigma^2}\right) u(y), \quad (11)$$

with $\mathbb{E}[y] = \sigma\sqrt{\pi/2}$ and $\sigma_y^2 = \text{var}(y) = \frac{4-\pi}{2}\sigma^2$.

Useful distributions

- ▶ Set $z = \sum_{i=1}^{2M} z_i^2$, $z_i \sim \mathcal{N}(0, 1)$ and $\{z_i\}$ i.i.d; as explained, z is distributed according to the Chi-squared distribution with $2M$ degrees of freedom.
- ▶ What is the distribution of $x = \theta z$ ($\theta > 0$)?
 - ▶ x can be viewed as the sum of M i.i.d. random variables distributed each according to the exponential distribution with parameter $\lambda = 1/(2\theta)$.
 - ▶ for any differentiable and invertible function $x = g(z)$, we do know that the new pdf can be found as follows:

$$f_x(x) = \frac{f_z(z)}{|g'(z)|} \Big|_{z=g^{-1}(x)} \quad (12)$$

and thus,

$$f_x(x) = \frac{f_z(z)}{\theta} \Big|_{z=x/\theta} = \frac{1}{(2\theta)^M \Gamma(M)} x^{(M-1)} \exp\left(-\frac{x}{2\theta}\right) u(x),$$

which corresponds to the Gamma distribution ($\Gamma(M, 2\theta)$), with parameters M , 2θ and $\mathbb{E}[x] = 2M\theta$, $\text{var}(x) = 4M\theta^2$.

Useful distributions

- ▶ Set $z = \sum_{i=1}^{2M} z_i^2$, $z_i \sim \mathcal{N}(0, 1)$ and $\{z_i\}$ i.i.d; as explained, z is distributed according to the Chi-squared distribution with $2M$ degrees of freedom.
- ▶ The pdf of $x = \theta z$ ($\theta > 0$) is the pdf of the sum of M i.i.d. exponentials:

$$f_x(x) = \frac{1}{(2\theta)^M \Gamma(M)} x^{(M-1)} \exp\left(-\frac{x}{2\theta}\right) u(x),$$

which corresponds, as shown, to the Gamma distribution $\Gamma(M, 2\theta)$, with parameters M , 2θ and $\mathbb{E}[x] = 2M\theta$, $\text{var}(x) = 4M\theta^2$.

- ▶ Other distributions of the exponential family to remember:
 - ▶ (discrete) Poisson: $\Pr(n) = \frac{\lambda^n}{n!} e^{-\lambda}$, $n \in \mathbb{N}$, $\mathbb{E}[n] = \lambda = \text{var}(n)$.
 - ▶ (continuous) Laplace: $f_x(x; \mu, \beta) = \frac{1}{2\beta} \exp\left(-\frac{|x-\mu|}{\beta}\right)$,
 $\mathbb{E}[x] = \mu$, $\text{var}(x) = 2\beta^2$.

Gaussian vectors (or jointly Gaussian random variables)

- ▶ Let $\mathbf{x} = [x_1 \ x_2 \ \dots \ x_M]^T$, where x_1, x_2, \dots, x_M are real. Vector \mathbf{x} is Gaussian, or equivalently, $x_1 \ x_2 \ \dots \ x_M$ are jointly Gaussian, if and only if the pdf of \mathbf{x} (or the joint pdf of $x_1 \ x_2 \ \dots \ x_M$) is given as follows:

Covariance form:

$$f_{\mathbf{x}}(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^M |\boldsymbol{\Sigma}|}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \mathbf{m})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \mathbf{m}) \right\} \quad (13)$$

denoted as $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \boldsymbol{\Sigma})$, with mean $\mathbf{m} \triangleq \mathbb{E}[\mathbf{x}]$, covariance matrix $\boldsymbol{\Sigma} \triangleq \mathbb{E} \left[(\mathbf{x} - \mathbf{m})(\mathbf{x} - \mathbf{m})^T \right]$ and $|\boldsymbol{\Sigma}|$ is the determinant of $\boldsymbol{\Sigma}$.

Gaussian vectors (or jointly Gaussian random variables)

- ▶ Let $\mathbf{x} = [x_1 \ x_2 \ \dots \ x_M]^T$, where x_1, x_2, \dots, x_M are real. Vector \mathbf{x} is Gaussian, or equivalently, $x_1 \ x_2 \ \dots \ x_M$ are jointly Gaussian, if and only if the pdf of \mathbf{x} (or the joint pdf of $x_1 \ x_2 \ \dots \ x_M$) is given as follows:

Information form:

$$f_{\mathbf{x}}(\mathbf{x}) \propto \exp \left\{ -\frac{1}{2} \mathbf{x}^T \mathbf{J} \mathbf{x} + \mathbf{h}^T \mathbf{x} \right\}, \quad (14)$$

denoted as $\mathbf{x} \sim \mathcal{N}^{-1}(\mathbf{h}, \mathbf{J})$, with potential vector $\mathbf{h} = \mathbf{J} \mathbf{m}$ and information (or precision) matrix $\mathbf{J} = \mathbf{\Sigma}^{-1}$.

Gaussian vector properties

- ▶ Let \mathbf{x} (real) Gaussian vector. The following hold:
 - ▶ Moment generating function (MGF) $M_{\mathbf{x}}(j\mathbf{u})$:

$$M_{\mathbf{x}}(j\mathbf{u}) \triangleq \mathbb{E} \left[e^{j\mathbf{u}^T \mathbf{x}} \right] = \exp \left\{ j\mathbf{u}^T \mathbf{m} - \frac{1}{2} \mathbf{u}^T \boldsymbol{\Sigma} \mathbf{u} \right\}. \quad (15)$$

- ▶ All linear combinations of elements of \mathbf{x} are scalar Gaussian random variables: $y = \mathbf{a}^T \mathbf{x}$ is Gaussian for all deterministic \mathbf{a} .
- ▶ There exists deterministic matrix \mathbf{A} , deterministic vector \mathbf{b} and random vector \mathbf{v} of i.i.d. $\mathcal{N}(0, 1)$ entries, such that $\mathbf{x} = \mathbf{A}\mathbf{v} + \mathbf{b}$.
- ▶ Affine transformation is also Gaussian, i.e., for any deterministic matrix \mathbf{A} and deterministic vector \mathbf{b} , random vector $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b}$ is Gaussian, according to $\mathcal{N}(\mathbf{A}\mathbf{m} + \mathbf{b}, \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}^T)$. [Simple proof from MGF]

Gaussian vector properties

- ▶ Let $\mathbf{x} = [\mathbf{y}^T \mathbf{z}^T]^T$ (real) Gaussian vector, where \mathbf{y} , \mathbf{z} are (real) vectors of appropriate dimensions. The following properties hold:
 - ▶ \mathbf{y} is Gaussian.
 - ▶ \mathbf{z} is Gaussian.
 - ▶ \mathbf{y} given \mathbf{z} is Gaussian.
 - ▶ \mathbf{z} given \mathbf{y} is Gaussian.
 - ▶ $\mathbb{E}[\mathbf{y}|\mathbf{z}] = \text{affine transformation of } \mathbf{z} \Rightarrow \text{Gaussian.}$
 - ▶ $\mathbb{E}[\mathbf{y} \mathbf{z}^T] = \mathbb{E}[\mathbf{y}] \mathbb{E}[\mathbf{z}]^T \Rightarrow \mathbf{y} \perp \mathbf{z}$, i.e., jointly Gaussian and uncorrelated results to independent!
- ▶ However, even if \mathbf{y} is Gaussian and \mathbf{z} is Gaussian, $\mathbf{x} = [\mathbf{y}^T \mathbf{z}^T]^T$ may not be Gaussian. In other words, \mathbf{y} and \mathbf{z} may not be necessarily *jointly* Gaussian!

- ▶ Counterexample: let x, y jointly Gaussian, zero mean, scalar random variables with joint pdf as follows:

$$p_{x,y}(x, y) = \frac{1}{2\pi\sigma_x\sigma_y} \exp\left\{-\frac{1}{2}\left(\frac{x^2}{\sigma_x^2} + \frac{y^2}{\sigma_y^2}\right)\right\}, \quad (16)$$

which corresponds to the Gaussian vector $[x \ y]^T$, with

$$\mathbf{m} = [0 \ 0]^T \text{ and } \mathbf{\Sigma} = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}.$$

- ▶ Clearly $p_y(y) = \int_{-\infty}^{+\infty} p_{x,y}(x, y)dx = 2 \int_0^{+\infty} p_{x,y}(x, y)dx$ corresponding to $\mathcal{N}(0, \sigma_y^2)$.
- ▶ Set the following non-Gaussian pdf:

$$\hat{p}_{x,y}(x, y) = \begin{cases} \frac{1}{\pi\sigma_x\sigma_y} \exp\left\{-\frac{1}{2}\left(\frac{x^2}{\sigma_x^2} + \frac{y^2}{\sigma_y^2}\right)\right\}, & \text{if } x \ y > 0, \\ 0, & \text{if } x \ y < 0. \end{cases}$$

In this case, for $y > 0$, $\int_{-\infty}^{+\infty} \hat{p}_{x,y}(x, y)dx = \int_0^{+\infty} \hat{p}_{x,y}(x, y)dx = 2 \int_0^{+\infty} \hat{p}_{x,y}(x, y)dx = p_y(y)$, i.e., Gaussian.

- [1] Bernard C. Levy, Principles of Signal Detection and Parameter Estimation, Springer 2008.
- [2] Class instructor notes.

Thank you!



Detection & Estimation Theory: Lecture 4

Prof. Aggelos Bletsas

Technical University of Crete
School of Electrical and Computer Engineering



European Union
European Social Fund

Operational Programme
**Human Resources Development,
Education and Lifelong Learning**

Co-financed by Greece and the European Union





“Support for the Internationalization of Higher Education, School of Electrical and Computer Engineering, Technical University of Crete” (MIS 5150766), under the call for proposals “Support of the Internationalization of Higher Education - Technical University of Crete” (EDULLL 153).

The project is co-financed by Greece and the European Union (European Social Fund- ESF) through the Operational Programme “Human Resources Development, Education and Lifelong Learning 2014-2020.



Operational Programme
**Human Resources Development,
Education and Lifelong Learning**
Co-financed by Greece and the European Union



- Probabilities that fully characterise a test: P_D vs P_F

- Neyman-Pearson Test
 - Derivation

Probability of Detection vs Probability of False Alarm

- ▶ There are two probability metrics for test $\delta(\mathbf{y})$ that fully characterise a binary hypothesis testing problem (as well as any binary classifier):
 - ▶ Probability of detection P_D :

$$P_D(\delta) \triangleq \int_{\mathcal{Y}_1} f_{\mathbf{y}|\mathbf{H}_1}(\mathbf{y}|\mathbf{H}_1) d\mathbf{y}. \quad (1)$$

Complementary to the above, Probability of a miss is defined as $P_M(\delta) \triangleq 1 - P_D(\delta) = \int_{\mathcal{Y}_0} f_{\mathbf{y}|\mathbf{H}_1}(\mathbf{y}|\mathbf{H}_1) d\mathbf{y}$.

- ▶ Probability of false alarm P_F :

$$P_F(\delta) \triangleq \int_{\mathcal{Y}_1} f_{\mathbf{y}|\mathbf{H}_0}(\mathbf{y}|\mathbf{H}_0) d\mathbf{y}. \quad (2)$$

- ▶ Notice that both are calculated over \mathcal{Y}_1 , i.e., for space of measurements where decision is $\delta(\mathbf{y}) = 1$.
- ▶ Think of a radar system: false alarm is when the radar reports an airplane is coming, when it is not.

Probability of Detection vs Probability of False Alarm

- Remember the Bayes conditional risk:

$$R(\delta(\mathbf{y})|H_j) = C_{1j} \Pr(\delta(\mathbf{y}) = 1|H_j) + C_{0j} \Pr(\delta(\mathbf{y}) = 0|H_j)$$

Thus,

$$\begin{aligned} R(\delta(\mathbf{y})|H_0) &= C_{10} \Pr(\delta(\mathbf{y}) = 1|H_0) + C_{00} \Pr(\delta(\mathbf{y}) = 0|H_0) \\ &= C_{10} P_F + C_{00} (1 - P_F). \end{aligned} \quad (3)$$

and similarly,

$$\begin{aligned} R(\delta(\mathbf{y})|H_1) &= C_{11} \Pr(\delta(\mathbf{y}) = 1|H_1) + C_{01} \Pr(\delta(\mathbf{y}) = 0|H_1) \\ &= C_{11} P_D + C_{01} (1 - P_D). \end{aligned} \quad (4)$$

- Thus, the (unconditional) Bayes risk of test δ is fully characterised by the pair (P_F, P_D) of specific test δ :

$$\begin{aligned} R(\delta) &= R(\delta(\mathbf{y})|H_0) \pi_0 + R(\delta(\mathbf{y})|H_1) \pi_1 \\ &= C_{00} \pi_0 + C_{01} \pi_1 + \pi_0 (C_{10} - C_{00}) P_F + \pi_1 (C_{11} - C_{01}) P_D \end{aligned} \quad (5)$$

Probability of Detection vs Probability of False Alarm

- ▶ The (unconditional) Bayes risk of test δ is fully characterised by the pair (P_F, P_D) of specific test δ :

$$\begin{aligned} R(\delta) &= R(\delta(\mathbf{y})|H_0) \pi_0 + R(\delta(\mathbf{y})|H_1) \pi_1 \\ &= C_{00} \pi_0 + C_{01} \pi_1 + \pi_0 (C_{10} - C_{00}) P_F + \pi_1 (C_{11} - C_{01}) P_D \end{aligned} \tag{6}$$

- ▶ Ideally, we would like to have a test (or a binary classifier) with $(P_F, P_D) = (0, 1)$. However, this is not feasible!
- ▶ Next lecture will offer the feasible pairs (P_F, P_D) , as well as properties of the boundary between feasible and non-feasible pairs for all tests!
- ▶ Boundary between feasible and non-feasible tests: receiver operating characteristic (ROC) [next lecture].

Neyman-Pearson Test

- ▶ Problem definition: among the tests that bound false alarm probability, find the test that maximises probability of detection. The formulation is given below:

$$\delta_{\text{NP}} = \arg \max_{\delta \in D_\alpha} P_D(\delta) \quad (7)$$

$$D_\alpha = \{\text{all tests } \delta : P_F(\delta) \leq \alpha\}$$

- ▶ The problem could be also formulated with bounded probability of detection and minimised probability of false alarm.

Neyman-Pearson Test Derivation

- ▶ Constrained optimization problem: we need to set Lagrangian and KKT condition(s).
 - ▶ Lagrangian $L(\cdot, \cdot)$ for Lagrange multiplier $\lambda \geq 0$:

$$\begin{aligned} L(\delta, \lambda) &= P_D(\delta) + \lambda(\alpha - P_F(\delta)) \\ &= \lambda\alpha + \int_{\mathcal{Y}_1} [f_{\mathbf{y}|\mathbf{H}_1}(\mathbf{y}|\mathbf{H}_1) - \lambda f_{\mathbf{y}|\mathbf{H}_0}(\mathbf{y}|\mathbf{H}_0)] d\mathbf{y} \quad (8) \end{aligned}$$

- ▶ KKT condition ($\lambda \geq 0$):

$$\lambda(\alpha - P_F(\delta)) = 0 \quad (9)$$

- ▶ Optimal test maximizes Lagrangian in Eq. (10) AND also satisfies KKT condition in Eq. (9).

Neyman-Pearson Test Derivation

- ▶ Maximization of Lagrangian $L(\delta, \lambda) (\lambda \geq 0)$:

$$L(\delta, \lambda) = \lambda \alpha + \int_{\mathcal{Y}_1} [f_{\mathbf{y}|\mathbf{H}_1}(\mathbf{y}|\mathbf{H}_1) - \lambda f_{\mathbf{y}|\mathbf{H}_0}(\mathbf{y}|\mathbf{H}_0)] d\mathbf{y} \quad (10)$$

- ▶ if $\mathbf{y} \in \mathcal{Y}_1$ then $[f_{\mathbf{y}|\mathbf{H}_1}(\mathbf{y}|\mathbf{H}_1) - \lambda f_{\mathbf{y}|\mathbf{H}_0}(\mathbf{y}|\mathbf{H}_0)] > 0$, otherwise the Lagrangian is not maximized.
 - ▶ Equivalently, if $\mathbf{y} \in \mathcal{Y}_0$, then $[f_{\mathbf{y}|\mathbf{H}_1}(\mathbf{y}|\mathbf{H}_1) - \lambda f_{\mathbf{y}|\mathbf{H}_0}(\mathbf{y}|\mathbf{H}_0)] < 0$ and that particular \mathbf{y} is not taken into account in Eq. (10).
 - ▶ for any \mathbf{y} with $[f_{\mathbf{y}|\mathbf{H}_1}(\mathbf{y}|\mathbf{H}_1) - \lambda f_{\mathbf{y}|\mathbf{H}_0}(\mathbf{y}|\mathbf{H}_0)] = 0$, what decision should we adopt?
- ▶ Thus, maximization of the Lagrangian results to testing the sign of $[f_{\mathbf{y}|\mathbf{H}_1}(\mathbf{y}|\mathbf{H}_1) - \lambda f_{\mathbf{y}|\mathbf{H}_0}(\mathbf{y}|\mathbf{H}_0)]$.

Neyman-Pearson Test Derivation

- ▶ Maximization of the Lagrangian results to testing the sign of $\left[f_{\mathbf{y}|\mathbf{H}_1}(\mathbf{y}|\mathbf{H}_1) - \lambda f_{\mathbf{y}|\mathbf{H}_0}(\mathbf{y}|\mathbf{H}_0) \right]$. Thus, optimal test for given $\lambda \geq 0$ follows:

$$\delta(\mathbf{y}) = \begin{cases} 1, & f_{\mathbf{y}|\mathbf{H}_1}(\mathbf{y}|\mathbf{H}_1) - \lambda f_{\mathbf{y}|\mathbf{H}_0}(\mathbf{y}|\mathbf{H}_0) > 0 \\ 0, & f_{\mathbf{y}|\mathbf{H}_1}(\mathbf{y}|\mathbf{H}_1) - \lambda f_{\mathbf{y}|\mathbf{H}_0}(\mathbf{y}|\mathbf{H}_0) < 0 \\ 0 \text{ or } 1 \text{ (TBD)}, & f_{\mathbf{y}|\mathbf{H}_1}(\mathbf{y}|\mathbf{H}_1) - \lambda f_{\mathbf{y}|\mathbf{H}_0}(\mathbf{y}|\mathbf{H}_0) = 0. \end{cases} \quad (11)$$

- ▶ Setting the likelihood ratio

$L(\mathbf{y}) \triangleq f_{\mathbf{y}|\mathbf{H}_1}(\mathbf{y}|\mathbf{H}_1)/f_{\mathbf{y}|\mathbf{H}_0}(\mathbf{y}|\mathbf{H}_0)$, Eq. (11) is equivalent to:

$$\delta(\mathbf{y}) = \begin{cases} 1, & L(\mathbf{y}) > \lambda, \\ 0, & L(\mathbf{y}) < \lambda, \\ 0 \text{ or } 1 \text{ (TBD)}, & L(\mathbf{y}) = \lambda. \end{cases} \quad (12)$$

Neyman-Pearson Test Derivation

- ▶ Maximization of the Lagrangian results to the following optimal test for given $\lambda \geq 0$:

$$\delta(\mathbf{y}) = \begin{cases} 1, & L(\mathbf{y}) > \lambda, \\ 0, & L(\mathbf{y}) < \lambda, \\ 0 \text{ or } 1 \text{ (TBD)}, & L(\mathbf{y}) = \lambda. \end{cases} \quad (13)$$

- ▶ We need to find out λ and decision for $L(\mathbf{y}) = \lambda$.
- ▶ We define the following conditional cumulative distribution function (cdf):

$$F_L(l|H_0) \triangleq \Pr(L(\mathbf{y}) \leq l|H_0). \quad (14)$$

- ▶ As any cdf, the above should be:
 1. right-continuous,
 2. non-decreasing with increasing l ,
 3. 0 for $l \rightarrow -\infty$ and
 4. 1 for $l \rightarrow +\infty$.

Neyman-Pearson Test Derivation

- ▶ Any cdf is "right-continuous": check figure below!

$$F_L(l|H_0) \triangleq \Pr(L(\mathbf{y}) \leq l|H_0). \quad (15)$$

- ▶ As any cdf, the above should be:
 1. right-continuous,
 2. non-decreasing with increasing l ,
 3. 0 for $l \rightarrow -\infty$ and
 4. 1 for $l \rightarrow +\infty$.

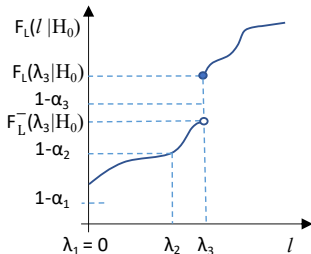
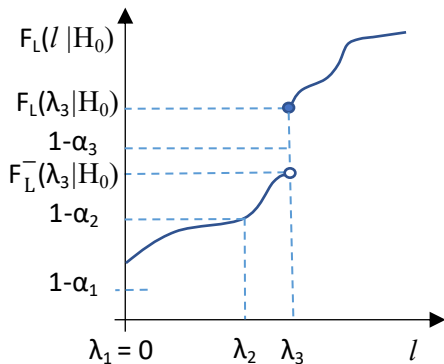


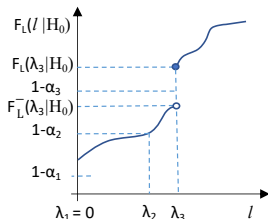
Figure 1: Example cdf with right-continuity.

Neyman-Pearson Test Derivation



- ▶ Three cases of likelihood ratio threshold λ occur, depending on value of $1 - \alpha$ vs $F_L(l = 0 | H_0)$.
- ▶ Reminder: α is the upper bound of P_F .

Neyman-Pearson Test Derivation

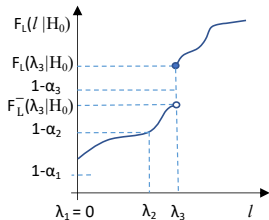


- ▶ Case I: $1 - \alpha < F_L(l = 0|H_0) \triangleq f_0 \Leftrightarrow 1 - f_0 < \alpha$:
 - ▶ set $\lambda = 0$ and $\delta(\mathbf{y}) = 0$ for $L(\mathbf{y}) = \lambda = 0$, i.e.,

$$\delta(\mathbf{y}) = \begin{cases} 1, & L(\mathbf{y}) > \lambda = 0, \\ 0, & L(\mathbf{y}) \leq \lambda = 0. \end{cases} \quad (16)$$

- ▶ KTT $\lambda(\alpha - P_F) = 0$ is satisfied for $\lambda = 0$.
- ▶ $P_F = 1 - \Pr(\delta(\mathbf{y}) = 0|H_0) = 1 - \Pr(L(\mathbf{y}) \leq 0|H_0) = 1 - f_0 < \alpha \Rightarrow$ probability of false alarm is bounded.

Neyman-Pearson Test Derivation

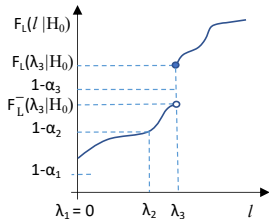


- ▶ Case II: $1 - \alpha \geq F_L(l=0|H_0) \triangleq f_0$ and λ^* is in the range of $F_L(l|H_0)$, i.e., there is λ^* such that $F_L(\lambda^*|H_0) = 1 - \alpha$.
 - ▶ set $\lambda = \lambda^*$ and $\delta(\mathbf{y}) = 0$ for $L(\mathbf{y}) = \lambda^*$, i.e.,

$$\delta(\mathbf{y}) = \begin{cases} 1, & L(\mathbf{y}) > \lambda^*, \\ 0, & L(\mathbf{y}) \leq \lambda^*. \end{cases} \quad (17)$$

- ▶ $P_F = 1 - \Pr(\delta(\mathbf{y}) = 0|H_0) = 1 - \Pr(L(\mathbf{y}) \leq \lambda^*|H_0) = 1 - (1 - \alpha) = \alpha$.
- ▶ KTT $\lambda(\alpha - P_F) = 0$ is satisfied for $\lambda = \lambda^*$.

Neyman-Pearson Test Derivation

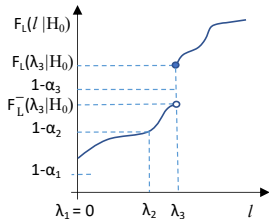


- ▶ Case III: $1 - \alpha \geq F_L(l = 0|H_0) \triangleq f_0$ and λ^* is NOT in the range of $F_L(l|H_0)$, i.e., $F_L^-(\lambda^*|H_0) < 1 - \alpha < F_L(\lambda^*|H_0)$.
 - ▶ set $\lambda = \lambda^*$ and $\delta(\mathbf{y}) = 0$ for $L(\mathbf{y}) = \lambda$, i.e.,

$$\delta(\mathbf{y}) \triangleq \delta_{L, \lambda^*}(\mathbf{y}) = \begin{cases} 1, & L(\mathbf{y}) > \lambda^*, \\ 0, & L(\mathbf{y}) \leq \lambda^*. \end{cases} \quad (18)$$

- ▶ $P_F = 1 - \Pr(\delta(\mathbf{y}) = 0|H_0) = 1 - \Pr(L(\mathbf{y}) \leq \lambda^*|H_0) = 1 - F_L(\lambda^*|H_0) < \alpha$.
- ▶ KTT $\lambda(\alpha - P_F) = 0$ is NOT satisfied!

Neyman-Pearson Test Derivation

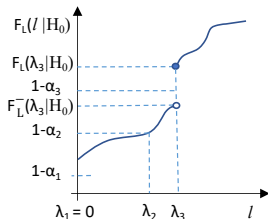


- ▶ Case III: $1 - \alpha \geq F_L(l = 0|H_0) \triangleq f_0$ and λ^* is NOT in the range of $F_L(l|H_0)$, i.e., $F_L^-(\lambda^*|H_0) < 1 - \alpha < F_L(\lambda^*|H_0)$.
 - ▶ set $\lambda = \lambda^*$ and $\delta(\mathbf{y}) = 1$ for $L(\mathbf{y}) = \lambda$, i.e.,

$$\delta(\mathbf{y}) \triangleq \delta_{U, \lambda^*}(\mathbf{y}) = \begin{cases} 1, & L(\mathbf{y}) \geq \lambda^*, \\ 0, & L(\mathbf{y}) < \lambda^*. \end{cases} \quad (19)$$

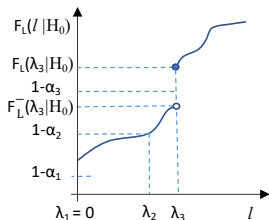
- ▶ $P_F = 1 - \Pr(\delta(\mathbf{y}) = 0|H_0) = 1 - \Pr(L(\mathbf{y}) < \lambda^*|H_0) = 1 - F_L^-(\lambda^*|H_0) > \alpha$.
- ▶ KTT $\lambda(\alpha - P_F) = 0$ is NOT satisfied!

Neyman-Pearson Test Derivation



- ▶ Case III: $1 - \alpha \geq F_L(l=0|H_0) \triangleq f_0$ and λ^* is NOT in the range of $F_L(l|H_0)$, i.e., $F_L^-(\lambda^*|H_0) < 1 - \alpha < F_L(\lambda^*|H_0)$.
So far:
 - ▶ Test $\delta_{L,\lambda^*}(\mathbf{y})$ with $P_F(\delta_{L,\lambda^*}) < \alpha$.
 - ▶ Test $\delta_{U,\lambda^*}(\mathbf{y})$ with $P_F(\delta_{U,\lambda^*}) > \alpha$.
 - ▶ KTT $\lambda(\alpha - P_F) = 0$ requires exactly $P_F = \alpha$.
- ▶ Solution?

Neyman-Pearson Test Derivation

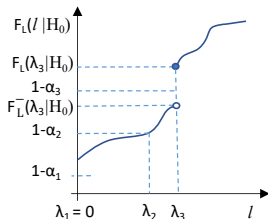


- ▶ Case III: $1 - \alpha \geq F_L(l = 0 | H_0) \triangleq f_0$ and λ^* is NOT in the range of $F_L(l | H_0)$, i.e., $F_L^-(\lambda^* | H_0) < 1 - \alpha < F_L(\lambda^* | H_0)$.
 - ▶ Solution: set $\lambda = \lambda^*$ and randomize decision for $L(\mathbf{y}) = \lambda$:

$$\delta(\mathbf{y}) = \begin{cases} \delta_{U, \lambda^*}(\mathbf{y}), & \text{with probability } \rho, \\ \delta_{L, \lambda^*}(\mathbf{y}), & \text{with probability } 1 - \rho, \end{cases} \quad (20)$$

- ▶ Set $0 < \rho < 1$ such that $P_F \equiv \alpha$ (and KKT is thus satisfied).

Neyman-Pearson Test Derivation



- ▶ Case III: $1 - \alpha \geq F_L(l=0|H_0) \triangleq f_0$ and λ^* is NOT in the range of $F_L(l|H_0)$, i.e., $F_L^-(\lambda^*|H_0) < 1 - \alpha < F_L(\lambda^*|H_0)$.
 - ▶ Solution: set $\lambda = \lambda^*$ and randomize decision for $L(\mathbf{y}) = \lambda$:

$$\delta(\mathbf{y}) = \begin{cases} \delta_{U,\lambda^*}(\mathbf{y}), & \text{with probability } \rho, \\ \delta_{L,\lambda^*}(\mathbf{y}), & \text{with probability } 1 - \rho, \end{cases} \quad (21)$$

$$\rho = \frac{F_L(\lambda^*|H_0) - (1 - \alpha)}{F_L(\lambda^*|H_0) - F_L^-(\lambda^*|H_0)}, \quad 0 < \rho < 1. \quad (22)$$

- ▶ Such $0 < \rho < 1$ guarantees $P_F \equiv \alpha$.

Neyman-Pearson Test Derivation

- ▶ Case III: $1 - \alpha \geq F_L(l=0|H_0) \triangleq f_0$ and λ^* is NOT in the range of $F_L(l|H_0)$, i.e., $F_L^-(\lambda^*|H_0) < 1 - \alpha < F_L(\lambda^*|H_0)$.
 - ▶ Solution: set $\lambda = \lambda^*$ and randomize decision for $L(\mathbf{y}) = \lambda$:

$$\delta(\mathbf{y}) = \begin{cases} \delta_{U,\lambda^*}(\mathbf{y}), & \text{with probability } \rho, \\ \delta_{L,\lambda^*}(\mathbf{y}), & \text{with probability } 1 - \rho, \end{cases} \quad (23)$$

$$\rho = \frac{F_L(\lambda^*|H_0) - (1 - \alpha)}{F_L(\lambda^*|H_0) - F_L^-(\lambda^*|H_0)}, \quad 0 < \rho < 1. \quad (24)$$

- ▶ Such ρ guarantees $P_F \equiv \alpha$.

Proof:

$$P_F = \rho P_F(\delta_{U,\lambda^*}) + (1 - \rho) P_F(\delta_{L,\lambda^*}) \quad (25)$$

$$= \rho (1 - F_L^-(\lambda^*|H_0)) + (1 - \rho) (1 - F_L(\lambda^*|H_0)) \quad (26)$$

$$= 1 - F_L(\lambda^*|H_0) + \rho (F_L(\lambda^*|H_0) - F_L^-(\lambda^*|H_0)) \quad (27)$$

$$= 1 - F_L(\lambda^*|H_0) + F_L(\lambda^*|H_0) - (1 - \alpha) \quad (28)$$

$$= \alpha \quad (29)$$

Neyman-Pearson Test

- ▶ Neyman-Pearson-optimal detector is likelihood ratio test (LRT)!
- ▶ As already mentioned, we could have minimized P_F subject to bounded P_D .
- ▶ Next lecture: feasible points (P_F, P_D) for any test!

- [1] Bernard C. Levy, Principles of Signal Detection and Parameter Estimation, Springer 2008.
- [2] Class instructor notes.

Thank you!



Detection & Estimation Theory: Lectures 5 & 6

Prof. Aggelos Bletsas

Technical University of Crete
School of Electrical and Computer Engineering



European Union
European Social Fund

Operational Programme
**Human Resources Development,
Education and Lifelong Learning**

Co-financed by Greece and the European Union





“Support for the Internationalization of Higher Education, School of Electrical and Computer Engineering, Technical University of Crete” (MIS 5150766), under the call for proposals “Support of the Internationalization of Higher Education - Technical University of Crete” (EDULLL 153).

The project is co-financed by Greece and the European Union (European Social Fund- ESF) through the Operational Programme “Human Resources Development, Education and Lifelong Learning 2014-2020.



Operational Programme
**Human Resources Development,
Education and Lifelong Learning**
Co-financed by Greece and the European Union



- Receiver Operating Characteristic (ROC)
- ROC Properties
- Remarks
- Examples

- ▶ Remember the two probability metrics that fully characterise a test $\delta(\mathbf{y})$ (as well as any binary classifier):
 - ▶ Prob. of detection P_D and prob. of false alarm P_F :

$$P_D(\delta) \triangleq \int_{\mathcal{Y}_1} f_{\mathbf{y}|\mathbf{H}_1}(\mathbf{y}|\mathbf{H}_1)d\mathbf{y}, P_F(\delta) \triangleq \int_{\mathcal{Y}_1} f_{\mathbf{y}|\mathbf{H}_0}(\mathbf{y}|\mathbf{H}_0)d\mathbf{y}.$$

- ▶ The Bayes risk of test δ is fully characterised by the pair (P_F, P_D) of specific test δ (previous lecture).
- ▶ Which pairs (P_F, P_D) are feasible?
- ▶ Boundary between feasible and non-feasible tests = receiver operating characteristic (ROC).
- ▶ ROC properties?

ROC definition

- ▶ Remember the two probability metrics that fully characterise a test $\delta(\mathbf{y})$ (as well as any binary classifier):
 - ▶ Prob. of detection P_D and prob. of false alarm P_F :

$$P_D(\delta) \triangleq \int_{\mathcal{Y}_1} f_{\mathbf{y}|\mathbf{H}_1}(\mathbf{y}|\mathbf{H}_1) d\mathbf{y}, \quad P_F(\delta) \triangleq \int_{\mathcal{Y}_1} f_{\mathbf{y}|\mathbf{H}_0}(\mathbf{y}|\mathbf{H}_0) d\mathbf{y}.$$

- ▶ Define likelihood ratio and conditional pdf of likelihood ratio $f_{L|\mathbf{H}_j}(l|\mathbf{H}_j)$:

$$L(\mathbf{y}) \triangleq \frac{f_{\mathbf{y}|\mathbf{H}_1}(\mathbf{y}|\mathbf{H}_1)}{f_{\mathbf{y}|\mathbf{H}_0}(\mathbf{y}|\mathbf{H}_0)}. \quad (1)$$

- ▶ For likelihood ratio test $L(\mathbf{y}) \stackrel{H_1}{\geq} \tau$, P_F, P_D can be redefined:

$$P_D(\tau) \triangleq \int_{\tau}^{+\infty} f_{L|\mathbf{H}_1}(l|\mathbf{H}_1) dl, \quad P_F(\tau) \triangleq \int_{\tau}^{+\infty} f_{L|\mathbf{H}_0}(l|\mathbf{H}_0) dl.$$

ROC Property 1

- ▶ Points $(0, 0)$ and $(1, 1)$ of (P_F, P_D) belong to ROC.

Proof:

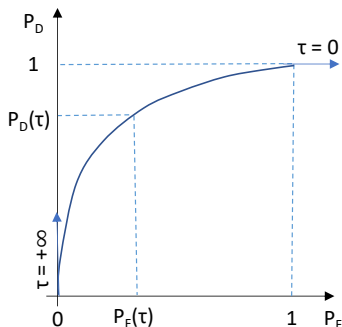
- ▶ Always select H_1 (or equivalently, set $\tau = 0$):

$$P_F(\tau = 0) = 1, P_D(\tau = 0) = 1. \quad (2)$$

- ▶ Always select H_0 (or equivalently, set $\tau = +\infty$):

$$P_F(\tau = +\infty) = 0, P_D(\tau = +\infty) = 0. \quad (3)$$

ROC Property 2



- Slope of ROC at $(P_F(\tau), P_D(\tau))$ is equal to τ , i.e.,

$$\frac{dP_D(\tau)}{dP_F(\tau)} = \tau.$$

Proof:

- $P_D(\delta) = \int_{\tau}^{+\infty} f_{L|H_1}(l|H_1) dy = 1 - \int_{-\infty}^{+\tau} f_{L|H_1}(l|H_1) dy \Rightarrow$

ROC Property 2

$$\blacktriangleright P_D(\tau) = \int_{\tau}^{+\infty} f_{L|H_1}(l|H_1) dl = 1 - \int_{-\infty}^{+\tau} f_{L|H_1}(l|H_1) dl \Rightarrow$$

$$\frac{dP_D}{d\tau}(\tau) = -f_{L|H_1}(\tau|H_1), \quad (4)$$

$$\text{similarly, } \frac{dP_F}{d\tau}(\tau) = -f_{L|H_0}(\tau|H_0). \quad (5)$$

From Eqs. (4), (5):

$$\frac{dP_D}{dP_F}(\tau) = \frac{f_{L|H_1}(\tau|H_1)}{f_{L|H_0}(\tau|H_0)}. \quad (6)$$

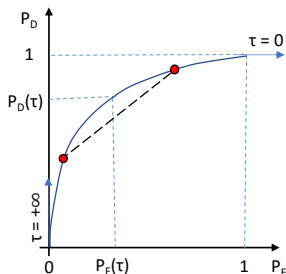
\blacktriangleright Recall that:

$$\begin{aligned} P_D(\tau) &= \int_{\tau}^{+\infty} f_{L|H_1}(l|H_1) dl = \int_{\mathcal{Y}_1} f_{\mathbf{y}|H_1}(\mathbf{y}|H_1) d\mathbf{y} = \int_{\mathcal{Y}_1} \frac{f_{\mathbf{y}|H_1}(\mathbf{y}|H_1)}{f_{\mathbf{y}|H_0}(\mathbf{y}|H_0)} f_{\mathbf{y}|H_0}(\mathbf{y}|H_0) d\mathbf{y} \\ &= \int_{\mathcal{Y}_1 = \{\mathbf{y}: L(\mathbf{y}) \geq \tau\}} L(\mathbf{y}) f_{\mathbf{y}|H_0}(\mathbf{y}|H_0) d\mathbf{y} = \int_{\tau}^{+\infty} l \cdot f_{L|H_0}(l|H_0) dl \Rightarrow \end{aligned} \quad (7)$$

$$\frac{dP_D}{d\tau}(\tau) = -\tau \cdot f_{L|H_0}(\tau|H_0) \stackrel{(4)}{\Rightarrow} f_{L|H_1}(\tau|H_1) = \tau f_{L|H_0}(\tau|H_0). \quad (8)$$

From Eqs. (6), (8), the proof is completed.

ROC Property 3



- ▶ The domain of feasible points (P_F, P_D) is convex.

