

Stochastic Signals and Systems

Contents

- 1 Probability Theory
- 2 Stochastic Processes
- 3 Parameter Estimation**
- 4 Signal Detection
- 5 Spectrum Analysis
- 6 Optimal Filtering

3	Parameter Estimation	3
3.1	Estimating Function and Estimator	4
3.2	Sufficient Statistic, Exponential Family	15
3.3	Linear Least Squares Estimation	29
3.4	Confidence Intervals	53
3.5	Cramer-Rao Lower Bound	57
3.6	Maximum Likelihood Estimation	75
3.7	Bayesian Estimation	80
3.7.1	Minimum Mean Square Error Estimation	84
3.7.2	Minimum Mean Absolute Error Estimation	91
3.7.3	Maximum A Posteriori Estimation	93
	References to Chapter 3	97

3 Parameter Estimation

The observations $(x_1, \dots, x_n)^T = \mathbf{x}$ are realizations of the random variables $(X_1, \dots, X_n)^T = \mathbf{X}$ with density $f_{\mathbf{x}}(\mathbf{x})$, which is element of a known set $\{f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}): \boldsymbol{\theta} \in \Omega\}$ and whose parameter vector $\boldsymbol{\theta} = (\theta_1, \dots, \theta_p)^T$ is unknown.

Problem:

For given observations \mathbf{x} we are looking for an estimate $\hat{\boldsymbol{\theta}} = (\hat{\theta}_1, \dots, \hat{\theta}_p)^T$ of $\boldsymbol{\theta}$ that depends on \mathbf{x} , i.e. $\hat{\boldsymbol{\theta}} = \hat{\boldsymbol{\theta}}(\mathbf{x})$.

In parameter estimation problems one distinguishes whether the parameter vector is a perfectly unknown quantity or whether the prior density of the parameter vector is supposed to be known.

In the latter case one interprets the parameter vector as random vector, whose posterior density is used for the subsequent statistical inference. This approach leads to Bayes estimators.

3.1 Estimating Function and Estimator

For determining θ_i , we denote the mapping

$$(x_1, \dots, x_n) \mapsto \hat{\theta}_i(x_1, \dots, x_n), \quad i = 1, \dots, p$$

a estimating function for θ_i and its function value for a given set of observations an estimate. Since x_1, \dots, x_n can be considered as random values the estimates are also random and therefore only approximate the true values θ_i .

We examine the accuracy properties of the estimates on the basis of the corresponding random variables.

