

# Ch. 11

## General Bayesian Estimators

# Introduction

## In Chapter 10 we:

- introduced the idea of a “a priori” information on  $\theta$ 
  - $\Rightarrow$  use “prior” pdf:  $p(\theta)$
- defined a new optimality criterion
  - $\Rightarrow$  Bayesian MSE
- showed the Bmse is minimized by  $E \{ \theta | \mathbf{x} \}$

called:

- “mean of posterior pdf”
- “conditional mean”

## In Chapter 11 we will:

- define a more general optimality criterion
  - $\Rightarrow$  leads to several different Bayesian approaches
  - $\Rightarrow$  includes Bmse as special case

*Why?* Provides flexibility in balancing:

- model,
- performance, and
- computations

# 11.3 Risk Functions

Previously we used Bmse as the Bayesian measure to minimize

$$Bmse = E\left\{\left(\theta - \hat{\theta}\right)^2\right\} \quad w.r.t. \quad p(\mathbf{x}, \theta)$$

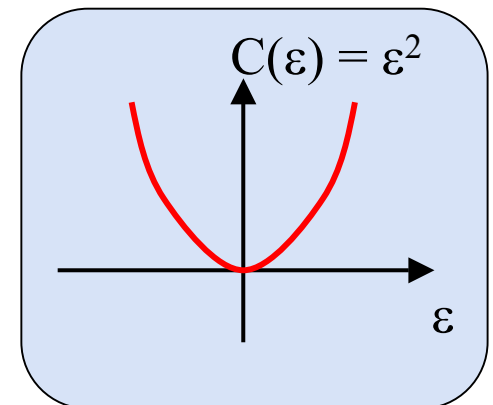
$$\theta - \hat{\theta} \triangleq \varepsilon$$

So, Bmse is... Expected value of square of error

Let's write this in a way that will allow us to generalize it.

Define a quadratic Cost Function:  $C(\varepsilon) = \varepsilon^2 = (\theta - \hat{\theta})^2$

Then we have that  $Bmse = E\{C(\varepsilon)\}$



**Why limit the cost function to just quadratic?**

# General Bayesian Criteria

1. Define a cost function:  $C(\varepsilon)$
2. Define Bayes Risk:  $\mathcal{R} = E\{C(\varepsilon)\}$  w.r.t.  $p(\mathbf{x}, \theta)$

$$\mathcal{R}(\hat{\theta}) = E\{C(\theta - \hat{\theta})\}$$

Depends on choice of estimator

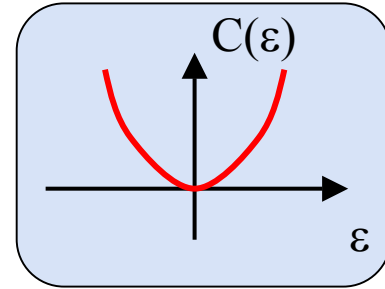
3. Minimize Bayes Risk w.r.t. estimate  $\hat{\theta}$

The choice of the cost function can be tailored to:

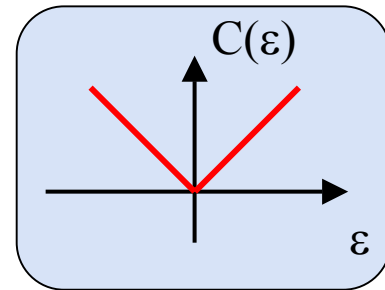
- Express importance of avoiding certain kinds of errors
- Yield desirable forms for estimates
  - e.g., easily computed
- Etc.