Proof:

- ▶ We need to show that for any two feasible points (P_{F1}, P_{D1}) , (P_{F2}, P_{D2}) , the line connecting them is also included in the feasible points.
- ▶ Such line is described by $\rho \in [0, 1]$:

$$P_F(\rho) = \rho P_{F1} + (1 - \rho) P_{F2}, \quad (9)$$

$$P_D(\rho) = \rho P_{D1} + (1 - \rho) P_{D2}. \quad (10)$$

ROC Property 3

- ▶ Such line is described by $\rho \in [0, 1]$:

$$P_F(\rho) = \rho P_{F1} + (1 - \rho) P_{F2}, \quad (11)$$

$$P_D(\rho) = \rho P_{D1} + (1 - \rho) P_{D2}. \quad (12)$$

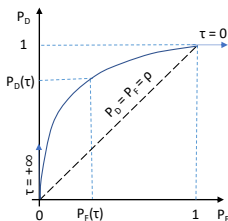
- ▶ Why? solve each of the above for ρ and equate: you will find out the line equation connecting the two points.
- ▶ Define the following randomized test that selects probabilistic between two tests:

$$\delta(\mathbf{y}) = \begin{cases} \delta_1(\mathbf{y}), & \text{with probability } \rho, \\ \delta_2(\mathbf{y}), & \text{with probability } 1 - \rho, \end{cases} \quad (13)$$

where $\delta_i(\mathbf{y})$ is the feasible test with P_{Fi}, P_{Di} , $i \in \{1, 2\}$.

- ▶ The test above achieves $P_F(\rho)$ and $P_D(\rho)$ given in Eqs. (11), (12).

ROC Property 3 Remarks

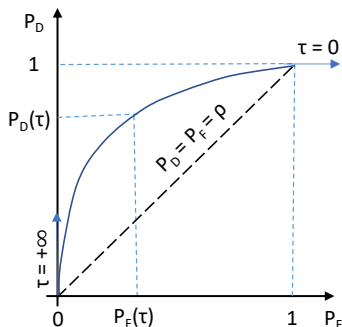


- Define the following randomized test that selects probabilistic between two hypothesis, *independently* of the measurements \mathbf{y} :

$$\delta(\mathbf{y}) = \begin{cases} H_1, & \text{with probability } \rho, \\ H_0, & \text{with probability } 1 - \rho, \end{cases} \quad (14)$$

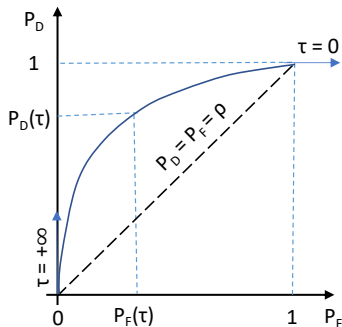
- The test above achieves $P_D(\rho) = \Pr(\mathcal{Y}_1|H_1) = \Pr(\mathcal{Y}_1) = \rho$ and $P_F(\rho) = \Pr(\mathcal{Y}_1|H_0) = \Pr(\mathcal{Y}_1) = \rho$.
- Thus, line $P_D = P_F = \rho$ belongs to the feasible points.

ROC Property 3 Remarks



- ▶ The domain of feasible points (P_F, P_D) is convex.
- ▶ Line $P_D = P_F = \rho$ belongs to the feasible points.
 - ▶ ...thus, domain of feasible points are located *below* the ROC curve!
- ▶ Domain of feasible points is convex and located below ROC.
 - ▶ ...thus, ROC curve is concave!

ROC Property 4

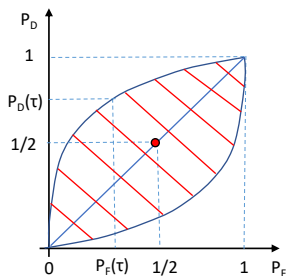


- ▶ For tests on the ROC curve, $P_F \leq P_D$.

Proof:

- ▶ ROC curve is concave.
- ▶ $(0, 0)$ and $(1, 1)$ belong to the ROC curve.
 - ▶ ...thus, $P_F \leq P_D$.

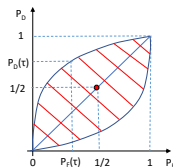
Remarks



- ▶ What about “bad” tests? Are all points with $P_D \leq P_F$ feasible? [NO!]
- ▶ Assume a “bad” test δ with specific $(P_F(\delta), P_D(\delta))$.
- ▶ Define the following test $\hat{\delta}$ that flips the decision:

$$\hat{\delta}(\mathbf{y}) = 1 - \delta(\mathbf{y}). \quad (15)$$

- ▶ $\hat{\delta}(\mathbf{y})$ achieves $(\widehat{P}_F, \widehat{P}_D) = (1 - P_F(\delta), 1 - P_D(\delta))$, i.e., region of feasible tests is symmetric around $(1/2, 1/2)$.



- Define the following test $\hat{\delta}$ that flips the decision:

$$\hat{\delta}(\mathbf{y}) = 1 - \delta(\mathbf{y}). \quad (16)$$

- $\hat{\delta}(\mathbf{y})$ achieves $(\widehat{P}_F, \widehat{P}_D) = (1 - P_F(\delta), 1 - P_D(\delta))$, i.e., region of feasible tests is symmetric around $(1/2, 1/2)$. Proof:

$$P_F(\delta) = \int_{\mathcal{Y}_1} f_{\mathbf{y}|\mathbf{H}_0}(\mathbf{y}|\mathbf{H}_0) d\mathbf{y} \Rightarrow \quad (17)$$

$$P_F(\hat{\delta} = 1 - \delta) = \int_{\mathcal{Y}_0} f_{\mathbf{y}|\mathbf{H}_0}(\mathbf{y}|\mathbf{H}_0) d\mathbf{y} = 1 - P_F(\delta). \quad (18)$$

...and similarly for $P_D(\hat{\delta})$.

Examples

- Assume binary hypothesis testing, with $\mathbf{y} \sim \mathcal{N}(\mathbf{m}_j, \mathbf{K}_j)$ under hypothesis H_j , $j \in \{0, 1\}$ and $\mathbf{y} \in \mathbb{R}^{N \times 1}$.

Preliminaries:

- For any matrix \mathbf{K} with inverse, $(\mathbf{K}^{-1})^T = (\mathbf{K}^T)^{-1}$.
- Thus, for any *symmetric* matrix \mathbf{K} with inverse, the inverse is also symmetric:
 $(\mathbf{K}^{-1})^T = \mathbf{K}^{-1}$.
- For any scalar z with $z = \mathbf{a}^T \mathbf{b}$, $z = \mathbf{b}^T \mathbf{a}$, where \mathbf{a} , \mathbf{b} vectors of the same dimension.

$$f_{\mathbf{y}|H_j}(\mathbf{y}|H_j) = \frac{1}{\sqrt{(2\pi)^N |\mathbf{K}_j|}} \exp\left(-\frac{1}{2}(\mathbf{y} - \mathbf{m}_j)^T \mathbf{K}_j^{-1}(\mathbf{y} - \mathbf{m}_j)\right) \Leftrightarrow \quad (19)$$

$$\frac{f_{\mathbf{y}|H_1}(\mathbf{y}|H_1)}{f_{\mathbf{y}|H_0}(\mathbf{y}|H_0)} \stackrel{H_1}{\geq} \tau \Leftrightarrow \quad (20)$$

$$\frac{1}{2}(\mathbf{y} - \mathbf{m}_0)^T \mathbf{K}_0^{-1}(\mathbf{y} - \mathbf{m}_0) - \frac{1}{2}(\mathbf{y} - \mathbf{m}_1)^T \mathbf{K}_1^{-1}(\mathbf{y} - \mathbf{m}_1) + \frac{1}{2} \ln\left(\frac{|\mathbf{K}_0|}{|\mathbf{K}_1|}\right) \stackrel{H_1}{\geq} \ln(\tau) \Leftrightarrow \quad (21)$$

...simple calculations exploiting 2. and 3. above [try them!]

$$\Leftrightarrow \underbrace{\frac{1}{2} \mathbf{y}^T (\mathbf{K}_0^{-1} - \mathbf{K}_1^{-1}) \mathbf{y} + \mathbf{y}^T (\mathbf{K}_1^{-1} \mathbf{m}_1 - \mathbf{K}_0^{-1} \mathbf{m}_0)}_{S(\mathbf{y})} \stackrel{H_1}{\geq} \eta, \quad (22)$$

$$\eta = \frac{1}{2} \mathbf{m}_1^T \mathbf{K}_1 \mathbf{m}_1 - \frac{1}{2} \mathbf{m}_0^T \mathbf{K}_0 \mathbf{m}_0 - \frac{1}{2} \ln\left(\frac{|\mathbf{K}_0|}{|\mathbf{K}_1|}\right) + \ln(\tau) \quad (23)$$

Examples

- Assume $\mathbf{K}_0 = \mathbf{K}_1 = \mathbf{K}$. In that case, the test is simplified as follows:

$$\mathbf{y}^T \mathbf{K}^{-1} \underbrace{(\mathbf{m}_1 - \mathbf{m}_0)}_{\Delta \mathbf{m}} = \mathbf{y}^T \mathbf{K}^{-1} \Delta \mathbf{m} \stackrel{H_1}{\geq} \eta \Leftrightarrow \quad (24)$$

$$S_s(\mathbf{y}) \triangleq \mathbf{y}^T \mathbf{K}^{-1} \Delta \mathbf{m} - \mathbf{m}_0^T \mathbf{K}^{-1} \Delta \mathbf{m} \stackrel{H_1}{\geq} \eta - \mathbf{m}_0^T \mathbf{K}^{-1} \Delta \mathbf{m} \triangleq \eta_s \Leftrightarrow \quad (25)$$

$$S_s(\mathbf{y}) = \Delta \mathbf{m}^T \mathbf{K}^{-1} \mathbf{y} - \Delta \mathbf{m}^T \mathbf{K}^{-1} \mathbf{m}_0 \stackrel{H_1}{\geq} \eta - \mathbf{m}_0^T \mathbf{K}^{-1} \Delta \mathbf{m} \triangleq \eta_s, \quad (26)$$

where we have used the fact that \mathbf{K}^{-1} is symmetric; $S_s(\mathbf{y})$ is the shifted sufficient statistic, which is affine transformation of a Gaussian vector, and thus, it is also Gaussian:

$$H_0 : S_s(\mathbf{y}) \sim \mathcal{N}\left(0, \Delta \mathbf{m}^T \mathbf{K}^{-1} \Delta \mathbf{m}\right) \quad (27)$$

$$H_1 : S_s(\mathbf{y}) \sim \mathcal{N}\left(\Delta \mathbf{m}^T \mathbf{K}^{-1} \Delta \mathbf{m}, \Delta \mathbf{m}^T \mathbf{K}^{-1} \Delta \mathbf{m}\right) \quad (28)$$

Notice that \mathbf{K}^{-1} (and \mathbf{K}) are positive definite, and thus, $\Delta \mathbf{m}^T \mathbf{K}^{-1} \Delta \mathbf{m} > 0$. We set $d^2 \triangleq \Delta \mathbf{m}^T \mathbf{K}^{-1} \Delta \mathbf{m}$. We also need the following definition of the (Gauss) Q-function $Q(x)$ and its properties:

$$Q(x) \triangleq \int_x^{+\infty} \frac{1}{2\pi} e^{-t^2/2} dt = 1 - \int_{-\infty}^x \frac{1}{2\pi} e^{-t^2/2} dt, \quad (29)$$

$$Q(-x) = 1 - Q(x), \quad \frac{dQ(x)}{dx} = -\frac{1}{2\pi} e^{-x^2/2}. \quad (30)$$

Examples

$$H_0 : S_s(\mathbf{y}) \sim \mathcal{N}\left(0, \Delta \mathbf{m}^T \mathbf{K}^{-1} \Delta \mathbf{m}\right) \equiv \mathcal{N}\left(0, d^2\right) \quad (31)$$

$$H_1 : S_s(\mathbf{y}) \sim \mathcal{N}\left(\Delta \mathbf{m}^T \mathbf{K}^{-1} \Delta \mathbf{m}, \Delta \mathbf{m}^T \mathbf{K}^{-1} \Delta \mathbf{m}\right) \equiv \mathcal{N}\left(d^2, d^2\right) \quad (32)$$

We can now calculate the basic probabilities for the test $S_s(\mathbf{y}) \stackrel{H_1}{\geq} \eta_s$:

$$P_D = \int_{\eta_s}^{+\infty} f_{S_s(\mathbf{y})|H_1}(s|H_1) ds = \int_{\eta_s}^{+\infty} \frac{1}{\sqrt{2\pi d^2}} \exp\left(-\frac{(s-d)^2}{2d^2}\right) ds \quad (33)$$

$$\stackrel{\frac{s-d}{d}=t}{=} \int_{\frac{\eta_s-d}{d}}^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt = Q\left(\frac{\eta_s-d}{d}\right) = 1 - Q\left(d - \frac{\eta_s}{d}\right), \quad (34)$$

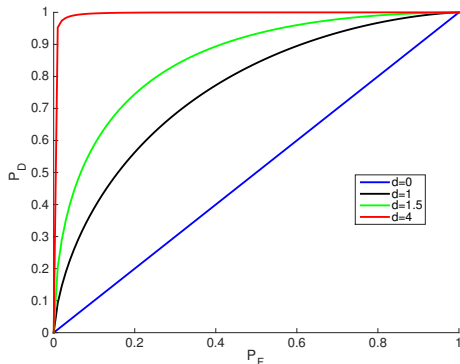
$$P_F = \int_{\eta_s}^{+\infty} f_{S_s(\mathbf{y})|H_0}(s|H_0) ds = \int_{\eta_s}^{+\infty} \frac{1}{\sqrt{2\pi d^2}} \exp\left(-\frac{s^2}{2d^2}\right) ds \quad (35)$$

$$\stackrel{\frac{s}{d}=t}{=} \int_{\frac{\eta_s}{d}}^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-t^2/2} dt = Q\left(\frac{\eta_s}{d}\right) \Rightarrow \frac{\eta_s}{d} = Q^{-1}(P_F). \quad (36)$$

$$(34), (36) \Rightarrow P_D = 1 - Q\left(d - Q^{-1}(P_F)\right). \quad (37)$$

$$P_D = 1 - Q\left(d - Q^{-1}(P_F)\right). \quad (38)$$

- We need d as large as possible! Why?



Whitening Procedure

- ▶ In many cases, it is useful to simplify the observation model with linear transformations. To this end, the eigendecomposition of positive-definite matrix \mathbf{K} is exploited:

$$\mathbf{K}\mathbf{P} = \mathbf{P}\mathbf{\Lambda} \Leftrightarrow \mathbf{K} \begin{bmatrix} | & | & \dots & | \\ \mathbf{p}_1 & \mathbf{p}_2 & \dots & \mathbf{p}_N \\ | & | & \dots & | \end{bmatrix} = \begin{bmatrix} | & | & \dots & | \\ \lambda_1 \mathbf{p}_1 & \lambda_2 \mathbf{p}_2 & \dots & \lambda_N \mathbf{p}_N \\ | & | & \dots & | \end{bmatrix} \quad (39)$$

$$= \begin{bmatrix} | & | & \dots & | \\ \mathbf{P}_1 & \mathbf{P}_2 & \dots & \mathbf{P}_N \\ | & | & \dots & | \end{bmatrix} \begin{bmatrix} \lambda_1 & 0 & \dots & 0 \\ 0 & \lambda_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \lambda_N \end{bmatrix} \Leftrightarrow \mathbf{K} = \mathbf{P} \mathbf{\Lambda} \mathbf{P}^T \quad (40)$$

with $\mathbf{K}\mathbf{p}_i = \lambda_i \mathbf{p}_i$, $i \in \{1, 2, \dots, N\}$, i.e., columns of \mathbf{P} are the eigenvectors of \mathbf{K} and $\{\lambda_i\}$ the corresponding eigenvalues, \mathbf{P} is orthogonal, i.e., $\mathbf{P}\mathbf{P}^T = \mathbf{P}^T\mathbf{P} = \mathbf{I}_N$ and $\mathbf{\Lambda}$ diagonal matrix, with main diagonal the positive eigenvalues of \mathbf{K} , i.e., $\mathbf{\Lambda} = \text{diag}[\lambda_1 \lambda_2 \dots \lambda_N]$.

- ▶ Set $\mathbf{\Lambda}^{-1/2} = \text{diag} \left[\sqrt{\lambda_1} \sqrt{\lambda_2} \dots \sqrt{\lambda_N} \right]$ and multiply \mathbf{y} by $\mathbf{\Lambda}^{-1/2} \mathbf{P}^T$:

Whitening Procedure

- ▶ Multiply \mathbf{y} by $\Lambda^{-1/2} \mathbf{P}^T$:

$$\underbrace{\Lambda^{-1/2} \mathbf{P}^T \mathbf{y}}_{\mathbf{y}_w} = \underbrace{\Lambda^{-1/2} \mathbf{P}^T \mathbf{m}_j}_{\boldsymbol{\mu}_j} + \underbrace{\Lambda^{-1/2} \mathbf{P}^T \mathbf{v}}_{\mathbf{v}_w}, \quad (41)$$

$$\Leftrightarrow \mathbf{y}_w = \boldsymbol{\mu}_j + \mathbf{v}_w, \quad (42)$$

with

$$\mathbb{E}[\mathbf{v}_w] = \Lambda^{-1/2} \mathbf{P}^T \mathbb{E}[\mathbf{v}] = \mathbf{0} \quad (43)$$

$$\mathbb{E}[\mathbf{v}_w \mathbf{v}_w^T] = \Lambda^{-1/2} \mathbf{P}^T \mathbf{P} \Lambda \mathbf{P}^T \mathbf{P} \Lambda^{-1/2} = \Lambda^{-1/2} \Lambda \Lambda^{-1/2} = \mathbf{I}_N. \quad (44)$$

- ▶ Thus, $\mathbf{v}_w \sim \mathcal{N}(\mathbf{0}, \mathbf{I}_N)$ and the detection problem is simplified in Eq. (42). Notice that

$$\delta\boldsymbol{\mu} \triangleq \boldsymbol{\mu}_1 - \boldsymbol{\mu}_0 = \Lambda^{-1/2} \mathbf{P}^T (\mathbf{m}_1 - \mathbf{m}_0) = \Lambda^{-1/2} \mathbf{P}^T \delta\mathbf{m} \quad (45)$$

$$\begin{aligned} d^2 &\triangleq \delta\mathbf{m}^T \mathbf{K}^{-1} \delta\mathbf{m} = \delta\mathbf{m}^T \mathbf{P} \Lambda^{-1} \mathbf{P}^T \delta\mathbf{m} = \delta\mathbf{m}^T \mathbf{P} \Lambda^{-1/2} \Lambda^{-1/2} \mathbf{P}^T \delta\mathbf{m} \equiv \delta\boldsymbol{\mu}^T \delta\boldsymbol{\mu} \\ &= \|\delta\boldsymbol{\mu}\|_2^2. \end{aligned} \quad (46)$$

- ▶ Another example:

$$\mathbf{y} = \boldsymbol{\mu}_j + \mathbf{v}, \quad (47)$$

with $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_N)$, $j \in \{0, 1\}$ and $\pi_0 = \pi_1 = 1/2$.
It can be easily shown that:

1. Minimum probability of error detection rule is the minimum distance rule:

$$\|\mathbf{y} - \boldsymbol{\mu}_0\|_2 \stackrel{H_1}{\geq} \|\mathbf{y} - \boldsymbol{\mu}_1\|_2$$

2. The probability of error of the above rule is given by:

$$\Pr(e) = Q\left(\frac{\|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0\|_2}{2\sigma}\right).$$

Proof: simply write the ML rule and exploit the fact that affine transformation of Gaussian vectors is also Gaussian.

- [1] Bernard C. Levy, Principles of Signal Detection and Parameter Estimation, Springer 2008.
- [2] Instructor notes.

Thank you!



Detection & Estimation Theory: Lectures 7 & 8

Prof. Aggelos Bletsas

Technical University of Crete
School of Electrical and Computer Engineering



European Union
European Social Fund

Operational Programme
**Human Resources Development,
Education and Lifelong Learning**

Co-financed by Greece and the European Union





“Support for the Internationalization of Higher Education, School of Electrical and Computer Engineering, Technical University of Crete” (MIS 5150766), under the call for proposals “Support of the Internationalization of Higher Education - Technical University of Crete” (EDULLL 153).

The project is co-financed by Greece and the European Union (European Social Fund- ESF) through the Operational Programme “Human Resources Development, Education and Lifelong Learning 2014-2020.



Operational Programme
Human Resources Development,
Education and Lifelong Learning
Co-financed by Greece and the European Union



- M-ary Hypothesis Testing
- Error Probability Bounds
- Example

M-ary Hypothesis Testing

- ▶ Assume M hypotheses $\{H_j\}$, $j \in \{0, 1, \dots, M-1\}$, with priors $\Pr(H_j) \triangleq \pi_j$.
 1. Observe \mathbf{y} and find out decision rule $\delta(\mathbf{y}) = j$, i.e., decide H_j for that specific \mathbf{y} .
 2. Equivalently, find out $\mathcal{Y}_i = \{\mathbf{y} : \delta(\mathbf{y}) = i\}$, with $\mathcal{Y} = \cup_{i=0}^{M-1} \mathcal{Y}_i$ and $\mathcal{Y}_i \cap \mathcal{Y}_j = \emptyset \forall i \neq j$.
- ▶ We will revert to Bayesian formulation. Remember Bayes Risk $R(\delta)$ and conditional Bayes Risk $R(\delta|H_j)$:

$$R(\delta) = \sum_{j=0}^{M-1} R(\delta|H_j)\pi_j \quad (1)$$

$$R(\delta|H_j) = \sum_{i=0}^{M-1} C_{ij} \Pr(\delta(\mathbf{y}) = i|H_j) = \sum_{i=0}^{M-1} C_{ij} \Pr(\mathcal{Y}_i|H_j)$$

$$\Leftrightarrow R(\delta) = \sum_{i=0}^{M-1} \sum_{j=0}^{M-1} C_{ij} \Pr(\mathcal{Y}_i|H_j) \pi_j, \quad (2)$$

where C_{ij} is the cost of deciding i when H_j holds.

M-ary Hypothesis Testing

► Bayesian formulation:

$$R(\delta) = \sum_{i=0}^{M-1} \sum_{j=0}^{M-1} C_{ij} \Pr(\mathcal{Y}_i | H_j) \pi_j \quad (3)$$

$$= \sum_{i=0}^{M-1} \sum_{j=0}^{M-1} C_{ij} \int_{\mathcal{Y}_i} f_{\mathbf{y}|H_j}(\mathbf{y}|H_j) \pi_j d\mathbf{y} \quad (4)$$

$$= \sum_{i=0}^{M-1} \int_{\mathcal{Y}_i} \sum_{j=0}^{M-1} C_{ij} f_{\mathbf{y}|H_j}(\mathbf{y}|H_j) \pi_j d\mathbf{y} \quad (5)$$

$$= \sum_{i=0}^{M-1} \int_{\mathcal{Y}_i} \sum_{j=0}^{M-1} C_{ij} \Pr(H_j | \mathbf{y}) f_{\mathbf{y}}(\mathbf{y}) d\mathbf{y} \quad (6)$$

$$= \sum_{i=0}^{M-1} \int_{\mathcal{Y}_i} f_{\mathbf{y}}(\mathbf{y}) \underbrace{\sum_{j=0}^{M-1} C_{ij} \Pr(H_j | \mathbf{y})}_{C_i(\mathbf{y})} d\mathbf{y} \quad (7)$$

M-ary Hypothesis Testing

- ▶ Bayesian formulation:

$$R(\delta) = \sum_{i=0}^{M-1} \int_{\mathcal{Y}_i} f_{\mathbf{y}}(\mathbf{y}) \underbrace{\sum_{j=0}^{M-1} C_{ij} \Pr(H_j|\mathbf{y})}_{C_i(\mathbf{y})} d\mathbf{y} \quad (8)$$

$$= \sum_{i=0}^{M-1} \int_{\mathcal{Y}_i} C_i(\mathbf{y}) f_{\mathbf{y}}(\mathbf{y}) d\mathbf{y}. \quad (9)$$

- ▶ From Eq. (9),

$$\delta_B(\mathbf{y}) = \arg \min_{i \in \{0,1,\dots,M-1\}} C_i(\mathbf{y}) \quad (10)$$

i.e., we select H_k if $C_k(\mathbf{y}) \leq C_i(\mathbf{y}), \forall i \in \{0, 1, \dots, M-1\}$.

min $\Pr(e)$ rule: MAP rule

- ▶ Set symmetric costs:

$$C_{ij} = 1 - \delta_{ij} = \begin{cases} 0, & i = j, \\ 1, & i \neq j. \end{cases} \quad (11)$$

In that case, min of prob. of error is equivalent to risk minimization (as in the binary case):

$$R(\delta|H_j) = \sum_{i=0}^{M-1} C_{ij} \Pr(\mathcal{Y}_i|H_j) = \sum_{i \neq j} \Pr(\mathcal{Y}_i|H_j) \quad (12)$$

$$\Rightarrow R(\delta|H_j) = 1 - \Pr(\mathcal{Y}_j|H_j) \equiv \Pr(e|H_j) \quad (13)$$

$$R(\delta) = \sum_{j=0}^{M-1} R(\delta|H_j) \pi_j = \sum_{j=0}^{M-1} \Pr(e|H_j) \pi_j \equiv \Pr(e) \quad (14)$$

$$C_i(\mathbf{y}) = \sum_{j=0}^{M-1} C_{ij} \Pr(H_j|\mathbf{y}) = \sum_{j \neq i} \Pr(H_j|\mathbf{y}) = 1 - \Pr(H_i|\mathbf{y}) \quad (15)$$

min $\Pr(e)$ rule: MAP rule

- ▶ Set symmetric costs:

$$C_{ij} = 1 - \delta_{ij} = \begin{cases} 0, & i = j, \\ 1, & i \neq j. \end{cases} \quad (16)$$

- ▶ Min prob. of error rule:

$$C_i(\mathbf{y}) = \sum_{j=0}^{M-1} C_{ij} \Pr(H_j|\mathbf{y}) = \sum_{j \neq i} \Pr(H_j|\mathbf{y}) = 1 - \Pr(H_i|\mathbf{y})$$

$$\delta_{MAP}(\mathbf{y}) = \arg \min_{i \in \{0,1,\dots,M-1\}} \{1 - \Pr(H_i|\mathbf{y})\} \quad (17)$$

$$= \arg \max_{i \in \{0,1,\dots,M-1\}} \Pr(H_i|\mathbf{y}), \text{ (MAP rule)} \quad (18)$$

$$= \arg \max_{i \in \{0,1,\dots,M-1\}} \frac{f_{\mathbf{y}|H_i}(\mathbf{y}|H_i) \pi_i}{f_{\mathbf{y}}(\mathbf{y})} \quad (19)$$

$$= \arg \max_{i \in \{0,1,\dots,M-1\}} f_{\mathbf{y}|H_i}(\mathbf{y}|H_i) \pi_i \quad (20)$$

min $\Pr(e)$ rule and equiprobable hypotheses: ML rule

- ▶ Set symmetric costs and equiprobable hypotheses ($\pi_j = 1/M$):

$$C_{ij} = 1 - \delta_{ij} = \begin{cases} 0, & i = j, \\ 1, & i \neq j. \end{cases} \quad (21)$$

- ▶ Min prob. of error rule:

$$C_i(\mathbf{y}) = \sum_{j=0}^{M-1} C_{ij} \Pr(H_j|\mathbf{y}) = \sum_{j \neq i} \Pr(H_j|\mathbf{y}) = 1 - \Pr(H_i|\mathbf{y})$$

$$\delta_{ML}(\mathbf{y}) = \arg \max_{i \in \{0,1,\dots,M-1\}} f_{\mathbf{y}|H_i}(\mathbf{y}|H_i) \cdot (\pi_i = 1/M) \quad (22)$$

$$= \arg \max_{i \in \{0,1,\dots,M-1\}} f_{\mathbf{y}|H_i}(\mathbf{y}|H_i) \text{ (ML rule)}. \quad (23)$$

- ▶ ...generalisation of the binary hypothesis testing case!

Example

- ▶ Under hypothesis H_i , with $i \in \{0, 1, \dots, M-1\}$, $\pi_i = 1/M$ and $\mathbf{y} = [y_1 \ y_2 \ \dots \ y_N]^T$:

$$\mathbf{y} \sim \mathcal{N}(\mathbf{m}_i, \mathbf{K}). \quad (24)$$

What is the minimum probability of error detection rule?

- ▶ Minimum probability of error detection rule for equiprobable hypotheses is the ML rule:

$$\delta_{ML}(\mathbf{y}) = \arg \max_{i \in \{0, 1, \dots, M-1\}} f_{\mathbf{y}|H_i}(\mathbf{y}|H_i) \quad (25)$$

$$= \arg \max_{i \in \{0, 1, \dots, M-1\}} \ln \left[f_{\mathbf{y}|H_i}(\mathbf{y}|H_i) \right] \quad (26)$$

$$= \arg \max_{i \in \{0, 1, \dots, M-1\}} \ln \left[-\frac{1}{2} \left((\mathbf{y} - \mathbf{m}_i)^T \mathbf{K}^{-1} (\mathbf{y} - \mathbf{m}_i) \right) - \frac{N}{2} \ln(2\pi) - \frac{\ln(|\mathbf{K}|)}{2} \right] \quad (27)$$

$$= \arg \min_{i \in \{0, 1, \dots, M-1\}} \ln \left[\left((\mathbf{y} - \mathbf{m}_i)^T \mathbf{K}^{-1} (\mathbf{y} - \mathbf{m}_i) \right) \right] \quad (28)$$

$$= \arg \min_{i \in \{0, 1, \dots, M-1\}} \left((\mathbf{y} - \mathbf{m}_i)^T \mathbf{K}^{-1} (\mathbf{y} - \mathbf{m}_i) \right) \quad (29)$$

- ▶ Note that through whitening $\mathbf{z} = \mathbf{\Lambda}^{-1/2} \mathbf{P}^T \mathbf{y}$, i.e., using $\mathbf{K} = \mathbf{P} \mathbf{\Lambda} \mathbf{P}^T$, under hypothesis H_i :

$$\mathbf{z} \sim \mathcal{N} \left(\mathbf{\Lambda}^{-1/2} \mathbf{P}^T \mathbf{m}_i, \mathbf{I}_N \right). \quad (30)$$

...analytical proof in the previous lectures.

Error Probability Bounds

- ▶ Set symmetric costs:

$$C_{ij} = 1 - \delta_{ij} = \begin{cases} 0, & i = j, \\ 1, & i \neq j. \end{cases} \quad (31)$$

As we have already seen:

$$R(\delta|H_j) = \sum_{i \neq j} \Pr(\mathcal{Y}_i|H_j) = 1 - \Pr(\mathcal{Y}_j|H_j) \equiv \Pr(e|H_j) \triangleq \Pr(\mathcal{Y}_j^c|H_j) \quad (32)$$

$$R(\delta) = \sum_{j=0}^{M-1} R(\delta|H_j)\pi_j = \sum_{j=0}^{M-1} \Pr(e|H_j)\pi_j \equiv \Pr(e) \quad (33)$$

- ▶ So, \mathcal{Y}_j^c is the region where hypothesis H_j is NOT selected. Define formally the following:

$$\mathcal{Y}_j^c = \bigcup_{k \neq j} \mathcal{E}_{kj}, \quad (34)$$

$$\mathcal{E}_{kj} = \left\{ \mathbf{y} : \Pr(H_k|\mathbf{y}) > \Pr(H_j|\mathbf{y}) \right\} \quad (35)$$

$$= \left\{ \mathbf{y} : \frac{f_{\mathbf{y}|H_k}(\mathbf{y}|H_k)}{f_{\mathbf{y}|H_j}(\mathbf{y}|H_j)} > \frac{\pi_j}{\pi_k} \right\}, \quad (36)$$

i.e., \mathcal{E}_{kj} is the region of $\{\mathbf{y}\}$'s where hypothesis H_k is preferred over H_j .

Error Probability Bounds

- ▶ The areas $\{\mathcal{E}_{kj}\}$ for given j usually overlap.
- ▶ It is often possible to find a subset of such areas, such that:

$$\mathcal{Y}_j^c = \bigcup_{k \in N(j)} \mathcal{E}_{kj}, \text{ with } N(j) \subset \{0, 1, \dots, M-1\} / j \quad (37)$$

i.e., $N(j)$ does not include element j .

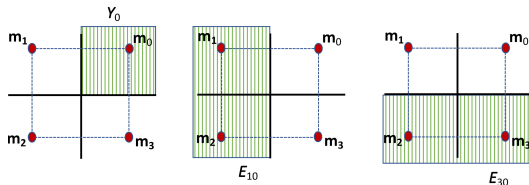
- ▶ Thus,

$$\Pr(e|H_j) \stackrel{\Delta}{=} \Pr(\mathcal{Y}_j^c|H_j) \leq \underbrace{\sum_{k \in N(j)} \Pr(\mathcal{E}_{kj}|H_j)}_{\text{Improved union bound}} \quad (38)$$

$$\max_{k \neq j} \Pr(\mathcal{E}_{kj}|H_j) \leq \Pr(\mathcal{Y}_j^c|H_j) \quad (39)$$

- ▶ The above bounds are usually simple to calculate!

Example



- Find error probability for optimal detection for the following:

$$H_j : \mathbf{y} = \mathbf{m}_j + \mathbf{v}, \quad (40)$$

with $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_2)$, $j \in \{0, 1, 2, 3\}$ and $\pi_j = 1/4$.

- MAP is simplified to ML, simplified to minimum distance. In addition:

$$\mathcal{Y}_0^c = \mathcal{E}_{10} \cup \mathcal{E}_{30} \quad (41)$$

$$d \triangleq \|\mathbf{m}_1 - \mathbf{m}_0\|_2 = \|\mathbf{m}_2 - \mathbf{m}_1\|_2 = \|\mathbf{m}_3 - \mathbf{m}_2\|_2 = \|\mathbf{m}_3 - \mathbf{m}_0\|_2 \quad (42)$$

$$\Pr(e|H_0) = \Pr(\mathcal{Y}_0^c|H_0) = \Pr(\mathcal{E}_{10} \cup \mathcal{E}_{30}|H_0) \leq \underbrace{\Pr(\mathcal{E}_{10}|H_0)}_{Q\left(\frac{d}{2\sigma}\right)} + \underbrace{\Pr(\mathcal{E}_{30}|H_0)}_{Q\left(\frac{d}{2\sigma}\right)} \quad (43)$$

$$\Leftrightarrow \Pr(e|H_0) \leq 2Q\left(\frac{d}{2\sigma}\right), \quad (44)$$

where the latter is due to minimum distance binary error detection in white Gaussian noise (reminder in next slide).

► Reminder:

$$H_j : \mathbf{y} = \boldsymbol{\mu}_j + \mathbf{v}, \quad (45)$$

with $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_N)$, $j \in \{0, 1\}$ and $\pi_0 = \pi_1 = 1/2$.

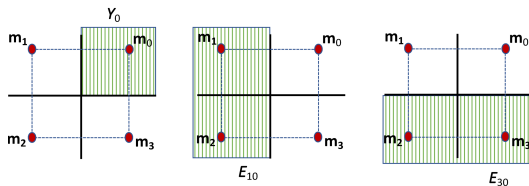
1. Minimum probability of error detection rule is the minimum distance rule:

$$\|\mathbf{y} - \boldsymbol{\mu}_0\|_2 \stackrel{H_1}{\geq} \|\mathbf{y} - \boldsymbol{\mu}_1\|_2$$

2. The probability of error of the above rule is given by:

$$\Pr(e) = Q\left(\frac{\|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0\|_2}{2\sigma}\right).$$

Example



► Error analysis:

$$\mathcal{Y}_0^c = \mathcal{E}_{10} \cup \mathcal{E}_{30} \quad (46)$$

$$d \triangleq \|\mathbf{m}_1 - \mathbf{m}_0\|_2 = \|\mathbf{m}_2 - \mathbf{m}_1\|_2 = \|\mathbf{m}_3 - \mathbf{m}_2\|_2 = \|\mathbf{m}_3 - \mathbf{m}_0\|_2 \quad (47)$$

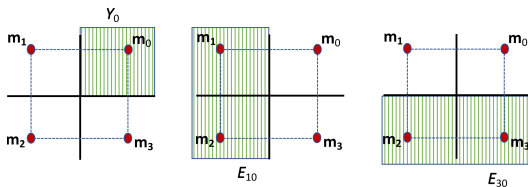
$$\Pr(e|\mathbf{H}_0) = \Pr(\mathcal{Y}_0^c|\mathbf{H}_0) = \Pr(\mathcal{E}_{10} \cup \mathcal{E}_{30}|\mathbf{H}_0) \leq \underbrace{\Pr(\mathcal{E}_{10}|\mathbf{H}_0)}_{Q\left(\frac{d}{2\sigma}\right)} + \underbrace{\Pr(\mathcal{E}_{30}|\mathbf{H}_0)}_{Q\left(\frac{d}{2\sigma}\right)} = 2Q\left(\frac{d}{2\sigma}\right)$$

$$Q\left(\frac{d}{2\sigma}\right) = \Pr(\mathcal{E}_{10}|\mathbf{H}_0) = \Pr(\mathcal{E}_{30}|\mathbf{H}_0) \leq \Pr(e|\mathbf{H}_0) \quad (48)$$

$$\Pr(e|\mathbf{H}_0) = \Pr(e), \text{ due to symmetry} \quad (49)$$

$$\Rightarrow Q\left(\frac{d}{2\sigma}\right) \leq \Pr(e) \leq 2Q\left(\frac{d}{2\sigma}\right), \quad (50)$$

Example



- ▶ Error analysis with improved union bound:

$$d \triangleq \|m_1 - m_0\|_2 = \|m_2 - m_1\|_2 = \|m_3 - m_2\|_2 = \|m_3 - m_0\|_2 \quad (51)$$

$$Q\left(\frac{d}{2\sigma}\right) \leq \Pr(e) \leq 2Q\left(\frac{d}{2\sigma}\right) \quad (52)$$

- ▶ Exact error analysis:

$$\Pr(e) = 1 - \left(1 - Q\left(\frac{d}{2\sigma}\right)\right)^2 \quad (53)$$

$$= 2Q\left(\frac{d}{2\sigma}\right) - \left[Q\left(\frac{d}{2\sigma}\right)\right]^2 \quad (54)$$

Thus, upper error probability bound is tight!

- [1] Bernard C. Levy, Principles of Signal Detection and Parameter Estimation, Springer 2008.
- [2] Instructor notes.

Thank you!



Detection & Estimation Theory: Lectures 9-10

Prof. Aggelos Bletsas

Technical University of Crete
School of Electrical and Computer Engineering



Operational Programme
**Human Resources Development,
Education and Lifelong Learning**
Co-financed by Greece and the European Union





“Support for the Internationalization of Higher Education, School of Electrical and Computer Engineering, Technical University of Crete” (MIS 5150766), under the call for proposals “Support of the Internationalization of Higher Education - Technical University of Crete” (EDULLL 153).

The project is co-financed by Greece and the European Union (European Social Fund- ESF) through the Operational Programme “Human Resources Development, Education and Lifelong Learning 2014-2020.



Operational Programme
**Human Resources Development,
Education and Lifelong Learning**
Co-financed by Greece and the European Union



- Bayesian Estimation: Problem Definition
 - Details on Problem Formulation

- Optimum Bayesian Estimator
 - Bayesian MSE Estimator
 - MSE Performance Evaluation
 - Bayesian MAE Estimator
 - Bayesian MAP Estimator

Parameter Estimation Theory: Bayesian Formulation

- ▶ Formulation for Bayesian estimation is similar to Bayesian detection!
 1. Observe $\mathbf{y} \in \mathbb{R}^n$ for estimation of parameter vector $\mathbf{x} \in \mathbb{R}^m$.
 2. Detection: find out decision rule $\delta(\mathbf{y}) = j$, where j is discrete.
 3. Estimation: find out estimate $\hat{\mathbf{x}}(\mathbf{y}) \in \mathbb{R}^m$.
- ▶ Formulation for Bayesian estimation requires the following:
 1. Observation model: $f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})$, i.e., conditional p.d.f. from measurements!
 2. Prior density: $f_{\mathbf{x}}(\mathbf{x})$, i.e., prior p.d.f. density of the unknown parameter. Notice that in the Bayesian formulation the unknown parameter is assumed random!
 3. Cost function: $C(\hat{\mathbf{x}}(\mathbf{y}), \mathbf{x}) \equiv C(\hat{\mathbf{x}}, \mathbf{x})$, i.e., the cost of estimating \mathbf{x} as $\hat{\mathbf{x}}(\mathbf{y})$.
- ▶ Derivations in Bayesian estimation proceed alongside similar lines to Bayesian detection!
...some details on the problem formulation follow...