Therefore, we define the random variable

$$\hat{\Theta}_i = \hat{\theta}_i(X_1, \dots, X_n), \quad i = 1, \dots, p,$$

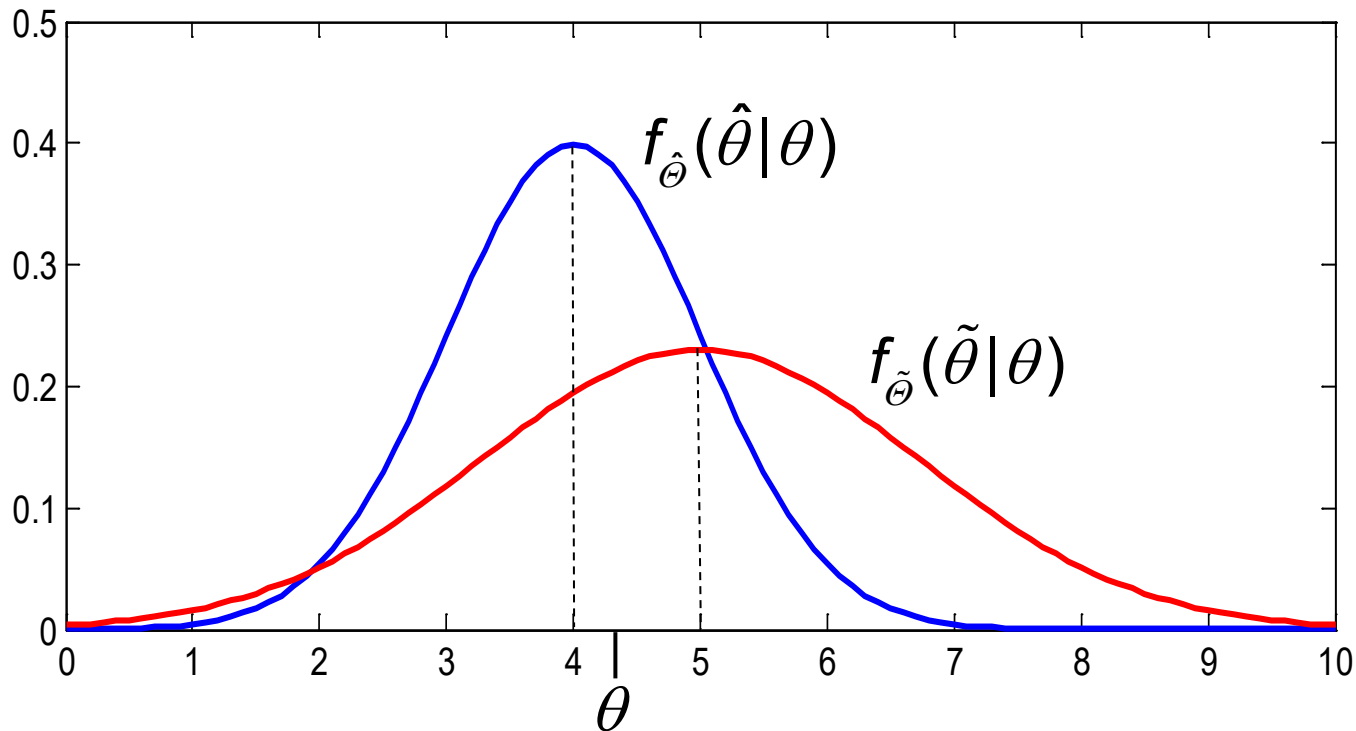
and denote it as estimator of θ_i .

To characterize the properties of the estimator

$$\hat{\Theta} = (\hat{\Theta}_1, \dots, \hat{\Theta}_p)^T = \hat{\Theta}(X_1, \dots, X_n)$$

completely, the density function $f_{\hat{\Theta}}(\hat{\Theta} | \Theta)$ that generally depends on Θ has to be determined, e.g. using the methods described in Chapter 1.7.

For simplicity let $p = 1$. Exemplarily, the densities of two alternative estimators $\hat{\Theta}$ and $\tilde{\Theta}$ for θ are depicted below.



The determination of the density function can be rather complicated. Therefore, one must be often satisfied with the determination of the moments of an estimator.

Bias:

The bias or systematic error is define by

$$\begin{aligned} \mathbf{b}(\hat{\Theta}) &= \left(b(\hat{\Theta}_1), \dots, b(\hat{\Theta}_p) \right)^T = \left(\mathbf{E}(\hat{\Theta}_1) - \theta_1, \dots, \mathbf{E}(\hat{\Theta}_p) - \theta_p \right)^T \\ &= \left(\mathbf{E}(\hat{\Theta}_1), \dots, \mathbf{E}(\hat{\Theta}_p) \right)^T - \left(\theta_1, \dots, \theta_p \right)^T = \mathbf{E}(\hat{\Theta}) - \theta, \end{aligned}$$

where

$$\mathbf{E}(\hat{\Theta}_i) = \int_{-\infty}^{\infty} \dots \int_{-\infty}^{\infty} \hat{\theta}_i(x_1, \dots, x_n) f_{X_1, \dots, X_n}(x_1, \dots, x_n | \theta) dx_1 \dots dx_n.$$

If $\mathbf{b}(\hat{\Theta}) = \mathbf{0}$ the estimator $\hat{\Theta}$ is said to be unbiased.