# Three Common Cost Functions

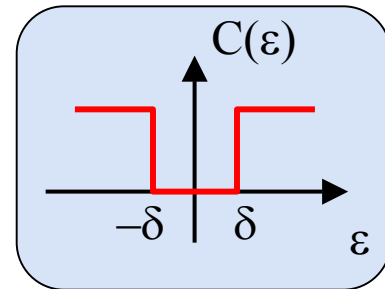
1. Quadratic:  $C(\varepsilon) = \varepsilon^2$



2. Absolute:  $C(\varepsilon) = |\varepsilon|$



3. Hit-or-Miss:  $C(\varepsilon) = \begin{cases} 0, & |\varepsilon| < \delta \\ 1, & |\varepsilon| \geq \delta \end{cases}$   
 $\delta > 0$  and small



# General Bayesian Estimators

Derive how to choose estimator to minimize the chosen risk:

$$\begin{aligned}\mathcal{R}(\hat{\theta}) &= E\{C(\theta - \hat{\theta})\} \\ &= \iint C(\theta - \hat{\theta}) \underbrace{p(x, \theta)}_{= p(\theta|x)p(x)} dx d\theta \\ &= \int \left[ \underbrace{\int C(\theta - \hat{\theta}) p(\theta|x) d\theta}_{\triangleq g(\hat{\theta})} \right] p(x) dx\end{aligned}$$

must minimize this for each x value

So... for a given desired cost function...

you have to find the form of the optimal estimator

# The Optimal Estimates for the Typical Costs

1. **Quadratic**:  $\mathcal{R}(\hat{\theta}) = E\left\{\left(\theta - \hat{\theta}\right)^2\right\} = Bmse(\hat{\theta})$

As we saw in Ch. 10

$$\hat{\theta} = E\{\theta | \mathbf{x}\}$$

= mean of  $p(\theta | \mathbf{x})$

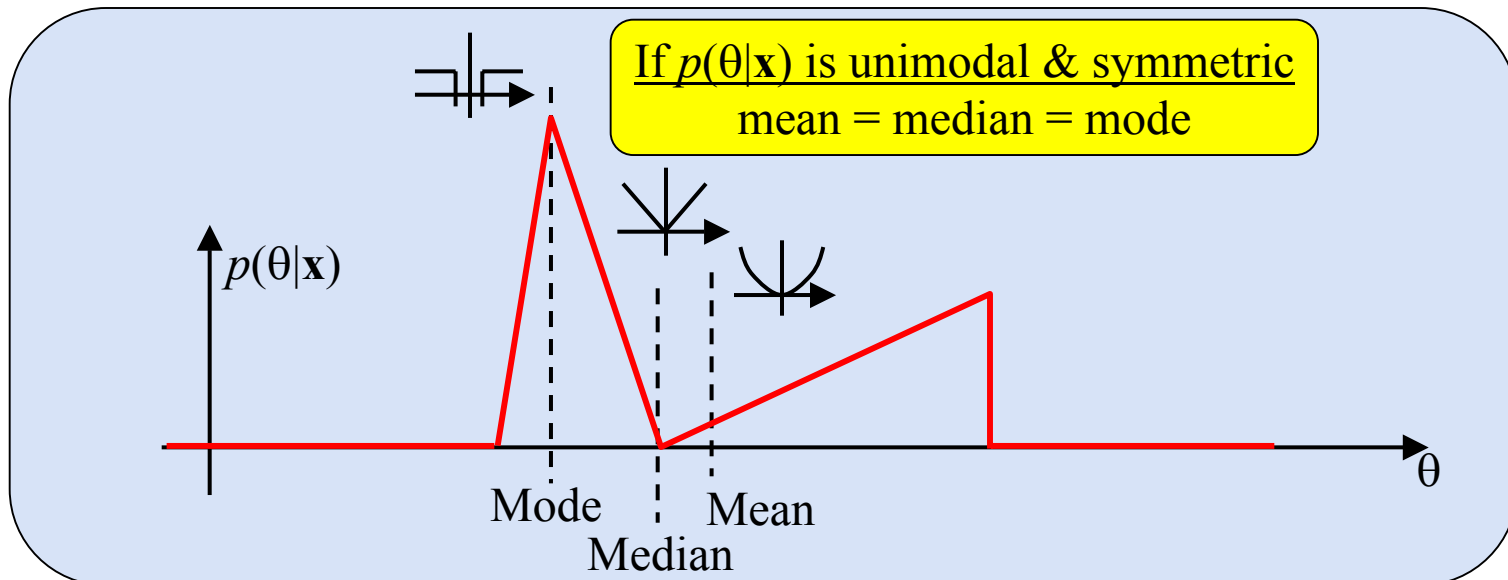
2. **Absolute**:  $\mathcal{R}(\hat{\theta}) = E\left\{|\theta - \hat{\theta}|\right\}$