Details on problem formulation

- ▶ Observation model $f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})$ is defined explicitly or indirectly through a measurement model.
- ▶ One example: $\mathbf{y} = \mathbf{h}(\mathbf{x}) + \mathbf{v}$, with $f_{\mathbf{v}}(\mathbf{v})$ known and $\mathbf{h}(\mathbf{x})$ is a deterministic vector function of \mathbf{x} ; assume $n = m$.
 $\Rightarrow f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) \stackrel{n=m}{=} f_{\mathbf{v}}(\mathbf{y} - \mathbf{h}(\mathbf{x}))$.

Proof:

$$\mathbf{J}_{n \times m}(\mathbf{x}, \mathbf{y}) = \begin{bmatrix} \frac{\partial y_1}{\partial x_1} & \frac{\partial y_1}{\partial x_2} & \cdots & \frac{\partial y_1}{\partial x_m} \\ \frac{\partial y_2}{\partial x_1} & \frac{\partial y_2}{\partial x_2} & \cdots & \frac{\partial y_2}{\partial x_m} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y_n}{\partial x_1} & \frac{\partial y_n}{\partial x_2} & \cdots & \frac{\partial y_n}{\partial x_m} \end{bmatrix} = \begin{bmatrix} \nabla y_1 \\ \nabla y_2 \\ \vdots \\ \nabla y_n \end{bmatrix}, \quad (1)$$

$$\text{where } \nabla f = \left[\frac{\partial f}{\partial x_1} \quad \frac{\partial f}{\partial x_2} \quad \cdots \quad \frac{\partial f}{\partial x_m} \right], \quad (2)$$

$$\mathbf{y} = \mathbf{h}(\mathbf{x}) + \mathbf{v} \Rightarrow$$

$$f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) \stackrel{n=m}{=} \frac{f_{\mathbf{v}}(\mathbf{v})}{\det(\mathbf{J}(\mathbf{x}, \mathbf{y}))} \Big|_{\mathbf{v}=\mathbf{y}-\mathbf{h}(\mathbf{x})} = \frac{f_{\mathbf{v}}(\mathbf{y} - \mathbf{h}(\mathbf{x}))}{\det(\mathbf{I}_n)} = f_{\mathbf{v}}(\mathbf{y} - \mathbf{h}(\mathbf{x})) \quad (3)$$

Details on problem formulation

- ▶ Prior density $f_{\mathbf{x}}(\mathbf{x})$ is known.
- ▶ ...unfortunately, prior density biases the estimator towards more probable values of \mathbf{x} , i.e., values of \mathbf{x} , where $f_{\mathbf{x}}(\mathbf{x})$ is larger.
- ▶ ...remember that we don't know anything about $f_{\mathbf{x}}(\mathbf{x})$ (i.e., it is random), apart from possible values.

Details on problem formulation

- ▶ Cost function: $C(\hat{\mathbf{x}}(\mathbf{y}), \mathbf{x}) \equiv C(\hat{\mathbf{x}}, \mathbf{x})$, i.e., the cost of estimating \mathbf{x} as $\hat{\mathbf{x}}(\mathbf{y})$.

$$C(\hat{\mathbf{x}}, \mathbf{x}) = C(\hat{\mathbf{x}} - \mathbf{x}) \equiv L(\mathbf{e}) \quad (4)$$

- ▶ Loss function $L(\mathbf{e})$ is a not decreasing function of error $\mathbf{e} \triangleq \hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x}$.
- ▶ 3 different versions of loss function are typically used:
 1. $L_{\text{MSE}}(\mathbf{e}) = \|\mathbf{e}\|_2^2$
 2. $L_{\text{MAE}}(\mathbf{e}) = \|\mathbf{e}\|_1$
 3. $L_\epsilon(\mathbf{e})$: notch function.

Details on problem formulation

► 3 different versions of loss function are typically used:

1. Euclidean norm or 2-norm squared:

$$L_{\text{MSE}}(\mathbf{e}) = \|\mathbf{e}\|_2^2 = \mathbf{e}^T \mathbf{e} = \sum_{i=1}^m e_i^2. \quad (5)$$

Due to the square, it penalises more larger errors; improbable instances of \mathbf{x} matter a lot; sensitive to modelling errors.

2. Sum norm or 1-norm:

$$L_{\text{MAE}}(\mathbf{e}) = \|\mathbf{e}\|_1 = \sum_{i=1}^m |e_i|. \quad (6)$$

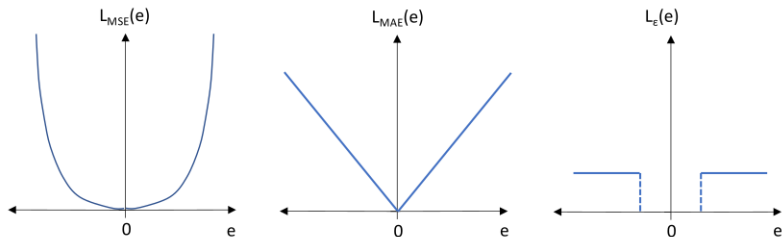
It weights equally the magnitude of all errors.

3. using the infinity norm:

$$L_{\epsilon}(\mathbf{e}) = \begin{cases} 0, & \text{if } \|\mathbf{e}\|_{\infty} < \epsilon \\ 1, & \text{otherwise,} \end{cases} \quad (7)$$

where infinity norm is given by $\|\mathbf{e}\|_{\infty} = \max_i |e_i|$; notch function cares about small errors and not about errors above ϵ .

Details on problem formulation: error (cost) functions



- ▶ The first two are convex, while the last one is non-convex (for scalar error).

Optimum Bayesian Estimator

- ▶ Cost is a function of random vectors, since both $\hat{\mathbf{x}}(\mathbf{y})$, \mathbf{x} are random.
- ▶ Bayesian objective: find $\hat{\mathbf{x}}(\mathbf{y}) = [x_1(\mathbf{y}) \ x_2(\mathbf{y}) \ \dots \ x_m(\mathbf{y})]^T$ by minimising the expected cost:

$\min \mathbb{E} [C(\hat{\mathbf{x}}(\mathbf{y}), \mathbf{x})]$, where

$$\mathbb{E} [C(\hat{\mathbf{x}}(\mathbf{y}), \mathbf{x})] = \int_{\mathbf{x}} \int_{\mathbf{y}} C(\hat{\mathbf{x}}(\mathbf{y}), \mathbf{x}) f_{\mathbf{x},\mathbf{y}}(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} \quad (8)$$

$$= \int_{\mathbf{y}} \left[\int_{\mathbf{x}} C(\hat{\mathbf{x}}(\mathbf{y}), \mathbf{x}) f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) d\mathbf{x} \right] f_{\mathbf{y}}(\mathbf{y}) d\mathbf{y}. \quad (9)$$

- ▶ Notice that posterior p.d.f. $f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y})$ and measurements p.d.f. $f_{\mathbf{y}}(\mathbf{y})$ can be (at least in principle) known:

$$f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) = \frac{f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) f_{\mathbf{x}}(\mathbf{x})}{f_{\mathbf{y}}(\mathbf{y})}, \quad f_{\mathbf{y}}(\mathbf{y}) = \int_{\mathbf{x}} f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) f_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} \quad (10)$$

Optimum Bayesian Estimator

$\min \mathbb{E} [C(\hat{\mathbf{x}}(\mathbf{y}), \mathbf{x})]$, where

$$\mathbb{E} [C(\hat{\mathbf{x}}(\mathbf{y}), \mathbf{x})] = \int_{\mathbf{y}} \left[\int_{\mathbf{x}} C(\hat{\mathbf{x}}(\mathbf{y}), \mathbf{x}) f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) d\mathbf{x} \right] f_{\mathbf{y}}(\mathbf{y}) d\mathbf{y}. \quad (11)$$

- ▶ Notice that since measurements p.d.f. $f_{\mathbf{y}}(\mathbf{y})$ is non-negative for each given \mathbf{y} , the term between brackets above is minimised for each given \mathbf{y} according to the following:

$$\hat{\mathbf{x}}(\mathbf{y}) = \arg \min_{\mathbf{x}} \int_{\mathbf{x}} C(\hat{\mathbf{x}}(\mathbf{y}), \mathbf{x}) f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) d\mathbf{x} \quad (12)$$

$$= \arg \min_{\mathbf{x}} \frac{1}{f_{\mathbf{y}}(\mathbf{y})} \int_{\mathbf{x}} C(\hat{\mathbf{x}}(\mathbf{y}), \mathbf{x}) f_{\mathbf{x},\mathbf{y}}(\mathbf{x}, \mathbf{y}) d\mathbf{x} \quad (13)$$

$$= \arg \min_{\mathbf{x}} \int_{\mathbf{x}} C(\hat{\mathbf{x}}(\mathbf{y}), \mathbf{x}) f_{\mathbf{x},\mathbf{y}}(\mathbf{x}, \mathbf{y}) d\mathbf{x} \quad (14)$$

$$= \arg \min_{\mathbf{x}} \int_{\mathbf{x}} C(\hat{\mathbf{x}}(\mathbf{y}), \mathbf{x}) f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) f_{\mathbf{x}}(\mathbf{x}) d\mathbf{x}. \quad (15)$$

Optimum Bayesian Estimator: MSE case

- ▶ For minimum square error (MSE) loss function:

$$C(\hat{\mathbf{x}}(\mathbf{y}), \mathbf{x}) = \underbrace{(\hat{\mathbf{x}} - \mathbf{x})}_{\mathbf{e}}^T (\hat{\mathbf{x}} - \mathbf{x}) = L_{\text{MSE}}(\mathbf{e}) = \|\mathbf{e}\|_2^2 \quad (16)$$

$$= \hat{\mathbf{x}}^T \hat{\mathbf{x}} - \hat{\mathbf{x}}^T \mathbf{x} - \mathbf{x}^T \hat{\mathbf{x}} + \mathbf{x}^T \mathbf{x} \quad (17)$$

- ▶ Since $\mathbf{a}^T \mathbf{b} = \mathbf{b}^T \mathbf{a}$, $\frac{\partial \mathbf{b}^T \mathbf{x}}{\partial \mathbf{x}} = \mathbf{b}$ and $\frac{\partial \mathbf{x}^T \mathbf{A} \mathbf{x}}{\partial \mathbf{x}} = (\mathbf{A} + \mathbf{A}^T) \mathbf{x}$,

$$\frac{\partial C(\hat{\mathbf{x}}(\mathbf{y}), \mathbf{x})}{\partial \hat{\mathbf{x}}} = 2(\mathbf{I} + \mathbf{I}^T) \hat{\mathbf{x}} - \mathbf{x} - \mathbf{x} = 2(\hat{\mathbf{x}} - \mathbf{x}) \quad (18)$$

- ▶ Denoting $\nabla_{\hat{\mathbf{x}}} = \left[\frac{\partial}{\partial \hat{x}_1} \quad \frac{\partial}{\partial \hat{x}_2} \quad \dots \quad \frac{\partial}{\partial \hat{x}_m} \right]^T = \frac{\partial}{\partial \hat{\mathbf{x}}}$, and the following (scalar) function of $\hat{\mathbf{x}}$ (from Eq. (12)),

$$J(\hat{\mathbf{x}}|\mathbf{y}) = \int_{\mathbf{x}} C(\hat{\mathbf{x}}(\mathbf{y}), \mathbf{x}) f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) d\mathbf{x} \quad (19)$$

$$\Rightarrow \nabla_{\hat{\mathbf{x}}} J(\hat{\mathbf{x}}|\mathbf{y}) = \int_{\mathbf{x}} \nabla_{\hat{\mathbf{x}}} C(\hat{\mathbf{x}}(\mathbf{y}), \mathbf{x}) f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) d\mathbf{x} = \mathbf{0} \quad (20)$$

$$= \int_{\mathbf{x}} 2(\hat{\mathbf{x}} - \mathbf{x}) f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) d\mathbf{x} = \mathbf{0} \quad (21)$$

- ▶ ...continued from previous page...

$$\int_{\mathbf{x}} 2(\hat{\mathbf{x}} - \mathbf{x}) f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) d\mathbf{x} = \mathbf{0} \quad (22)$$

$$\Leftrightarrow \int_{\mathbf{x}} \hat{\mathbf{x}} f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) d\mathbf{x} = \int_{\mathbf{x}} \mathbf{x} f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) d\mathbf{x} \quad (23)$$

$$\Leftrightarrow \hat{\mathbf{x}} \int_{\mathbf{x}} f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) d\mathbf{x} = \int_{\mathbf{x}} \mathbf{x} f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) d\mathbf{x} \triangleq \mathbb{E}[\mathbf{x}|\mathbf{y}] \quad (24)$$

$$\Leftrightarrow \hat{\mathbf{x}}(\mathbf{y})_{\text{MSE}} = \int_{\mathbf{x}} \mathbf{x} f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) d\mathbf{x} \triangleq \mathbb{E}[\mathbf{x}|\mathbf{y}]. \quad (25)$$

- ▶ Bayesian MSE estimator is the conditional (on \mathbf{y}) mean!

MSE Performance Evaluation

- ▶ First, conditional mean square error is calculated:

$$J(\hat{\mathbf{x}}(\mathbf{y})|\mathbf{y}) \stackrel{\hat{\mathbf{x}}(\mathbf{y}) \equiv \hat{\mathbf{x}}}{=} \int_{\mathbf{x}} (\hat{\mathbf{x}} - \mathbf{x})^T (\hat{\mathbf{x}} - \mathbf{x}) f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) d\mathbf{x} \quad (26)$$

$$\stackrel{\Delta}{=} \mathbb{E}_{\mathbf{x}|Y=\mathbf{y}} [(\hat{\mathbf{x}} - \mathbf{x})^T (\hat{\mathbf{x}} - \mathbf{x}) | Y = \mathbf{y}] \quad (27)$$

$$\stackrel{(*)}{=} \mathbb{E}_{\mathbf{x}|Y=\mathbf{y}} [\text{Trace} \{(\hat{\mathbf{x}} - \mathbf{x})^T (\hat{\mathbf{x}} - \mathbf{x})\} | Y = \mathbf{y}] \quad (28)$$

$$\stackrel{(**)}{=} \text{Trace} \{ \mathbb{E}_{\mathbf{x}|Y=\mathbf{y}} [(\hat{\mathbf{x}} - \mathbf{x})(\hat{\mathbf{x}} - \mathbf{x})^T | Y = \mathbf{y}] \} \quad (29)$$

$$= \text{Trace} \{ \mathbf{K}_{\mathbf{x}|Y=\mathbf{y}} \}, \quad (30)$$

where $\text{Trace}(\mathbf{A}\mathbf{B}) = \text{Trace}(\mathbf{B}\mathbf{A})$ property was used in (*) and $\mathbb{E}[\text{Trace}\{\cdot\}] = \text{Trace}\{\mathbb{E}[\cdot]\}$ property in (**), and

$$\mathbf{K}_{\mathbf{x}|Y=\mathbf{y}} \equiv \mathbf{K}_{\mathbf{x}|Y=\mathbf{y}}(\mathbf{y}) \stackrel{\Delta}{=} \mathbb{E}_{\mathbf{x}|Y=\mathbf{y}} [(\hat{\mathbf{x}} - \mathbf{x})(\hat{\mathbf{x}} - \mathbf{x})^T | Y = \mathbf{y}] \quad (31)$$

$$= \int_{\mathbf{x}} (\mathbb{E}[\mathbf{x}|\mathbf{y}] - \mathbf{x})(\mathbb{E}[\mathbf{x}|\mathbf{y}] - \mathbf{x})^T f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) d\mathbf{x} \quad (32)$$

$$= \int_{\mathbf{x}} (\mathbf{x} - \mathbb{E}[\mathbf{x}|\mathbf{y}])(\mathbf{x} - \mathbb{E}[\mathbf{x}|\mathbf{y}])^T f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) d\mathbf{x}. \quad (33)$$

MSE Performance Evaluation

- ▶ Then, (unconditional) minimum mean square error (MMSE) is calculated:

$$\text{MMSE} \triangleq \mathbb{E}_{\mathbf{x}, \mathbf{y}} [(\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x})^T (\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x})] \quad (34)$$

$$= \mathbb{E}_{\mathbf{x}, \mathbf{y}} [\text{Trace} \{(\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x})^T (\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x})\}] \quad (35)$$

$$= \text{Trace} \{ \mathbb{E}_{\mathbf{x}, \mathbf{y}} [(\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x})(\hat{\mathbf{x}}(\mathbf{y}) - \mathbf{x})^T] \} \quad (36)$$

$$= \text{Trace} \{ \mathbf{K}_E \} = \text{Trace} \{ \mathbf{K}_{X/Y} \}, \quad (37)$$

where $\text{Trace}(\mathbf{A}\mathbf{B}) = \text{Trace}(\mathbf{B}\mathbf{A})$ and $\mathbb{E}[\text{Trace}\{\cdot\}] = \text{Trace}\{\mathbb{E}[\cdot]\}$ properties were again used and \mathbf{K}_E follows:

$$\mathbf{K}_E \triangleq \int_{\mathbf{y}} \int_{\mathbf{x}} (\mathbf{x} - \mathbb{E}[\mathbf{x}|\mathbf{y}])(\mathbf{x} - \mathbb{E}[\mathbf{x}|\mathbf{y}])^T f_{\mathbf{x}, \mathbf{y}}(\mathbf{x}, \mathbf{y}) \, d\mathbf{x} \, d\mathbf{y} \quad (38)$$

$$= \int_{\mathbf{y}} \int_{\mathbf{x}} (\mathbf{x} - \mathbb{E}[\mathbf{x}|\mathbf{y}])(\mathbf{x} - \mathbb{E}[\mathbf{x}|\mathbf{y}])^T f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) f_{\mathbf{y}}(\mathbf{y}) \, d\mathbf{x} \, d\mathbf{y}$$

$$= \int_{\mathbf{y}} \left[\int_{\mathbf{x}} (\mathbf{x} - \mathbb{E}[\mathbf{x}|\mathbf{y}])(\mathbf{x} - \mathbb{E}[\mathbf{x}|\mathbf{y}])^T f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) \, d\mathbf{x} \right] f_{\mathbf{y}}(\mathbf{y}) \, d\mathbf{y}$$

$$\stackrel{\text{Eq. (33)}}{=} \int_{\mathbf{y}} \mathbf{K}_{\mathbf{x}|Y=\mathbf{y}} f_{\mathbf{y}}(\mathbf{y}) \, d\mathbf{y}. \quad (39)$$

Optimum Bayesian Estimator: MAE case

- ▶ For minimum absolute error (MAE) loss function:

$$C(\hat{\mathbf{x}}(\mathbf{y}), \mathbf{x}) = \|\hat{\mathbf{x}} - \mathbf{x}\|_1 = |\hat{x}_1 - x_1| + |\hat{x}_2 - x_2| + \dots + |\hat{x}_m - x_m|$$
$$J(\hat{\mathbf{x}}|\mathbf{y}) = \int_{\mathbf{x}} \|\hat{\mathbf{x}} - \mathbf{x}\|_1 f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) d\mathbf{x} \quad (40)$$

- ▶ Denoting $\text{sign}(z) = +1$ if $z \geq 0$ and $\text{sign}(z) = -1$, otherwise, the following are calculated:

$$\frac{\partial J(\hat{\mathbf{x}}|\mathbf{y})}{\partial \hat{x}_i} = \int_{\mathbf{x}} \text{sign}(\hat{x}_i - x_i) f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) d\mathbf{x} \quad (41)$$

$$= \int_{x_i} \text{sign}(\hat{x}_i - x_i) \times$$
$$\left[\int_{x_1} \int_{x_2} \dots \int_{x_{i-1}} \int_{x_{i+1}} \dots \int_{x_m} f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) dx_1 dx_2 \dots dx_{i-1} dx_{i+1} \dots dx_m \right] dx_i$$
$$= \int_{x_i} \text{sign}(\hat{x}_i - x_i) f_{x_i|\mathbf{y}}(x_i|\mathbf{y}) dx_i. \quad (42)$$

Optimum Bayesian Estimator: MAE case

- ▶ Continuing from previous slide,

$$\frac{\partial J(\hat{\mathbf{x}}|\mathbf{y})}{\partial \hat{x}_i} = 0 \Leftrightarrow \int_{x_i} \text{sign}(\hat{x}_i - x_i) f_{x_i|\mathbf{y}}(x_i|\mathbf{y}) dx_i = 0 \Leftrightarrow \quad (43)$$

$$\int_{-\infty}^{\hat{x}_i} \text{sign}(\hat{x}_i - x_i) f_{x_i|\mathbf{y}}(x_i|\mathbf{y}) dx_i + \int_{\hat{x}_i}^{+\infty} \text{sign}(\hat{x}_i - x_i) f_{x_i|\mathbf{y}}(x_i|\mathbf{y}) dx_i = 0 \Leftrightarrow \quad (44)$$

$$\int_{-\infty}^{\hat{x}_i} f_{x_i|\mathbf{y}}(x_i|\mathbf{y}) dx_i - \int_{\hat{x}_i}^{+\infty} f_{x_i|\mathbf{y}}(x_i|\mathbf{y}) dx_i = 0 \Leftrightarrow \quad (45)$$

$$I_1 \triangleq \int_{-\infty}^{\hat{x}_i} f_{x_i|\mathbf{y}}(x_i|\mathbf{y}) dx_i = \int_{\hat{x}_i}^{+\infty} f_{x_i|\mathbf{y}}(x_i|\mathbf{y}) dx_i \triangleq I_2 \quad (46)$$

- ▶ Since $I_1 + I_2 = 1$ and from above $I_1 = I_2$, the Bayesian MAE estimate is the median of the posterior density:

$$\int_{-\infty}^{\hat{x}_i^{\text{MAE}}} f_{x_i|\mathbf{y}}(x_i|\mathbf{y}) dx_i = \int_{\hat{x}_i^{\text{MAE}}}^{+\infty} f_{x_i|\mathbf{y}}(x_i|\mathbf{y}) dx_i = 1/2. \quad (47)$$

- ▶ Thus, the i -th entry of $\hat{\mathbf{x}}_{\text{MAE}}(\mathbf{y})$ is the median of the posterior density $f_{x_i|\mathbf{y}}(x_i|\mathbf{y})$:

$$\int_{-\infty}^{\hat{x}_i^{\text{MAE}}} f_{x_i|\mathbf{y}}(x_i|\mathbf{y}) dx_i = \int_{\hat{x}_i^{\text{MAE}}^{+\infty}} f_{x_i|\mathbf{y}}(x_i|\mathbf{y}) dx_i = 1/2. \quad (48)$$

Optimum Bayesian Estimator: Notch Cost Function

- For notch cost function $L_\epsilon(\mathbf{e} = \hat{\mathbf{x}} - \mathbf{x}) = \begin{cases} 0, & \text{if } \|\mathbf{e}\|_\infty < \epsilon \\ 1, & \text{otherwise,} \end{cases}$,
where $\|\mathbf{e}\|_\infty = \max_i |e_i|$, assuming $\mathbf{e} = [e_1 \ e_2 \ \dots \ e_m]^T$. Thus,

$$\begin{aligned} J(\hat{\mathbf{x}}|\mathbf{y}) &= \int_{\mathbf{x}} L_\epsilon(\hat{\mathbf{x}} - \mathbf{x}) f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) d\mathbf{x} = \int_{\|\hat{\mathbf{x}}-\mathbf{x}\|_\infty \geq \epsilon} f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) d\mathbf{x} \\ &= 1 - \int_{\|\hat{\mathbf{x}}-\mathbf{x}\|_\infty < \epsilon} f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) d\mathbf{x} \end{aligned} \quad (49)$$

- For $\epsilon \rightarrow 0^+$, minimization of $J(\cdot)$ above is equivalent to maximizing the following:

$$\int_{\|\hat{\mathbf{x}}-\mathbf{x}\|_\infty < \epsilon} f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) d\mathbf{x} \approx (2\epsilon)^m f_{\mathbf{x}|\mathbf{y}}(\hat{\mathbf{x}}|\mathbf{y}) \quad (50)$$

- In other words, for small ϵ ,

$$\hat{\mathbf{x}}(\mathbf{y}) \equiv \hat{\mathbf{x}}_{\text{MAP}}(\mathbf{y}) = \arg \max_{\mathbf{x}} f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}). \quad (51)$$

Optimum Bayesian Estimator: Notch Cost Function

- ▶ For notch cost function $L_\epsilon(\mathbf{e} = \hat{\mathbf{x}} - \mathbf{x})$ and small ϵ ,

$$\hat{\mathbf{x}}(\mathbf{y}) \equiv \hat{\mathbf{x}}_{\text{MAP}}(\mathbf{y}) = \arg \max_{\mathbf{x}} f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) \quad (52)$$

- ▶ ...the Bayesian estimate becomes the mode (i.e., maximum) of the posterior density (MAP estimate).
- ▶ It makes sense if the posterior density $f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y})$ has a single dominant peak, or multiple peaks of the same size.
- ▶ Example: for jointly Gaussian \mathbf{x}, \mathbf{y} , conditional mean, median and mode coincide, i.e.,

$$\hat{\mathbf{x}}_{\text{MSE}}(\mathbf{y}) = \hat{\mathbf{x}}_{\text{MAE}}(\mathbf{y}) = \hat{\mathbf{x}}_{\text{MAP}}(\mathbf{y}) = \mathbb{E}[\mathbf{x}|\mathbf{y}].$$

- [1] Bernard C. Levy, Principles of Signal Detection and Parameter Estimation, Springer 2008.
- [2] Instructor notes.

Thank you!



Detection & Estimation Theory: Lectures 11-12

Prof. Aggelos Bletsas

Technical University of Crete
School of Electrical and Computer Engineering



Operational Programme
**Human Resources Development,
Education and Lifelong Learning**
Co-financed by Greece and the European Union





“Support for the Internationalization of Higher Education, School of Electrical and Computer Engineering, Technical University of Crete” (MIS 5150766), under the call for proposals “Support of the Internationalization of Higher Education - Technical University of Crete” (EDULLL 153).

The project is co-financed by Greece and the European Union (European Social Fund- ESF) through the Operational Programme “Human Resources Development, Education and Lifelong Learning 2014-2020.



Operational Programme
**Human Resources Development,
Education and Lifelong Learning**
Co-financed by Greece and the European Union



- Examples of Bayesian Estimation
- Properties of Bayesian MSE Estimator
- MMSE Estimation in Linear Gaussian Systems

Bayesian Estimation Example 1

- ▶ The following (exponential) p.d.f.s are given:
 $f_{y|x}(y|x) = xe^{-xy}u(y)$, $f(x) = ae^{-ax}u(x)$
- ▶ Need to compute $f_y(y)$ to compute $f_{x|y}(x|y)$:

$$\begin{aligned}f(x, y) &= f(y|x)f(x) = axe^{-(a+y)x}u(x)u(y) \Rightarrow \\f(y) &= \int_0^{+\infty} axe^{-(a+y)x}u(y)dx = \frac{a}{(y+a)^2}u(y) \\ \Rightarrow f(x|y) &= \frac{f_{x,y}(x, y)}{f(y)} = (y+a)^2 xe^{-(y+a)x}u(x)\end{aligned}$$

- ▶ MSE estimator:

$$\hat{x}_{\text{MSE}}(y) = \int_{-\infty}^{+\infty} xf(x|y)dx = \dots = \frac{2}{y+a}$$

Bayesian Estimation Example 1

- MAE estimator:

$$\begin{aligned}\frac{1}{2} &= \int_{-\infty}^{\hat{x}} f(x|y) dx = (y+a)^2 \int_0^{\hat{x}} x e^{-(y+a)x} dx = \dots = \\ &= [1 + (a+y)\hat{x}] e^{-(a+y)\hat{x}} \stackrel{c=(a+y)\hat{x}}{=} (1+c)e^{-c} \Rightarrow \\ e^c &= 2(1+c) \Rightarrow c = \ln[2(1+c)] \Rightarrow c \approx 1.68 \Rightarrow \\ \hat{x}_{\text{MAE}}(y) &= \frac{c}{a+y} = \frac{1.68}{y+a}\end{aligned}$$

- MAP estimator:

$$\begin{aligned}\frac{\partial f_{x|y}(x|y)}{\partial x} &= (y+a)^2 e^{-(y+a)x} + (y+a)^2 x e^{-(y+a)x} (-(y+a)) = 0 \Rightarrow \\ (y+a)^2 e^{-(y+a)x} (1 - (y+a)x) &= 0 \Rightarrow \\ \hat{x}_{\text{MAP}} &= \frac{1}{y+a}\end{aligned}$$

Bayesian Estimation Example 2

- $\mathbf{x} \in \mathbb{R}^m, \mathbf{y} \in \mathbb{R}^n$ jointly Gaussian $\Leftrightarrow \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N}(\mathbf{m}, \mathbf{K})$

$$\begin{aligned} \mathbf{m} \triangleq \begin{bmatrix} \mathbb{E}[\mathbf{x}] \\ \mathbb{E}[\mathbf{y}] \end{bmatrix} &= \begin{bmatrix} \mathbf{m}_x \\ \mathbf{m}_y \end{bmatrix}, \mathbf{K} \triangleq \mathbb{E} \left\{ \begin{bmatrix} \mathbf{x} - \mathbf{m}_x \\ \mathbf{y} - \mathbf{m}_y \end{bmatrix} \begin{bmatrix} (\mathbf{x} - \mathbf{m}_x)^T & (\mathbf{y} - \mathbf{m}_y)^T \end{bmatrix} \right\} \\ &= \begin{bmatrix} \mathbf{K}_X & \mathbf{K}_{XY} \\ \mathbf{K}_{YX} & \mathbf{K}_Y \end{bmatrix}, \mathbf{K}_{XY} = \mathbf{K}_{YX}^T \end{aligned}$$

- $\mathbf{x} \in \mathbb{R}^m, \mathbf{y} \in \mathbb{R}^n$ jointly Gaussian \Rightarrow

$$f_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{m}_{X|Y}, \mathbf{K}_{X|Y})$$

$$\mathbf{m}_{X|Y} = \mathbf{m}_x + \mathbf{K}_{XY} \mathbf{K}_Y^{-1} (\mathbf{y} - \mathbf{m}_y)$$

$$\mathbf{K}_{X|Y} = \mathbf{K}_X - \mathbf{K}_{XY} \mathbf{K}_Y^{-1} \mathbf{K}_{YX}$$

$$\hat{\mathbf{x}}_{\text{MSE}}(\mathbf{y}) = \mathbb{E}[\mathbf{x}|\mathbf{y}] \equiv \mathbf{m}_{X|Y} = \mathbf{m}_x + \mathbf{K}_{XY} \mathbf{K}_Y^{-1} (\mathbf{y} - \mathbf{m}_y)$$

$$\text{MMSE} \triangleq \mathbb{E}[\|\mathbf{x} - \hat{\mathbf{x}}_{\text{MSE}}\|_2^2] = \text{Trace}(\mathbf{K}_{X|Y})$$

Bayesian Estimation Example 2

- ▶ The three estimates coincide (conditional mean, median, maximum):

$$\hat{\mathbf{x}}_{\text{MSE}}(\mathbf{y}) = \hat{\mathbf{x}}_{\text{MAE}}(\mathbf{y}) = \hat{\mathbf{x}}_{\text{MAP}}(\mathbf{y}) = \mathbf{m}_x + \mathbf{K}_{XY}\mathbf{K}_Y^{-1}(\mathbf{y} - \mathbf{m}_y).$$

1 Orthogonality Property:

$$\mathbb{E} \left[(\mathbf{x} - \mathbb{E}[\mathbf{x}|\mathbf{y}]) g(\mathbf{y}) \right] = \mathbf{0} \quad (1)$$

- ▶ The above states that the expected value of the inner product of the error vector with any function of the measurements is always zero.
- ▶ This property is just an expression of the fact that the conditional mean $\mathbb{E}[\mathbf{x}|\mathbf{y}]$ extracts all the information in \mathbf{y} that can be used to reduce the MSE.
- ▶ Notice that $g(\mathbf{y})$ can be scalar or (row) vector.

Properties of MSE Estimator

Proof.

$$\begin{aligned} \text{1st method: } \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}|\mathbf{y}]) g(\mathbf{y})] &= \mathbb{E}[\mathbf{x} g(\mathbf{y})] - \mathbb{E}[\mathbb{E}[\mathbf{x}|\mathbf{y}]g(\mathbf{y})] \\ &= \mathbb{E}[\mathbf{x} g(\mathbf{y})] - \mathbb{E}_Y \left[\mathbb{E}_{X|Y} [\mathbf{x} g(\mathbf{y})|\mathbf{y}] \right] \\ &\stackrel{(*)}{=} \mathbb{E}[\mathbf{x} g(\mathbf{y})] - \mathbb{E}[\mathbf{x} g(\mathbf{y})] = \mathbf{0} \end{aligned}$$

where at step (*) the law of iterated expectation was used:

$$\begin{aligned} E[h(x, y)] &= E_Y E_{X|Y} [h(x, y)|y] \Leftrightarrow \\ \int \int h(x, y) f(x, y) dx dy &= \int_y \int_x h(x, y) f(x|y) dx f(y) dy \end{aligned}$$

$$\text{2nd method: } \mathbb{E}_{X,Y} [\mathbf{x} g(\mathbf{y})] = \mathbb{E}_Y \left\{ \mathbb{E}_{X|Y} [\mathbf{x} g(\mathbf{y})|\mathbf{y}] \right\} = \mathbb{E}_Y \left[\mathbb{E}_{X|Y} [\mathbf{x}|\mathbf{y}] g(\mathbf{y}) \right]$$

Properties of MSE Estimator

2 Uniqueness Property:

$\mathbb{E}[\mathbf{x}|\mathbf{y}]$ is the **unique** vector function $\in \mathbb{R}^m$ that adheres to the orthogonality property.

Proof.

Suppose that $\mathbf{h}(\mathbf{y})$ is another function with $\mathbb{E}[(\mathbf{x} - \mathbf{h}(\mathbf{y})) g(\mathbf{y})] = \mathbf{0}$ for all functions $g(\cdot)$. Then, it follows:

$$\begin{aligned}\mathbb{E} [\|\mathbb{E}[\mathbf{x}|\mathbf{y}] - \mathbf{h}(\mathbf{y})\|_2^2] &= \mathbb{E} \left[\underbrace{(\mathbb{E}[\mathbf{x}|\mathbf{y}] - \mathbf{h}(\mathbf{y}))}_{g(\mathbf{y})}^T (\mathbb{E}[\mathbf{x}|\mathbf{y}] - \mathbf{x} + \mathbf{x} - \mathbf{h}(\mathbf{y})) \right] \\ &= \mathbb{E} [g^T(\mathbf{y}) (\mathbb{E}[\mathbf{x}|\mathbf{y}] - \mathbf{x}) + g^T(\mathbf{y}) (\mathbf{x} - \mathbf{h}(\mathbf{y}))] \\ &= \mathbf{0} + \mathbf{0} \text{ (due to orthogonality principle)} \Rightarrow \\ &\mathbb{E}[\mathbf{x}|\mathbf{y}] = \mathbf{h}(\mathbf{y}),\end{aligned}$$

since at the last step, the expected value of a non-negative random variable is zero only when the variable is always zero. \square

Properties of MSE Estimator

3 Variance reduction:

if $\mathbf{K}_X = \mathbb{E}[(\mathbf{x} - \mathbf{m}_x)(\mathbf{x} - \mathbf{m}_x)^T]$

and

$$\mathbf{K}_E = \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}|\mathbf{y}])(\mathbf{x} - \mathbb{E}[\mathbf{x}|\mathbf{y}])^T] \equiv \mathbf{K}_{X|Y}$$

then

- $\mathbf{K}_E \leq \mathbf{K}_X$ i.e., $\mathbf{K}_X - \mathbf{K}_E$ is positive semi-definite,
- $\mathbf{K}_E = \mathbf{K}_x$ if and only if (iff) $\mathbf{m}_x = \mathbb{E}[\mathbf{x}|\mathbf{y}]$ i.e., knowledge of the observation \mathbf{y} does not improve the estimate of \mathbf{x} .

Properties of MSE Estimator

Proof.

$\mathbf{x} - \mathbf{m}_x = \mathbf{x} - \mathbb{E}[\mathbf{x}|\mathbf{y}] + \mathbb{E}[\mathbf{x}|\mathbf{y}] - \mathbf{m}_x$, therefore

$$\begin{aligned}\mathbf{K}_X &= \mathbf{K}_E + \mathbb{E} \left[(\mathbf{x} - \mathbb{E}[\mathbf{x}|\mathbf{y}]) (\mathbb{E}[\mathbf{x}|\mathbf{y}] - \mathbf{m}_x)^T \right] \\ &\quad + \mathbb{E} \left[(\mathbb{E}[\mathbf{x}|\mathbf{y}] - \mathbf{m}_x) (\mathbf{x} - \mathbb{E}[\mathbf{x}|\mathbf{y}])^T \right] \\ &\quad + \mathbb{E} \left[\underbrace{(\mathbb{E}[\mathbf{x}|\mathbf{y}] - \mathbf{m}_x) (\mathbb{E}[\mathbf{x}|\mathbf{y}] - \mathbf{m}_x)^T}_{\Delta(\mathbf{y})} \right] \Leftrightarrow\end{aligned}$$

$$\begin{aligned}\mathbf{K}_X &= \mathbf{K}_E + \mathbb{E} \left[(\mathbf{x} - \mathbb{E}[\mathbf{x}|\mathbf{y}]) \Delta(\mathbf{y})^T \right] + \mathbb{E} \left[\Delta(\mathbf{y}) (\mathbf{x} - \mathbb{E}[\mathbf{x}|\mathbf{y}])^T \right] \\ &\quad + \mathbb{E} \left[\underbrace{\Delta(\mathbf{y}) \Delta(\mathbf{y})^T}_{\mathbf{K}_\Delta} \right] = \mathbf{K}_E + \mathbf{K}_\Delta \Leftrightarrow\end{aligned}$$

$$\begin{aligned}\mathbf{K}_X - \mathbf{K}_E &= \mathbf{K}_\Delta = \mathbb{E}[\Delta(\mathbf{y}) \Delta(\mathbf{y})^T] \geq 0 \text{ since} \\ \mathbf{z}^T \mathbb{E} [\Delta(\mathbf{y}) \Delta(\mathbf{y})^T] \mathbf{z} &= \mathbb{E} \left[\underbrace{\mathbf{z}^T \Delta \mathbf{y} \Delta \mathbf{y}^T \mathbf{z}}_{\mathbf{z}_0^T} \right] = \mathbb{E} [\|\mathbf{z}_0\|_2^2] \geq 0\end{aligned}$$

□

Properties of MSE Estimator

...proof continued.

- $\mathbf{K}_\Delta = \mathbf{0}$ (all elements zero) \Rightarrow $\text{Trace}(\mathbf{K}_\Delta) = 0$ and the following holds:

$$\begin{aligned}\text{Trace}(\mathbf{K}_\Delta) &= \text{Trace} [\mathbb{E}[\Delta \Delta^T]] = \mathbb{E}[\text{Trace}(\Delta \Delta^T)] \\ &= \mathbb{E}[\Delta^T \Delta] = \mathbb{E}[\|\Delta\|_2^2] = 0 \Rightarrow \\ \Delta \equiv \Delta(\mathbf{y}) &= \mathbf{0} \Rightarrow \mathbb{E}[\mathbf{x}|\mathbf{y}] = \mathbf{m}_x\end{aligned}$$

- For the other direction, i.e., $\Delta(\mathbf{y}) = \mathbf{0} \Rightarrow \mathbf{K}_\Delta = \mathbf{0}$ the proof is trivial.