Covariance:

The covariance matrix is define by

$$\begin{aligned} \text{Cov}(\hat{\Theta}) &= \left[\text{Cov}(\hat{\Theta}_i, \hat{\Theta}_j) \right]_{i,j=1,\dots,p} \\ &= \left[\text{E} \left(\left(\hat{\Theta}_i - \text{E}(\hat{\Theta}_i) \right) \left(\hat{\Theta}_j - \text{E}(\hat{\Theta}_j) \right) \right) \right]_{i,j=1,\dots,p} \\ &= \text{E} \left(\left(\hat{\Theta} - \text{E}(\hat{\Theta}) \right) \left(\hat{\Theta} - \text{E}(\hat{\Theta}) \right)^T \right), \end{aligned}$$

where the diagonal elements represent the variances

$$\begin{aligned} \text{Var}(\hat{\Theta}_i) &= \text{Cov}(\hat{\Theta}_i, \hat{\Theta}_i) \\ &= \text{E}(\hat{\Theta}_i - \text{E} \hat{\Theta}_i)^2 = \text{E}(\hat{\Theta}_i^2) - (\text{E} \hat{\Theta}_i)^2 \end{aligned}$$

for $i = 1, \dots, p$.

Mean Square Error:

The matrix-valued mean square error is define by

$$\begin{aligned}\text{MSE}(\hat{\Theta}) &= \left[\text{E} \left((\hat{\Theta}_i - \theta_i)(\hat{\Theta}_j - \theta_j) \right) \right]_{i,j=1,\dots,p} \\ &= \text{E} \left((\hat{\Theta} - \theta)(\hat{\Theta} - \theta)^T \right) \\ &= \text{Cov}(\hat{\Theta}) + \mathbf{b}(\hat{\Theta})\mathbf{b}(\hat{\Theta})^T,\end{aligned}$$

where the diagonal elements represent the mean square errors of the individual components $\hat{\Theta}_i$ given by

$$\text{MSE}(\hat{\Theta}_i) = \text{E} \left((\hat{\Theta}_i - \theta_i)^2 \right) = \text{Var}(\hat{\Theta}_i) + b(\hat{\Theta}_i)^2$$

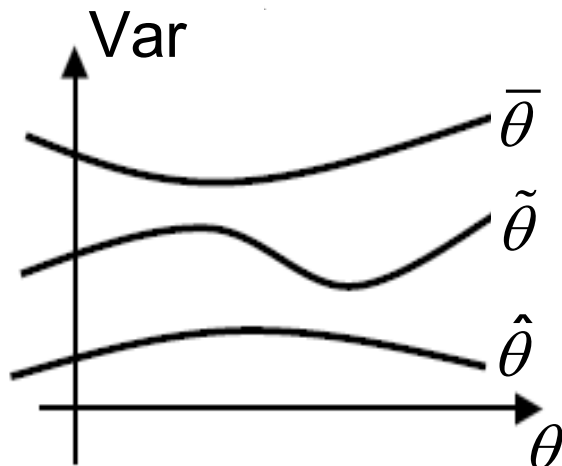
for $i = 1, \dots, p$.

Minimum Variance Unbiased (MVU) Estimator:

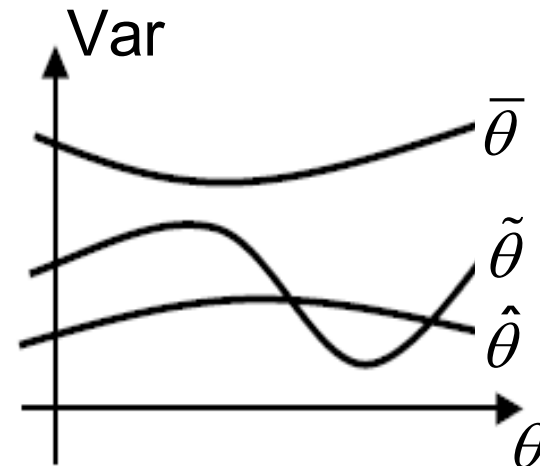
An unbiased estimator $\hat{\Theta}$ is MVU if for any unbiased estimator $\tilde{\Theta}$ the following inequality holds.

$$\mathbf{a}^T \text{Cov}(\hat{\Theta}) \mathbf{a} \leq \mathbf{a}^T \text{Cov}(\tilde{\Theta}) \mathbf{a} \quad \forall \mathbf{a} \in \mathbb{R}^p, \quad \forall \boldsymbol{\theta} \in \Omega$$

MVU exists



no MVU exists



Note:

To emphasize that the variance is smallest for all $\boldsymbol{\theta} \in \Omega$ an MVU estimator is sometimes called uniformly minimum variance unbiased UMVU estimator.