$$\hat{\theta} = \text{median of } p(\theta | \mathbf{x})$$

3. **Hit-or-Miss**:

$$\hat{\theta} = \text{mode of } p(\theta | \mathbf{x})$$

“Maximum A Posteriori”  
or MAP



# Derivation for Absolute Cost Function

Writing out the function to be minimized gives:

$$\begin{aligned}g(\hat{\theta}) &= \int_{-\infty}^{\infty} |\theta - \hat{\theta}| p(\theta | \mathbf{x}) d\theta \\ &= \underbrace{\int_{-\infty}^{\hat{\theta}} (\hat{\theta} - \theta) p(\theta | \mathbf{x}) d\theta}_{\text{region where } |\theta - \hat{\theta}| = \hat{\theta} - \theta} + \underbrace{\int_{\hat{\theta}}^{\infty} (\theta - \hat{\theta}) p(\theta | \mathbf{x}) d\theta}_{\text{region where } |\theta - \hat{\theta}| = \theta - \hat{\theta}}\end{aligned}$$

Now set  $\frac{\partial g(\hat{\theta})}{\partial \hat{\theta}} = 0$  and use Leibnitz's rule for  $\frac{\partial}{\partial u} \int_{\phi_1(u)}^{\phi_2(u)} h(u, v) dv$

$$\Rightarrow \int_{-\infty}^{\hat{\theta}} p(\theta | \mathbf{x}) d\theta - \int_{\hat{\theta}}^{\infty} p(\theta | \mathbf{x}) d\theta = 0$$

which is satisfied if... (area to the left) = (area to the right)

$\Rightarrow$  Median of conditional PDF



# Derivation for Hit-or-Miss Cost Function

Writing out the function to be minimized gives:

$$\begin{aligned}g(\hat{\theta}) &= \int_{-\infty}^{\infty} C(\theta - \hat{\theta}) p(\theta | \mathbf{x}) d\theta \\&= \int_{-\infty}^{\hat{\theta}-\delta} 1 \cdot p(\theta | \mathbf{x}) d\theta + \int_{\hat{\theta}+\delta}^{\infty} 1 \cdot p(\theta | \mathbf{x}) d\theta \\&= 1 - \int_{\hat{\theta}-\delta}^{\hat{\theta}+\delta} p(\theta | \mathbf{x}) d\theta\end{aligned}$$

Almost all the probability  
= 1 – left out

Maximize this integral

So... center the integral around peak of integrand  
 $\Rightarrow$  Mode of conditional PDF

# 11.4 MMSE Estimators

We've already seen the solution for the scalar parameter case

$$\begin{aligned}\hat{\theta} &= E\{\theta | \mathbf{x}\} \\ &= \text{mean of } p(\theta | \mathbf{x})\end{aligned}$$

Here we'll look at:

- Extension to the vector parameter case
- Analysis of Useful Properties

# Vector MMSE Estimator

The criterion is... minimize the MSE for each component

$$\text{Vector Parameter: } \boldsymbol{\theta} = [\theta_1 \quad \theta_2 \quad \cdots \quad \theta_p]^T$$

$$\text{Vector Estimate: } \hat{\boldsymbol{\theta}} = [\hat{\theta}_1 \quad \hat{\theta}_2 \quad \cdots \quad \hat{\theta}_p]^T$$

is chosen to minimize each of the MSE elements:

$$E\{(\theta_i - \hat{\theta}_i)^2\} = \int (\theta_i - \hat{\theta}_i)^2 p(\mathbf{x}, \theta_i) d\mathbf{x} d\theta_i$$

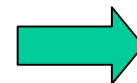
=  $p(\mathbf{x}, \boldsymbol{\theta})$  integrated over all other  $\theta_j$ 's