□

MMSE Estimation in Linear Gaussian Systems

- ▶ Let $\mathbf{x} \in \mathbb{R}^{D_x}$ to be estimated and $\mathbf{y} \in \mathbb{R}^{D_y}$ with

$$p(\mathbf{x}) \sim \mathcal{N}(\boldsymbol{\mu}_x, \boldsymbol{\Sigma}_x), p(\mathbf{y}|\mathbf{x}) \sim \mathcal{N}(\mathbf{A}\mathbf{x} + \mathbf{b}, \boldsymbol{\Sigma}_y)$$

where \mathbf{A} (deterministic) $D_y \times D_x$ real matrix, $\mathbf{b} \in \mathbb{R}^{D_y}$ and $\boldsymbol{\mu}_y = \mathbb{E}[\mathbf{y}] = \mathbf{A}\boldsymbol{\mu}_x + \mathbf{b}$.

Then, $p(\mathbf{x}|\mathbf{y}) \sim \mathcal{N}(\boldsymbol{\mu}_{x|y}, \boldsymbol{\Sigma}_{x|y})$ with

$$\boldsymbol{\mu}_{x|y} = \boldsymbol{\Sigma}_{x|y} [\mathbf{A}^T \boldsymbol{\Sigma}_y^{-1} (\mathbf{y} - \mathbf{b}) + \boldsymbol{\Sigma}_x^{-1} \boldsymbol{\mu}_x]$$

and

$$\boldsymbol{\Sigma}_{x|y} = (\boldsymbol{\Sigma}_x^{-1} + \mathbf{A}^T \boldsymbol{\Sigma}_y^{-1} \mathbf{A})^{-1}$$

- ▶ Thus,

$$\hat{\mathbf{x}}_{\text{MMSE}}(\mathbf{y}) \equiv \mathbb{E}[\mathbf{x}|\mathbf{y}] \equiv \boldsymbol{\mu}_{x|y} = \boldsymbol{\Sigma}_{x|y} [\mathbf{A}^T \boldsymbol{\Sigma}_y^{-1} (\mathbf{y} - \mathbf{b}) + \boldsymbol{\Sigma}_x^{-1} \boldsymbol{\mu}_x].$$

MMSE Estimation in Linear Gaussian Systems

Proof (1/4):

$$\begin{aligned} p(\mathbf{x}, \mathbf{y}) &= p(\mathbf{y}|\mathbf{x})p(\mathbf{x}) \Rightarrow \\ \log p(\mathbf{x}, \mathbf{y}) &= \log p(\mathbf{y}|\mathbf{x}) + \log p(\mathbf{x}) = \\ &= -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_x)^T \boldsymbol{\Sigma}_x^{-1}(\mathbf{x} - \boldsymbol{\mu}_x) - \frac{1}{2}(\mathbf{y} - \mathbf{Ax} - \mathbf{b})^T \boldsymbol{\Sigma}_y^{-1}(\mathbf{y} - \mathbf{Ax} - \mathbf{b}) + \\ &+ \text{constant terms} \\ &= -\frac{1}{2}\mathbf{x}^T \boldsymbol{\Sigma}_x^{-1}\mathbf{x} - \frac{1}{2}\mathbf{y}^T \boldsymbol{\Sigma}_y^{-1}\mathbf{y} - \frac{1}{2}(\mathbf{Ax})^T \boldsymbol{\Sigma}_y^{-1}(\mathbf{Ax}) + \mathbf{y}^T \boldsymbol{\Sigma}_y^{-1}(\mathbf{Ax}) + \\ &+ \text{linear terms} + \text{constant terms} \\ &= -\frac{1}{2}\mathbf{x}^T \boldsymbol{\Sigma}_x^{-1}\mathbf{x} - \frac{1}{2}\mathbf{x}^T \mathbf{A}^T \boldsymbol{\Sigma}_y^{-1} \mathbf{Ax} - \frac{1}{2}\mathbf{y}^T \boldsymbol{\Sigma}_y^{-1}\mathbf{y} + \mathbf{y}^T \boldsymbol{\Sigma}_y^{-1}(\mathbf{Ax}) + \\ &+ \text{linear terms} + \text{constant terms} \end{aligned}$$

MMSE Estimation in Linear Gaussian Systems

Continue proof (2/4):

$$\begin{aligned} &= -\frac{1}{2} \begin{bmatrix} \mathbf{x}^T & \mathbf{y}^T \end{bmatrix} \begin{bmatrix} \Sigma_x^{-1} + \mathbf{A}^T \Sigma_y^{-1} \mathbf{A} & -\mathbf{A}^T \Sigma_y^{-1} \\ -\Sigma_y^{-1} \mathbf{A} & \Sigma_y^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} + \text{constant terms} \\ &= -\frac{1}{2} \begin{bmatrix} \mathbf{x}^T & \mathbf{y}^T \end{bmatrix} \Sigma^{-1} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \Rightarrow \\ \Sigma^{-1} &= \begin{bmatrix} \Sigma_x^{-1} + \mathbf{A}^T \Sigma_y^{-1} \mathbf{A} & -\mathbf{A}^T \Sigma_y^{-1} \\ -\Sigma_y^{-1} \mathbf{A} & \Sigma_y^{-1} \end{bmatrix} \triangleq \Lambda = \begin{bmatrix} \Lambda_{xx} & \Lambda_{xy} \\ \Lambda_{yx} & \Lambda_{yy} \end{bmatrix} \quad (2) \end{aligned}$$

Useful (for the proof) Theorem

Continue proof (3/4): The following theorem will be utilized; its proof will be given in the problem sets and can be found in various textbooks, e.g., Chapter 4 in *"Machine Learning, a Probabilistic Perspective"* by Kevin Murphy):

- Assume $\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}$ Gaussian vector, i.e. $\mathbf{x}_1, \mathbf{x}_2$ jointly Gaussians, with $\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}$, $\boldsymbol{\Sigma} \triangleq \mathbb{E} [(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T] = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}$, $\boldsymbol{\Lambda} \triangleq \boldsymbol{\Sigma}^{-1} = \begin{bmatrix} \boldsymbol{\Lambda}_{11} & \boldsymbol{\Lambda}_{12} \\ \boldsymbol{\Lambda}_{21} & \boldsymbol{\Lambda}_{22} \end{bmatrix}$ (\mathbf{A}^{**})

Then:

$$\begin{aligned} p(\mathbf{x}_1 | \mathbf{x}_2) &= \mathcal{N}(\boldsymbol{\mu}_{1|2}, \boldsymbol{\Sigma}_{1|2}), \quad p(\mathbf{x}_1) = \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{11}), \quad p(\mathbf{x}_2) = \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22}), \\ \boldsymbol{\mu}_{1|2} &= \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{x}_2 - \boldsymbol{\mu}_2) \\ &= \boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_{11}^{-1} \boldsymbol{\Lambda}_{12} (\mathbf{x}_2 - \boldsymbol{\mu}_2) = \boldsymbol{\Lambda}_{11}^{-1} [\boldsymbol{\Lambda}_{11} \boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_{12} (\mathbf{x}_2 - \boldsymbol{\mu}_2)] \\ &= \boldsymbol{\Sigma}_{1|2} [\boldsymbol{\Lambda}_{11} \boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_{12} (\mathbf{x}_2 - \boldsymbol{\mu}_2)] \end{aligned} \tag{3}$$

$$\boldsymbol{\Sigma}_{1|2} = \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} = \boldsymbol{\Lambda}_{11}^{-1} \tag{4}$$

MMSE Estimation in Linear Gaussian Systems

Continue proof (4/4).

Thus, from (\mathbf{A}^{**}) and Eq. (4),

$$\Sigma_{x|y} = \Lambda_{xx}^{-1} = (\Sigma_x^{-1} + \mathbf{A}^T \Sigma_y^{-1} \mathbf{A})^{-1}$$

From from (\mathbf{A}^{**}) , Eq. (3), and $\boldsymbol{\mu}_y = \mathbf{A}\boldsymbol{\mu}_x + \mathbf{b}$,

$$\begin{aligned}\boldsymbol{\mu}_{x|y} &= \Sigma_{x|y} [\Sigma_{x|y}^{-1} \boldsymbol{\mu}_x - \Lambda_{xy} (\mathbf{y} - \boldsymbol{\mu}_y)] \\ &= \Sigma_{x|y} [(\Sigma_x^{-1} + \mathbf{A}^T \Sigma_y^{-1} \mathbf{A}) \boldsymbol{\mu}_x + \mathbf{A}^T \Sigma_y^{-1} (\mathbf{y} - \mathbf{A}\boldsymbol{\mu}_x - \mathbf{b})] \\ &= \Sigma_{x|y} [\Sigma_x^{-1} \boldsymbol{\mu}_x + \mathbf{A}^T \Sigma_y^{-1} \mathbf{A} \boldsymbol{\mu}_x + \mathbf{A}^T \Sigma_y^{-1} (\mathbf{y} - \mathbf{b}) - \mathbf{A}^T \Sigma_y^{-1} \mathbf{A} \boldsymbol{\mu}_x] \\ &= \Sigma_{x|y} [\Sigma_x^{-1} \boldsymbol{\mu}_x + \mathbf{A}^T \Sigma_y^{-1} (\mathbf{y} - \mathbf{b})]\end{aligned}$$

□

- Important Remark: for the above $p(\mathbf{x})$ and $p(\mathbf{y})$, the following can be also shown:

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{A}\boldsymbol{\mu}_x + \mathbf{b}, \Sigma_y + \mathbf{A}\Sigma_x\mathbf{A}^T)$$

MMSE in Linear Gaussian Systems: Example

$$f(\mathbf{y}_i|\mathbf{x}) = \mathcal{N}(\mathbf{x}, \Sigma_y), f(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_0, \Sigma_0)$$

- We observe $\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_N$, which are i.i.d.
- We would like to estimate \mathbf{x} based on \mathbf{y}_0 , where

$$\mathbf{y}_0 = \frac{1}{N} \sum_{i=1}^N \mathbf{y}_i \quad (5)$$

- Notice that $\mathbf{y}_0 \sim \mathcal{N}(\mathbf{x}, \frac{1}{N} \Sigma_y)$. This can be easily shown by the fact that affine transformation of a Gaussian vector is again a Gaussian vector¹ and the fact that:

$$\mathbf{y}_0 = \frac{1}{N} \underbrace{\begin{bmatrix} 1 & 1 & \dots & 1 \end{bmatrix}}_B \underbrace{\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_N \end{bmatrix}}_y$$

¹if $\mathbf{x} \sim \mathcal{N}(\mathbf{m}, \Sigma)$ then $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b} \sim \mathcal{N}(\mathbf{A}\mathbf{m} + \mathbf{b}, \mathbf{A}\Sigma\mathbf{A}^T)$

MMSE in Linear Gaussian Systems: Example

- Thus, $f(\mathbf{y}_0|\mathbf{x}) = \mathcal{N}(\mathbf{x}, \frac{1}{N}\boldsymbol{\Sigma}_y)$, i.e.,
 $\mathbf{A}\mathbf{x} + \mathbf{b} = \mathbf{x} \Rightarrow \mathbf{A} = \mathbf{I}, \mathbf{b} = \mathbf{0}$ and $f(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$.
- With direct application of the above theorem,
 $f(\mathbf{x}|\mathbf{y}_0) = \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ where

$$\boldsymbol{\mu} = \boldsymbol{\Sigma} \left[\mathbf{I}^T \left(\frac{1}{N} \boldsymbol{\Sigma}_y \right)^{-1} (\mathbf{y}_0 - \mathbf{0}) + \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0 \right]$$

and

$$\boldsymbol{\Sigma} = \left(\boldsymbol{\Sigma}_0^{-1} + \left(\frac{1}{N} \boldsymbol{\Sigma}_y \right)^{-1} \right)^{-1} = \left(\boldsymbol{\Sigma}_0^{-1} + N \boldsymbol{\Sigma}_y^{-1} \right)^{-1}$$

Therefore,

$$\hat{\mathbf{x}}(\mathbf{y})_{\text{MSE}} \equiv \boldsymbol{\mu} = \left(\boldsymbol{\Sigma}_0^{-1} + N \boldsymbol{\Sigma}_y^{-1} \right)^{-1} \left(N \boldsymbol{\Sigma}_y^{-1} \mathbf{y}_0 + \boldsymbol{\Sigma}_0^{-1} \boldsymbol{\mu}_0 \right).$$

Thank you!



Detection & Estimation Theory: Lecture 13

Prof. Aggelos Bletsas

Technical University of Crete
School of Electrical and Computer Engineering



Operational Programme
**Human Resources Development,
Education and Lifelong Learning**
Co-financed by Greece and the European Union





“Support for the Internationalization of Higher Education, School of Electrical and Computer Engineering, Technical University of Crete” (MIS 5150766), under the call for proposals “Support of the Internationalization of Higher Education - Technical University of Crete” (EDULLL 153).

The project is co-financed by Greece and the European Union (European Social Fund- ESF) through the Operational Programme “Human Resources Development, Education and Lifelong Learning 2014-2020.



Operational Programme
**Human Resources Development,
Education and Lifelong Learning**
Co-financed by Greece and the European Union



- Bayesian Linear Estimators
 - Derivation

- Remarks
 - Proof on Remark 4

- Linear MSE Estimator Example

Bayesian Linear Estimators: Problem Formulation

- ▶ Assume $\mathbf{m}_x = \mathbb{E}[\mathbf{x}]$, $\mathbf{m}_y = \mathbb{E}[\mathbf{y}]$ and the joint covariance matrix $\mathbf{K} = \mathbb{E} \left[\begin{bmatrix} \mathbf{x} - \mathbf{m}_x \\ \mathbf{y} - \mathbf{m}_y \end{bmatrix} \begin{bmatrix} (\mathbf{x} - \mathbf{m}_x)^T & (\mathbf{y} - \mathbf{m}_y)^T \end{bmatrix} \right] = \begin{bmatrix} \mathbf{K}_X & \mathbf{K}_{XY} \\ \mathbf{K}_{YX} & \mathbf{K}_Y \end{bmatrix}$ are known
- ▶ $\mathbf{K} > 0$ (\mathbf{K}^{-1} exists); otherwise a non-trivial linear combination in vector $\mathbf{y} - \mathbf{m}_y$ exists, so we could replace observation \mathbf{y} by a vector of smaller dimension!
- ▶ $f(\mathbf{y}|\mathbf{x})$ and $f(\mathbf{x})$ are UNKNOWN!
- ▶ $\mathbf{y} \in \mathbb{R}^n, \mathbf{x} \in \mathbb{R}^m$
- ▶ we are looking for $\hat{\mathbf{x}}_L(\mathbf{y}) = \mathbf{A}\mathbf{y} + \mathbf{b}$ that minimizes $\mathbb{E}[\|\mathbf{e}\|_2^2]$ (MSE) with the following:
 1. \mathbf{A} a $m \times n$ matrix
 2. $\mathbf{b} \in \mathbb{R}^m$
 3. $\mathbf{e} = \mathbf{x} - \hat{\mathbf{x}}_L(\mathbf{y}) = \mathbf{x} - \mathbf{A}\mathbf{y} - \mathbf{b}$
 4. $\mathbb{E}[\mathbf{e}] \equiv \mathbf{m}_e = \mathbf{m}_x - \mathbf{A}\mathbf{m}_y - \mathbf{b}$

Bayesian Linear Estimators

- ▶ Notice that

$$\begin{aligned}\hat{\mathbf{x}}(\mathbf{y}) = \mathbb{E}[\mathbf{x}|\mathbf{y}] &\Rightarrow \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}|\mathbf{y}]) \cdot \mathbf{g}(\mathbf{y})] = 0 \quad (\text{orthogonality property}) \\ &\Rightarrow \mathbb{E}[(\mathbf{x} - \mathbb{E}[\mathbf{x}|\mathbf{y}])] = 0 \quad (g(\mathbf{y}) = 1) \\ &\Rightarrow \mathbb{E}[\mathbf{e}] = 0\end{aligned}$$

- ▶ For $\hat{\mathbf{x}}_L(\mathbf{y}) = \mathbf{A}\mathbf{y} + \mathbf{b} \Rightarrow \mathbb{E}[\mathbf{e}] = \mathbf{m}_x - \mathbf{A}\mathbf{m}_y - \mathbf{b} \neq 0$

- ▶ Thus,

$$\mathbf{K}_E \triangleq \mathbb{E}[(\mathbf{e} - \mathbf{m}_e)(\mathbf{e} - \mathbf{m}_e)^T] \quad (1)$$

$$\mathbf{e} - \mathbf{m}_e = (\mathbf{x} - \mathbf{A}\mathbf{y} - \mathbf{b}) - (\mathbf{m}_x - \mathbf{A}\mathbf{m}_y - \mathbf{b}) \quad (2)$$

$$= [\mathbf{I}_m \quad -\mathbf{A}] \begin{bmatrix} \mathbf{x} - \mathbf{m}_x \\ \mathbf{y} - \mathbf{m}_y \end{bmatrix} \quad (3)$$

- ▶ From Eq. (1) and Eq. (3):

$$\mathbf{K}_E = [\mathbf{I}_m \quad -\mathbf{A}] \begin{bmatrix} \mathbf{K}_X & \mathbf{K}_{XY} \\ \mathbf{K}_{YX} & \mathbf{K}_Y \end{bmatrix} \begin{bmatrix} \mathbf{I}_m \\ -\mathbf{A}^T \end{bmatrix}$$

Bayesian Linear Estimators

- ▶ MMSE:

$$\mathbb{E} [\|\mathbf{e}\|_2^2] = \mathbb{E} [\mathbf{e}^T \mathbf{e}] = \|\mathbf{e}\|_2^2 - \|\mathbf{m}_e\|_2^2 + \|\mathbf{m}_e\|_2^2 \quad (4)$$

$$= \mathbb{E} [(\mathbf{e} - \mathbf{m}_e)^T (\mathbf{e} - \mathbf{m}_e)] + \|\mathbf{m}_e\|_2^2 \quad (5)$$

$$= \text{Trace}(\mathbf{K}_E) + \|\mathbf{m}_e\|_2^2 \quad (6)$$

- ▶ \mathbf{K}_E depends on \mathbf{A} , \mathbf{m}_e depends on \mathbf{A} and \mathbf{b}
- ▶ Set $\mathbf{m}_e = 0 \Rightarrow \mathbf{m}_x - \mathbf{A}\mathbf{m}_y - \mathbf{b} = 0 \Rightarrow$

$$\mathbf{b} = \mathbf{m}_x - \mathbf{A}\mathbf{m}_y \quad (7)$$

- ▶ Now we need to find \mathbf{A} . We work as follows:

$$\begin{aligned} \mathbf{K}_E &= \begin{bmatrix} \mathbf{I}_m & -\mathbf{A} \end{bmatrix} \begin{bmatrix} \mathbf{K}_X & \mathbf{K}_{XY} \\ \mathbf{K}_{YX} & \mathbf{K}_Y \end{bmatrix} \begin{bmatrix} \mathbf{I}_m \\ -\mathbf{A}^T \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{K}_X - \mathbf{A}\mathbf{K}_{YX} & \mathbf{K}_{XY} - \mathbf{A}\mathbf{K}_Y \end{bmatrix} \begin{bmatrix} \mathbf{I}_m \\ -\mathbf{A}^T \end{bmatrix} \\ &= \mathbf{K}_X - \mathbf{A}\mathbf{K}_{YX} - \mathbf{K}_{XY}\mathbf{A}^T + \mathbf{A}\mathbf{K}_Y\mathbf{A}^T \\ &= \mathbf{K}_{XY}\mathbf{K}_Y^{-1}\mathbf{K}_{YX} - \mathbf{K}_{XY}\mathbf{A}^T - \mathbf{A}\mathbf{K}_{YX} + \mathbf{A}\mathbf{K}_Y\mathbf{A}^T + \mathbf{S} \\ &= (\mathbf{K}_{XY} - \mathbf{A}\mathbf{K}_Y) \left((\mathbf{K}_{XY}\mathbf{K}_Y^{-1})^T - \mathbf{A}^T \right) + \mathbf{S} \\ &= (\mathbf{K}_{XY}\mathbf{K}_Y^{-1} - \mathbf{A}) \mathbf{K}_Y (\mathbf{K}_{XY}\mathbf{K}_Y^{-1} - \mathbf{A})^T + \mathbf{S} \end{aligned}$$

Bayesian Linear Estimators

- ▶ Schur complement $\mathbf{S} \triangleq \mathbf{K}_X - \mathbf{K}_{XY}\mathbf{K}_Y^{-1}\mathbf{K}_{YX}$
 - ▶ Schur complement of \mathbf{K}_Y in \mathbf{K} plays a role in evaluating the determinant and the inverse of block matrices.
 - ▶ Schur complement is constant and known and does not depend on \mathbf{A}

Thus, $\text{Trace}(\mathbf{K}_E) =$

$\text{Trace} \left((\mathbf{K}_{XY}\mathbf{K}_Y^{-1} - \mathbf{A}) \mathbf{K}_Y (\mathbf{K}_{XY}\mathbf{K}_Y^{-1} - \mathbf{A})^T \right) + \text{Trace}(\mathbf{S})$,
since $\text{Trace}(\mathbf{A} + \mathbf{B}) = \text{Trace}(\mathbf{A}) + \text{Trace}(\mathbf{B})$.

- ▶ $\text{Trace} \left((\mathbf{K}_{XY}\mathbf{K}_Y^{-1} - \mathbf{A}) \mathbf{K}_Y (\mathbf{K}_{XY}\mathbf{K}_Y^{-1} - \mathbf{A})^T \right) \geq 0$, since $(\mathbf{K}_{XY}\mathbf{K}_Y^{-1} - \mathbf{A}) \mathbf{K}_Y (\mathbf{K}_{XY}\mathbf{K}_Y^{-1} - \mathbf{A})^T$ is positive semi-definite.
Thus, $\text{Trace} \left((\mathbf{K}_{XY}\mathbf{K}_Y^{-1} - \mathbf{A}) \mathbf{K}_Y (\mathbf{K}_{XY}\mathbf{K}_Y^{-1} - \mathbf{A})^T \right) = 0 \Leftrightarrow \mathbf{K}_{XY}\mathbf{K}_Y^{-1} = \mathbf{A}$.
- ▶ Thus, $\text{Tr}(\mathbf{K}_E)$ is minimized iff

$$\mathbf{K}_{XY}\mathbf{K}_Y^{-1} = \mathbf{A} \quad (8)$$

- ▶ From Eq. (7) and Eq. (8), $\hat{\mathbf{x}}_L(\mathbf{y}) = \mathbf{m}_x + \mathbf{K}_{XY}\mathbf{K}_Y^{-1}(\mathbf{y} - \mathbf{m}_y)$.

Remarks

1. $\mathbb{E}[\mathbf{e}] = \mathbf{m}_e = \mathbf{0}$
2. $\hat{\mathbf{x}}_L(\mathbf{y}) = \mathbf{m}_x + \mathbf{K}_{XY}\mathbf{K}_Y^{-1}(\mathbf{y} - \mathbf{m}_y) = \hat{\mathbf{x}}_{\text{MSE}}(\mathbf{y})$
for \mathbf{x}, \mathbf{y} jointly Gaussians

$$\begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{m}_x \\ \mathbf{m}_y \end{bmatrix}, \begin{bmatrix} \mathbf{K}_X & \mathbf{K}_{XY} \\ \mathbf{K}_{YX} & \mathbf{K}_Y \end{bmatrix} \right)$$

Linear estimate and MSE estimate coincide for the Gaussian case, i.e., all MSE estimates will necessarily be linear.

3. As promised, linear-least-square estimate $\hat{\mathbf{x}}_L(\mathbf{y})$ requires knowledge of first and second moments and not knowledge of $f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}), f_{\mathbf{x}}(\mathbf{x})$.
4. Orthogonality property holds only for linear functions of \mathbf{y} :

$$\mathbb{E}[(\mathbf{x} - \hat{\mathbf{x}}_L(\mathbf{y})) \cdot \mathbf{g}(\mathbf{y})] = \mathbf{0}$$

for all linear functions $\mathbf{g}(\mathbf{y}) = g_0 + \mathbf{y}^T \mathbf{g}_1$ where g_0 a real scalar and \mathbf{g}_1 a vector in \mathbb{R}^n .

Proof.

▶ Estimation error: $\mathbf{x} - \hat{\mathbf{x}}_L(\mathbf{y}) = \begin{bmatrix} \mathbf{I}_m & -\mathbf{K}_{XY}\mathbf{K}_Y^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{x} - \mathbf{m}_x \\ \mathbf{y} - \mathbf{m}_y \end{bmatrix}$

▶ $\mathbf{g}(\mathbf{y}) = g_0 + \mathbf{m}_y^T \mathbf{g}_1 + (\mathbf{y} - \mathbf{m}_y)^T \mathbf{g}_1$ where $g_0 + \mathbf{m}_y^T \mathbf{g}_1$ is a constant term.

▶ Therefore,

$$\begin{aligned} \mathbb{E}[(\mathbf{x} - \hat{\mathbf{x}}_L(\mathbf{y})) \mathbf{g}(\mathbf{y})] &= \mathbb{E}[(\mathbf{x} - \hat{\mathbf{x}}_L(\mathbf{y})) (\mathbf{y} - \mathbf{m}_y)^T \mathbf{g}_1] \\ &= \begin{bmatrix} \mathbf{I}_m & -\mathbf{K}_{XY}\mathbf{K}_Y^{-1} \end{bmatrix} \mathbb{E} \left[\begin{bmatrix} \mathbf{x} - \mathbf{m}_x \\ \mathbf{y} - \mathbf{m}_y \end{bmatrix} \begin{bmatrix} \mathbf{y} & -\mathbf{m}_y \end{bmatrix}^T \right] \mathbf{g}_1 \\ &= \begin{bmatrix} \mathbf{I}_m & -\mathbf{K}_{XY}\mathbf{K}_Y^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{K}_{XY} \\ \mathbf{K}_Y \end{bmatrix} \mathbf{g}_1 \\ &\quad \begin{matrix} m \times (m+n) & (m+n) \times n \\ n \times 1 \end{matrix} \\ &= (\mathbf{K}_{XY} \quad -\mathbf{K}_{XY}) \mathbf{g}_1 \\ &\quad \begin{matrix} m \times n & n \times 1 \end{matrix} \\ &= \mathbf{0} \end{aligned}$$

□

Proof.

5. $\hat{\mathbf{x}}_L(\mathbf{y})$ is the unique linear estimator that adheres to the orthogonality property.

▶ Suppose that $\mathbf{h}(\mathbf{y})$ is another linear estimator with the property $\mathbb{E}[(\mathbf{x} - \mathbf{h}(\mathbf{y})) \cdot \mathbf{g}(\mathbf{y})] = \mathbf{0}$.

▶ Thus

$$\begin{aligned}\mathbb{E}[\|\hat{\mathbf{x}}_L(\mathbf{y}) - \mathbf{h}(\mathbf{y})\|_2^2] &= \mathbb{E}[\mathbf{g}^T(\mathbf{y})(\hat{\mathbf{x}}_L(\mathbf{y}) - \mathbf{x} + \mathbf{x} - \mathbf{h}(\mathbf{y}))] \\ &= \mathbb{E}[\mathbf{g}^T(\mathbf{y})(\hat{\mathbf{x}}_L(\mathbf{y}) - \mathbf{x})] + \mathbb{E}[\mathbf{g}^T(\mathbf{y})(\mathbf{x} - \mathbf{h}(\mathbf{y}))] = \mathbf{0} + \mathbf{0} \\ &= \mathbf{0},\end{aligned}$$

since $\mathbf{g}(\mathbf{y}) = \hat{\mathbf{x}}_L(\mathbf{y}) - \mathbf{h}(\mathbf{y})$ is linear (because $\hat{\mathbf{x}}_L(\mathbf{y}), \mathbf{h}(\mathbf{y})$ are linear in \mathbf{y}).

□

Example

- ▶ Assume $\mathbf{y} = \mathbf{H}\mathbf{x} + \mathbf{v}$
 1. $\mathbf{x} \in \mathbb{R}^m$, $\mathbf{y} \in \mathbb{R}^n$
 2. \mathbf{v} uncorrelated with \mathbf{x}
 3. $\mathbb{E}[\mathbf{v}] = \mathbf{0}$ and $\mathbb{E}[\mathbf{v}\mathbf{v}^T] = \mathbf{R} > \mathbf{0}$, i.e., \mathbf{R} is positive-definite
 4. \mathbf{H} (known) constant matrix
- ▶ We need to find the $\hat{\mathbf{x}}_L(\mathbf{y})$:
 - ▶ $\mathbf{m}_y = \mathbf{H}\mathbf{m}_x$
 - ▶ $\mathbf{K}_{YX} = \mathbb{E}[(\mathbf{y} - \mathbf{m}_y)(\mathbf{x} - \mathbf{m}_x)^T]$
$$= \mathbb{E}[(\mathbf{H}(\mathbf{x} - \mathbf{m}_x) + \mathbf{v})(\mathbf{x} - \mathbf{m}_x)^T]$$
$$= \mathbf{H}\mathbf{K}_X$$
since $\mathbb{E}[\mathbf{v}] = \mathbf{0}$ and \mathbf{v}, \mathbf{x} are uncorrelated
 - ▶ $\mathbf{K}_Y = \mathbb{E}[(\mathbf{y} - \mathbf{m}_y)(\mathbf{y} - \mathbf{m}_y)^T]$
$$= \mathbb{E}[(\mathbf{H}(\mathbf{x} - \mathbf{m}_x) + \mathbf{v})(\mathbf{H}(\mathbf{x} - \mathbf{m}_x) + \mathbf{v})^T]$$
$$= \mathbf{H}\mathbf{K}_X\mathbf{H}^T + \mathbf{R}$$

Example

► Therefore,

► $\hat{\mathbf{x}}_L(\mathbf{y}) = \mathbf{m}_x + \mathbf{K}_X \mathbf{H}^T (\mathbf{H} \mathbf{K}_X \mathbf{H}^T + \mathbf{R})^{-1} (\mathbf{y} - \mathbf{H} \mathbf{m}_x)$

► $\mathbf{K}_L = \mathbf{K}_X - \mathbf{K}_X \mathbf{H}^T (\mathbf{H} \mathbf{K}_X \mathbf{H}^T + \mathbf{R})^{-1} \mathbf{H} \mathbf{K}_X \equiv \mathbf{K}_E$

► We can show that

$$\mathbf{K}_L^{-1} = \mathbf{K}_X^{-1} + \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H} \quad (\text{A})$$

i.e. $\mathbf{K}_L = (\mathbf{K}_X^{-1} + \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H})^{-1}$

("Sherman-Morrison-Woodbury" identity)

Proof:

► $(\mathbf{A} + \mathbf{BCD})^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{B} (\mathbf{C}^{-1} + \mathbf{DA}^{-1} \mathbf{B})^{-1} \mathbf{DA}^{-1}$

► $(\mathbf{K}_L)^{-1} = \mathbf{K}_X - \mathbf{K}_X \mathbf{H}^T (\mathbf{R} + \mathbf{H} \mathbf{K}_X \mathbf{H}^T)^{-1} \mathbf{H} \mathbf{K}_X,$

where $\mathbf{A} = \mathbf{K}_X^{-1}, \mathbf{B} = \mathbf{H}^T, \mathbf{C} = \mathbf{R}^{-1}, \mathbf{D} = \mathbf{H}$

Example

► Denote $\mathbf{G} \triangleq \mathbf{K}_X \mathbf{H}^T (\mathbf{H} \mathbf{K}_X \mathbf{H}^T + \mathbf{R})^{-1}$ (α)

► It is true that $\mathbf{G} = \mathbf{K}_L \mathbf{H}^T \mathbf{R}^{-1}$ (β)

► Proof:

1. $\mathbf{K}_L^{-1} \cdot \mathbf{G} \cdot (\mathbf{H} \mathbf{K}_X \mathbf{H}^T + \mathbf{R}) \stackrel{(\beta)}{=} \mathbf{H}^T \mathbf{R}^{-1} (\mathbf{H} \mathbf{K}_X \mathbf{H}^T + \mathbf{R})$ (6)

2. $\mathbf{K}_L^{-1} \cdot \mathbf{G} \cdot (\mathbf{H} \mathbf{K}_X \mathbf{H}^T + \mathbf{R}) \stackrel{(\alpha)}{=} (\mathbf{K}_X^{-1} + \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H}) \mathbf{K}_X \mathbf{H}^T$ (7)
which is true

3. (6) = (7) after simple manipulations...

► Also

$$\begin{aligned} \mathbf{K}_L^{-1} \hat{\mathbf{x}}_L(\mathbf{y}) &= (\mathbf{K}_X^{-1} + \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H}) (\mathbf{m}_x + \mathbf{G} \cdot (\mathbf{y} - \mathbf{H} \mathbf{m}_x)) \\ &= \mathbf{K}_X^{-1} \mathbf{m}_x + \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H} \mathbf{m}_x + \mathbf{H}^T \mathbf{R}^{-1} (\mathbf{y} - \mathbf{H} \mathbf{m}_x) \\ &= \mathbf{K}_X^{-1} \mathbf{m}_x + \mathbf{H}^T \mathbf{R}^{-1} \mathbf{y} \quad (\text{B}) \end{aligned}$$

► (A), (B) are used in the derivation of Kalman Filter.

Example

► Since

1. $\mathbf{G} \equiv \mathbf{K}_L \mathbf{H}^T \mathbf{R}^{-1}$

2. $\hat{\mathbf{x}}_L(\mathbf{y}) = \mathbf{m}_x + \underbrace{\mathbf{K}_X \mathbf{H}^T (\mathbf{H} \mathbf{K}_X \mathbf{H}^T + \mathbf{R})^{-1}}_{\mathbf{G}} (\mathbf{y} - \mathbf{H} \mathbf{m}_x)$

3. $\mathbf{K}_L \equiv (\mathbf{K}_X^{-1} + \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H})^{-1}$

► Then $\hat{\mathbf{x}}_L(\mathbf{y}) = \mathbf{m}_x + (\mathbf{K}_X^{-1} + \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{R}^{-1} (\mathbf{y} - \mathbf{H} \mathbf{m}_x)$

► Notice that the above expression is regularly used in various textbooks.

Thank you!



Detection & Estimation Theory: Lectures 14-16

Prof. Aggelos Bletsas

Technical University of Crete
School of Electrical and Computer Engineering



Operational Programme
**Human Resources Development,
Education and Lifelong Learning**
Co-financed by Greece and the European Union





“Support for the Internationalization of Higher Education, School of Electrical and Computer Engineering, Technical University of Crete” (MIS 5150766), under the call for proposals “Support of the Internationalization of Higher Education - Technical University of Crete” (EDULLL 153).

The project is co-financed by Greece and the European Union (European Social Fund- ESF) through the Operational Programme “Human Resources Development, Education and Lifelong Learning 2014-2020.



Operational Programme
**Human Resources Development,
Education and Lifelong Learning**
Co-financed by Greece and the European Union



- Estimation of Non-random Parameters
 - ML Estimation Examples
- Performance: Cramer-Rao Bound
- Cramer-Rao Bound Derivation
- Existence of Efficient Estimator

Estimation of Non-random Parameters

- ▶ Alternative view: $f_{\mathbf{x}}(\mathbf{x})$ not available, $\mathbf{x} \in \mathbb{R}^m$ is viewed as unknown and non-random!
 1. Likelihood function $f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})$ of measurements vector \mathbf{y} when the parameter vector is \mathbf{x} .
 2. One simple solution: maximum likelihood (ML) estimate!

$$\hat{\mathbf{x}}_{\text{ML}}(\mathbf{y}) = \arg \max_{\mathbf{x} \in \mathbb{R}^m} f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) \quad (1)$$

$$= \arg \max_{\mathbf{x} \in \mathbb{R}^m} \ln f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) \quad (2)$$

...the latter (logarithmic) is convenient for p.d.f. in the exponential family (Poisson, Exponential, Gaussian):

$$f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) = \exp(\mathbf{x}^T \mathbf{s}(\mathbf{y}) - \mathbf{t}(\mathbf{x}))$$

- ▶ Bias: $\mathbf{b} \triangleq \mathbb{E}[\mathbf{x} - \hat{\mathbf{x}}(\mathbf{y})] = \mathbf{x} - \mathbb{E}[\hat{\mathbf{x}}(\mathbf{y})]$
- ▶ Bias of ML estimate may not be zero.

Estimation of Non-random Parameters: Bias

- ▶ Bias: $\mathbf{b} \triangleq \mathbb{E}[\mathbf{x} - \hat{\mathbf{x}}(\mathbf{y})] = \mathbf{x} - \mathbb{E}[\hat{\mathbf{x}}(\mathbf{y})]$
- ▶ ...is the expected value of the error.
- ▶ ...is a weak metric, since it does not ensure that for a single measurement vector \mathbf{y} the estimate will offer the true parameter vector \mathbf{x} .

Estimation of Non-random Parameters: Examples

- ▶ Assume $\mathbf{y} \sim \mathcal{N}(\mathbf{A}\mathbf{s}, \sigma^2 \mathbf{I}_N)$, where $\mathbf{y} \in \mathbb{R}^N$ and \mathbf{s}, σ^2 known.
- ▶ $\hat{A}_{\text{ML}}(\mathbf{y})$?

$$\ln [f_{\mathbf{y}|\mathbf{A}}(\mathbf{y}|\mathbf{A})] = -\frac{1}{2\sigma^2} \|\mathbf{y} - \mathbf{A}\mathbf{s}\|_2^2 + \ln \left[\frac{1}{\sqrt{(2\pi\sigma^2)^N}} \right] \quad (3)$$

$$\Rightarrow \arg \max_{A \in \mathbb{R}} \ln f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) = \arg \min_A \|\mathbf{y} - \mathbf{A}\mathbf{s}\|_2^2 \quad (4)$$

$$= \arg \min_A (\mathbf{y} - \mathbf{A}\mathbf{s})^T (\mathbf{y} - \mathbf{A}\mathbf{s}) \quad (5)$$

$$= \arg \min_A (\|\mathbf{s}\|_2^2 A^2 - 2\mathbf{s}^T \mathbf{y} A + \|\mathbf{y}\|_2^2) \quad (6)$$

$$= \frac{\mathbf{s}^T \mathbf{y}}{\|\mathbf{s}\|_2^2} \quad (7)$$

$$\Rightarrow \hat{A}_{\text{ML}}(\mathbf{y}) = \frac{\mathbf{s}^T \mathbf{y}}{\|\mathbf{s}\|_2^2}. \quad (8)$$

Estimation of Non-random Parameters: Examples

- ▶ Assume $\mathbf{y} = [y_1 \ y_2 \ \dots \ y_N]^T$ of N independent and identically distributed (i.i.d.) r.v.'s $\{y_k\}$, with $y_k \sim \mathcal{N}(m, v)$.
- ▶ Estimation problems:
 1. Estimate m with v known.
 2. Estimate v with m known.
 3. Estimate m, v .
- ▶ Check bias of estimate(s).

Estimation of Non-random Parameters: Examples

- ▶ Assume $\mathbf{y} = [y_1 \ y_2 \ \dots \ y_N]^T$ of N independent and identically distributed (i.i.d.) r.v.'s $\{y_k\}$, with $y_k \sim \mathcal{N}(m, v)$.
- ▶ Estimation problem:
 - 1 Estimate m with v known. Check bias of estimate.