Minimum Mean Square Error (MMSE) Estimator:

An estimator $\hat{\boldsymbol{\Theta}}$ provides an MMSE if for any estimator $\tilde{\boldsymbol{\Theta}}$ the following inequality holds.

$$\mathbf{a}^T \text{MSE}(\hat{\boldsymbol{\Theta}}) \mathbf{a} \leq \mathbf{a}^T \text{MSE}(\tilde{\boldsymbol{\Theta}}) \mathbf{a} \quad \forall \mathbf{a} \in \mathbb{R}^p, \quad \forall \boldsymbol{\theta} \in \Omega$$

Linear Estimator:

An estimator $\hat{\boldsymbol{\Theta}}$ is called linear if it can be expressed as a linear function of the observations, i.e.

$$\hat{\boldsymbol{\Theta}} = \mathbf{A} \mathbf{X} \quad \text{with} \quad \mathbf{X} = (X_1, \dots, X_n) \quad \text{and} \quad \mathbf{A} = (a_{i,j})_{i=1, \dots, p; j=1, \dots, n}$$

Exercise 3.1-1:
MVU and MMSE variance estimator

Consistency:

Often one is interested in the behavior of an estimator if the number of observations grows.

An estimator is said to be

- strongly consistent if $\hat{\Theta}_n$ converges with probability 1 towards Θ , i.e.

$$\hat{\Theta}_n = \hat{\Theta}(X_1, \dots, X_n) \xrightarrow[n \rightarrow \infty]{a.s.} \Theta$$

- mean square consistent if $\hat{\Theta}_n$ converges in mean square sense towards Θ , i.e.

$$\hat{\Theta}_n = \hat{\Theta}(X_1, \dots, X_n) \xrightarrow[n \rightarrow \infty]{m.s.} \Theta$$

- consistent if $\hat{\Theta}_n$ converges in probability towards Θ , i.e.

$$\hat{\Theta}_n = \hat{\Theta}(X_1, \dots, X_n) \xrightarrow[n \rightarrow \infty]{P} \Theta$$

Asymptotic Normality:

In certain cases one can show that an estimator is asymptotically normally distributed such that

$$\lim_{n \rightarrow \infty} \sqrt{n}(\hat{\Theta}_n - \boldsymbol{\theta}) \sim \mathcal{N}_p(\mathbf{0}, \boldsymbol{\Sigma}).$$

Consequently, for large n the density function of $\hat{\Theta}_n$ can be approximated by

$$f_{\hat{\Theta}_n}(\hat{\boldsymbol{\theta}} | \boldsymbol{\theta}) = \frac{n^{p/2}}{(2\pi)^{p/2} \sqrt{\det(\boldsymbol{\Sigma})}} \cdot \exp \left\{ -\frac{n}{2} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta})^T \boldsymbol{\Sigma}^{-1} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}) \right\}.$$

3.2 Sufficient Statistic, Exponential Families

A function $\mathbf{T} = \mathbf{t}(\mathbf{X})$ that is only depending on the observation model \mathbf{X} is called a statistic.

Let $\mathbf{X} = (X_1, \dots, X_n)^T$ be a model of a binary sequence with probability p and $1-p$ of observing a one and a zero respectively. Furthermore, we suppose that the X_1, \dots, X_n are independent.