From the scalar case we know the solution is:

$$\hat{\theta}_i = \int \theta_i p(\mathbf{x}, \theta_i) d\mathbf{x} d\theta_i$$

$$= \int \cdots \int \theta_i p(\mathbf{x}, \theta_1, \dots, \theta_p) d\mathbf{x} d\theta_1 \cdots d\theta_p$$

$$= \int \int \theta_i p(\mathbf{x}, \boldsymbol{\theta}) d\mathbf{x} d\boldsymbol{\theta}$$



$$\hat{\theta}_i = E\{\theta_i | \mathbf{x}\}$$

So... putting all these into a vector gives:

$$\begin{aligned}\hat{\boldsymbol{\theta}} &= \begin{bmatrix} \hat{\theta}_1 & \hat{\theta}_2 & \dots & \hat{\theta}_p \end{bmatrix}^T \\ &= \begin{bmatrix} E\{\theta_1 | \mathbf{x}\} & E\{\theta_2 | \mathbf{x}\} & \dots & E\{\theta_p | \mathbf{x}\} \end{bmatrix}^T \\ &= E\left\{ \begin{bmatrix} \theta_1 & \theta_2 & \dots & \theta_p \end{bmatrix}^T | \mathbf{x} \right\}\end{aligned}$$



$$\hat{\boldsymbol{\theta}} = E\{\boldsymbol{\theta} | \mathbf{x}\}$$

**Vector MMSE Estimate  
= Vector Conditional Mean**

Similarly...  $Bmse(\hat{\theta}_i) = \int [C_{\boldsymbol{\theta}|\mathbf{x}}]_{ii} p(\mathbf{x}) d\mathbf{x} \quad i = 1, \dots, p$

where  $C_{\boldsymbol{\theta}|\mathbf{x}} = E_{\boldsymbol{\theta}|\mathbf{x}} \left\{ [\boldsymbol{\theta} - E\{\boldsymbol{\theta} | \mathbf{x}\}] [\boldsymbol{\theta} - E\{\boldsymbol{\theta} | \mathbf{x}\}]^T \right\}$

## Ex. 11.1 Bayesian Fourier Analysis

Signal model is:  $x[n] = a\cos(2\pi f_o n) + b\sin(2\pi f_o n) + w[n]$

$$\boldsymbol{\theta} = \begin{bmatrix} a \\ b \end{bmatrix} \sim N(\mathbf{0}, \sigma_\theta^2 \mathbf{I})$$

AWGN  
w/ zero mean and  $\sigma^2$

$\boldsymbol{\theta}$  and  $w[n]$  are independent for each  $n$

This is a common propagation model called Rayleigh Fading

Write in matrix form:  $\mathbf{x} = \mathbf{H}\boldsymbol{\theta} + \mathbf{w}$  Bayesian Linear Model

$$\mathbf{H} = \begin{bmatrix} \uparrow & \uparrow \\ \text{cosine} & \text{sine} \\ \downarrow & \downarrow \end{bmatrix}$$

Results from Ch. 10 show that

$$\hat{\boldsymbol{\theta}} = E\{\boldsymbol{\theta} | \mathbf{x}\} = \left[ \frac{1}{\sigma_\theta^2} \mathbf{I} + \frac{\mathbf{H}^T \mathbf{H}}{\sigma^2} \right]^{-1} \frac{\mathbf{H}^T \mathbf{x}}{\sigma^2} \quad \mathbf{C}_{\boldsymbol{\theta}|\mathbf{x}} = \left[ \frac{1}{\sigma_\theta^2} \mathbf{I} + \frac{\mathbf{H}^T \mathbf{H}}{\sigma^2} \right]^{-1}$$

For  $f_o$  chosen such that  $\mathbf{H}$  has orthogonal columns then

$$\hat{\boldsymbol{\theta}} = E\{\boldsymbol{\theta} | \mathbf{x}\} = \begin{bmatrix} \frac{1}{\sigma^2} \\ \frac{1}{\sigma_\theta^2} + \frac{1}{\sigma^2} \end{bmatrix} \mathbf{H}^T \mathbf{x} \quad \begin{matrix} \hat{a} = \beta \left[ \frac{2}{N} \sum_{n=0}^{N-1} x[n] \cos(2\pi f_o n) \right] \\ \hat{b} = \beta \left[ \frac{2}{N} \sum_{n=0}^{N-1} x[n] \sin(2\pi f_o n) \right] \end{matrix} \quad \beta = \frac{1}{1 + \frac{2\sigma^2}{N\sigma_\theta^2}}$$