Solution: ...this is the previous example, with $A \equiv m$, $\mathbf{s} = [1 \ 1 \ \dots \ 1]^T$ and $\sigma^2 = v$, Thus,

$$\hat{m}_{\text{ML}}(\mathbf{y}) = \frac{\mathbf{s}^T \mathbf{y}}{N} = \frac{\sum_{k=1}^N y_k}{N} \quad (9)$$

$$\mathbb{E}[\hat{m}_{\text{ML}}(\mathbf{y})] = \frac{Nm}{N} = m \text{ (unbiased estimate)} \quad (10)$$

Estimation of Non-random Parameters: Examples

- ▶ Estimation problem:

2 Estimate v with m known. Check bias of estimate.

$$\text{Solution: } \ln [f_{\mathbf{y}|v}(\mathbf{y}|v)] = -\frac{N}{2} \ln(2\pi v) - \frac{1}{2v} \sum_{k=1}^N (y_k - m)^2$$

$$\frac{d}{dv} \ln [f_{\mathbf{y}|v}(\mathbf{y}|v)] = 0 \Rightarrow \frac{1}{2v} \left[-N + \frac{1}{v} \sum_{k=1}^N (y_k - m)^2 \right] = 0 \quad (11)$$

$$\Rightarrow \hat{v}_{\text{ML}} = \frac{1}{N} \sum_{k=1}^N (y_k - m)^2 \quad (12)$$

- ▶ Need to make sure that $\frac{d^2}{dv^2} \ln [f_{\mathbf{y}|v}(\mathbf{y}|v)] < 0$ at $v = \hat{v}_{\text{ML}}$:

$$\frac{d^2}{dv^2} \ln [f_{\mathbf{y}|v}(\mathbf{y}|v)] = \dots = \frac{1}{v^2} \left(\frac{N}{2} - \frac{1}{v} \sum_{k=1}^N (y_k - m)^2 \right) \quad (13)$$

$$\underset{v=\hat{v}_{\text{ML}}}{=} -\frac{1}{\hat{v}_{\text{ML}}^2} \frac{N}{2} < 0 \quad (14)$$

Estimation of Non-random Parameters: Examples

- ▶ Assume $\mathbf{y} = [y_1 \ y_2 \ \dots \ y_N]^T$ of N independent and identically distributed (i.i.d.) r.v.'s $\{y_k\}$, with $y_k \sim \mathcal{N}(m, v)$.
- ▶ Estimation problem:
 - 2 Estimate v with m known. Check bias of estimate.

$$\hat{v}_{\text{ML}} = \frac{1}{N} \sum_{k=1}^N (y_k - m)^2 \quad (15)$$

$$\mathbb{E} [\hat{v}_{\text{ML}}] = \frac{1}{N} \sum_{k=1}^N \mathbb{E} [(y_k - m)^2] = \frac{Nv}{N} = v \text{ (unbiased estimate)}. \quad (16)$$

Estimation of Non-random Parameters: Examples

- ▶ Estimation problem:

3 Estimate m, v . Check bias of estimates.

$$\frac{\partial}{\partial m} \ln [f_{\mathbf{y}|m,v}(\mathbf{y}|m,v)] = 0 \Rightarrow \frac{1}{v} \sum_{k=1}^N (y_k - m) = 0 \quad (17)$$

$$\Rightarrow \hat{m}_{\text{ML}} = \frac{1}{N} \sum_{k=1}^N y_k \quad (18)$$

$$\frac{\partial}{\partial v} \ln [f_{\mathbf{y}|m,v}(\mathbf{y}|m,v)] = 0 \Rightarrow \frac{1}{2v} \left[-N + \frac{1}{v} \sum_{k=1}^N (y_k - m)^2 \right] = 0 \quad (19)$$

$$\Rightarrow \hat{v}_{\text{ML}} = \frac{1}{N} \sum_{k=1}^N (y_k - \hat{m}_{\text{ML}})^2 \quad (20)$$

- ▶ How do we know that the above maximise the log-likelihood?

Estimation of Non-random Parameters: Examples

- ▶ Estimation problem:

3 Estimate m, v . Check bias of estimates.

$$\hat{m}_{\text{ML}} = \frac{1}{N} \sum_{k=1}^N y_k \quad (21)$$

$$\hat{v}_{\text{ML}} = \frac{1}{N} \sum_{k=1}^N (y_k - \hat{m}_{\text{ML}})^2 \quad (22)$$

$$g(m, v) \triangleq \ln [f_{\mathbf{y}|m,v}(\mathbf{y}|m, v)] \quad (23)$$

$$\mathbf{H}_g(m, v) = \begin{bmatrix} \frac{\partial^2 g(m, v)}{\partial m^2} & \frac{\partial^2 g(m, v)}{\partial m \partial v} \\ \frac{\partial^2 g(m, v)}{\partial v \partial m} & \frac{\partial^2 g(m, v)}{\partial v^2} \end{bmatrix} \quad (24)$$

- ▶ Need to check that the Hessian matrix \mathbf{H}_g on $g(m, v)$ for $m = \hat{m}_{\text{ML}}$ and $v = \hat{v}_{\text{ML}}$ is non-negative definite (left as an exercise for the reader).

Estimation of Non-random Parameters: Examples

- ▶ Estimation problem:

3 Estimate m, v . Check bias of estimates.

$$\mathbb{E} [\hat{v}_{\text{ML}}] = \mathbb{E} \left[\frac{1}{N} \sum_{k=1}^N (y_k - \hat{m}_{\text{ML}})^2 \right] \quad (25)$$

$$= \mathbb{E} \left[\frac{1}{N} \sum_{k=1}^N y_k^2 \right] + \mathbb{E} \left[\frac{1}{N} \sum_{k=1}^N \hat{m}_{\text{ML}}^2 \right] - \mathbb{E} \left[\frac{2}{N} \sum_{k=1}^N (\hat{m}_{\text{ML}} y_k) \right] \quad (26)$$

$$= \dots = \frac{N-1}{N} v \neq v \text{ (biased estimate)} \quad (27)$$

- ▶ That is why numerical packages utilise the following, non-ML, unbiased variance estimate:

$$\frac{1}{N-1} \sum_{k=1}^N (y_k - \hat{m}_{\text{ML}})^2$$

Performance of Non-random Parameter Estimation

- ▶ Apart from bias, we need the mean square error (MSE) and the corresponding error matrix:

$$\mathbf{C}_E = \mathbb{E} [(\mathbf{x} - \hat{\mathbf{x}}(\mathbf{y})) (\mathbf{x} - \hat{\mathbf{x}}(\mathbf{y}))^T] \quad (28)$$

$$\text{MSE} = \text{Trace}(\mathbf{C}_E) = \mathbb{E} [\|\mathbf{x} - \hat{\mathbf{x}}(\mathbf{y})\|_2^2] \quad (29)$$

- ▶ Denote the gradient (vector) $\nabla_x = [\frac{\partial}{\partial x_1} \quad \frac{\partial}{\partial x_2} \dots \frac{\partial}{\partial x_m}]^T$ and the Hessian (matrix) $\nabla_x \nabla_x^T$.

Non-random Parameter Estimation: Cramer-Rao Bound

- ▶ $m \times m$ Fisher Information matrix $\mathbf{J}(\mathbf{x})$ characterises the information in $\mathbf{y} \in \mathbb{R}^n$ about the parameter vector $\mathbf{x} \in \mathbb{R}^m$:

$$\mathbf{J}(\mathbf{x}) \triangleq \mathbb{E}_{\mathbf{y}} \left[\left[\nabla_{\mathbf{x}} \ln f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) \right] \left[\nabla_{\mathbf{x}} \ln f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) \right]^T \right] \quad (30)$$

- ▶ It turns out that

$$\mathbf{J}(\mathbf{x}) = -\mathbb{E}_{\mathbf{y}} \left[\underbrace{\nabla_{\mathbf{x}} \nabla_{\mathbf{x}}^T \ln f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})}_{\text{Hessian of the log-likelihood}} \right] \quad (31)$$

Theorem

For any **unbiased** estimator $\hat{\mathbf{x}}(\mathbf{y})$, the MSE is bounded from the Cramer-Rao bound, which stems from the diagonal elements of the inverse Fisher Information matrix:

$$\mathbb{E} [\|\mathbf{x}_i - \hat{\mathbf{x}}(\mathbf{y})_i\|_2^2] \geq [\mathbf{J}^{-1}(\mathbf{x})]_{ii}, \quad (32)$$

where $\mathbf{a}_i, \mathbf{A}_{ii}$ denotes the i -th element and i -th diagonal element of vector \mathbf{a} and matrix \mathbf{A} , respectively.

Cramer-Rao Bound Example

- Assume $\mathbf{y} = [y_1 \ y_2 \ \dots \ y_N]^T$, with $\{y_k\}$ i.i.d. and $y_k \sim \mathcal{N}(m, v)$.

$$\ln [f_{\mathbf{y}|m,v}(\mathbf{y}|m, v)] = -\frac{N}{2} \ln(2\pi v) - \frac{1}{2v} \sum_{k=1}^N (y_k - m)^2$$

$$\frac{\partial}{\partial m} \ln [f_{\mathbf{y}|m,v}(\mathbf{y}|m, v)] = +\frac{1}{v} \sum_{k=1}^N (y_k - m) \quad (33)$$

$$\frac{\partial^2}{\partial m^2} \ln [f_{\mathbf{y}|m,v}(\mathbf{y}|m, v)] = -\frac{N}{v} \quad (34)$$

$$\frac{\partial}{\partial v} \ln [f_{\mathbf{y}|m,v}(\mathbf{y}|m, v)] = -\frac{N}{2v} + \frac{1}{2v^2} \sum_{k=1}^N (y_k - m)^2 \quad (35)$$

$$\frac{\partial^2}{\partial v^2} \ln [f_{\mathbf{y}|m,v}(\mathbf{y}|m, v)] = +\frac{N}{2v^2} - \frac{1}{v^3} \sum_{k=1}^N (y_k - m)^2 \quad (36)$$

$$\frac{\partial^2}{\partial m \partial v} \ln [f_{\mathbf{y}|m,v}(\mathbf{y}|m, v)] = -\frac{1}{v^2} \sum_{k=1}^N (y_k - m) \quad (37)$$

Cramer-Rao Bound Example

► ...will use the Hessian version of \mathbf{J} . Thus,

$$- \mathbb{E}_{\mathbf{y}} \left[\frac{\partial^2}{\partial m^2} \ln [f_{\mathbf{y}|m,v}(\mathbf{y}|m,v)] \right] = - \mathbb{E}_{\mathbf{y}} \left[-\frac{N}{v} \right] = \frac{N}{v} \quad (38)$$

$$\begin{aligned} - \mathbb{E}_{\mathbf{y}} \left[\frac{\partial^2}{\partial v^2} \ln [f_{\mathbf{y}|m,v}(\mathbf{y}|m,v)] \right] &= - \mathbb{E}_{\mathbf{y}} \left[+\frac{N}{2v^2} - \frac{1}{v^3} \sum_{k=1}^N (y_k - m)^2 \right] = \\ &= +\frac{N}{2v^2} \end{aligned} \quad (39)$$

$$- \mathbb{E}_{\mathbf{y}} \left[\frac{\partial^2}{\partial m \partial v} \ln [f_{\mathbf{y}|m,v}(\mathbf{y}|m,v)] \right] = - \mathbb{E}_{\mathbf{y}} \left[-\frac{1}{v^2} \sum_{k=1}^N (y_k - m) \right] = 0 \quad (40)$$

Cramer-Rao Bound Example

- ▶ Now, Fisher Information matrix and its inverse can be calculated:

$$\mathbf{J}(m, v) = \begin{bmatrix} \frac{N}{v} & 0 \\ 0 & \frac{N}{2v^2} \end{bmatrix} \Leftrightarrow \mathbf{J}^{-1}(m, v) = \begin{bmatrix} \frac{v}{N} & 0 \\ 0 & \frac{2v^2}{N} \end{bmatrix} \quad (41)$$

- ▶ Therefore, for any unbiased estimate of m, v ,

$$\mathbb{E} [(m - \hat{m}(\mathbf{y}))^2] \geq \frac{v}{N} \quad (42)$$

$$\mathbb{E} [(v - \hat{v}(\mathbf{y}))^2] \geq \frac{2v^2}{N} \quad (43)$$

Schur Complement Properties

- ▶ Before offering the derivation, we first list some basic properties. For any symmetric matrix \mathbf{M} of the following form:¹

$$\mathbf{M} = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^T & \mathbf{C} \end{bmatrix} \quad (44)$$

- ▶ If \mathbf{C} is invertible then:

1. $\mathbf{M} > \mathbf{0}$ iff $\mathbf{C} > \mathbf{0}$ and $\mathbf{A} - \mathbf{B}\mathbf{C}^{-1}\mathbf{B}^T > \mathbf{0}$
2. For $\mathbf{C} > \mathbf{0}$: $\mathbf{M} \geq \mathbf{0} \Leftrightarrow \mathbf{A} - \mathbf{B}\mathbf{C}^{-1}\mathbf{B}^T \geq \mathbf{0}$
3. Schur complement

$$\mathbf{M}|\mathbf{C} \triangleq \mathbf{A} - \mathbf{B}\mathbf{C}^{-1}\mathbf{B}^T \Rightarrow \det(\mathbf{M}) = \det(\mathbf{M}|\mathbf{C}) \det(\mathbf{C})$$

- ▶ If \mathbf{A} is invertible then:

1. $\mathbf{M} > \mathbf{0}$ iff $\mathbf{A} > \mathbf{0}$ and $\mathbf{C} - \mathbf{B}^T\mathbf{A}^{-1}\mathbf{B} > \mathbf{0}$
2. For $\mathbf{A} > \mathbf{0}$: $\mathbf{M} \geq \mathbf{0} \Leftrightarrow \mathbf{C} - \mathbf{B}^T\mathbf{A}^{-1}\mathbf{B} \geq \mathbf{0}$
3. Schur complement

$$\mathbf{M}|\mathbf{A} \triangleq \mathbf{C} - \mathbf{B}^T\mathbf{A}^{-1}\mathbf{B} \Rightarrow \det(\mathbf{M}) = \det(\mathbf{M}|\mathbf{A}) \det(\mathbf{A})$$

¹matrix inequality in the positive semi-definite sense:

$$\mathbf{A} \geq \mathbf{0} \Leftrightarrow \mathbf{z}^T \mathbf{A} \mathbf{z} \geq 0, \forall \mathbf{z}$$

Cramer-Rao Bound Derivation

- Proof: First, we show the two equivalent forms of the Fisher $m \times m$ matrix $\mathbf{J}(\mathbf{x})$:

$$\int f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})d\mathbf{y} = 1 \quad (45)$$

$$\nabla_x \ln [f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})] = \frac{1}{f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})} \nabla_x f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) \Rightarrow \quad (46)$$

$$\nabla_x^T f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) = \nabla_x^T \ln [f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})] f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) \quad (47)$$

$$\stackrel{(45)}{\Rightarrow} \int \nabla_x^T f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})d\mathbf{y} \stackrel{(47)}{=} \int \nabla_x^T \ln [f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})] f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})d\mathbf{y} = \mathbf{0}$$

$$\stackrel{\nabla_x}{\Rightarrow} \int \nabla_x \nabla_x^T \ln [f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})] f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})d\mathbf{y} + \int \nabla_x \ln [f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})] \nabla_x^T \ln [f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})] f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})d\mathbf{y} = \mathbf{0} \quad (48)$$

$$\begin{aligned} \Rightarrow \mathbf{J}_{ij}(\mathbf{x}) &= \mathbb{E}_{\mathbf{y}} \left[\frac{\partial}{\partial x_i} \ln [f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})] \frac{\partial}{\partial x_j} \ln [f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})] \right] \\ &= -\mathbb{E}_{\mathbf{y}} \left[\frac{\partial^2}{\partial x_i \partial x_j} \ln [f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})] \right] \end{aligned} \quad (49)$$

Cramer-Rao Bound Derivation

- ▶ Next, we define the $2m \times 1$ vector \mathbf{z} , corresponding positive semi-definite matrix \mathbf{C}_z and bias $\mathbf{b}(\mathbf{x})$:

$$\mathbf{z} = \begin{bmatrix} \underbrace{\mathbf{x} - \hat{\mathbf{x}}(\mathbf{y})}_{\mathbf{e}} - \mathbf{b}(\mathbf{x}) \\ \nabla_x \ln f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) \end{bmatrix}, \mathbf{C}_z = \begin{bmatrix} \mathbf{C}_{11} & \mathbf{C}_{12} \\ \mathbf{C}_{21} & \mathbf{C}_{22} \end{bmatrix} = \mathbb{E} [\mathbf{z}\mathbf{z}^T] \geq 0 \quad (50)$$

- ▶ The following hold:

1. $\mathbf{C}_{11} = \mathbb{E} [(\mathbf{e} - \mathbf{x})(\mathbf{e} - \mathbf{x})^T] = \mathbf{C}_E - \mathbf{b}(\mathbf{x})\mathbf{b}^T(\mathbf{x})$.
2. $\mathbf{C}_{22} = \mathbf{J}(\mathbf{x})$.
3. $\mathbf{C}_{12} = \mathbf{C}_{21}^T = \mathbb{E} [(\mathbf{x} - \hat{\mathbf{x}} - \mathbf{b}(\mathbf{x}))\nabla_x^T \ln f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})] = -\mathbf{I}_m + \nabla_x^T \mathbf{b}(\mathbf{x})$.

- ▶ The proof for 3. above follows.

Cramer-Rao Bound Derivation

► Next, we show $\mathbf{C}_{12} = \mathbf{C}_{21}^T = -\mathbf{I}_m + \nabla_x^T \mathbf{b}(\mathbf{x})$.

$$\mathbf{C}_{12} = \mathbf{C}_{21}^T = \mathbb{E} [(\mathbf{x} - \hat{\mathbf{x}} - \mathbf{b}(\mathbf{x})) \nabla_x^T \ln f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})] \quad (51)$$

$$= \int_{\mathbf{y}} (\mathbf{x} - \hat{\mathbf{x}} - \mathbf{b}(\mathbf{x})) \underbrace{\nabla_x^T \ln [f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})] f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})}_{\nabla_x^T f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})} d\mathbf{y} \quad (52)$$

$$\mathbb{E} [(\mathbf{x} - \hat{\mathbf{x}}(\mathbf{y}) - \mathbf{b}(\mathbf{x}))] = \mathbf{0} \quad (53)$$

$$\Rightarrow \int_{\mathbf{y}} (\mathbf{x} - \hat{\mathbf{x}}(\mathbf{y}) - \mathbf{b}(\mathbf{x})) f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) d\mathbf{y} = \mathbf{0} \quad (54)$$

$$\stackrel{\nabla_x^T}{\Rightarrow} \underbrace{\int_{\mathbf{y}} (\mathbf{x} - \hat{\mathbf{x}}(\mathbf{y}) - \mathbf{b}(\mathbf{x})) \nabla_x^T f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) d\mathbf{y}}_{\mathbf{C}_{12}} +$$

$$+ \int_{\mathbf{y}} (\mathbf{I}_m - \nabla_x^T \mathbf{b}(\mathbf{x})) f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) d\mathbf{y} = \mathbf{0} \quad (55)$$

$$\Rightarrow \mathbf{C}_{12} = -\mathbf{I}_m + \nabla_x^T \mathbf{b}(\mathbf{x}) \quad (56)$$

Cramer-Rao Bound Derivation

► In summary:

1. $\mathbf{C}_{11} = \mathbb{E}[(\mathbf{e} - \mathbf{x})(\mathbf{e} - \mathbf{x})^T] = \mathbf{C}_E - \mathbf{b}(\mathbf{x})\mathbf{b}^T(\mathbf{x})$.
2. $\mathbf{C}_{22} = \mathbf{J}(\mathbf{x})$.
3. $\mathbf{C}_{12} = -\mathbf{I}_m + \nabla_x^T \mathbf{b}(\mathbf{x})$.

$$\mathbf{z} = \begin{bmatrix} \underbrace{\mathbf{x} - \hat{\mathbf{x}}(\mathbf{y})}_{\mathbf{e}} - \mathbf{b}(\mathbf{x}) \\ \nabla_x \ln f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) \end{bmatrix}, \mathbf{C}_z = \begin{bmatrix} \mathbf{C}_{11} & \mathbf{C}_{12} \\ \mathbf{C}_{21} & \mathbf{C}_{22} \end{bmatrix} = \mathbb{E}[\mathbf{z}\mathbf{z}^T] \geq 0 \quad (57)$$

- We assume that positive semi-definite $\mathbf{J}(\mathbf{x})$ is invertible, i.e., it is positive definite. We also know that \mathbf{C}_z is positive semi-definite.
- Thus, the Schur complement is also positive semi-definite:

$$\mathbf{C}_{11} - \mathbf{C}_{12}\mathbf{C}_{22}^{-1}\mathbf{C}_{21}^T \geq 0 \quad (58)$$

Cramer-Rao Bound Derivation

► In summary:

1. $\mathbf{C}_{11} = \mathbb{E}[(\mathbf{e} - \mathbf{x})(\mathbf{e} - \mathbf{x})^T] = \mathbf{C}_E - \mathbf{b}(\mathbf{x})\mathbf{b}^T(\mathbf{x})$.
2. $\mathbf{C}_{22} = \mathbf{J}(\mathbf{x})$.
3. $\mathbf{C}_{12} = -\mathbf{I}_m + \nabla_x^T \mathbf{b}(\mathbf{x})$.

$$\mathbf{C}_{11} - \mathbf{C}_{12}\mathbf{C}_{22}^{-1}\mathbf{C}_{21}^T \geq \mathbf{0} \Rightarrow \quad (59)$$

$$\mathbf{C}_E - \mathbf{b}(\mathbf{x})\mathbf{b}^T(\mathbf{x}) - (\mathbf{I}_m - \nabla_x^T \mathbf{b}(\mathbf{x})) \mathbf{J}^{-1}(\mathbf{x}) (\mathbf{I}_m - \nabla_x^T \mathbf{b}(\mathbf{x}))^T \geq \mathbf{0} \quad (60)$$

► For unbiased estimator, i.e., $\mathbf{b}(\mathbf{x}) = \mathbf{0}$, the above is simplified to:

$$\mathbf{C}_E - \mathbf{J}^{-1}(\mathbf{x}) \geq \mathbf{0}$$

which completes the proof, if we consider that the diagonal elements of a positive semi-definite matrix are non-negative. ■

Existence of Efficient Estimator - Some Properties

- ▶ From Cramer-Rao proof, the following positive semi-definite matrix was utilized:

$$\mathbf{C}_z = \begin{bmatrix} \mathbf{C}_{11} = \mathbf{C}_E & \mathbf{C}_{12} \\ \mathbf{C}_{21} & \mathbf{C}_{22} = \mathbf{J}(\mathbf{x}) \end{bmatrix} = \mathbb{E} [\mathbf{z}\mathbf{z}^T] \geq 0 \Rightarrow \quad (61)$$

$$\mathbf{C}_z^{-1} = \begin{bmatrix} \mathbf{S}^{-1} & \mathbf{B} \\ \mathbf{C} & \mathbf{D} \end{bmatrix}, \quad \mathbf{z} = \begin{bmatrix} \mathbf{x} - \hat{\mathbf{x}}(\mathbf{y}) - \mathbf{b} \\ \nabla_x \ln f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) \end{bmatrix}, \quad (62)$$

- ▶ The following properties hold:

1. $\mathbf{S} = \mathbf{C}_{11} - \mathbf{C}_{12}\mathbf{C}_{22}^{-1}\mathbf{C}_{21}$.
2. $\mathbf{C}_z > 0 \Rightarrow \mathbf{S} > 0$.
3. $\mathbf{C}_z \geq 0 \Rightarrow \mathbf{S} \geq 0$.
4. $\mathbf{S} = \mathbf{0} \Rightarrow 2m \times 2m$ matrix \mathbf{C}_z is of rank m .²

² $\mathbf{A} \geq 0$ in the positive-semi definite sense, i.e., $\mathbf{A} \geq 0 \Leftrightarrow \mathbf{z}^T \mathbf{A} \mathbf{z} \geq 0$

- ▶ From property 3 in the previous slide, it follows for any biased estimator, with bias $\mathbf{b}(\mathbf{x})$:

$$\mathbf{S} = \mathbf{C}_E - \mathbf{b}(\mathbf{x})\mathbf{b}(\mathbf{x})^T - (\mathbf{I}_m - \nabla_x^T \mathbf{b}(\mathbf{x})) \mathbf{J}^{-1} (\mathbf{I}_m - \nabla_x^T \mathbf{b}(\mathbf{x}))^T \geq 0 \quad (63)$$

- ▶ For an unbiased estimator (i.e., $\mathbf{b}(\mathbf{x}) = \mathbf{0}$), the above leads to the Cramer-Rao bound:

$$\mathbf{S} \geq \mathbf{0} \Leftrightarrow \mathbf{C}_E \geq \mathbf{J}^{-1}. \quad (64)$$

- ▶ Efficient estimator is the unbiased estimator for which $\mathbf{C}_E \equiv \mathbf{J}^{-1}$.
- ▶ Thus, for an efficient estimator it holds that $\mathbf{S} = \mathbf{0}$; from property 4 in the previous slide, matrix \mathbf{C}_z is of rank m for an efficient estimator.

Efficient Estimator

- ▶ For an efficient estimator it holds that $\mathbf{S} = \mathbf{0}$; from property 4 in the previous slide, matrix \mathbf{C}_z is of rank m for an efficient estimator.
- ▶ ...the above means that m rows can be written as a linear combination of the other m rows; having in mind the definition of \mathbf{C}_z and the definition of \mathbf{z} , the above can be stated as follows:

$$\mathbf{x} - \hat{\mathbf{x}}(\mathbf{y}) = \mathbf{M} \nabla_x \ln f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) \Rightarrow \quad (65)$$

$$(\mathbf{x} - \hat{\mathbf{x}}(\mathbf{y})) \nabla_x^T \ln f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) = \mathbf{M} \nabla_x \ln f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) \nabla_x^T \ln f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) \Rightarrow$$

$$\mathbb{E}_{\mathbf{y}} [(\mathbf{x} - \hat{\mathbf{x}}(\mathbf{y})) \nabla_x^T \ln f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})] =$$

$$\mathbf{M} \mathbb{E}_{\mathbf{y}} [\nabla_x \ln f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) \nabla_x^T \ln f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x})]$$

$$\Rightarrow \mathbf{C}_{12} = \mathbf{M} \mathbf{J}$$

$$\Rightarrow \mathbf{M} = \mathbf{C}_{12} \mathbf{J}^{-1} = (-\mathbf{I}_m + \nabla_x^T \mathbf{b}(\mathbf{x})) \mathbf{J}^{-1} = -\mathbf{J}^{-1} \quad (66)$$

$$\Rightarrow \hat{\mathbf{x}}(\mathbf{y}) = \mathbf{x} - \mathbf{M} \nabla_x \ln f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) \quad (67)$$

$$\Rightarrow \hat{\mathbf{x}}(\mathbf{y}) = \mathbf{x} + \mathbf{J}^{-1} \nabla_x \ln f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) \quad (68)$$

Efficient Estimator Existence

- ▶ Thus, an efficient estimator has the following form:

$$\Rightarrow \hat{\mathbf{x}}(\mathbf{y}) = \mathbf{x} + \mathbf{J}^{-1} \nabla_x \ln f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) \quad (69)$$

- ▶ The left-hand side (LHS) is independent of \mathbf{x} ; thus, an efficient estimator exists iff the right-hand side (RHS) of the above equation is independent of \mathbf{x} .
- ▶ Notice that $\nabla_x \ln f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}|\mathbf{x}) = \mathbf{0} \Rightarrow \hat{\mathbf{x}}(\mathbf{y}) = \mathbf{x}$.
- ▶ Remarks follow:
 1. If an efficient exists, it must be a stationary point of the likelihood function; if there is only one such point, it must be the ML estimator.
 2. If the likelihood function has a single maximum and the estimator is efficient, it must be the ML estimator.
 3. ...the above does not mean that all ML estimators are efficient... they may not be!
- ▶ ...will see an example at next lecture.

- [1] Bernard C. Levy, Principles of Signal Detection and Parameter Estimation, Springer 2008.
- [2] Instructor notes.

Thank you!



Detection & Estimation Theory: Lectures 17-18

Prof. Aggelos Bletsas

Technical University of Crete
School of Electrical and Computer Engineering



Operational Programme
**Human Resources Development,
Education and Lifelong Learning**
Co-financed by Greece and the European Union





“Support for the Internationalization of Higher Education, School of Electrical and Computer Engineering, Technical University of Crete” (MIS 5150766), under the call for proposals “Support of the Internationalization of Higher Education - Technical University of Crete” (EDULLL 153).

The project is co-financed by Greece and the European Union (European Social Fund- ESF) through the Operational Programme “Human Resources Development, Education and Lifelong Learning 2014-2020.



Operational Programme
**Human Resources Development,
Education and Lifelong Learning**
Co-financed by Greece and the European Union



- ML Estimates and 1-1 Functions
- Phase Estimation Example
 - Geometric interpretation
- Sufficient Statistic
- UMVU Estimator
- Complete Sufficient Statistic
- The RBLT Theorem (& Proof)

ML Estimates and 1-1 Functions

- Suppose $\hat{\mathbf{x}}_{\text{ML}}(\mathbf{y})$ and $\mathbf{z} = \mathbf{g}(\mathbf{x})$ with $\mathbf{x} = \mathbf{g}^{-1}(\mathbf{z})$, i.e., $\mathbf{g}(\mathbf{x}) = \mathbf{z}$ is a “1-1” mapping.
- Thus, $f_{\mathbf{y}}(\mathbf{y}|\mathbf{z}) = f_{\mathbf{y}}(\mathbf{y}|\underbrace{\mathbf{g}^{-1}(\mathbf{z})}_{\mathbf{x}})$
- ▶ Then if $\hat{\mathbf{x}}_{\text{ML}} = \mathbf{g}^{-1}(\hat{\mathbf{z}}) \Rightarrow \hat{\mathbf{z}}_{\text{ML}} = \mathbf{g}(\hat{\mathbf{x}}_{\text{ML}})$.
- ▶ However, the transformation does not preserve unbiasedness or efficiency.

Phase estimation example

Now let's examine a phase estimation example:

$$\mathbf{y} = \begin{bmatrix} y_c \\ y_s \end{bmatrix} = A \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix} + \mathbf{v}, \quad \mathbf{v} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_2),$$

where σ and A are known and θ is unknown.

$$\text{Thus, } \mathbf{y} \sim N\left(\overbrace{A \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix}}^{\boldsymbol{\mu}}, \sigma^2 \mathbf{I}_2\right).$$

So,

$$\begin{aligned} f_{\mathbf{y}}(\mathbf{y}|\theta) &= \frac{1}{\sqrt{(2\pi)^2 \sigma^4}} e^{-\frac{1}{2\sigma^2} \|\mathbf{y} - \boldsymbol{\mu}\|^2} \\ \Rightarrow \ln f_{\mathbf{y}}(\mathbf{y}|\theta) &= -\frac{1}{2} \ln (2\pi)^2 \sigma^4 - \frac{1}{2\sigma^2} [(y_c - A \cos \theta)^2 + (y_s - A \sin \theta)^2] \\ \Rightarrow L(\mathbf{y}) &\triangleq \ln f_{\mathbf{y}}(\mathbf{y}|\theta) \\ &= -\ln 2\pi\sigma^2 - \frac{1}{2\sigma^2} (y_c^2 + y_s^2 + A^2) + \frac{1}{\sigma^2} (Ay_c \cos \theta + Ay_s \sin \theta) \end{aligned}$$

Phase estimation example (cont.)

$$\begin{aligned}\frac{\partial L}{\partial \theta} &= \frac{A}{\sigma^2}(-y_c \sin \theta + y_s \cos \theta) = 0 \\ \Rightarrow y_c \sin \theta &= y_s \cos \theta \Rightarrow \tan \theta = \frac{y_s}{y_c} \\ &\Rightarrow \hat{\theta}_{\text{ML}} = \tan^{-1} \frac{y_s}{y_c}\end{aligned}$$

$$\begin{aligned}\frac{\partial^2 L}{\partial \theta^2} &= \frac{A}{\sigma^2}(-y_c \cos \theta - y_s \sin \theta) \\ &= -\frac{A}{\sigma^2}(y_c \cos \theta + y_s \sin \theta)\end{aligned}$$

$$\begin{aligned}d(\theta) &= -\mathbb{E}\left[\frac{\partial^2 L}{\partial \theta^2}\right] = \frac{A}{\sigma^2}\mathbb{E}[y_c \cos \theta + y_s \sin \theta] \\ &= \frac{A}{\sigma^2}\mathbb{E}[A \cos^2 \theta + A \sin^2 \theta] = \frac{A^2}{\sigma^2}\end{aligned}$$

Phase estimation example (cont.)

$$\mathbb{E}[(\theta - \hat{\theta}(\mathbf{y}))^2] \geq \mathbf{J}^{-1}(\theta) = \frac{\sigma^2}{A^2} \simeq \frac{1}{\text{SNR}}$$

- $\mathbb{E}[\hat{\theta}_{\text{ML}}(\mathbf{y})] = ?$
- Is $\hat{\theta}_{\text{ML}}(\mathbf{y})$ efficient?

$$\text{Set } \left. \begin{array}{l} y_c = r \cos \phi \\ y_s = r \sin \phi \end{array} \right\} \hat{\theta}_{\text{ML}} = \tan^{-1} \frac{y_s}{y_c} = \phi$$

$$\left. \begin{array}{l} r > 0, \quad y_c^2 + y_s^2 = r^2 \\ \tan^{-1} \frac{y_s}{y_c} = \phi \end{array} \right\} \begin{bmatrix} r \\ \phi \end{bmatrix} \leftarrow \begin{bmatrix} y_c \\ y_s \end{bmatrix}$$

Phase estimation example (cont.)

$$\text{Jacobian} = \begin{bmatrix} \nabla_{\mathbf{y}}^T r \\ \nabla_{\mathbf{y}}^T \phi \end{bmatrix} = \begin{bmatrix} \frac{\partial r}{\partial y_c} & \frac{\partial r}{\partial y_s} \\ \frac{\partial \phi}{\partial y_c} & \frac{\partial \phi}{\partial y_s} \end{bmatrix} = \begin{bmatrix} \frac{y_c}{\sqrt{y_c^2 + y_s^2}} & \frac{y_s}{\sqrt{y_c^2 + y_s^2}} \\ -y_s \frac{1}{y_c^2} & \frac{1}{y_c} \end{bmatrix}$$

So,

$$|\det(\text{Jacobian})| = \left| \frac{1}{\sqrt{y_c^2 + y_s^2} \cdot \left(1 + \left(\frac{y_s}{y_c}\right)^2\right)} + \frac{\frac{y_s^2}{y_c^2}}{\left(1 + \left(\frac{y_s}{y_c}\right)^2\right) \cdot \sqrt{y_c^2 + y_s^2}} \right| = \frac{1}{\sqrt{y_c^2 + y_s^2}}$$

Phase estimation example (cont.)

$$\begin{aligned}\text{Hence, } f_{r,\phi}(r, \phi|\theta) &= \frac{f_{\mathbf{y}}(\mathbf{y})}{|\text{Jacobian}|} \Bigg|_{\substack{y_c = r \cos \phi \\ y_s = r \sin \phi}} \\ &= \frac{1}{\sqrt{(2\pi)^2(\sigma^2)^2}} e^{-\frac{1}{2\sigma^2} \|\mathbf{y} - \boldsymbol{\mu}\|^2} \\ &= \frac{1}{\sqrt{y_c^2 + y_s^2}} \\ &= \frac{\sqrt{y_c^2 + y_s^2}}{2\pi\sigma^2} e^{-\frac{1}{2\sigma^2} [y_c^2 + y_s^2 + A^2 - 2Ay_c \cos \theta - 2Ay_s \sin \theta]} \\ &= \frac{r}{2\pi\sigma^2} \underbrace{e^{-\frac{1}{2\sigma^2} [r^2 + A^2 - 2Ar \cos(\phi - \theta)]}}_{\text{even function of } \phi - \theta}\end{aligned}$$

$$\text{Thus, } \mathbb{E}[\hat{\theta}_{\text{ML}}] = \mathbb{E}[\phi] = \theta + \mathbb{E}[\cancel{\phi - \theta}] \overset{\theta}{=} \theta$$

Therefore $\hat{\theta}_{\text{ML}}$ is unbiased.

Phase estimation example (cont.)

Remember that

$$\hat{\mathbf{x}}(\mathbf{y}) = \mathbf{x} + \mathbf{J}^{-1}(\mathbf{x}) \nabla_{\mathbf{x}} \ln(f_{\mathbf{y}}(\mathbf{y}|\mathbf{x}))$$

unbiased efficient estimator exists if and only if the RHS does not depend on \mathbf{x} .

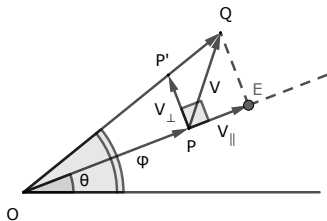
$$\hat{\mathbf{x}}(\mathbf{y}) = \theta + \frac{\sigma^2}{A^2} \frac{A}{\sigma^2} (-y_c \sin \theta + y_s \cos \theta) = \theta + \frac{r}{A} \sin(\phi - \theta)$$

As long as $A \simeq r \rightarrow \phi \simeq \theta$ and since $\hat{\theta}_{\text{ML}}(\mathbf{y}) = \phi \simeq \theta$ then,
 $\theta + \frac{r}{A}(\phi - \theta) = \phi = \hat{\theta}_{\text{ML}}(\mathbf{y})$

Thus,

$$\mathbb{E}[(\phi - \theta)^2] = \mathbb{E}[(\hat{\theta}_{\text{ML}} - \theta)^2] \simeq \mathbf{J}^{-1} = \frac{\sigma^2}{A^2} = \frac{1}{\text{SNR}}$$

Geometric interpretation



- $|\overrightarrow{OP}| = A$, $\overrightarrow{OQ} = \mathbf{y}$
- $\mathbf{V} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_2)$
- $\mathbf{V} = \mathbf{V}_{\parallel} + \mathbf{V}_{\perp}$ (orthogonal vector subspace)
- $|\mathbf{V}_{\perp}|$ is $\mathcal{N}(0, \sigma^2)$

$$\sin(\phi - \theta) \simeq \phi - \theta \simeq \frac{|\mathbf{V}_{\perp}|}{A} \sim \mathcal{N}(0, \frac{\sigma^2}{A^2})$$

$$\Rightarrow \mathbb{E}[\phi - \theta] = 0 \quad \text{and} \quad \mathbb{E}[(\phi - \theta)^2] = \frac{\sigma^2}{A^2}$$

Thus, efficient estimator for high SNR.

Sufficient Statistic

- \mathbf{y} is a sufficient statistic of \mathbf{x} if $f_{\mathbf{y}|\mathbf{s}}$ is independent of \mathbf{x} i.e., all information about \mathbf{x} has been “squeezed” in $f_{\mathbf{s}}(\mathbf{s}|\mathbf{x})$ and there is no leftover information about \mathbf{x} that could be extracted from $f_{\mathbf{y}|\mathbf{s}}$, which means that the latter is independent of \mathbf{x} .
- In practice, sufficient statistic $\mathbf{s}(\mathbf{y})$ can be directly found if $f_{\mathbf{y}}(\mathbf{y}|\mathbf{x})$ belongs to the exponential class of densities:

$$f_{\mathbf{y}}(\mathbf{y}|\mathbf{x}) = u(\mathbf{y}) \cdot \exp[\mathbf{x}^T \mathbf{s}(\mathbf{y}) - t(\mathbf{t})],$$

which includes discrete Poisson, Exponential and Gaussian distributions as special cases.