Thus, our model can be described by the probabilities

$$\begin{aligned} P(\mathbf{X} = \mathbf{x} \mid p) &= \prod_{k=1}^n p^{x_k} (1-p)^{1-x_k} \\ &= p^{\sum_{k=1}^n x_k} (1-p)^{n-\sum_{k=1}^n x_k} = p^{t(\mathbf{x})} (1-p)^{n-t(\mathbf{x})} \end{aligned}$$

Now the following question arises. Do we collect all information in terms of inference about p by recording only

$$t(\mathbf{x}) = \sum_{k=1}^n x_k$$

The statistical understanding about the collection of all information is quantified in the following definition.

Definition:

A statistic $\mathbf{T} = \mathbf{t}(\mathbf{X})$ is called sufficient for the parameter $\boldsymbol{\theta}$ if the conditional distribution of \mathbf{X} given $\mathbf{T} = \mathbf{t}(\mathbf{x})$ is independent of $\boldsymbol{\theta}$ for all \mathbf{t} , i.e.

$$F_{\mathbf{x}}(\mathbf{x} | \mathbf{T} = \mathbf{t}(\mathbf{x}); \boldsymbol{\theta}) = F_{\mathbf{x}}(\mathbf{x} | \mathbf{T} = \mathbf{t}(\mathbf{x})).$$

Hence, \mathbf{T} contains "all information" about $\boldsymbol{\theta}$ included in \mathbf{x} .

Exercise 3.2-1:

Sufficiency of $t(\mathbf{x}) = \sum_{k=1}^n x_k$

Because the conditional distribution has to be determined a direct evaluation of sufficiency is usually difficult.

Fortunately, the following theorem exist whose conditions can be verified easily.

Theorem: (*factorization theorem for densities*)

A necessary and sufficient condition for a statistic $\mathbf{T} = \mathbf{t}(\mathbf{X})$ to be sufficient is that there exist non-negative functions $g(\mathbf{t} | \boldsymbol{\theta})$ and $h(\mathbf{x})$ such that $f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta})$ satisfies

$$f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}) = g(\mathbf{t}(\mathbf{x}) | \boldsymbol{\theta}) \cdot h(\mathbf{x}).$$

Exercise 3.2-2:
Proof of the Theorem

Exercise 3.2-3:
Sufficient statistic for mean and variance

Minimal Sufficient Statistic

A sufficient statistic \mathbf{T} is said to be minimal if of all sufficient statistics it provides the greatest possible reduction of data, i.e. if for any sufficient statistic \mathbf{T}' there exists a function \mathbf{s} such that $\mathbf{T} = \mathbf{s}(\mathbf{T}')$.

Complete Sufficient Statistic

A sufficient statistic \mathbf{T} is said to be complete (or unique) if the condition

$$E_{\theta}(\mathbf{f}(\mathbf{T})) = \mathbf{0} \quad \forall \theta \in \Omega$$

implies $\mathbf{f}(\mathbf{T}) = \mathbf{0}$ with probability 1 for all θ .

Note: Completeness ensures minimality

Exercise 3.2-4:
Minimal and complete sufficient statistics

Uniqueness of unbiased estimating function

Completeness implies that there is just one estimating function of the sufficient statistic that provides an unbiased estimator of $\boldsymbol{\theta}$.

Let \mathbf{T} be a complete sufficient statistic and \mathbf{f}_1 and \mathbf{f}_2 be two functions such that

$$E_{\boldsymbol{\theta}}(\mathbf{f}_1(\mathbf{T})) = E_{\boldsymbol{\theta}}(\mathbf{f}_2(\mathbf{T})) = \boldsymbol{\theta} \quad \forall \boldsymbol{\theta} \in \Omega.$$

Then

$$E_{\boldsymbol{\theta}}(\mathbf{f}_1(\mathbf{T}) - \mathbf{f}_2(\mathbf{T})) = E_{\boldsymbol{\theta}}(\mathbf{f}(\mathbf{T})) = \mathbf{0} \quad \forall \boldsymbol{\theta} \in \Omega$$

and due to the completeness of \mathbf{T}

$$\mathbf{f}(\mathbf{T}) = \mathbf{0} \Rightarrow \mathbf{f}_1(\mathbf{T}) = \mathbf{f}_2(\mathbf{T}) \quad \text{with probability 1} \quad \forall \boldsymbol{\theta} \in \Omega$$