Fourier Coefficients in the Brackets

Recall: Same form as classical result, except there  $\beta = 1$

*Note:*  $\beta \approx 1$  if  $\sigma_\theta^2 \gg 2\sigma^2/N$

$\Rightarrow$  if prior knowledge is poor, this degrades to classical

# Impact of Poor Prior Knowledge

Conclusion: For poor prior knowledge in Bayesian Linear Model  
MMSE Est.  $\rightarrow$  MVU Est.

Can see this holds in general: Recall that

$$\hat{\boldsymbol{\theta}} = E\{\boldsymbol{\theta} | \mathbf{x}\} = \boldsymbol{\mu}_{\boldsymbol{\theta}} + \left[ \mathbf{C}_{\boldsymbol{\theta}}^{-1} + \mathbf{H}^T \mathbf{C}_{\mathbf{w}}^{-1} \mathbf{H} \right]^{-1} \mathbf{H}^T \mathbf{C}_{\mathbf{w}}^{-1} [\mathbf{x} + \mathbf{H} \boldsymbol{\mu}_{\boldsymbol{\theta}}]$$

For no prior information:  $\mathbf{C}_{\boldsymbol{\theta}}^{-1} \rightarrow \mathbf{0}$  and  $\boldsymbol{\mu}_{\boldsymbol{\theta}} \rightarrow \mathbf{0}$

$$\hat{\boldsymbol{\theta}} \rightarrow \underbrace{\left[ \mathbf{H}^T \mathbf{C}_{\mathbf{w}}^{-1} \mathbf{H} \right]^{-1} \mathbf{H}^T \mathbf{C}_{\mathbf{w}}^{-1} \mathbf{x}}_{\text{MVUE for General Linear Model}}$$

MVUE for General Linear Model

# Useful Properties of MMSE Est.

Will be used for  
Kalman Filter

## 1. Commutes over affine mappings:

If we have  $\alpha = A\theta + \mathbf{b}$  then  $\hat{\alpha} = A\hat{\theta} + \mathbf{b}$

## 2. Additive Property for independent data sets

Assume  $\theta$ ,  $\mathbf{x}_1$ ,  $\mathbf{x}_2$  are jointly Gaussian w/  $\mathbf{x}_1$  and  $\mathbf{x}_2$  independent

$$\hat{\theta} = E\{\theta\} + C_{\theta x_1} C_{x_1}^{-1} [\mathbf{x}_1 - E\{\mathbf{x}_1\}] + C_{\theta x_2} C_{x_2}^{-1} [\mathbf{x}_2 - E\{\mathbf{x}_2\}]$$

a priori Estimate

Update due to  $\mathbf{x}_1$

Update due to  $\mathbf{x}_2$

Proof: Let  $\mathbf{x} = [\mathbf{x}_1^T \ \mathbf{x}_2^T]^T$ . The jointly Gaussian assumption gives:

$$\hat{\theta} = E\{\theta\} + C_{\theta x} C_x^{-1} [\mathbf{x} - E\{\mathbf{x}\}]$$

Indep.  $\Rightarrow$  Block Diagonal

$$= E\{\theta\} + \begin{bmatrix} C_{\theta x_1} & C_{\theta x_2} \end{bmatrix} \begin{bmatrix} C_{x_1}^{-1} & 0 \\ 0 & C_{x_2}^{-1} \end{bmatrix} \begin{bmatrix} \mathbf{x}_1 - E\{\mathbf{x}_1\} \\ \mathbf{x}_2 - E\{\mathbf{x}_2\} \end{bmatrix}$$

Simplify to  
get the result

## 3. Jointly Gaussian case leads to a linear estimator: $\hat{\theta} = P\mathbf{x} + \mathbf{m}$