Example 1

- iid $\{y_k\}$'s, $y_k \sim \mathcal{N}(m, u)$, $\mathbf{y} = [y_1 \ y_2 \ \dots \ y_N]^\top$

$$\begin{aligned} f_{\mathbf{y}}(m, u) &= \prod_{k=1}^N \frac{1}{\sqrt{2\pi u}} \cdot e^{-\frac{1}{2u}(y_k - m)^2} \\ &= \left(\frac{1}{\sqrt{2\pi u}} \right)^N \cdot e^{-\frac{1}{2u} \sum_{k=1}^N (y_k - m)^2} \\ &= \frac{1}{(2\pi u)^{\frac{N}{2}}} \cdot e^{-\frac{1}{2u} \left[\sum_{k=1}^N y_k^2 + Nm^2 - 2m \sum_{k=1}^N y_k \right]} \\ &= \frac{1}{(2\pi u)^{\frac{N}{2}}} \cdot e^{-\frac{Nm^2}{2u}} \cdot e^{-\frac{1}{2u} \sum_{k=1}^N y_k^2 + \frac{m}{u} \sum_{k=1}^N y_k} \end{aligned}$$

Example 1 (cont.)

$$\begin{aligned} f_{\mathbf{y}}(m, u) &= \frac{1}{(2\pi u)^{\frac{N}{2}}} \cdot e^{-\frac{Nm^2}{2u}} \cdot e^{-\frac{1}{2u} \sum_{k=1}^N y_k^2 + \frac{m}{u} \sum_{k=1}^N y_k} \\ &= \frac{1}{(2\pi u)^{\frac{N}{2}}} \cdot e^{-\frac{Nm^2}{2u}} \cdot \exp\left(\underbrace{\begin{bmatrix} N\frac{m}{u} & -\frac{N}{2u} \end{bmatrix}}_{\mathbf{x}^\top} \underbrace{\begin{bmatrix} \frac{1}{N} \sum_{k=1}^N y_k \\ \frac{1}{N} \sum_{k=1}^N y_k^2 \end{bmatrix}}_{\mathbf{s}(\mathbf{y})} \right) \end{aligned}$$

$\mathbf{s}(\mathbf{y}) = \begin{bmatrix} s_1 \\ s_2 \end{bmatrix}$

i.e. $\mathbf{s}(\mathbf{y})$ is a sufficient statistic for estimating parameter \mathbf{x} .

Example 1 (cont.)

Sufficient statistic requires the definition of the unknown parameter

- if $u = \text{unknown}$ \rightarrow sufficient statistic is $s_1(\mathbf{y})$
- if $m = \text{unknown}$ \rightarrow sufficient statistic is $s_2(\mathbf{y})$

Notice that we have shown that:

$$\begin{aligned}\hat{m}_{\text{ML}} &= \frac{\sum y_k}{N} = s_1(\mathbf{y}) \\ \hat{u}_{\text{ML}} &= \frac{\sum (y_k - \hat{m}_{\text{ML}})^2}{N} = \frac{\sum y_k^2}{N} + \hat{m}_{\text{ML}}^2 - \frac{2\hat{m}_{\text{ML}}}{N} \sum y_k \\ &= \frac{\sum y_k^2}{N} - \hat{m}_{\text{ML}}^2 = s_2(\mathbf{y}) - s_1^2(\mathbf{y})\end{aligned}$$

Example 2

- iid $\{y_k\}$'s, $y_k \sim$ exponential with known parameter $1/\theta$,
 $\mathbf{y} = [y_1 \quad y_2 \quad \dots \quad y_N]^\top$

$$f_{\mathbf{y}|\theta} = \left(\frac{1}{\theta}\right)^N \cdot e^{-\frac{1}{\theta} \sum y_k} \prod_{k=1}^N u(y_k) \Rightarrow \begin{cases} x = \frac{1}{\theta} \\ s(\mathbf{y}) = \sum y_k \end{cases}$$

$$L(\mathbf{y}|\theta) = \ln [f(\mathbf{y}|\theta)] = -N \ln \theta - \frac{s(\mathbf{y})}{\theta}, \quad y_k \geq \emptyset$$

$$\frac{\partial}{\partial \theta} L(\mathbf{y}|\theta) = -\frac{N}{\theta} + \frac{s(\mathbf{y})}{\theta^2} = 0 \Rightarrow \hat{\theta}_{\text{ML}}(\mathbf{y}) = \frac{s(\mathbf{y})}{N}$$

$$\frac{\partial^2}{\partial \theta^2} L(\mathbf{y}|\theta) = \frac{N}{\theta^2} - \frac{2s(\mathbf{y})}{\theta^3} = \frac{1}{\theta^2} \left(N - \frac{2s(\mathbf{y})}{\theta} \right) = \frac{N}{\theta^2} \left(1 - \frac{2\hat{\theta}}{\theta} \right)$$

Example 2 (cont.)

$$\mathbb{E} [\hat{\theta}_{\text{ML}}(\mathbf{y})] = \frac{1}{N} \sum \mathbb{E}[y_k] = \frac{N\theta}{N} = \theta, \text{ (unbiased estimate)}$$

$$J(\theta) = -\mathbb{E} \left[\frac{\partial^2}{\partial \theta^2} L(\mathbf{y}|\theta) \right] = -\frac{N}{\theta^2} + \frac{2}{\theta^3} \mathbb{E}[s(\mathbf{y})] = -\frac{N}{\theta^2} + \frac{2}{\theta^3} N\theta = \frac{N}{\theta^2}$$

$$\text{Thus } \mathbb{E}[(\theta - \hat{\theta}_{\text{ML}})^2] \geq J^{-1}(\theta) = \frac{\theta^2}{N}$$

$$\begin{aligned} \mathbb{E}[(\theta - \hat{\theta}_{\text{ML}})^2] &= \theta^2 + \mathbb{E}[\hat{\theta}_{\text{ML}}^2] - 2\theta \mathbb{E}[\hat{\theta}_{\text{ML}}] \\ &= \mathbb{E}[\hat{\theta}_{\text{ML}}^2] - \theta^2 = \frac{1}{N^2} \mathbb{E}[(\sum y_i)^2] - \theta^2 \\ &= \frac{1}{N^2} \left[N \cdot \mathbb{E}[y_i^2] + 2\mathbb{E}^2[y_i] \cdot \binom{N}{2} \right] - \theta^2 \end{aligned}$$

Example 2 (cont.)

Since $\mathbb{E}[y_i] = \theta = 1/\lambda$

and $\mathbb{E}[y_i^2] - \mathbb{E}^2[y_i] = 1/\lambda^2 = \theta^2 \Rightarrow \mathbb{E}[y_i^2] = 2\theta^2$:

$$\begin{aligned}\mathbb{E}[(\theta - \hat{\theta}_{\text{ML}})^2] &= \frac{1}{N^2} \left[N \cdot 2\theta + 2\theta \cdot \frac{N(N-1)}{2} \right] - \theta^2 \\ &= \frac{1}{N} \left(2\theta^2 + (N-1)\theta^2 \right) - \theta^2 = \frac{\theta^2}{N} \equiv J^{-1}(\theta)\end{aligned}$$

Thus, in this case, the ML unbiased estimate is efficient and $s(\mathbf{y})$ for the parameter θ is $s(\mathbf{y}) = \sum y_k$.

Discussion on unbiased estimates

Set $J(\hat{\mathbf{x}}, \mathbf{x}) = \mathbb{E} \left[\|\mathbf{x} - \hat{\mathbf{x}}(\mathbf{y})\|_2^2 \right] = \text{MSE}$

“Uniform minimum variance unbiased estimate $\hat{\mathbf{x}}_{\text{UMVUE}}(\mathbf{y})$ ”

$J(\hat{\mathbf{x}}_{\text{UMVUE}}, \mathbf{x}) \leq J(\hat{\mathbf{x}}, \mathbf{x})$ for all other unbiased $\hat{\mathbf{x}}$ estimate.

- How do we find it?
 - search for ML estimate, see if it is unbiased and see if it is efficient.
 - if that approach fails, look for complete, sufficient statistic, as well as an unbiased estimator $\check{\mathbf{x}}(\mathbf{y})$.
- Apply RBLS theorem: if $\mathbf{s}(\mathbf{y})$ is complete sufficient statistic, then the estimate $\hat{\mathbf{x}}(\mathbf{s})$ (stemming from RBLS) is a UMVUE of \mathbf{x} .

Complete Sufficient Statistic

- What is a complete sufficient statistic?
- ▶ Let $\mathbf{s}(\mathbf{y})$ is a sufficient statistic for parameter \mathbf{x} .
- ▶ \mathbf{s} is complete if $\mathbb{E}[\mathbf{h}(\mathbf{s})] = \mathbf{0} \Leftrightarrow \mathbf{h}(\mathbf{s}) = \mathbf{0} \Leftrightarrow$
there is at most one unbiased
estimator of \mathbf{x} depending on \mathbf{s} only.

Note: if $\mathbf{h}(\mathbf{s}) = \mathbf{0} \Rightarrow \mathbb{E}[\mathbf{h}(\mathbf{s})] = \mathbf{0}$ is trivial.

Obviously, $\mathbb{E}[\mathbf{h}(\mathbf{s})] = \mathbf{0} \Rightarrow \mathbf{h}(\mathbf{s}) = \mathbf{0}$ is non-trivial

Complete Sufficient Statistic

- How do we find sufficient statistics which are complete?
In general it is hard.
- However, for $f_{\mathbf{y}}(\mathbf{y}|\mathbf{x}) = u(\mathbf{y}) \cdot \exp[\mathbf{x}^\top \mathbf{s}(\mathbf{y}) - t(\mathbf{x})]$, $\mathbf{s}(\mathbf{y})$ is complete sufficient statistic!
- ...the above includes Poisson, Exponential and Gaussian.

Rao-Blackwell-Lehmann-Sheffe (RBLs) Theorem

The Rao-Blackwell-Lehmann-Sheffe Theorem states that for an unbiased estimate $\check{\mathbf{x}}(\mathbf{y})$ of \mathbf{x} and a sufficient statistic $\mathbf{s}(\mathbf{y})$, the estimate can be improved:

$$\text{If } \mathbb{E}[\check{\mathbf{x}}(\mathbf{y})] = \mathbf{x}, \tag{1}$$

then $\hat{\mathbf{x}}(\mathbf{s}) = \mathbb{E}[\check{\mathbf{x}}(\mathbf{y})|\mathbf{s}]$ is unbiased

with $\hat{\mathbf{K}}(\mathbf{x}) \leq \check{\mathbf{K}}(\mathbf{x})$ i.e., their differ. is positive semi-definite (2)

$$\text{and } \check{\mathbf{K}}(\mathbf{x}) = \mathbb{E}[(\mathbf{x} - \check{\mathbf{x}})(\mathbf{x} - \check{\mathbf{x}})^\top]$$

$$\hat{\mathbf{K}}(\mathbf{x}) = \mathbb{E}[(\mathbf{x} - \hat{\mathbf{x}})(\mathbf{x} - \hat{\mathbf{x}})^\top]$$

- If \mathbf{s} is complete then $\hat{\mathbf{x}}(\mathbf{s})$ is a uniformly minimum-variance unbiased estimator (UMVUE).

Proof of RBLT Theorem

Proof of RBLT Theorem:

Let's start from the last one and assume that (1), (2) are true:

- If $\hat{\mathbf{s}}$ is complete, then there is at most one unbiased estimate of \mathbf{x} that depends on \mathbf{s} : $\hat{\mathbf{x}}(\mathbf{s})$.
- Suppose that there is a second $\hat{\mathbf{x}}_2(\mathbf{y})$ that achieves smaller $\hat{\mathbf{K}}_2(\mathbf{x}) < \hat{\mathbf{K}}(\mathbf{x})$.
- If we condition on \mathbf{s} , then we must get $\hat{\mathbf{x}}(\mathbf{s})$ with $\hat{\mathbf{K}}(\mathbf{x}) \leq \hat{\mathbf{K}}_2(\mathbf{x})$ which is a contradiction.
- Thus, there is no other estimator that minimizes the mean squared error, meaning that $\hat{\mathbf{x}}(\mathbf{s})$ is UMVUE.

Proof of RBLT Theorem (cont.)

Now let's prove that $\mathbb{E}_{\mathbf{s}}[\hat{\mathbf{x}}(\mathbf{s})] = \mathbf{x}$

$$\mathbb{E}_{\mathbf{s}}[\hat{\mathbf{x}}(\mathbf{s})] = \mathbb{E}_{\mathbf{s}} \left[\mathbb{E}_{\mathbf{y}|\mathbf{s}}[\check{\mathbf{x}}(\mathbf{y})|\mathbf{s}] \right] \triangleq \mathbb{E}[\check{\mathbf{x}}] = \mathbf{x}$$

law of iterated/repeated expectation

Finally we need to prove that $\hat{\mathbf{K}}_{\mathbf{x}} \leq \check{\mathbf{K}}_{\mathbf{x}}$:

$$\begin{aligned} \check{\mathbf{K}}_{\mathbf{x}} &= \mathbb{E}[(\mathbf{x} - \check{\mathbf{x}})(\mathbf{x} - \check{\mathbf{x}})^{\top}] \\ &= \mathbb{E}[(\mathbf{x} - \hat{\mathbf{x}} + \hat{\mathbf{x}} - \check{\mathbf{x}})(\mathbf{x} - \hat{\mathbf{x}} + \hat{\mathbf{x}} - \check{\mathbf{x}})^{\top}] \\ &= \hat{\mathbf{K}}_{\mathbf{x}} + \mathbb{E} \left[\underbrace{(\mathbf{x} - \hat{\mathbf{x}})(\hat{\mathbf{x}} - \check{\mathbf{x}})^{\top}}_{\mathbf{g}(\mathbf{y})} \right] + \mathbb{E} \left[\underbrace{(\hat{\mathbf{x}} - \check{\mathbf{x}})(\mathbf{x} - \hat{\mathbf{x}})^{\top}}_{\mathbf{g}(\mathbf{y})} \right] + \\ &\quad + \mathbb{E}[(\hat{\mathbf{x}} - \check{\mathbf{x}})(\hat{\mathbf{x}} - \check{\mathbf{x}})] \end{aligned} \tag{3}$$

* $\hat{\mathbf{x}} = \mathbb{E}[\check{\mathbf{x}}(\mathbf{y})|\mathbf{s}]$ is the MSE estimate of $\check{\mathbf{x}}(\mathbf{y})$ given \mathbf{s} thus $\mathbf{x} - \hat{\mathbf{x}}(\mathbf{y})$ is orthogonal to any function of \mathbf{s} (or \mathbf{y}).

Proof of RBLS Theorem (cont.)

Proof.

Thus from (3) we are left with:

$$\begin{aligned}\check{\mathbf{K}}_{\mathbf{x}} &= \hat{\mathbf{K}}_{\mathbf{x}} + \mathbb{E}[(\hat{\mathbf{x}} - \check{\mathbf{x}})(\hat{\mathbf{x}} - \check{\mathbf{x}})] \\ \Rightarrow \check{\mathbf{K}}_{\mathbf{x}} - \hat{\mathbf{K}}_{\mathbf{x}} &= \text{covariance matrix} \\ \Rightarrow \check{\mathbf{K}}_{\mathbf{x}} - \hat{\mathbf{K}}_{\mathbf{x}} &\geq 0 \text{ i.e., } \check{\mathbf{K}}_{\mathbf{x}} - \hat{\mathbf{K}}_{\mathbf{x}} \text{ is positive semi-definite.}\end{aligned}$$

□

Example

- iid s $y_k = [y_1 \ y_2 \ \dots \ y_N]^T$, $y_k \sim \frac{1}{\theta} \cdot e^{-\frac{y_k}{\theta}}$

$$f(\mathbf{y}|\theta) = \frac{1}{\theta^N} \exp\left(-\frac{s(\mathbf{y})}{\theta}\right) \prod_{k=1}^N u(y_k), \quad s(\mathbf{y}) = \sum_{k=1}^N y_k$$

$$\hat{\theta}_{\text{ML}}(\mathbf{y}) = \frac{s(\mathbf{y})}{N}$$

$$\check{\theta}(\mathbf{y}) = y_1 \text{ since } \mathbb{E}[y_1] = \theta \quad (\text{unbiased estimate})$$

$$s(\mathbf{y}) = \sum y_k = \underline{\text{complete}} \text{ sufficient statistic}$$

Thus, according to RBLS,

$$\hat{\theta}(s) = \mathbb{E}_{\mathbf{y}} [\mathbf{x}(\check{\theta})|\mathbf{s}] = \mathbb{E}_{y_1} [y_1|\mathbf{s}] \quad \text{is an UMVUE}$$

$$= \int y_1 \cdot f_{y_1|\mathbf{s}}(y_1|\mathbf{s}) dy_1$$

Example (cont.)

- Thus, I need to find the $f_{y_1|s}$, $s = \sum_{k=1}^N y_k$

$$\text{Set } \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_{N-1} \\ s \end{bmatrix}}_{\tilde{\mathbf{y}}} = \underbrace{\begin{bmatrix} 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 \\ \vdots & & & & \vdots \\ 0 & 0 & \dots & 1 & 0 \\ 1 & 1 & \dots & 1 & 1 \end{bmatrix}}_{\mathbf{A}} \underbrace{\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_{N-1} \\ y_N \end{bmatrix}}_{\mathbf{y}}$$

$$\text{Jacobian} = \begin{bmatrix} \nabla_{\mathbf{y}}^T y_1 \\ \nabla_{\mathbf{y}}^T y_2 \\ \vdots \\ \nabla_{\mathbf{y}}^T y_{N-1} \\ \nabla_{\mathbf{y}}^T s \end{bmatrix} = \begin{bmatrix} 1 & 0 & \dots & 0 & 0 \\ 0 & 1 & \dots & 0 & 0 \\ \vdots & & & & \vdots \\ 0 & 0 & \dots & 1 & 0 \\ 1 & 1 & \dots & 1 & 1 \end{bmatrix} = \mathbf{A}$$

* $\mathbf{A} \rightarrow$ upper diagonal $\Rightarrow \det(\mathbf{A}) =$ product of diagonal elements, and $\det(\mathbf{A}) = 1$

Example (cont.)

$$\nabla_{\mathbf{y}} = \begin{bmatrix} \frac{\partial}{\partial y_1} \\ \frac{\partial}{\partial y_2} \\ \vdots \\ \frac{\partial}{\partial y_N} \end{bmatrix} \quad \text{thus} \quad f_{\tilde{\mathbf{y}}}(\tilde{\mathbf{y}}) \triangleq f_{y_1, y_2, \dots, y_{N-1}, s}(\tilde{\mathbf{y}})$$
$$= \frac{f_{\mathbf{y}}(\mathbf{y})}{|\det(\text{Jacobian})|} \Big|_{\mathbf{y}=\mathbf{A}^{-1} \cdot \tilde{\mathbf{y}}}$$

$$s = \sum_{k=1}^{N-1} y_k + y_n \Rightarrow y_n = s - \sum_{k=1}^{N-1} y_k$$

$$\Rightarrow f_{y_1, y_2, \dots, y_{N-1}, s}(\tilde{\mathbf{y}}) = \frac{1}{\theta^N} \cdot \exp\left(-\frac{s}{\theta}\right) \left(\prod_{k=1}^{N-1} u(y_k)\right) u\left(s - \sum_{k=1}^{N-1} y_k\right)$$

Example (cont.)

$$\begin{aligned}f_{y_1, s} &= \int f_{y_1, y_2, \dots, y_{N-1}, s}(\tilde{\mathbf{y}}) dy_2 dy_3 \dots dy_{N-1} \\&= \frac{1}{\theta^N} \exp\left(-\frac{s}{\theta}\right) \int \prod_{k=1}^{N-2} u(y_k) \left[\int u(y_{N-1}) u\left(s - \sum_{k=1}^{N-1} y_k\right) dy_{N-1} \right] dy_2 dy_3 \dots dy_{N-2} \\&\stackrel{*}{=} \frac{1}{\theta^N} \exp\left(-\frac{s}{\theta}\right) \int \prod_{k=1}^{N-2} u(y_k) \cdot u\left(s - \sum_{k=1}^{N-2} y_k\right) \cdot \left(s - \sum_{k=1}^{N-2} y_k\right) dy_2 dy_3 \dots dy_{N-2} \\&= \dots = \frac{1}{\theta^N} \exp\left(-\frac{s}{\theta}\right) \frac{(s - y_1)^{N-2}}{(N-2)!} u(s - y_1) u(y_1)\end{aligned}$$

$s(y) : \sum_{k=1}^N y_k =$ sum of N identically distributed exponentials with parameters $\frac{1}{\theta}$ each $\Rightarrow s : \text{Gamma distribution with}$

$$f(s) = \Gamma(N, \theta) = s^{N-1} \frac{\exp(-s/\theta)}{\Gamma(N) \cdot \theta^N}, \quad \mathbb{E}[s] = N \cdot \theta, \quad \text{Var}(s) = N \cdot \theta^2$$

* $y_{N-1} \geq 0, s - \sum_{k=1}^{N-1} y_k \geq 0 \Rightarrow s - \sum_{k=1}^{N-2} y_k \geq y_{N-1} \geq 0$

Example (cont.)

- Alternatively, we could integrate $f_{y_1,s}$

$$\begin{aligned}\int f_{y_1,s}(y_1, s) dy_1 &= \frac{1}{\theta^N} e^{-s/\theta} \frac{1}{(N-2)!} \int_0^s (s-y_1)^{N-2} dy_1 = \\ &= \frac{1}{\theta^N} e^{-s/\theta} \frac{1}{(N-2)!} \left[-\frac{(s-y_1)^{N-1}}{N-1} \right]_0^s \\ &= \frac{1}{\theta^N} e^{-s/\theta} \frac{1}{(N-1)!} s^{N-1} \\ &= \frac{1}{\theta^N} e^{-s/\theta} \frac{1}{\Gamma(N)} s^{N-1}\end{aligned}$$

Example (cont.)

- Thus, $f_{y_1|s} = \frac{f_{y_1,s}}{f_s} = \frac{\cancel{\frac{1}{\theta^N} e^{-s/\theta}} \frac{(s-y_1)^{N-2}}{(N-2)!} u(s-y_1)u(y_1)}{\cancel{\frac{1}{\theta^N} \frac{1}{(N-1)!} s^{N-1} e^{-s/\theta}}}$
$$= \frac{(N-1)}{s^{N-1}} (s-y_1)^{N-2} u(s-y_1)u(y_1)$$

- and $\mathbb{E}[y_1|s] = \int y_1 \cdot f_{y_1|s}(y_1|s) dy_1$
$$= \int_0^s y_1 \frac{(N-1)}{s^{N-1}} (s-y_1)^{N-2} dy_1 = \dots = \frac{s}{N}$$

Thus, the RBLs procedure provides $\hat{\theta}(s) = \frac{s}{N} = \frac{\sum y_k}{N} \equiv \hat{\theta}_{\text{ML}}$, which can be shown to be both unbiased and efficient!

s complete \Rightarrow at most one unbiased estimate which is a function of s only.

Example (cont.)

- Trick : find one unbiased estimate which is a function of s only, provided that s is also a complete sufficient statistic. That will be the UMVUE!

Below, we offer an example where transformation offers ML estimate, without preserving unbiased property.

$$\text{if } \hat{\mathbf{x}}(\mathbf{y}) = \hat{\mathbf{x}}_{\text{ML}}(\mathbf{y}), \mathbf{z} = \mathbf{g}(\mathbf{x}) \quad 1-1 \quad \left[\mathbf{x} = \mathbf{g}^{-1}(\mathbf{z}) \right]$$

then $\hat{\mathbf{z}}_{\text{ML}}(\mathbf{y}) = \mathbf{g}(\hat{\mathbf{x}}(\mathbf{y}))$ “ML estimate is
parameterization independent”

Example (cont.)

$$\hat{\theta}_{\text{ML}} = \frac{s(\mathbf{y})}{N} = \frac{\sum y_k}{N}$$

$$\theta = \frac{1}{\lambda} \Rightarrow \lambda = \frac{1}{\theta} \Rightarrow \hat{\lambda}_{\text{ML}}(\mathbf{y}) = \frac{1}{\hat{\theta}_{\text{ML}}(\mathbf{y})} = \frac{N}{s}$$

$$\begin{aligned}\mathbb{E}\left[\frac{N}{s}\right] &= \int \frac{N}{s} f_s(s) ds = \frac{N}{\Gamma(N)\theta^N} \int_{-\infty}^{+\infty} s^{N-2} e^{-s/\theta} ds \\ &= \frac{N}{N-1} \frac{1}{\theta} \int \frac{1}{\theta^{N-1}} \frac{1}{\Gamma(N-1)} s^{N-2} e^{-s/\theta} ds \\ &= \frac{N}{N-1} \frac{1}{\theta} = \frac{N}{N-1} \lambda, \text{ thus the estimator is biased}\end{aligned}$$

Example (cont.)

- Set $\hat{\lambda}_0 = \frac{N-1}{N} \hat{\lambda}_{\text{ML}} = \frac{N-1}{s}$, $\mathbb{E}[\hat{\lambda}_0] = \lambda$ (unbiased)
- $\hat{\lambda}_0(s)$ is a function of s (which is a complete sufficient statistic) only, thus $\hat{\lambda}_0(s) = \hat{\lambda}_{\text{UMVUE}}$
- It can be easily shown that:

$$\left. \begin{array}{l} \text{CRLB: } \lambda^2/N \\ \mathbb{E}[(\lambda - \lambda_{\text{UMVUE}})^2] = \lambda^2/(N-2) \\ \mathbb{E}[(x - \hat{\lambda}_{\text{ML}})^2] = \frac{N+2}{(N-1)(N-2)} \lambda^2 \end{array} \right\} \begin{array}{l} \mathbb{E}[(\lambda - \lambda_{\text{ML}})^2] > \mathbb{E}[(\lambda - \lambda_{\text{UMVUE}})^2] > \\ > \text{CRLB} \end{array}$$

Important Remark: the conditional mean given the complete sufficient statistic should always give the same estimator.

Bernard C. Levy, Principles of Signal Detection and Parameter Estimation, Springer 2008.

Thank you!



Detection & Estimation Theory: Lecture 19

Prof. Aggelos Bletsas

Technical University of Crete
School of Electrical and Computer Engineering



Operational Programme
**Human Resources Development,
Education and Lifelong Learning**
Co-financed by Greece and the European Union





“Support for the Internationalization of Higher Education, School of Electrical and Computer Engineering, Technical University of Crete” (MIS 5150766), under the call for proposals “Support of the Internationalization of Higher Education - Technical University of Crete” (EDULLL 153).

The project is co-financed by Greece and the European Union (European Social Fund- ESF) through the Operational Programme “Human Resources Development, Education and Lifelong Learning 2014-2020.



Operational Programme
**Human Resources Development,
Education and Lifelong Learning**
Co-financed by Greece and the European Union



BLUE (Best Linear Unbiased Estimator)

- Problem Definition
 - Derivation
- Example 1
- Remarks
- Example 2
- Example 3

BLUE: Problem Definition

- ▶ Suppose that we want to estimate parameter vector $\boldsymbol{\theta}_{(p \times 1)}$ based on measurements $\mathbf{y}_{(N \times 1)}$ with linear estimator:

$$\hat{\boldsymbol{\theta}} = \mathbf{A}_{(p \times N)} \cdot \mathbf{y} \quad (1)$$

- ▶ We require unbiased estimator:

$$\mathbb{E}[\hat{\boldsymbol{\theta}}] = \mathbf{A} \cdot \mathbb{E}[\mathbf{y}] = \boldsymbol{\theta} \quad (2)$$

that can be achieved if and only if $\mathbb{E}[\mathbf{y}] = \mathbf{H}_{(N \times p)} \cdot \boldsymbol{\theta}_{(p \times 1)} \Rightarrow$

$$\mathbf{A}_{(p \times N)} \cdot \mathbf{H}_{(N \times p)} = \mathbf{I}_p \quad (3)$$

- ▶ $\text{rank}(\mathbf{I}_p) = p = \min(\text{rank}(\mathbf{A}), \text{rank}(\mathbf{H})) = \min(N, p)$, thus $N \geq p$

- ▶ Set $\mathbf{A} = \begin{bmatrix} \mathbf{a}_1^T \\ \mathbf{a}_2^T \\ \vdots \\ \mathbf{a}_p^T \end{bmatrix}$ (4) and $\mathbf{H} = [\mathbf{h}_1 \quad \mathbf{h}_2 \quad \cdots \quad \mathbf{h}_p]$ (5)

▶ from (1) and (4) \Rightarrow
$$\hat{\theta}_i = \mathbf{a}_i^T \cdot \mathbf{y} \quad (6)$$

▶ from (4), (5) and (3) $\Rightarrow \mathbf{a}_i^T \mathbf{h}_j = \delta_{ij}$

▶ $\text{var}(\hat{\theta}_i) = \mathbb{E} \left[\left(\hat{\theta}_i - \mathbb{E}[\hat{\theta}_i] \right)^2 \right] = \mathbb{E} \left[\left(\mathbf{a}_i^T \mathbf{y} - \mathbf{a}_i^T \mathbb{E}[\mathbf{y}] \right)^2 \right]$
 $= \mathbb{E} \left[\left[\mathbf{a}_i^T (\mathbf{y} - \mathbb{E}[\mathbf{y}]) \right]^2 \right]$
 $= \mathbb{E} \left[\mathbf{a}_i^T (\mathbf{y} - \mathbb{E}[\mathbf{y}]) (\mathbf{y} - \mathbb{E}[\mathbf{y}])^T \mathbf{a}_i \right] = \mathbf{a}_i^T \mathbf{K}_y \mathbf{a}_i$
 $\Rightarrow \text{var}(\hat{\theta}_i) = \mathbf{a}_i^T \mathbf{K}_y \mathbf{a}_i \quad (7)$

BLUE Derivation

- ▶ Minimize $\text{var}(\hat{\theta}_i) = \mathbf{a}_i^T \mathbf{K}_y \mathbf{a}_i$, for $i = 1, 2, \dots, p$ subject to the constraints $\mathbf{a}_i^T \mathbf{h}_j = \delta_{ij}$, $i, j \in \{1, 2, \dots, p\}$
- ▶ We have p constraints for each \mathbf{a}_i . Since each \mathbf{a}_i is free to assume any value, independently of the others, we actually have p separate minimization problems linked only by the constraints:

$$J_i = \mathbf{a}_i^T \mathbf{K}_y \mathbf{a}_i + \sum_{j=1}^p \left(\lambda_j^{(i)} (\mathbf{a}_i^T \mathbf{h}_j - \delta_{ij}) \right), \boldsymbol{\lambda}_i = \left[\lambda_1^{(i)} \lambda_2^{(i)} \dots \lambda_p^{(i)} \right]^T$$

$$\frac{\partial J_i}{\partial \mathbf{a}_i} = 2\mathbf{K}_y \mathbf{a}_i + \sum_{j=1}^p \lambda_j^{(i)} \mathbf{h}_j = 2\mathbf{K}_y \mathbf{a}_i + \mathbf{H} \boldsymbol{\lambda}_i = \mathbf{0}$$

$$\Rightarrow \frac{\partial J_i}{\partial \mathbf{a}_i} = \mathbf{0} \Rightarrow \mathbf{a}_i = -\frac{1}{2} \mathbf{K}_y^{-1} \mathbf{H} \boldsymbol{\lambda}_i \quad (8)$$

BLUE Derivation

- ▶ To find $\boldsymbol{\lambda}_i$, we need to exploit the constraints:

$$\mathbf{a}_i^T \cdot \mathbf{h}_j = \mathbf{h}_j^T \mathbf{a}_i = \delta_{ij}, \quad j = 1, 2, \dots, p$$

where $\delta_{ij} = 1$ for $i = j$ and $\delta_{ij} = 0$ for $i \neq j$.

- ▶ from (5) and the above $\Rightarrow \mathbf{H}^T \cdot \mathbf{a}_i = \mathbf{e}_i$, which is a vector with all zeros apart from position i , where it is one:

$$\begin{aligned} \mathbf{H}^T \cdot \mathbf{a}_i = \mathbf{e}_i &\Rightarrow \mathbf{H}^T \cdot \left(-\frac{1}{2} \mathbf{K}_y^{-1} \mathbf{H} \boldsymbol{\lambda}_i \right) = \mathbf{e}_i \\ &\Rightarrow -\frac{1}{2} \boldsymbol{\lambda}_i = \underbrace{\left(\mathbf{H}^T \mathbf{K}_y^{-1} \mathbf{H} \right)^{-1}}_{\text{assuming invertibility of } \mathbf{H}^T \mathbf{K}_y^{-1} \mathbf{H}} \end{aligned} \quad (9)$$

- ▶ from (8) and (9) \Rightarrow

$$\mathbf{a}_{i_{opt}} = \mathbf{K}_y^{-1} \mathbf{H} \left(\mathbf{H}^T \mathbf{K}_y^{-1} \mathbf{H} \right)^{-1} \mathbf{e}_i$$

$$\hat{\boldsymbol{\theta}} = \mathbf{A} \cdot \mathbf{y} = \begin{bmatrix} \mathbf{a}_1^T \\ \mathbf{a}_2^T \\ \vdots \\ \mathbf{a}_p^T \end{bmatrix} \cdot \mathbf{y} = \underbrace{\begin{bmatrix} \mathbf{e}_1^T \\ \mathbf{e}_2^T \\ \vdots \\ \mathbf{e}_p^T \end{bmatrix}}_{\mathbf{I}_p} \cdot \left(\mathbf{H}^T \mathbf{K}_y^{-1} \mathbf{H} \right)^{-1} \mathbf{H}^T \mathbf{K}_y^{-1} \cdot \mathbf{y}$$
$$\Rightarrow \hat{\boldsymbol{\theta}} = \left(\mathbf{H}^T \mathbf{K}_y^{-1} \mathbf{H} \right)^{-1} \mathbf{H}^T \mathbf{K}_y^{-1} \cdot \mathbf{y}$$

Example 1

- ▶ $\mathbf{y} = \mathbf{H}\boldsymbol{\theta} + \mathbf{w}$, \mathbf{w} has zero mean $\mathbb{E}[\mathbf{w}] = \mathbf{0}$ and $\mathbb{E}[\mathbf{w}\mathbf{w}^T] = \mathbf{K}_y$, so $\mathbb{E}[\mathbf{y}] = \mathbf{H}\boldsymbol{\theta}$
- ▶ $\hat{\boldsymbol{\theta}} - \mathbb{E}[\hat{\boldsymbol{\theta}}] = (\mathbf{H}^T \mathbf{K}_y^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{K}_y^{-1} \mathbf{y} - \mathbb{E}[\hat{\boldsymbol{\theta}}]$

$$\begin{aligned} &= \cancel{(\mathbf{H}^T \mathbf{K}_y^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{K}_y^{-1} \mathbf{H} \boldsymbol{\theta}} - \mathbb{E}[\hat{\boldsymbol{\theta}}] + \boldsymbol{\theta} \\ &\quad + (\mathbf{H}^T \mathbf{K}_y^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{K}_y^{-1} \mathbf{w} \\ &= (\mathbf{H}^T \mathbf{K}_y^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{K}_y^{-1} \mathbf{w} \end{aligned}$$

- ▶ Thus,

$$\begin{aligned} &\mathbb{E} \left[\left(\hat{\boldsymbol{\theta}} - \mathbb{E}[\hat{\boldsymbol{\theta}}] \right) \left(\hat{\boldsymbol{\theta}} - \mathbb{E}[\hat{\boldsymbol{\theta}}] \right)^T \right] = \\ &= (\mathbf{H}^T \mathbf{K}_y^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{K}_y^{-1} \underbrace{\mathbf{K}_y \mathbf{K}_y^{-1}}_{\mathbf{I}} \mathbf{H} (\mathbf{H}^T \mathbf{K}_y^{-1} \mathbf{H})^{-1} \\ &= (\mathbf{H}^T \mathbf{K}_y^{-1} \mathbf{H})^{-1} \equiv \mathbf{C}_{\hat{\boldsymbol{\theta}}} \\ &\quad \Rightarrow \text{var}(\hat{\theta}_i) = \left[(\mathbf{H}^T \mathbf{K}_y^{-1} \mathbf{H})^{-1} \right]_{ii} \end{aligned} \tag{10}$$

Example 1

- ▶ This could be also seen from

$$\hat{\theta}_i = \mathbf{a}_{i_{opt}}^T \cdot \mathbf{y} = \mathbf{e}_i^T \left(\mathbf{H}^T \mathbf{K}_y^{-1} \mathbf{H} \right)^{-1} \mathbf{H}^T \mathbf{K}_y^{-1} \cdot \mathbf{y}$$

and

$$\begin{aligned} \text{var} \hat{\theta}_i &= \mathbf{a}_{i_{opt}}^T \mathbf{K}_y \mathbf{a}_{i_{opt}} \\ &= \mathbf{e}_i^T \left(\mathbf{H}^T \mathbf{K}_y^{-1} \mathbf{H} \right)^{-1} \mathbf{H}^T \mathbf{K}_y^{-1} \mathbf{K}_y \mathbf{K}_y^{-1} \mathbf{H} \left(\mathbf{H}^T \mathbf{K}_y^{-1} \mathbf{H} \right)^{-1} \mathbf{e}_i \\ &= \mathbf{e}_i^T \left(\mathbf{H}^T \mathbf{K}_y^{-1} \mathbf{H} \right)^{-1} \mathbf{e}_i \\ &\stackrel{(7)}{=} \left[\left(\mathbf{H}^T \mathbf{K}_y^{-1} \mathbf{H} \right)^{-1} \right]_{ii} \end{aligned}$$

- ▶ *Remark 1:* MVUE for linear Gaussian Case \equiv BLUE, if $\mathbf{y} = \mathbf{H}\boldsymbol{\theta} + \mathbf{w}$, $\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$, then $\hat{\boldsymbol{\theta}} = \left(\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H}\right)^{-1} \mathbf{H}^T \mathbf{C}^{-1} \mathbf{y} \equiv$ BLUE is also UMVUE
- ▶ *Remark 2 (Gauss-Markov Theorem):* if the data are of the general linear model form $\mathbf{y} = \mathbf{H}\boldsymbol{\theta} + \mathbf{w}$, where \mathbf{H} is a known $N \times p$ matrix, $\boldsymbol{\theta}$ is a $p \times 1$ vector of parameters to be estimated, and \mathbf{w} is a $N \times 1$ noise vector with zero mean and covariance \mathbf{C}^1 , then BLUE of $\boldsymbol{\theta}$ is

$$\hat{\boldsymbol{\theta}} = \left(\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H}\right)^{-1} \mathbf{H}^T \mathbf{C}^{-1} \mathbf{y} \quad \text{and} \quad \mathbf{C}_{\hat{\boldsymbol{\theta}}} = \left(\mathbf{H}^T \mathbf{C}^{-1} \mathbf{H}\right)^{-1}$$

¹ N should be greater or equal than p ($N \geq p$)

Example 2

- $y[n] = A + w[n]$, $n = 0, 1, \dots, N - 1$: $w[n]$ white noise with variance σ^2 (not necessarily Gaussian)

$$\mathbb{E}[w[n]] = 0 \quad \text{and} \quad \mathbb{E}[w[n]w[n+m]] = \sigma^2\delta[m]$$

$$\mathbf{y} = \underbrace{\begin{bmatrix} y[0] \\ y[1] \\ \vdots \\ y[N-1] \end{bmatrix}}_{\mathbf{y}} = \underbrace{\begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}}_{\mathbf{H}_{N \times 1}} \cdot A = \underbrace{\begin{bmatrix} w[0] \\ w[1] \\ \vdots \\ w[N-1] \end{bmatrix}}_{\mathbf{w}}$$

$$\mathbb{E}[\mathbf{y}] = \underbrace{\begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix}}_{\mathbf{H}} \cdot A, \quad \mathbb{E}[\mathbf{w}] = 0 \quad \text{and} \quad \mathbb{E}[\mathbf{w}\mathbf{w}^T] = \sigma^2\mathbf{I}_N$$

Example 2

► Thus,

$$\begin{aligned}\hat{\theta} &= \left(\frac{N}{\sigma^2}\right)^{-1} \cdot \frac{1}{\sigma^2} \sum_{n=0}^{N-1} y[n] \\ &= \frac{1}{N} \sum_{n=1}^{N-1} y[n] = \hat{A}\end{aligned}$$

$$\mathbb{E} [A - \hat{A}] = \left(\frac{1}{\sigma^2} N\right)^{-1} = \frac{\sigma^2}{N}$$

Example 3

- ▶ We had seen that for y_k i.i.d. $\sim \underbrace{\mathcal{N}(\cancel{m}, \hat{u})}_{m \text{ known}} = \mathcal{N}(0, u)$

$$\hat{u}_{\text{ML}} = \frac{1}{N} \sum_{k=1}^N y_k^2 \quad \text{and} \quad \text{CRLB} = \frac{2u^2}{N}$$

- ▶ Is \hat{u}_{ML} efficient?
- ▶ What is the BLUE estimate of u ?