Theorem: (Rao-Blackwell)

Let $\tilde{\Theta}$ be an unbiased estimator of θ and $\mathbf{T} = \mathbf{t}(\mathbf{X})$ be a sufficient statistic for θ . Then the estimator defined by

$$\hat{\Theta} = E(\tilde{\Theta} | \mathbf{T})$$

is unbiased and improves on $\tilde{\Theta}$ as follows.

$$\mathbf{a}^T \text{Cov}(\hat{\Theta}) \mathbf{a} \leq \mathbf{a}^T \text{Cov}(\tilde{\Theta}) \mathbf{a} \quad \forall \mathbf{a} \in \mathbb{R}^p$$

Theorem: (Lehmann-Scheffe)

If in addition to the assumptions employed for the Rao-Blackwell theorem the sufficient statistic is complete, then $\hat{\Theta} = E(\tilde{\Theta} | \mathbf{T})$ is unique MVU estimator of θ .

Exercise 3.2-5:
Proof of the Theorems

Exponential Families

A family $\{F_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta})\}$ of distributions is forming a k -dimensional exponential family if the distributions $F_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta})$ have densities of the form

$$f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}) = h(\mathbf{x}) \cdot \exp\left(\sum_{i=1}^k \xi_i(\boldsymbol{\theta}) t_i(\mathbf{x}) - B(\boldsymbol{\theta})\right).$$

Frequently, it is more convenient to use the ξ_i as the parameters and write the density in the canonical form

$$f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\xi}) = h(\mathbf{x}) \cdot \exp\left(\sum_{i=1}^k \xi_i t_i(\mathbf{x}) - A(\boldsymbol{\xi})\right).$$

Applying the factorization theorem for densities one can easily observe that

$$\mathbf{T} = (T_1, \dots, T_k)^T = (t_1(\mathbf{X}), \dots, t_k(\mathbf{X}))^T$$

constitutes a sufficient statistic for the exponential family.

Note:

The parameter space $\Xi \subset \mathbb{R}^k$ of the natural parameter vector $\boldsymbol{\xi} = (\xi_1, \dots, \xi_k)^T$ is convex.

If the exponential family is of full rank, i.e. the parameter space Ξ contains a k -dimensional rectangle, then $\mathbf{T} = (T_1, \dots, T_k)^T$ is complete.

For exponential families one can claim, that the integral

$$\int_{\mathbb{R}^n} f_{\mathbf{X}}(\mathbf{x} | \boldsymbol{\xi}) d\mathbf{x} = \int_{\mathbb{R}^n} h(\mathbf{x}) \cdot \exp\left(\sum_{i=1}^k \xi_i t_i(\mathbf{x}) - A(\boldsymbol{\xi})\right) d\mathbf{x} = 1$$

has derivatives of all orders with respect to the ξ_i which can be obtained by differentiating under the integral sign.

Exploiting the properties of the integral we can find

$$E(T_i) = E(t_i(\mathbf{X})) = \frac{\partial}{\partial \xi_i} A(\boldsymbol{\xi})$$

and

$$\text{Cov}(T_i, T_j) = \text{Cov}(t_i(\mathbf{X}), t_j(\mathbf{X})) = \frac{\partial^2}{\partial \xi_i \partial \xi_j} A(\boldsymbol{\xi})$$

Exercise 3.2-6:

Rayleigh distribution, mean and variance of the natural sufficient statistic

3.3 Linear Least Squares Estimation

Consider the linear model

$$\mathbf{X} = \mathbf{H}\boldsymbol{\theta} + \mathbf{Z} \quad \text{with} \quad X_i = h_{i1}\theta_1 + \dots + h_{ip}\theta_p + Z_i, \quad i = 1, \dots, n,$$

where \mathbf{X} and \mathbf{Z} are $n \times 1$ vectors modeling the measurements and the measurement noise, respectively. Furthermore, \mathbf{H} denotes a known $n \times p$ matrix and $\boldsymbol{\theta}$ the $p \times 1$ parameter vector that has to be estimated.