Example 3

▶ $\mathbb{E}[\hat{u}_{\text{ML}}] = \frac{1}{N}Nu = u$ unbiased

▶
$$\begin{aligned}\mathbb{E}[(\hat{u}_{\text{ML}} - u)^2] &= \mathbb{E}[\hat{u}_{\text{ML}}^2] - u^2 = \frac{1}{N^2} \left(\sum_{k=1}^N y_k^2 \right)^2 - u^2 \\ &= \frac{1}{N^2} \left(N\mathbb{E}[y_k^4] + \binom{N}{2} 2u^2 \right) - u^2 \\ &= \frac{1}{N^2} \left(N \cdot 3u^2 + \frac{N!}{2!(N-2)!} 2u^2 \right) - u^2 \\ &= \frac{3u^2}{N} + \frac{N-1}{N}u^2 - u^2 = \frac{2u^2}{N} \equiv \text{CRLB}\end{aligned}$$

▶ Thus, $\hat{u}_{\text{ML}}(\mathbf{y}) = \frac{1}{N} \sum_{k=1}^N y_k^2$ is efficient.

$${}^2 y \sim \mathcal{N}(0, u), \sigma_u = \sqrt{u} \Rightarrow$$

$$\mathbb{E}[y^n] = \begin{cases} (\sigma_u)^n \cdot 1 \cdot 3 \cdots (n-1), & n \text{ even,} \\ 0, & n \text{ odd,} \end{cases}$$

Example 3 - BLUE estimate

- ▶ Search for BLUE of u :

$$\hat{u}_{\text{BLUE}}(\mathbf{y}) = \mathbf{a}^T \mathbf{y} = \sum_{k=1}^N a_k y_k$$

$$\mathbb{E}[\hat{u}_{\text{BLUE}}] = \sum a_k \mathbb{E}[y_k] = \emptyset \neq u \quad \text{Biased}$$

- ▶ However, we can use $z_k = y_k^2$ (data transformation) and test BLUE on the transformed data:³

$$\hat{u}_{\text{BLUE}}(\mathbf{z}) = \sum a_k z_k = \sum_{k=1}^N a_k y_k^2 \Rightarrow \quad (11)$$

$$\mathbb{E}[\hat{u}_{\text{BLUE}}] = \underbrace{\sum_{k=1}^N a_k}_{\sum_{k=1}^N a_k = 1} u = u \quad (12)$$

³ z_k iid, $\sum a_k = 1 \Rightarrow a_k = \frac{1}{N}$

Bernard C. Levy, Principles of Signal Detection and Parameter Estimation, Springer 2008.

Thank you!



Detection & Estimation Theory: Lectures 20-21

Prof. Aggelos Bletsas

Technical University of Crete
School of Electrical and Computer Engineering



Operational Programme
**Human Resources Development,
Education and Lifelong Learning**
Co-financed by Greece and the European Union





“Support for the Internationalization of Higher Education, School of Electrical and Computer Engineering, Technical University of Crete” (MIS 5150766), under the call for proposals “Support of the Internationalization of Higher Education - Technical University of Crete” (EDULLL 153).

The project is co-financed by Greece and the European Union (European Social Fund- ESF) through the Operational Programme “Human Resources Development, Education and Lifelong Learning 2014-2020.



Operational Programme
**Human Resources Development,
Education and Lifelong Learning**
Co-financed by Greece and the European Union



- Introduction to Composite Hypothesis Testing and UMP/GLRT
- GLRT: Examples and Properties
- Asymptotic Optimality of the GLRT

Composite Hypothesis Testing

- ▶ Composite Hypothesis Testing: Problem of both Detection and Estimation!
- ▶ Problem definition:
 1. $\mathcal{H}_0 : \mathbf{y} \sim f_{\mathbf{y}}(\mathbf{y}|\mathbf{x}, \mathcal{H}_0), \mathbf{x} \in \mathcal{X}_0$
 $\mathcal{H}_1 : \mathbf{y} \sim f_{\mathbf{y}}(\mathbf{y}|\mathbf{x}, \mathcal{H}_1), \mathbf{x} \in \mathcal{X}_1$
 2. \mathbf{x} defines \mathcal{H}_j ; if $\mathcal{X}_0 \equiv \mathcal{X}_1$ there would be no way to distinguish between $\mathcal{H}_0, \mathcal{H}_1 \Rightarrow$ no way to estimate \mathbf{x}
- ▶ Example: $y[k] = As[k] + v[k], k \in \{1, \dots, N\}$
 1. A unknown, v WGN with $v \sim \mathcal{N}(0, \sigma^2)$, σ^2 unknown
 2.
$$\mathbf{y} = \begin{bmatrix} y[1] \\ y[2] \\ \vdots \\ y[N] \end{bmatrix} = A \begin{bmatrix} s[1] \\ s[2] \\ \vdots \\ s[N] \end{bmatrix} + \begin{bmatrix} v[1] \\ v[2] \\ \vdots \\ v[N] \end{bmatrix}$$
 3. $\mathcal{H}_0: A = 0, \sigma^2$ unknown $\Rightarrow \mathcal{X}_0 = \{(A, \sigma^2) : A = 0\}$
 $\mathcal{H}_1: A \neq 0, \sigma^2$ unknown $\Rightarrow \mathcal{X}_1 = \{(A, \sigma^2) : A \neq 0\}$
 4. σ^2 is common to both hypothesis, thus σ^2 does not play a role in determining which hypothesis holds (i.e., it is a "nuisance parameter").

Composite Hypothesis Testing

- ▶ Remark: In the above example $\mathcal{X}_0 \cap \mathcal{X}_1 = \emptyset$, even though σ^2 is common - remember $\mathcal{X} = (A, \sigma^2)$. If under $\mathcal{X}_0, \mathcal{X}_1$ the observation distribution is the same, then the detection problem cannot be solved, unless $\mathcal{X}_0 \cup \mathcal{X}_1 = \emptyset$.
- ▶ Special Case: \mathcal{H}_1 is composite but \mathcal{H}_0 is "simple". This means that the domain \mathcal{X}_0 reduces to a single point $\mathbf{x}_0 \rightarrow$ easier analysis than the case where both hypotheses are composite.
- ▶ Example: The previous example with σ^2 known instead
 1. $\mathcal{H}_0: \mathcal{X}_0: (A, \sigma^2): A = 0$ and σ^2 fixed and known
 2. set $\mathcal{Y} = \mathcal{Y}_0 \cup \mathcal{Y}_1$ ($\mathcal{Y}_0 \cap \mathcal{Y}_1 = \emptyset$)
 3. $\mathcal{Y}_j = \mathbf{y} : \delta(\mathbf{y}) = j, j \in 0, 1, \delta(\cdot)$ decision rule
 4. Probability of detection:¹

$$P_D(\delta, \mathbf{x}) = \Pr(\delta = 1 | \mathbf{x}, \mathcal{H}_1) = \int_{\mathcal{Y}_1} f(\mathbf{y} | \mathbf{x}, \mathcal{H}_1) d\mathbf{y}$$

5. Probability of false alarm:¹

$$P_F(\delta, \mathbf{x}) = \Pr(\delta = 1 | \mathbf{x}, \mathcal{H}_0) = \int_{\mathcal{Y}_1} f(\mathbf{y} | \mathbf{x}, \mathcal{H}_0) d\mathbf{y}$$

¹function of \mathbf{x}

Composite Hypothesis Testing

- ▶ When $P_D(\delta, \mathbf{x})$ is viewed as a function of \mathbf{x} , it is called the "power of the test"
- ▶ Neyman–Pearson approach is followed (even though Bayesian approach is also possible):
 1. set upper bound for probability of false alarm
 2. $\max_{x \in \mathcal{X}_0} P_F(S, \mathbf{x}) \geq a$ (1)
 3. a is called the size of the test
 4. Then, among all tests δ obeying Eq. (1), we say that δ_{UMP} is a uniformly most powerful (UMP) test if it satisfies:

$$P_D(\delta, \mathbf{x}) \leq P_D(\delta_{\text{UMP}}, \mathbf{x})$$

for all $\mathbf{x} \in \mathcal{X}_1$.

- I. very strong property - rarely we find UMP
- II. UMP test δ_{UMP} cannot depend on \mathbf{x}
- III. if \mathbf{x} is viewed as being fixed, δ_{UMP} must be the optimum test in the sense of Neyman-Pearson tests ($\max P_D$ for bounded P_F), so it must take the form of a LRT, possibly involving randomization.

Composite Hypothesis Testing - UMP

- ▶ Thus, from II) and III) we need to find LRT and then try to transform it in such a way that the parameter vector \mathbf{x} disappears from the test statistic. Then the threshold of the test is computed in such a way that the P_F upper bound is satisfied. If that is possible, then a UMP test exists!
- ▶ Example: The previous example rewritten:
 - ▶ $\mathcal{H}_1: y[k] = As[k] + w[k], A > 0$ unknown, $w[k] \sim \mathcal{N}(0, \sigma^2)$ (WGN)
 - ▶ $\mathcal{H}_0: y[k] = w[k], A \neq 0$
 - ▶ $k \in 1, 2, \dots, N$
 - ▶ Thus,
 - ▶ $\mathcal{H}_1: \mathbf{y} = A\mathbf{s} + \mathbf{w}, A > 0, \mathbf{w} \sim \mathcal{N}(0, \sigma^2 I_N)$
 - ▶ $\mathcal{H}_0: \mathbf{y} = \mathbf{w}$
 - ▶ case I: σ^2 known
 - ▶ $\mathcal{X}_1 = \{A > 0\}$ (or $A < 0$)(one-sided test), $\mathcal{X}_0 = \{A = 0\}$ "simple"

Composite Hypothesis Testing - UMP

▶ LRT: $L(\mathbf{y}|A) = \frac{f(\mathbf{y}|A>0, \mathcal{H}_1)}{f(\mathbf{y}|A=0, \mathcal{H}_0)} = \frac{f(\mathbf{y}|A>0)}{f(\mathbf{y}|A=0)} \stackrel{\mathcal{H}_1}{\geq} \tau^1$

▶ $f(\mathbf{y}|A > 0) = \frac{1}{\sqrt{(2\pi)^N (\sigma^2)^N}} e^{-\frac{1}{2\sigma^2} \|\mathbf{y} - A\mathbf{s}\|^2}$

▶ $\|\mathbf{y} - A\mathbf{s}\|^2 = (\mathbf{y} - A\mathbf{s})^T (\mathbf{y} - A\mathbf{s})$
 $= (\mathbf{y}^T - A\mathbf{s}^T) (\mathbf{y} - A\mathbf{s})$
 $= \|\mathbf{y}\|^2 - A\mathbf{y}^T \mathbf{s} - A\mathbf{s}^T \mathbf{y} + A^2 \|\mathbf{s}\|^2$
 $= \|\mathbf{y}\|^2 - 2A\mathbf{s}^T \mathbf{y} + A^2 \|\mathbf{s}\|^2$

▶ set $\|\mathbf{s}\|^2 = \mathbf{s}^T \mathbf{s} = \sum_{k=1}^N s^2[k] = E$

▶ $L(\mathbf{y}|A) = e^{-\frac{1}{2\sigma^2} (-2A\mathbf{s}^T \mathbf{y} + A^2 \|\mathbf{s}\|^2)} \Rightarrow$

$$\ln L(\mathbf{y}|A) = \frac{A}{\sigma^2} \mathbf{s}^T \mathbf{y} - \frac{A^2}{2\sigma^2} E \Rightarrow$$

$$\frac{A}{\sigma^2} \mathbf{s}^T \mathbf{y} - \frac{A^2 E}{2\sigma^2} \stackrel{\mathcal{H}_1}{\geq} \ln \tau$$

¹yet to be specified

Composite Hypothesis Testing - UMP

- ▶ Remember that we need a test, which does not depend on the unknown parameter(s). Thus, we need to get rid of A .

- ▶ $A > 0 \Rightarrow \frac{1}{\sigma^2} \mathbf{s}^T \mathbf{y} - \frac{AE}{2\sigma^2} \stackrel{\mathcal{H}_1}{\geq} \frac{1}{A} \ln \tau \Rightarrow$
 $\frac{\mathbf{s}^T \mathbf{y}}{\sqrt{E}} - \frac{A\sqrt{E}}{2} \stackrel{\mathcal{H}_1}{\geq} \frac{\sigma^2}{A\sqrt{E}} \ln \tau \Rightarrow$
 $\frac{\mathbf{s}^T \mathbf{y}}{\sqrt{E}} \stackrel{\mathcal{H}_1}{\geq} \frac{A\sqrt{E}}{2} + \frac{\sigma^2}{A\sqrt{E}} \ln \tau \triangleq \eta$

- ▶ $s(\mathbf{y}) = \frac{\mathbf{s}^T \mathbf{y}}{\sqrt{E}}$ is Gaussian

- ▶ $\mathbb{E}[s(\mathbf{y})] = 0$ under \mathcal{H}_0 and

$$\mathbb{E}[s(\mathbf{y})] = \frac{\mathbf{s}^T A \mathbf{s}}{\sqrt{E}} = \frac{A \|\mathbf{s}\|^2}{\sqrt{E}} = \frac{AE}{\sqrt{E}} = A\sqrt{E} \text{ under } \mathcal{H}_1$$

- ▶ $\text{Var}(s(\mathbf{y})) = \mathbb{E}[\frac{1}{E} \mathbf{s}^T \mathbf{y} \mathbf{y}^T \mathbf{s}] = \frac{1}{E} [\mathbf{s}^T \sigma^2 I_N \mathbf{s}] = \frac{E}{E} \sigma^2 = \sigma^2$

- ▶ Thus $s(\mathbf{y}) \sim \mathcal{N}(0, \sigma^2)$ under \mathcal{H}_0 and $s(\mathbf{y}) \sim \mathcal{N}(A\sqrt{E}, \sigma^2)$ under \mathcal{H}_1

Composite Hypothesis Testing - UMP

- ▶ $s(\mathbf{y})$ independent of A , need to calculate η .
- ▶ $P_F(\delta = 1 | \mathcal{H}_0, A = 0) = \Pr(s(\mathbf{y}) \geq \eta | \mathcal{H}_0, A = 0)$

$$\begin{aligned} &= \int_y^{+\infty} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}s^2} ds \\ &= \int_{\frac{\eta}{\sigma}}^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}t^2} dt \\ &= Q\left(\frac{\eta}{\sigma}\right) \end{aligned}$$

where $t = \frac{s}{\sigma} \Rightarrow dt = \frac{1}{\sigma} ds$

- ▶ Be careful: we don't need $\max_{\mathbf{x} \in \mathcal{X}_0} P_F$ since $\mathbf{x} = \mathbf{x}_0$ (simple)
- ▶ $Q\left(\frac{\eta}{\sigma}\right) = a \Rightarrow \frac{\eta}{\sigma} = Q^{-1}(a) \Rightarrow \eta = \sigma \cdot Q^{-1}(a)$
- ▶ Thus the test $s(\mathbf{y}) \stackrel{\mathcal{H}_1}{\geq} \eta$ is UMP!

Composite Hypothesis Testing - UMP

- ▶ The power of the test can be calculated as follows:

$$P_D(A) = \Pr(s(\mathbf{y}) \geq \eta | A > 0, \mathcal{H}_1) = \int_y^{+\infty} \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(s-A\sqrt{E})^2} ds$$

- ▶ set $\frac{s-A\sqrt{E}}{\sigma} = t$

$$\begin{aligned} \text{▶ } P_D(A) &= Q\left(\frac{y - A\sqrt{E}}{\sigma}\right) \\ &= 1 - Q\left(\frac{A\sqrt{E} - \eta}{\sigma}\right) \\ &= 1 - Q\left(\frac{A\sqrt{E}}{\sigma} - Q^{-1}(a)\right) \end{aligned}$$

which is monotone increasing with A .

- ▶ Remark: We managed to find UMP since we managed to get rid of the dependence of $s(\mathbf{y})$ from A . The same would be possible for $A < 0$. But it could be impossible for $A \neq 0$, since the order of the inequality would be unknown. Thus, UMP test exists only for the one-sided test $A > 0$ (or $A < 0$).

Composite Hypothesis Testing - UMP

- ▶ case II: σ^2 unknown
- ▶ In that case hypothesis \mathcal{H}_0 is also composite since $(A, \sigma^2) = (0, \sigma^2)$ (σ^2 unknown)
- ▶ The LRT derivation still stands. How do we select η ?
- ▶ one approach: set $\eta = +\infty \Rightarrow P_D = 0$ (not very good)
- ▶ second approach: $\sigma_L^2 \leq \sigma^2 \leq \sigma_U^2$ given that $P_F = Q\left(\frac{\eta}{\sigma}\right) = a$
 $P_F = Q\left(\frac{\eta}{\sigma}\right) \leq Q\left(\frac{\eta}{\sigma_U}\right) = a$, since $Q(x)$ is increasing with decreasing x .
- ▶ Thus $\eta = \sigma_U \cdot Q^{-1}(a)$ satisfies $P_F \leq a$ and thus UMP still exists!

Composite Hypothesis Testing - GLRT

- ▶ UMP is rarely found. We revert to generalised likelihood ratio test (GLRT).
- ▶ GLRT: Suboptimal technique in general, even though it can provide UMP tests in special cases.
- ▶ $\mathcal{H}_0: \mathbf{y} \sim f(\mathbf{y}|\mathbf{x}, \mathcal{H}_0), \mathbf{x} \in \mathcal{X}_0$
- ▶ $\mathcal{H}_1: \mathbf{y} \sim f(\mathbf{y}|\mathbf{x}, \mathcal{H}_1), \mathbf{x} \in \mathcal{X}_1$
- ▶ $L_G(\mathbf{y}) = \frac{\max_{\mathbf{x} \in \mathcal{X}_1} f(\mathbf{y}|\mathbf{x}, \mathcal{H}_1)}{\max_{\mathbf{x} \in \mathcal{X}_0} f(\mathbf{y}|\mathbf{x}, \mathcal{H}_0)} = \frac{f(\mathbf{y}|\hat{\mathbf{x}}_1, \mathcal{H}_1)}{f(\mathbf{y}|\hat{\mathbf{x}}_0, \mathcal{H}_0)}$,
 $\hat{\mathbf{x}}_i = \arg \max_{\mathbf{x}_i \in \mathcal{X}_i} f(\mathbf{y}|\mathbf{x}, \mathcal{H}_i)$
- ▶ GLR is obtained by replacing the unknown parameter vector \mathbf{x} by its estimate!
- ▶ $L_G(\mathbf{y}) \stackrel{\mathcal{H}_1}{\geq} \tau$ and then we select τ by the size of the test:

$$\max_{\mathbf{x} \in \mathcal{X}} Pr(L_G(\mathbf{y}) \geq \tau | \mathcal{H}_0, \mathbf{x}) \leq a,$$

where a is the size of the test.

Composite Hypothesis Testing

Summary - Composite Hypothesis Testing with 3 approaches:

- ▶ UMP(τ) (to be explained later)
- ▶ GLRT
- ▶ "Frequentist approach": treat \mathbf{x} as a random vector rather than a constant (i.e. non-random) parameter:
 1. $f(\mathbf{y}|\mathcal{H}_i) = \int f(\mathbf{y}|\mathbf{x}, \mathcal{H}_i) f(\mathbf{x}|\mathcal{H}_i) d\mathbf{x}$
 2. Notice that $f(\mathbf{y}|\mathcal{H}_i) = \mathbb{E}_{\mathbf{x}|\mathcal{H}_i}[f(\mathbf{y}|\mathbf{x}, \mathcal{H}_i)]$ and
$$L = \frac{\mathbb{E}_{\mathbf{x}|\mathcal{H}_1}[f(\mathbf{y}|\mathbf{x}, \mathcal{H}_1)]}{\mathbb{E}_{\mathbf{x}|\mathcal{H}_0}[f(\mathbf{y}|\mathbf{x}, \mathcal{H}_0)]}$$
 3. In other words, set $f(\mathbf{x}|\mathcal{H}_i)$, calculate the above and then use detection theory.

Example

- ▶ $\mathbf{y} = \begin{bmatrix} y_c \\ y_s \end{bmatrix} = A \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix} + \mathbf{v}$, $\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_2)$
- ▶ Incoherent Detection
 - ▶ σ^2 known, A , θ unknown
 - ▶ $\mathcal{H}_0: A = 0$
 - ▶ $\mathcal{H}_1: A \neq 0$
- ▶ Polar Coordinates:
 1. $y_c = r \cos \phi$
 2. $y_s = r \sin \phi$
 3. $r = \sqrt{y_c^2 + y_s^2}$
 4. $\phi = \tan^{-1} \frac{y_s}{y_c}$
- ▶ $f(r, \phi | A, \theta) = \frac{r}{2\pi\sigma^2} e^{-\frac{1}{2\sigma^2}(A^2 + r^2 - 2Ar \cos(\phi - \theta))}$
(we have showed this in a previous lecture).

Example

$$\blacktriangleright \mathcal{H}_1: \mathbf{y} = A \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix} + \mathbf{v}$$

$$\begin{aligned} \blacktriangleright f(\mathbf{y}|A, \theta) &= \mathcal{N} \left(A \begin{bmatrix} \cos \theta \\ \sin \theta \end{bmatrix}, \sigma^2 \mathbf{I}_2 \right) \\ &= \frac{1}{\sqrt{(2\pi)^2 \sigma^4}} e^{-\frac{1}{2\sigma^2} [(y_c - A \cos \theta)^2 + (y_s - A \sin \theta)^2]} \\ &= \frac{1}{2\pi\sigma^2} e^{-\frac{1}{2\sigma^2} (y_c^2 + y_s^2 + A^2 - 2Ay_c \cos \theta - 2Ay_s \sin \theta)} \end{aligned}$$

$$\blacktriangleright \ln[f(\mathbf{y}|A, \theta, \mathcal{H}_1)] = -\frac{A}{2\sigma^2} + \frac{A}{\sigma^2} (y_c \cos \theta + y_s \sin \theta) + \underbrace{c(\mathbf{y})}_{\text{not dependent on } A, \theta}$$

$$\begin{aligned} \blacktriangleright \frac{\partial}{\partial A} \ln[f(\mathbf{y}|A, \theta, \mathcal{H}_1)] &= -\frac{A}{\sigma^2} + \frac{y_c \cos \theta + y_s \sin \theta}{\sigma^2} = 0 \Rightarrow \\ A &= y_c \cos \theta + y_s \sin \theta \end{aligned}$$

Example

$$\blacktriangleright \frac{\partial}{\partial \theta} \ln[f(\mathbf{y}|A, \theta, \mathcal{H}_1)] = \frac{A}{\sigma^2} (-y_c \sin \theta + y_s \cos \theta) = 0 \Rightarrow$$

$$y_c \sin \theta = y_s \cos \theta \Rightarrow$$

$$\tan \theta = \frac{y_s}{y_c} \Rightarrow$$

$$\hat{\theta}_{ML} = \tan^{-1} \frac{y_s}{y_c} \equiv \phi$$

$$\blacktriangleright \text{set } \theta = \hat{\theta}_{ML} \Rightarrow \hat{A}_{ML} = \frac{y_c^2}{r} + \frac{y_s^2}{r} = \frac{r^2}{r} = r$$

$$\blacktriangleright \hat{\theta}_{ML} = \phi, \hat{A}_{ML} = r$$

$$\begin{aligned} \blacktriangleright \text{Thus, } L_G(\mathbf{y}) &= \frac{f(\mathbf{y}|\hat{A}, \hat{\theta}, \mathcal{H}_1)}{f(\mathbf{y}|\hat{A}, \hat{\theta}, \mathcal{H}_0)} \\ &= \frac{\frac{r}{2\pi\sigma^2} e^{-\frac{r^2+r^2}{2\sigma^2}} e^{\frac{r^2}{\sigma^2} \cos(\phi-\phi)}}{\frac{r}{2\pi\sigma^2} e^{-\frac{r^2}{2\sigma^2}}} \\ &= e^{\frac{r^2}{2\sigma^2}} = e^{\frac{1}{2\sigma^2}(y_c^2+y_s^2)} \end{aligned}$$

Example

$$\begin{aligned} \blacktriangleright L_G(\mathbf{y}) &\stackrel{\mathcal{H}_1}{\geq} \tau \Rightarrow \\ \frac{1}{2\sigma^2} r^2 &\stackrel{\mathcal{H}_1}{\geq} \ln \tau \Rightarrow \\ r^2 &\stackrel{\mathcal{H}_1}{\geq} 2\sigma^2 \ln \tau \Rightarrow \\ r &\stackrel{\mathcal{H}_1}{\geq} \sigma \sqrt{2 \ln \tau} = \eta \end{aligned}$$

▶ Thus,

$$\Pr(r \geq \eta | \mathcal{H}_0, A = 0) = \int_{\eta}^{+\infty} f(r | A = 0, \mathcal{H}_0) dr \quad (2)$$

▶ From previous results, $f(r, \phi | A = 0, \mathcal{H}_0) = \frac{r}{2\pi\sigma^2} e^{-\frac{r^2}{2\sigma^2}}$

▶ Thus,

$$f(r | A = 0, \mathcal{H}_0) = \int_0^{2\pi} f(r, \phi | A = 0, \mathcal{H}_0) d\phi = \frac{r}{\sigma^2} e^{-\frac{r^2}{2\sigma^2}} \quad (3)$$

Example

- ▶ From (2) and (3):

$$\begin{aligned}\Pr(r \geq \eta | \mathcal{H}_0, A = 0) &= \int_{\eta}^{+\infty} \frac{r}{\sigma^2} e^{-\frac{r^2}{2\sigma^2}} dr \\ &= \left[-e^{-\frac{r^2}{2\sigma^2}} \right]_{\eta}^{+\infty} \\ &= e^{-\frac{\eta^2}{2\sigma^2}}\end{aligned}$$

- ▶ $\Pr(r \geq \eta | \mathcal{H}_0) = a = e^{-\frac{\eta^2}{2\sigma^2}} \Rightarrow$

$$\ln a = -\frac{\eta^2}{2\sigma^2} \Rightarrow$$

$$-2\sigma^2 \ln a = \eta^2 \Rightarrow$$

$$-\sigma^2 \ln \frac{1}{a} = \eta^2 \Rightarrow$$

$$\eta = \underbrace{\sigma \sqrt{2 \ln \frac{1}{a}}}_{\text{fully defined}}$$

Example

- ▶ Power $P_D = \int_y^{+\infty} f(r|A, \theta, \mathcal{H}_1) dr$



$$\begin{aligned} f(r|A, \theta, \mathcal{H}_1) &= \int_0^{2\pi} f(r, \phi|A, \theta, \mathcal{H}_1) d\phi \\ &= \int_0^{2\pi} \frac{r}{2\pi\sigma^2} e^{-\frac{r^2+A^2}{2\sigma^2}} e^{\frac{Ar}{\sigma^2} \cos(\phi-\theta)} d\phi \\ &= \frac{r}{\sigma^2} e^{-\frac{r^2+A^2}{2\sigma^2}} \frac{1}{2\pi} \int_0^{2\pi} e^{\frac{Ar}{\sigma^2} \cos(\phi-\theta)} d\phi \end{aligned} \quad (4)$$

- ▶ Set $I_0(z) = \frac{1}{2\pi} \int_0^{2\pi} e^{z \cos(\psi)} d\psi$
- ▶ Modified Bessel function of zero-th order $I_0(\cdot)$ (monotone increasing function of $z > 0$)



$$\begin{aligned} \frac{1}{2\pi} \int_0^{2\pi} e^{\frac{Ar}{\sigma^2} \cos(\phi-\theta)} d\phi &= \frac{1}{2\pi} \int_{-\theta}^{2\pi-\theta} e^{\frac{Ar}{\sigma^2} \cos(\psi)} d\psi \\ &= \frac{1}{2\pi} \int_0^{2\pi} e^{\frac{Ar}{\sigma^2} \cos(\psi)} d\psi \end{aligned} \quad (5)$$

Example

- ▶ From Eqs. (4) and (5): $f(r|A, \theta, \mathcal{H}_1) = \frac{r}{\sigma^2} e^{-\frac{r^2+A^2}{2\sigma^2}} I_0\left(\frac{Ar}{\sigma^2}\right)$ independent of θ .
- ▶ This was expected since $y_c = r \cos \phi$, $y_s = r \sin \phi$, $r = \sqrt{y_c^2 + y_s^2}$. Sufficient statistic r is rotation-invariant - the whole detection problem is rotation-invariant.
- ▶ $P_D = \Pr(r > \eta | \mathcal{H}_1, A)$
 $\equiv \Pr(r > \eta | A, \mathcal{H}_1)$
 $= \int_{\eta}^{+\infty} \frac{r}{\sigma^2} e^{-\frac{r^2+A^2}{2\sigma^2}} I_0\left(\frac{rA}{\sigma^2}\right) dr$
- ▶ Marcum's Q Function: $Q_M(a, \beta) \triangleq \int_{\beta}^{+\infty} z e^{-\frac{z^2+a^2}{2}} I_0(az) dz$
 a^2 is called the non-centrality parameter.
- ▶ set $z = \frac{r}{\sigma}$
- ▶ $\Pr(r > \eta | A) = \int_{\frac{\eta}{\sigma}}^{+\infty} z e^{-\frac{z^2+(\frac{A}{\sigma})^2}{2}} I_0\left(z \frac{A}{\sigma}\right) dz = Q_M\left(\frac{A}{\sigma}, \frac{\eta}{\sigma}\right)$.

Asymptotic Optimality of the GLRT I

- ▶ For continuous p.d.f. of the form $f(\mathbf{y}|\mathbf{x}) = u(\mathbf{y}) \cdot e^{[\mathbf{x}^T \mathbf{s}(\mathbf{y}) - t(\mathbf{x})]}$ (exponential family)
- ▶ For discrete p.m.f. of the form $\Pr(y_k = i) = p_i, i \in \{1, 2, \dots, k\}$ (multinomial distribution)
 $\mathbf{y} = [y_1 \ y_2 \ \dots \ y_N]^T$

GLRT:

1. P_F (probability of false alarm) has a guaranteed asymptotic exponential decay rate of η :

$$- \lim_{N \rightarrow +\infty} \frac{1}{N} \ln P_F(\delta_G, N, \mathbf{x}_0) \geq \eta$$

Asymptotic Optimality of the GLRT II

2. Among all tests that guarantee that the size of the test decays asymptotically at a rate greater or equal to η , the GLRT maximizes the asymptotic rate of decay of the probability of miss P_M :

$$-\lim_{N \rightarrow +\infty} \frac{1}{N} \ln P_M(\delta_G, N, \mathbf{x}_1) \geq -\lim_{N \rightarrow +\infty} \frac{1}{N} \ln P_M(\delta, N, \mathbf{x}_1)$$

Thank you!



Detection & Estimation Theory: Lectures 22-23

Prof. Aggelos Bletsas

Technical University of Crete
School of Electrical and Computer Engineering



European Union
European Social Fund

Operational Programme
**Human Resources Development,
Education and Lifelong Learning**

Co-financed by Greece and the European Union





“Support for the Internationalization of Higher Education, School of Electrical and Computer Engineering, Technical University of Crete” (MIS 5150766), under the call for proposals “Support of the Internationalization of Higher Education - Technical University of Crete” (EDULLL 153).

The project is co-financed by Greece and the European Union (European Social Fund- ESF) through the Operational Programme “Human Resources Development, Education and Lifelong Learning 2014-2020.



Operational Programme
**Human Resources Development,
Education and Lifelong Learning**
Co-financed by Greece and the European Union



Kalman Filtering

- Gauss-Markov Process
- Useful Theorem
- Kalman Filter Derivation
- Remarks

Gauss-Markov Process

- ▶ *First order Gauss-Markov process:*

$$x[n] = ax[n-1] + u[n], \quad n \geq 0$$

where $u[n]$ is (zero-mean) white Gaussian noise (WGN) with variance σ_u^2 , $x[-1] \sim \mathcal{N}(\mu_s, \sigma_s^2)$, and $x[-1]$ independent of $u[n]$, for all n .

- ▶ Are $x[0]$, $x[1]$, \dots , $x[n]$ correlated or not? \rightarrow *ANSWER:* of course they are!

$$\begin{aligned}x[0] &= ax[-1] + u[0] \\x[1] &= ax[0] + u[1] = a(ax[-1] + u[0]) + u[1] \\&= a^2x[-1] + au[0] + u[1]\end{aligned}$$

\vdots

$$x[n] = a^{n+1}x[-1] + \sum_{k=0}^n a^k u[n-k]$$

- ▶ $\mathbb{E}[x[n]] = a^{n+1}\mathbb{E}[x[-1]] = a^{n+1}\mu_s$ (depends on time, i.e., non-stationary).

► x Covariance

$$\begin{aligned}\mathbf{C}_s[m, n] &= \mathbb{E} [(x[m] - \mathbb{E}[x[m]]) (x[n] - \mathbb{E}[x[n]])] \\ &= \mathbb{E} \left[\left\{ a^{m+1} (x[-1] - \mu_s) + \sum_{k=0}^m a^k u[m-k] \right\} \cdot \right. \\ &\quad \left. \cdot \left\{ a^{n+1} (x[-1] - \mu_s) + \sum_{l=0}^n a^l u[n-l] \right\} \right] \\ &\stackrel{1}{=} a^{n+m+2} \sigma_s^2 + \sum_{k=0}^m \sum_{l=0}^n a^{k+l} \mathbb{E}[u[m-k] u[n-l]] \\ &\stackrel{2,3}{=} a^{n+m+2} \sigma_s^2 + \sum_{k=0}^m \sum_{l=0}^n a^{k+l} \sigma_u^2 \delta[l - (n - m + k)]\end{aligned}$$

¹ $x[-1]$ independent with $u[n]$

² $m - k - (n - l) = m - k - n + l = l - (n - m + k)$

³ Kronecker δ : $\delta[u] = 1$, when $u = 0$ and 0 , when $u \neq 0$

Gauss-Markov Process

- ▶ We assume $m \geq n$ and $0 \leq l \leq n$, $l = n - m + k$, then
- $$n - m + k \geq 0 \Rightarrow k \geq m - n \quad \text{and} \quad n - m + k \leq n \Rightarrow k \leq m$$
- ▶ Then,

$$\begin{aligned} \mathbf{C}_s[m, n] &= a^{n+m+2} \sigma_s^2 + \sum_{k=m-n}^m a^{n-m+2k} \sigma_u^2 \\ &\stackrel{4}{=} a^{n+m+2} \sigma_s^2 + \sigma_u^2 \sum_{k'=0}^n a^{2k'+m-n} \\ &= a^{n+m+2} \sigma_s^2 + a^{m-n} \sigma_u^2 \sum_{k=0}^n a^{2k} \end{aligned}$$

▶ *Properties:*

- clearly not WSS since depends on time (i.e., n or $n + m$)
- heavily correlated $|a| \rightarrow 1$
- heavily uncorrelated $|a| \rightarrow 0$

$${}^4k' = k - (m - n) \Rightarrow 2k = 2k' + 2(m - n)$$

Gauss-Markov Process

- ▶ for $n > m \Rightarrow \mathbf{C}_s[m, n] = \mathbf{C}_s[n, m]$
- ▶ However, Gauss-Markov process for $n \rightarrow +\infty$:

$$\mathbb{E}[x[n]] = a^{n+1} \underbrace{\mu_s}_{\mathbb{E}[x[-1]]} \xrightarrow{n \rightarrow +\infty} \emptyset \quad \text{iff}^5 \quad |a| < 1$$

$$\mathbf{C}_s[m, n] \xrightarrow{n \rightarrow +\infty} \sigma_u^2 a^{m-n} \sum_{k=0}^n a^{2k} = \sigma_u^2 a^{m-n} \frac{1}{1-a^2} \quad \text{for } |a| < 1$$

$$\Rightarrow \mathbf{C}_s[m, n] = \mathbf{C}_s[k = m-n] = \underbrace{\frac{\sigma_u^2}{1-a^2} a^k}_{\text{auto-correlation function}}, \quad k \geq 0 \quad (\text{AR}(1) \text{ process})$$

- ▶ if $\frac{\sigma_u^2}{1-a^2} = \sigma_s^2$ and $\mu_s = \emptyset$ then the above becomes wide-sense stationary (WSS) for $n \rightarrow +\infty$

⁵if and only if

- ▶ Gauss-Markov Process⁶: mean and variance can be obtained recursively:

$$\mathbb{E}[x[n]] = a\mathbb{E}[x[n-1]] + \mathbb{E}[u[n]] = a\mathbb{E}[x[n-1]]$$

$$\begin{aligned}\text{var}[x[n]] &= \mathbb{E} \left[(x[n] - \mathbb{E}[x[n]])^2 \right] \\ &= \mathbb{E} \left[(ax[n-1] + u[n] - a\mathbb{E}[x[n-1]])^2 \right] \\ &= \mathbb{E} \left[\{a(x[n-1] - \mathbb{E}[x[n-1]]) + u[n]\}^2 \right] \\ &= a^2 \text{var}(x[n-1]) + \sigma_u^2\end{aligned}$$

since $u[n]$ has zero mean, and $x[n-1]$ depends on $x[-1]$ and $u[0], u[1], \dots, u[n-1]$ which are independent of $u[n]$!