The measurement noise is supposed to be statistically characterized by $E(\mathbf{Z}) = \mathbf{0}$ and $\text{Cov}(\mathbf{Z}) = E(\mathbf{Z}\mathbf{Z}^T) = \sigma_Z^2 \mathbf{I}$.

Hence, the measurement model \mathbf{X} possess the mean vector $E(\mathbf{X}) = \mathbf{H}\boldsymbol{\theta}$ and covariance matrix $\text{Cov}(\mathbf{X}) = \sigma_Z^2 \mathbf{I}$.

Exercise 3.3-1:

Model examples,

- *polynomial curve fitting,*
- *amplitude and phase estimation of sinusoids*
- *FIR filter identification*

Now, the least squares criterion can be expressed by

$$q(\boldsymbol{\theta}) = \sum_{i=1}^n \left(x_i - (h_{i1}\theta_1 + \dots + h_{ip}\theta_p) \right)^2$$
$$= (\mathbf{x} - \mathbf{H}\boldsymbol{\theta})^T (\mathbf{x} - \mathbf{H}\boldsymbol{\theta}) = \mathbf{x}^T \mathbf{x} - 2\boldsymbol{\theta}^T \mathbf{H}^T \mathbf{x} + \boldsymbol{\theta}^T \mathbf{H}^T \mathbf{H} \boldsymbol{\theta}.$$

Differentiating with respect to $\boldsymbol{\theta}$ by using the identities

$$\nabla_{\boldsymbol{\theta}} \boldsymbol{\theta}^T \mathbf{A} \mathbf{a} = \mathbf{A} \mathbf{a}, \quad \nabla_{\boldsymbol{\theta}} \boldsymbol{\theta}^T \mathbf{A} \boldsymbol{\theta} = (\mathbf{A} + \mathbf{A}^T) \boldsymbol{\theta}$$

results after equating to zero in the following so-called normal equation system

$$\mathbf{H}^T \mathbf{H} \hat{\boldsymbol{\theta}} = \mathbf{H}^T \mathbf{x}.$$

A solution of the normal equation system is given by

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

where $(\mathbf{H}^T \mathbf{H})^-$ denotes a generalized inverse.

- 1) If the $\text{rank}(\mathbf{H}) = p$, the number of unknown parameters, then $(\mathbf{H}^T \mathbf{H})^- = (\mathbf{H}^T \mathbf{H})^{-1}$ is the ordinary inverse.
- 2) Let $\text{rank}(\mathbf{H}) = M < p$, i.e. either $n < p$ or the columns of \mathbf{H} are linearly dependent, then

$$(\mathbf{H}^T \mathbf{H})^- = (\mathbf{H}^T \mathbf{H})^+ = \sum_{m=1}^M \frac{1}{\lambda_m} \mathbf{u}_m \mathbf{u}_m^T$$

is the Moore-Penrose inverse, where $\lambda_1, \dots, \lambda_M$ and $\mathbf{u}_1, \dots, \mathbf{u}_M$ denote the non-zero eigenvalues and corresponding eigenvectors of $\mathbf{H}^T \mathbf{H}$, respectively.

Supplement:

A generalized inverse of \mathbf{A} is defined by the property

$$\mathbf{A}\mathbf{A}^-\mathbf{A} = \mathbf{A}.$$

It is not unique. The Moore-Penrose inverse of \mathbf{A} is defined by the properties

$$\begin{aligned} \mathbf{A}\mathbf{A}^+\mathbf{A} &= \mathbf{A}, & \mathbf{A}^+\mathbf{A}\mathbf{A}^+ &= \mathbf{A}^+, \\ (\mathbf{A}\mathbf{A}^+)^T &= \mathbf{A}\mathbf{A}^+ & (\mathbf{A}^+\mathbf{A})^T &= \mathbf{A}^+\mathbf{A}. \end{aligned}$$

It is unique and provides the minimum length solution of the linear equation system $\mathbf{A}\mathbf{x} = \mathbf{b}$.