⁶ $x[n] = ax[n-1] + u[n]$

Theorem

- If $\boldsymbol{\theta}$ has zero mean $\mathbb{E}[\boldsymbol{\theta}] = \mathbf{0}$ and $\boldsymbol{\theta}, \mathbf{y} = \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix}$ are jointly Gaussians, with $\mathbf{y}_1, \mathbf{y}_2$ uncorrelated, then

$$\mathbb{E}[\boldsymbol{\theta}|\mathbf{y}] \triangleq \mathbb{E}[\boldsymbol{\theta}|\mathbf{y}_1, \mathbf{y}_2] = \mathbb{E}[\boldsymbol{\theta}|\mathbf{y}_1] + \mathbb{E}[\boldsymbol{\theta}|\mathbf{y}_2]$$

- *Proof:* since $\boldsymbol{\theta}, \mathbf{y}$ are jointly Gaussians the MSE coincides with the linear least square estimate (LLSE):

$$\hat{\boldsymbol{\theta}} = \mathbb{E}[\boldsymbol{\theta}|\mathbf{y}] = \mathbb{E}[\boldsymbol{\theta}] + \mathbf{C}_{\boldsymbol{\theta}\mathbf{y}}\mathbf{C}_{\mathbf{y}}^{-1}(\mathbf{y} - \mathbb{E}[\mathbf{y}])$$

$$\begin{aligned} \mathbf{C}_{\mathbf{y}} &= \mathbb{E} \left[\begin{bmatrix} \mathbf{y}_1 - \mathbb{E}[\mathbf{y}_1] \\ \mathbf{y}_2 - \mathbb{E}[\mathbf{y}_2] \end{bmatrix} \begin{bmatrix} (\mathbf{y}_1 - \mathbb{E}[\mathbf{y}_1])^T & (\mathbf{y}_2 - \mathbb{E}[\mathbf{y}_2])^T \end{bmatrix} \right] \\ &\Rightarrow \mathbf{C}_{\mathbf{y}} = \begin{bmatrix} \mathbf{C}_{\mathbf{y}_1} & \mathbf{C}_{12} \\ \mathbf{C}_{12}^T & \mathbf{C}_{\mathbf{y}_2} \end{bmatrix} \end{aligned} \quad (1)$$

Theorem

- ▶ Notice that $\mathbf{y}_1, \mathbf{y}_2$ are uncorrelated:

$$\begin{aligned}\mathbf{C}_{12} &= \mathbb{E} \left[\left[(\mathbf{y}_1 - \mathbb{E}[\mathbf{y}_1]) \quad (\mathbf{y}_2 - \mathbb{E}[\mathbf{y}_2])^T \right] \right] \\ &= \mathbb{E} [\mathbf{y}_1 - \mathbb{E}[\mathbf{y}_1]] \mathbb{E} \left[(\mathbf{y}_2 - \mathbb{E}[\mathbf{y}_2])^T \right] = \mathbf{0} \cdot \mathbf{0}^T = \mathbf{0} \quad (2)\end{aligned}$$

- ▶ from (1) and (2)

$$\begin{aligned}\mathbf{C}_{\mathbf{y}} &= \begin{bmatrix} \mathbf{C}_{\mathbf{y}_1} & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_{\mathbf{y}_2} \end{bmatrix} \Rightarrow \mathbf{C}_{\mathbf{y}}^{-1} = \begin{bmatrix} \mathbf{C}_{\mathbf{y}_1}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_{\mathbf{y}_2}^{-1} \end{bmatrix} \\ \mathbf{C}_{\theta_{\mathbf{y}}} &\stackrel{\mathbb{E}[\theta]=0}{=} \mathbb{E} \left[\theta \cdot \begin{bmatrix} \mathbf{y}_1 - \mathbb{E}[\mathbf{y}_1] \\ \mathbf{y}_2 - \mathbb{E}[\mathbf{y}_2] \end{bmatrix}^T \right] = \begin{bmatrix} \mathbf{C}_{\theta_{\mathbf{y}_1}} & \mathbf{C}_{\theta_{\mathbf{y}_2}} \end{bmatrix}\end{aligned}$$

► Thus,

$$\begin{aligned}\hat{\boldsymbol{\theta}} &= \cancel{\mathbb{E}[\boldsymbol{\theta}]} + \mathbf{0} + \begin{bmatrix} \mathbf{C}_{\boldsymbol{\theta}\mathbf{y}_1} & \mathbf{C}_{\boldsymbol{\theta}\mathbf{y}_2} \end{bmatrix} \begin{bmatrix} \mathbf{C}_{\mathbf{y}_1}^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_{\mathbf{y}_2}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{y}_1 - \mathbb{E}[\mathbf{y}_1] \\ \mathbf{y}_2 - \mathbb{E}[\mathbf{y}_2] \end{bmatrix} \\ &= \mathbf{C}_{\boldsymbol{\theta}\mathbf{y}_1} \mathbf{C}_{\mathbf{y}_1}^{-1} (\mathbf{y}_1 - \mathbb{E}[\mathbf{y}_1]) + \mathbf{C}_{\boldsymbol{\theta}\mathbf{y}_2} \mathbf{C}_{\mathbf{y}_2}^{-1} (\mathbf{y}_2 - \mathbb{E}[\mathbf{y}_2]) \\ &= \mathbb{E}[\boldsymbol{\theta}|\mathbf{y}_1] + \mathbb{E}[\boldsymbol{\theta}|\mathbf{y}_2] \quad \blacksquare\end{aligned}$$

Derivation of (scalar) Kalman filter

- ▶ Random parameter to be estimated, according to a Gauss-Markov Process:

$$\text{(state equation)} \quad x[n] = ax[n-1] + u[n]$$

- $x[-1]$ independent of $u[n] \forall n$
- $u[n]$ WGN (zero mean) with variance σ_u^2

and

$$\text{(observation equation)} \quad y[n] = x[n] + w[n]$$

- $u[n]$ zero mean with independent samples and variance σ_u^2 independent of time (WGN)
- $w[n]$ zero mean Gaussian noise with independent samples and $\mathbb{E}[w[n]^2] = \sigma_n^2$ (depends on time)

- ▶ $x[-1] \sim \mathcal{N}(\mu_s, \sigma_s^2) \Rightarrow \mathbb{E}[x[n]] = a^{n+1} \mu_s = \emptyset$

Derivation of Kalman filter

- ▶ *Problem:* we wish to estimate $x[n]$ base on the observations $\{y[0], y[1], y[2], \dots, y[n]\}$ or filter $\{y[n]\}$ to produce $\hat{x}[n]$:

$$\hat{x}[n|y[0], y[1], y[2], \dots, y[m]] \triangleq \hat{x}[n|m]$$

- ▶ Optimality criterion: Bayesian MSE: $\mathbb{E} [(x[n] - \hat{x}[n|n])^2]$

$$\text{MSE: } \hat{x}[n|n] = \mathbb{E} [x[n]|y[0], y[1], y[2], \dots, y[n]]$$

$x[n], y[0], y[1], y[2], \dots, y[n]$ are linear combinations of $x[-1]$ plus Gaussian noise $\Rightarrow x[n], y[0], y[1], y[2], \dots, y[n]$ are jointly Gaussians!

Derivation of Kalman filter

- ▶ Thus $\text{MSE} \equiv \text{LLSE} \Rightarrow \hat{x}[n|n] = \cancel{\mathbb{E}[x[n]]} + \mathbf{C}_{xy} \mathbf{C}_y^{-1} \mathbf{y}$ (notice that $\mathbb{E}[\mathbf{y}] = \mathbf{0}$ with $\mathbf{y} = [y[0], y[1], y[2], \dots, y[n]]^T$)
- ▶ If the Gaussian assumption is not valid, then the above holds for the "optimal" LLS estimator. Returning to the original problem:
- ▶ We need to find a sequential estimator: if $\{y[n]\}$ were uncorrelated then

$$\begin{aligned}\hat{x}[n|n] &= \mathbb{E}[x[n]|y[0], y[1], y[2], \dots, y[n]] \\ &= \mathbb{E}[x[n]|y[0], y[1], y[2], \dots, y[n-1]] + \mathbb{E}[x[n]|y[n]] \\ &= \hat{x}[n|n-1] + \mathbb{E}[x[n]|y[n]]\end{aligned}$$

however $\{y[n]\}$ are NOT uncorrelated, thus a different approach is needed.

Derivation of Kalman filter

- ▶ Set $\tilde{y}[n] = y[n] - \hat{y}[n|n-1]$, where
 $\hat{y}[n|n-1] = \mathbb{E}[y[n]|y[0], \dots, y[n-1]]$ (MSE estimate)
- ▶ Thus $\tilde{y}[n] = e$, which \perp to the data $y[0], \dots, y[n-1]$
(orthogonality principle)
- ▶ Set $\mathbf{y}[n-1] = [y[0] \quad y[1] \quad \dots \quad y[n-1]]$
- ▶ Thus,

$$y[n] = \tilde{y}[n] + \hat{y}[n|n-1] = \tilde{y}[n] + \sum_{k=0}^{n-1} a_k y[k] \quad (3)$$

since $\{y[k]\}$ jointly Gaussians and thus MSE \equiv LLS
(linear).

- ▶ $\mathbf{y}[n-1], \tilde{y}[n]$ can give $y[n]$ through (3). Thus $\mathbf{y}[n-1], y[n]$
are equivalent to $\mathbf{y}[n-1], \tilde{y}[n]$:

$$\hat{x}[n|n] = \mathbb{E}[x[n]|\mathbf{y}[n-1], y[n]] = \mathbb{E}[x[n]|\mathbf{y}[n-1], \tilde{y}[n]]$$

Derivation of Kalman filter

- ▶ \tilde{y} is error in observation: the trick is to predict $y[n]$, which you already know!
- ▶ $\tilde{y}[n]$ uncorrelated with the observation data $\mathbf{y}[n-1]$, due to orthogonality principle.
- ▶ Thus,

$$\hat{x}[n|n] = \underbrace{\mathbb{E}[x[n]|\mathbf{y}[n-1]]}_{\hat{x}[n|n-1]} + \mathbb{E}[x[n]|\tilde{y}[n]] \quad (4)$$

$$\begin{aligned}\hat{x}[n|n-1] &= \mathbb{E}[x[n]|\mathbf{y}[n-1]] \\ &= \mathbb{E}[ax[n-1] + u[n]|\mathbf{y}[n-1]] \\ &= a\mathbb{E}[x[n-1]|\mathbf{y}[n-1]] \\ &= a\hat{x}[n-1|n-1]\end{aligned} \quad (5)$$

- ▶ $u[n]$ is independent of $\{w[n]\}$, $x[-1]$ and $u[n-1], u[n-2], \dots, u[0]$

Derivation of Kalman filter

- ▶ Thus, $\mathbb{E}[u[n]|\mathbf{y}[n-1]] = \mathbb{E}[u[n]] = 0$
- ▶ So far:

$$\hat{x}[n|n] = \underbrace{a\hat{x}[n-1|n-1]}_{\hat{x}[n|n-1]} + \mathbb{E}[x[n]|\tilde{y}[n]] \quad (6)$$

- $\mathbb{E}[x[n]|\tilde{y}[n]] = \cancel{\mathbb{E}[x[n]]} + \mathbf{C}_{x\tilde{y}}\mathbf{C}_{\tilde{y}}^{-1}\tilde{y}[n]$
 - $\mathbb{E}[\tilde{y}[n]] = 0$
 - $\mathbf{C}_{x\tilde{y}} = \mathbb{E}[x[n]\tilde{y}[n]]$
 - $\mathbf{C}_{\tilde{y}} = \mathbb{E}[(\tilde{y}[n])^2]$
- ▶ Thus,

$$\begin{aligned}\mathbb{E}[x[n]|\tilde{y}[n]] &= \frac{\mathbb{E}[x[n]\tilde{y}[n]]}{\mathbb{E}[(\tilde{y}[n])^2]}\tilde{y}[n] = \underbrace{\frac{\mathbb{E}[x[n]\tilde{y}[n]]}{\mathbb{E}[(\tilde{y}[n])^2]}}_{K[n]}(y[n] - \hat{y}[n|n-1]) \\ &= K[n](y[n] - \hat{y}[n|n-1])\end{aligned}$$

Derivation of Kalman filter

▶ $\hat{y}[n|n-1] \triangleq \mathbb{E}[y[n]|\mathbf{y}[n-1]]$

$$\stackrel{7}{=} \underbrace{\mathbb{E}[x[n]|\mathbf{y}[n-1]]}_{\hat{x}[n|n-1]} + \mathbb{E}[w[n]|\mathbf{y}[n-1]] \xrightarrow{\emptyset}$$

▶ $\mathbb{E}[x[n]|\tilde{y}[n]] = K[n](y[n] - \hat{y}[n|n-1])$
 $= K[n](y[n] - \hat{x}[n|n-1])$

▶ and thus from (6),

$$\hat{x}[n|n] = \underbrace{\hat{x}[n|n-1]}_{a\hat{x}[n-1|n-1]} + K[n](y[n] - \hat{x}[n|n-1])$$

▶ it remains to calculate the gain $K[n]$ in a recursive manner:

$$K[n] \triangleq \frac{\mathbb{E}[x[n](y[n] - \hat{x}[n|n-1])]}{\mathbb{E}[(y[n] - \hat{x}[n|n-1])^2]}$$

⁷ $y[n] = x[n] + w[n]$ and $w[n]$ independent to $y[0], \dots, y[n-1]$

Derivation of Kalman filter

We observe the following:

$$\begin{aligned} & \blacktriangleright \mathbb{E} [x[n] (y[n] - \hat{x}[n|n-1])] \\ & = \mathbb{E} [(x[n] - \hat{x}[n|n-1]) (y[n] - \hat{x}[n|n-1])] \stackrel{\Delta}{=} M(n|n-1) \\ & \text{since} \end{aligned}$$

$$\mathbb{E} \left[\underbrace{\hat{x}[n|n-1]}_{\substack{\text{linear combination of} \\ y[0], y[1], \dots, y[n-1]}} \underbrace{(y[n] - \hat{x}[n|n-1])}_e \right] \stackrel{1}{=} 0$$

\blacktriangleright and

$$\mathbb{E} [w[n] (x[n] - \hat{x}[n|n-1])] = 0$$

since $w[n]$ independent (and thus, uncorrelated) to $y[n-1], \dots, y[0]$.

¹ $y[n]=x[n]+w[n]$

Derivation of Kalman filter

► Thus,

$$\begin{aligned} K[n] &= \frac{\mathbb{E} [(x[n] - \hat{x}[n|n-1]) (x[n] - \hat{x}[n|n-1])]}{\mathbb{E} [(x[n] - \hat{x}[n|n-1] + w[n])^2]} \\ &\stackrel{8}{\Rightarrow} K[n] = \frac{\mathbb{E} [(x[n] - \hat{x}[n|n-1])^2]}{\mathbb{E} [(x[n] - \hat{x}[n|n-1])^2] + \sigma_n^2} \\ &\Rightarrow K[n] = \frac{M[n|n-1]}{M[n|n-1] + \sigma_n^2} \end{aligned} \quad (7)$$

► ...need recursion:

$$\begin{aligned} M[n|n-1] &\triangleq \mathbb{E} [(x[n] - \hat{x}[n|n-1])^2] \\ &= \mathbb{E} [(ax[n-1] + u[n] - \hat{x}[n|n-1])^2] \\ &= \mathbb{E} [\{a(x[n-1] - \hat{x}[n-1|n-1]) + u[n]\}^2] \\ &= a^2 \mathbb{E} [(x[n-1] - \hat{x}[n-1|n-1])^2] + \sigma_u^2 \end{aligned}$$

⁸numerator: MSE when $\mathbf{y}[n-1]$ is used instead of $\mathbf{y}[n]$

Derivation of Kalman filter

- ▶ ...where the last equality is due to the fact that $u[n]$ is independent of $\underbrace{x[0], x[1], \dots, x[n-1]}_{\text{depend on } x[-1], u[0], \dots, u[n-1]}$ and

$$\begin{aligned}\mathbf{y}[n-1] &= \begin{bmatrix} y[0] & y[1] & \dots & y[n-1] \end{bmatrix} \\ &= \begin{bmatrix} x[0] + w[0] & x[1] + w[1] & \dots & x[n-1] + w[n-1] \end{bmatrix} \\ &\Rightarrow \mathbb{E}[(x[n-1] - \hat{x}[n-1|n-1])u[n]] = 0\end{aligned}$$

- ▶ Thus,

$$\begin{aligned}M[n|n-1] &= a^2 \mathbb{E}[(x[n-1] - \hat{x}[n-1|n-1])^2] + \sigma_u^2 \\ &= a^2 M[n-1|n-1] + \sigma_u^2\end{aligned}$$

Derivation of Kalman filter

- ▶ Now we require a recursion for $M[n|n]$:⁹

$$\begin{aligned}M[n|n] &= \mathbb{E} \left[(x[n] - \hat{x}[n|n])^2 \right] \\&= \mathbb{E} \left[\underbrace{\{x[n] - \hat{x}[n|n-1] - K[n](y[n] - \hat{x}[n|n-1])\}^2}_{\text{constant as a function of } n} \right] \\&= \underbrace{\mathbb{E} \left[(x[n] - \hat{x}[n|n-1])^2 \right]}_{M[n|n-1]} \\&\quad - 2K[n] \underbrace{\mathbb{E} \left[(x[n] - \hat{x}[n|n-1]) (y[n] - \hat{x}[n|n-1]) \right]}_{\text{from (7), numerator of } K[n] \Rightarrow M[n|n-1]} \\&\quad + K^2[n] \underbrace{\mathbb{E} \left[(y[n] - \hat{x}[n|n-1])^2 \right]}_{\text{denominator of } K[n] \Rightarrow \frac{M[n|n-1]}{K[n]}}\end{aligned}$$

⁹ $\hat{x}[n|n] = \hat{x}[n|n-1] + K[n](y[n] - \hat{x}[n|n-1])$

► Thus,

$$\begin{aligned}M[n|n] &= M[n|n-1] - 2K[n]M[n|n-1] + K[n]M[n|n-1] \\ &= (1 - K[n]) M[n|n-1].\end{aligned}$$

⁹ $\hat{x}[n|n] = \hat{x}[n|n-1] + K[n](y[n] - \hat{x}[n|n-1])$

Derivation of Kalman filter: Summary

- ▶ *Prediction:*

$$\hat{x}[n|n-1] = a\hat{x}[n-1|n-1]$$

- ▶ *Prediction MSE:*

$$M[n|n-1] = a^2M[n-1|n-1] + \sigma_u^2$$

- ▶ *Kalman Gain:*

$$K[n] = \frac{M[n|n-1]}{\sigma_n^2 + M[n|n-1]}$$

- ▶ *Correction:*

$$\hat{x}[n|n] = \hat{x}[n|n-1] + K[n](y[n] - \hat{x}[n|n-1])$$

- ▶ *Minimum MSE:*

$$M[n|n] = (1 - K[n])M[n|n-1]$$

- ▶ ...the same for $\mu_s \neq \emptyset$

- ▶ *Initialization:*

$$\hat{x}[-1|-1] = \mathbb{E}[x[-1]] = \mu_s \quad \text{and} \quad M[-1|-1] = \sigma_s^2.$$

Remarks

▶ Remark 0:

Derived equations hold for $\mu_s \neq \emptyset$ too.

▶ Remark 1:

- LLS estimates: $\text{LLS}(\boldsymbol{\theta}|\mathbf{y}_1, \mathbf{y}_2) = \text{LLS}(\boldsymbol{\theta}|\mathbf{y}_1) + \text{LLS}(\boldsymbol{\theta}|\mathbf{y}_2)$, where $\mathbf{y}_1, \mathbf{y}_2$ uncorrelated
- orthogonality principle: $\mathbf{e} \perp$ linear combination of the data thus, Kalman filter is the "optimal" (in MSE sense) recursive LINEAR estimator
- if Gaussian statistics are employed \Rightarrow Kalman is the "optimal" (in MSE sense) estimator!

▶ Remark 2:

we used Gauss-Markov process for the parameter to be estimated

- (state equation) $x[n] = ax[n-1] + u[n]$
- (observation equation) $y[n] = x[n] + w[n]$
- $\mathbb{E}[(w[n])^2] = \sigma_n^2 \rightarrow$ function of n and $\mathbb{E}[x[n]] = a^{n+1}\mathbb{E}[x[-1]] = a^{n+1}\mu_s$

thus, Kalman filter holds for non-WSS processes (we haven't seen that so far)!

► Remark 3:

Set $a = 1$ and $\sigma_u^2 = 0 \Rightarrow x[n] = x[n - 1]$ (prediction: last estimate of $x[n]$).

In that case:

$$\hat{x}[n|n - 1] = \hat{x}[n - 1]$$

$$M[n|n - 1] = M[n - 1]$$

$$\hat{x}[n] = \hat{x}[n - 1] + K[n](y[n] - \hat{x}[n - 1])$$

$$K[n] = \frac{M[n - 1]}{M[n - 1] + \sigma_n^2}$$

$$M[n] = (1 - K[n])M[n - 1]$$

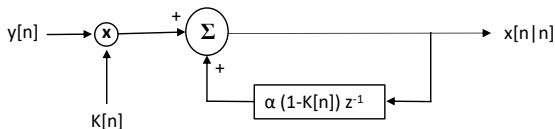
...can omit the prediction stage of Kalman.

Remarks

► Remark 4:

Kalman filter is a time varying linear filter:

$$\begin{aligned}\hat{x}[n|n] &= a\hat{x}[n-1|n-1] + K[n](y[n] - \underbrace{a\hat{x}[n-1|n-1]}_{\hat{x}[n|n-1]}) \\ &= \underbrace{K[n]}_{\text{time varying-constant}} y[n] + \underbrace{a(1-K[n])}_{\text{time varying-constant}} \hat{x}[n-1|n-1]\end{aligned}$$



- For $n \rightarrow +\infty$ the filter becomes time-invariant (steady-state).

Remarks

▶ Remark 5:

no-matrix inversion is needed (true only here). For the vector Kalman filter, this is not true.

▶ Remark 6:

Minimum prediction MSE:

$$\begin{aligned}M[n|n-1] &= \mathbb{E} [(x[n] - \hat{x}[n|n-1])^2] \\ &= a^2 M[n-1|n-1] + \sigma_u^2\end{aligned}$$

Kalman Gain:

$$K[n] = \frac{M[n|n-1]}{\sigma_n^2 + M[n|n-1]}$$

Minimum MSE:

$$M[n|n] = (1-K[n])M[n|n-1] \triangleq (1-K[n])\mathbb{E} [(x[n] - \hat{x}[n|n-1])^2]$$

Thus, $M[n|n]$ can be computed independently of the observation data $\{y[n]\} \Rightarrow$ can be computed offline!

▶ Remark 7:

Kalman filter is a filter: transient response and steady-state response. At steady-state (or $n \rightarrow +\infty$) it can be proved that $M[n|n-1] > M[n-1|n-1]$
Thus, error increases at prediction stage and decreases at correction stage ($K[n] < 1$)

$$M[n|n] = (1 - K[n])M[n|n-1] < M[n|n-1]$$

▶ Remark 8:

Infinite-length causal Wiener filter:

$$\hat{x}[n] = \sum_{k=0}^{+\infty} h[k]y[n-k]$$

solved through Wiener-Hopf equations.

► Remark 8 continued:

(A) We showed that Gauss-Markov $x[n]$ as $n \rightarrow +\infty$ could become WSS

(B) if $\mathbb{E} [(w[n])^2] = \sigma_n^2 = \sigma^2$ (time independent)

if (A), (B) hold then Kalman \equiv Wiener

► Remark 9:

steady-state: time invariant filter for conditions (A), (B)

$$K[n] = K[\infty]$$

$$\begin{aligned}\hat{x}[n|n] &= \hat{x}[n|n-1] + K[n](y[n] - \hat{x}[n|n-1]) \\ &= a\hat{x}[n-1|n-1] + K[n](y[n] - a\hat{x}[n-1|n-1]) \\ &= a(1 - K[\infty])\hat{x}[n-1|n-1] + K[\infty]y[n]\end{aligned}$$

- ▶ Remark 9 continued:

$$\hat{x}[n|n] - a(1 - K[\infty])\hat{x}[n-1|n-1] = K[\infty]y[n]$$

$$\Rightarrow \hat{X}(z) - a(1 - K[\infty])\hat{X}(z)z^{-1} = K[\infty]Y(z)$$

$$\Rightarrow H(z) = \frac{K[\infty]}{1 - a(1 - K[\infty])z^{-1}} = H(z = j\omega) = H(z = j2\pi f)$$

- ▶ Remark 10:

Same properties for vector Kalman filter (apart from matrix invertibility).

- ▶ Remark 10:
Equations of the vector Kalman filter can be found in any estimation theory textbook.
- ▶ Derivation of the (scalar or vector) Kalman filter equations can be performed in an elegant, simplified way, using (modern) inference theory!
- ▶ Kalman filter = Gaussian Belief Propagation in HMMs!
- ▶ Pls take the graduate course on *Probabilistic Graphical Models and Inference Algorithms* to see this.

Thank you!



Detection & Estimation Theory: Lecture 24

Prof. Aggelos Bletsas

Technical University of Crete
School of Electrical and Computer Engineering



European Union
European Social Fund

Operational Programme
**Human Resources Development,
Education and Lifelong Learning**

Co-financed by Greece and the European Union





“Support for the Internationalization of Higher Education, School of Electrical and Computer Engineering, Technical University of Crete” (MIS 5150766), under the call for proposals “Support of the Internationalization of Higher Education - Technical University of Crete” (EDULLL 153).

The project is co-financed by Greece and the European Union (European Social Fund- ESF) through the Operational Programme “Human Resources Development, Education and Lifelong Learning 2014-2020.



Operational Programme
**Human Resources Development,
Education and Lifelong Learning**
Co-financed by Greece and the European Union



Important Sampling / Particle Filters

- Problem Definition & Basic Assumptions
- Prediction/Correction Equations
- Particle Filtering Derivation
 - Importance Sampling
- Remarks

Problem Definition

We limit discussion on scalar case - same reasoning holds for the vector case.

- ▶ denote as $p_{n|m} \equiv p_{x_n|y_0,y_1,\dots,y_m}(x_n|y_0,y_1,\dots,y_m)$
- ▶ again, we follow the same Detection/Estimation notation:
 - ▶ y_m is the m-th measurement/observation
 - ▶ x_l is the l-th random variable to be estimated ("hidden state")
- ▶ General problem:
estimate x_0, x_1, \dots, x_n given observations y_0, y_1, \dots, y_n , i.e.,

$$\text{find } p_{x_i|y_0,y_1,\dots,y_n}(x_i|y_0,y_1,\dots,y_n), \underline{0 \leq i \leq n}$$

Prediction/Correction Equations (& Assumptions)

General equations, written in an iterative manner.

► *Prediction Equation:*

$$\begin{aligned} p_{n+1|n}(x_{n+1}|y_0, y_1, \dots, y_n) &= \int_{x_n} p(x_{n+1}, x_n|y_0, y_1, \dots, y_n) dx_n \\ &= \int_{x_n} p(x_{n+1}|x_n, y_0, y_1, \dots, y_n) p(x_n|y_0, y_1, \dots, y_n) dx_n \\ &\stackrel{1}{=} \int_{x_n} p(x_{n+1}|x_n) p(x_n|y_0, y_1, \dots, y_n) dx_n \\ &= \int_{x_n} p(x_{n+1}|x_n) \underbrace{p_{n|n}(x_n|y_0, y_1, \dots, y_n)}_{\text{previous iteration}} dx_n \end{aligned}$$

¹ $x_{n+1} \perp y_0, y_1, \dots, y_n | x_n$

Prediction/Correction Equations (& Assumptions)

- *Update/Correction Equation:*

$$\begin{aligned} p_{n+1|n+1}(x_{n+1}|y_0, y_1, \dots, y_{n+1}) &= \frac{p(x_{n+1}, y_{n+1}|y_0, y_1, \dots, y_n)}{p(y_{n+1}|y_0, y_1, \dots, y_n)} \\ &= \frac{p(y_{n+1}|x_{n+1}, y_0, y_1, \dots, y_n)p(x_{n+1}|y_0, y_1, \dots, y_n)}{\underbrace{p(y_{n+1}|y_0, y_1, \dots, y_n)}_z} \\ &= \frac{1}{z} p(y_{n+1}|x_{n+1}) p_{n+1|n}(x_{n+1}|y_0, y_1, \dots, y_n) \end{aligned}$$

- Notice that the above prediction/correction equations hold for:

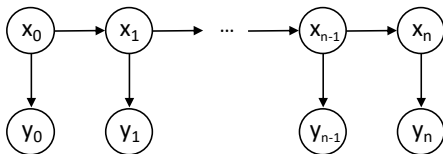
$$x_{n+1} \perp y_0, y_1, \dots, y_n | x_n \quad (1)$$

$$y_{n+1} \perp y_0, y_1, \dots, y_n | x_{n+1} \quad (2)$$

² $y_{n+1} \perp y_0, y_1, \dots, y_n | x_{n+1}$

Particle Filter for HMMs with General Continuous

- ▶ Eqs. (1), (2) are satisfied by hidden Markov model (HMM):



- ▶ HMM satisfies the following:

$$p(x_n | x_{n-1}, x_{n-2}, \dots, x_0) = p(x_n | x_{n-1}) \quad (3)$$

- ▶ Thus, one can produce $(x_0(s), x_1(s), \dots, x_n(s))$ that adhere to $p_{x_0, x_1, \dots, x_n}(x_0^n) \triangleq p_{x_0, x_1, \dots, x_n}(x_0, x_1, \dots, x_n)$ as follows:

$$x_0(s) \sim p_{x_0}(\cdot), \quad x_i(s) \sim p_{x_i | x_{i-1}}(\cdot | x_{i-1}(s)), \quad i = 1, \dots, n \quad (4)$$

- ▶ Note that due to (3),

$$p_{x_0, x_1, \dots, x_n}(x_0^n) = p_{x_0}(x_0) \prod_{i=1}^n p_{x_i | x_{i-1}}(x_i | x_{i-1})$$

Particle Filter Derivation

- ▶ HMM also satisfies the following:

$$p(y_n | x_n, x_{n-1}, y_{n-1}, x_{n-2}, y_{n-2}, \dots, x_0, y_0) = p(y_n | x_n)$$

thus,

$$\begin{aligned} p(y_n, y_{n-1}, \dots, y_0 | x_n, x_{n-1}, \dots, x_0) &= \\ &= p(y_n | y_{n-1}, \dots, y_0, x_n, x_{n-1}, \dots, x_0) \\ &\quad \cdot p(y_{n-1}, \dots, y_0 | x_n, x_{n-1}, \dots, x_0) \\ &\stackrel{(6)}{=} p(y_n | x_n) p(y_{n-1}, \dots, y_0 | x_n, x_{n-1}, \dots, x_0) \end{aligned}$$

Working inductively,

$$p(y_n, y_{n-1}, \dots, y_0 | x_n, x_{n-1}, \dots, x_0) = \prod_{i=0}^n p_{y_i | x_i}(y_i | x_i) \quad (5)$$

Particle Filter Derivation

- ▶ x_{n-1} separates y_{n-1}, \dots, y_0 from x_n and thus (6):

$$y_{n-1}, \dots, y_0 \perp x_n | x_{n-1}, x_{n-2}, \dots, x_0 \quad (6)$$

- ▶ In HMM-particle filtering, we are given $p(x_i | x_{i-1})$ and $p(y_i | x_i)$, $0 \leq i \leq n$
- ▶ Thus,³

$$\begin{aligned} p_{x_0, x_1, \dots, x_n | y_0, y_1, \dots, y_n}(x_0^n | y_0^n) &= \frac{p(y_0^n | x_0^n) p(x_0^n)}{p(y_0^n)} \\ &= \frac{p_{x_0, x_1, \dots, x_n}(x_0^n) \prod_{i=0}^n p_{y_i | x_i}(y_i | x_i)}{\underbrace{p_{y_0, y_1, \dots, y_n}(y_0^n)}}_{z \rightarrow \text{unknown (i.e., hard to compute) constant}} \end{aligned} \quad (7)$$

³notation: $x_0^n = x_0, x_1, \dots, x_n$ and $y_0^n = y_0, y_1, \dots, y_n$

Particle Filter Idea

- ▶ Particle filtering idea:

- 1) Produce samples that adhere to

$p_{x_0, x_1, \dots, x_n | y_0, y_1, \dots, y_n}(x_0^n | y_0^n)$ without knowing z .

- 2) These samples can be used to estimate any $\mathbb{E}_{p_{x_0^n | y_0^n}} [f(x_0^n)]$,

for any function $f(\cdot)$ as if one perfectly knew the p.d.f.

$p_{x_0^n | y_0^n}(\cdot | y_0^n)$.

- ▶ We do not know z ; can only produce as many samples $\{(x_0(s), x_1(s), \dots, x_n(s))\}$, $s = 1, \dots, \mathcal{S}$ according to known

$$\prod_{i=0}^n p_{y_i | x_i}(y_i | x_i)$$

- ▶ Particle filtering is a a case of non-parametric modeling, since we do not know the posterior p.d.f. in closed form $p_{x_0^n | y_0^n}(\cdot | y_0^n)$.

Example

- ▶ In localization research, we need the conditional mean (i.e., MMSE estimator):

$$\mathbb{E}_{p_{x_0^n | y_0^n}} [f(x_0^n)] \equiv \mathbb{E} [x_n | y_0, y_1, \dots, y_n]$$

Importance Sampling

- ▶ Theory to solve the above problem \Rightarrow *Importance Sampling*:

$$\mu(x) = \frac{q(x)}{z}$$

\leftarrow known function
 \leftarrow unknown function

We want samples from $\mu(\cdot)$ so that we can compute the following:

$$\mathbb{E}_{\mu} [f(x)] \triangleq \int f(x)\mu(x)dx$$

Important Sampling Algorithm

► *Important Sampling algorithm:*

- 1) Produce samples $x(1), x(2), \dots, x(s), \dots, x(\mathcal{S})$ from known distribution $v(x)$, called the "proposal" distribution;
- 2) Compute as many weights as the samples, according to

$$w(s) = \frac{q(x(s))}{v(x(s))}, \quad s = 1, 2, \dots, \mathcal{S}$$

Calculate $\hat{E}(\mathcal{S})$ instead of $\mathbb{E}_\mu [f(x)]$:

$$\mathbb{E}_\mu [f(x)] \rightarrow \hat{E}(\mathcal{S}) = \frac{\frac{1}{\mathcal{S}} \sum_{s=1}^{\mathcal{S}} w(s) f(x(s))}{\frac{1}{\mathcal{S}} \sum_{s=1}^{\mathcal{S}} w(s)}$$

► *Definition:* support of function $p(x)$ ($\text{supp}(p)$)

$$\text{supp}(p) = \{x : p(x) > 0\}$$

Important Sampling Algorithm: Theorem

- *Theorem 1:* Let $\text{supp}(\mu) \subseteq \text{supp}(v)$. Then for $\mathcal{S} \rightarrow +\infty$

$$\hat{E}(\mathcal{S}) \rightarrow \mathbb{E}_\mu [f(x)], \quad \text{with probability 1}$$

- *Proof:*

$$\lim_{\mathcal{S} \rightarrow +\infty} \frac{1}{\mathcal{S}} \sum_{s=1}^{\mathcal{S}} w(s) f(x(s)) \xrightarrow[\text{Large numbers}]{\text{strong Law of}} \underbrace{\mathbb{E}_v [w \cdot f(x)]}_{\substack{\text{expected value in terms of} \\ v(x) \text{ because we have samples from} \\ \text{the "proposal" pdf } v(x)}}$$

$$\begin{aligned} \mathbb{E}_v [w \cdot f(x)] &= \mathbb{E}_v \left[\frac{q(x)}{v(x)} f(x) \right] = \int_{\text{supp}(v)} \frac{q(x)}{v(x)} f(x) v(x) dx = \\ &= \int_{\text{supp}(v)} q(x) f(x) dx \stackrel{4}{=} \int_{\text{supp}(\mu)} \mu(x) f(x) dx = z \mathbb{E}_\mu [f(x)] \quad (8) \end{aligned}$$

⁴ $q(x) = 0 \forall x \notin \text{supp}(\mu)$

Important Sampling Algorithm: Theorem

- *Proof continued:*
similarly,

$$\begin{aligned} \lim_{S \rightarrow +\infty} \frac{1}{S} \sum_{s=1}^S w(s) &\xrightarrow[\text{Large numbers}]{\text{strong Law of}} \mathbb{E}_v[w] = \int_{\text{supp}(v)} \frac{q(x)}{v(x)} v(x) dx \stackrel{5}{=} \\ &= \int_{\text{supp}(\mu)} q(x) dx = \\ &= z \int_{\text{supp}(\mu)} \mu(x) dx \stackrel{1}{=} z \end{aligned} \tag{9}$$

from (8) and (9) \Rightarrow proof is completed. ■

⁵ $q(x) = 0 \forall x \notin \text{supp}(\mu)$

Remarks

- ▶ It is worth noting that as long as $\text{supp}(\mu)$ is contained in $\text{supp}(v)$, the estimation converges irrespective of the choice of v .

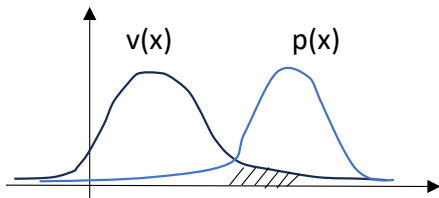


Figure 1: Area where proposal distribution obtains small values, as opposed to true distribution; that amplifies the number of required samples.

- ▶ However, the choice of v determines the variance of the estimator and hence the number of samples \mathcal{S} required to obtain good estimation.

- ▶ Need to estimate $\mathbb{E}_{p_{x_0^n|y_0^n}} [f(x_0^n)|y_0^n]$

- ▶ Set

$$q(x_0^n) = p_{x_0, x_1, \dots, x_n}(x_0^n) \prod_{i=0}^n p_{y_i|x_i}(y_i|x_i)$$

- ▶ Set

$$v(x_0^n) = \text{prior} = p_{x_0, x_1, \dots, x_n}(x_0^n)$$

- ▶ For any given $(x_0, x_1, \dots, x_n) \equiv x_0^n$ sample, calculate the corresponding weight as follows:

$$w(x_0^n) \equiv w_0^n = \frac{q(x_0^n)}{v(x_0^n)} \stackrel{(7)}{=} \prod_{i=0}^n p_{y_i|x_i}(y_i|x_i)$$

HMM-PF Summary

► Thus, for given $y_0, y_1, \dots, y_n = y_0^n$

1) Obtain \mathcal{S} iid samples $x_0^n(s)$, $1 \leq s \leq \mathcal{S}$ according to

$$p_{x_0, x_1, \dots, x_n}(\cdot)$$

2) Compute \mathcal{S} corresponding weights

$$w_0^n(s) = \prod_{i=0}^n p(y_i | x_i(s)) \quad \forall s, 1 \leq s \leq \mathcal{S}$$

3) Output estimation of $\mathbb{E}_{p_{x_0^n | y_0^n}} [f(x_0, x_1, \dots, x_n) | y_0, y_1, \dots, y_n]$

$$\text{as } \frac{\sum_{s=1}^{\mathcal{S}} w_0^n(s) f(x_0^n(s))}{\sum_{s=1}^{\mathcal{S}} w_0^n(s)}.$$

Remark 1

- ▶ The above estimator can be implemented in a sequential fashion, exploiting the HMM properties:
 - ▶ for given s in $\{1, 2, \dots, \mathcal{S}\}$
sample $x_{k+1}(s) \sim p_{x_{k+1}|x_k}(\cdot|x_k(s))$
 - ▶ compute

$$\underbrace{w_0^{k+1}(s)}_{\substack{\text{corresponding} \\ \text{weight for} \\ \text{particle } x_{k+1}(s)}} = \underbrace{w_0^{k+1}(s)}_{\substack{\text{weight for} \\ \text{sample } x_k(s)}} p_{y_{k+1}|x_{k+1}}(y_{k+1}|x_{k+1}(s))$$

- ▶ repeat for all k up to n and all s up to \mathcal{S}
- ▶ re-sample the weights/particles before you proceed to next $k + 1$ (particle depletion problem)

Remark 2

- ▶ Can we extend particle filtering to other probabilistic graphical models (PGM), other than HMMs?
- ▶ *Answer:* YES, using mixture of Gaussians and producing samples based on specific Markov chain Monte Carlo (MCMC) techniques: say $p(x) = \frac{q(x)}{z}$ ($z \rightarrow$ unknown constant)

You can craft a Markov chain (MC) that produces samples according to $p(x)$, even though the MC was built using $q(x) \rightarrow$ Metropolis-Hastings technique (Gibbs sampling is a special case).

- ▶ ...you can take graduate class *Probabilistic Graphical Models and Inference Algorithms* to see the above!

Thank you!