The Moore-Penrose inverse \mathbf{A}^+ can be determined by employing the singular value decomposition of \mathbf{A} .

Exercise 3.3-2:
Normal equation system and generalized inverse

Properties of the least squares estimator

The mean vector is given by

$$\begin{aligned}
 E(\hat{\boldsymbol{\Theta}}) &= (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T E(\mathbf{X}) = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{H} \boldsymbol{\theta} \\
 &= \begin{cases} \boldsymbol{\theta} & \text{rank}(\mathbf{H}) = p \\ \mathbf{V}_1 \mathbf{V}_1^T \boldsymbol{\theta} = (\mathbf{I} - \mathbf{V}_2 \mathbf{V}_2^T) \boldsymbol{\theta} & \text{rank}(\mathbf{H}) = M < p \end{cases}
 \end{aligned}$$

where the matrices \mathbf{V}_1 and \mathbf{V}_2 can be derived from the singular value decomposition

$$\mathbf{H}^T \mathbf{H} = \mathbf{V} \text{diag}(\lambda_1, \dots, \lambda_M, 0, \dots, 0) \mathbf{V}^T$$

with

$$\mathbf{V} = (\mathbf{V}_1, \mathbf{V}_2), \quad \mathbf{V}_1 = (\mathbf{v}_1, \dots, \mathbf{v}_M) \quad \mathbf{V}_2 = (\mathbf{v}_{M+1}, \dots, \mathbf{v}_p).$$

The covariance matrix of the least squares estimator

$$\mathbf{C}_{\hat{\Theta}\hat{\Theta}} = \text{Cov}(\hat{\Theta}) = E\left(\left(\hat{\Theta} - E(\hat{\Theta})\right)\left(\hat{\Theta} - E(\hat{\Theta})\right)^T\right)$$

results after exploiting

$$\hat{\Theta} = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{X} \quad \text{and} \quad E(\hat{\Theta}) = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{H} \boldsymbol{\theta}$$

in the expression

$$\begin{aligned} \mathbf{C}_{\hat{\Theta}\hat{\Theta}} &= E\left(\left(\mathbf{H}^T \mathbf{H}\right)^{-1} \mathbf{H}^T (\mathbf{X} - \mathbf{H}\boldsymbol{\theta})(\mathbf{X} - \mathbf{H}\boldsymbol{\theta})^T \mathbf{H} \left(\mathbf{H}^T \mathbf{H}\right)^{-1}\right) \\ &= \left(\mathbf{H}^T \mathbf{H}\right)^{-1} \mathbf{H}^T E\left((\mathbf{X} - \mathbf{H}\boldsymbol{\theta})(\mathbf{X} - \mathbf{H}\boldsymbol{\theta})^T\right) \mathbf{H} \left(\mathbf{H}^T \mathbf{H}\right)^{-1} \\ &= \left(\mathbf{H}^T \mathbf{H}\right)^{-1} \mathbf{H}^T (\sigma_z^2 \mathbf{I}) \mathbf{H} \left(\mathbf{H}^T \mathbf{H}\right)^{-1} = \sigma_z^2 \left(\mathbf{H}^T \mathbf{H}\right)^{-1}. \end{aligned}$$

Exercise 3.3-3:
LSE for $p = 1$ and sample mean

Theorem: (Gauß-Markov)

Given the model

$$\mathbf{X} = \mathbf{H}\boldsymbol{\theta} + \mathbf{Z} \quad \text{with} \quad X_i = h_{i1}\theta_1 + \dots + h_{ip}\theta_p + Z_i, \quad i = 1, \dots, n,$$

where

$$\mathbf{E}(\mathbf{X}) = \mathbf{H}\boldsymbol{\theta}, \quad \mathbf{Cov}(\mathbf{X}) = \mathbf{Cov}(\mathbf{Z}) = \sigma_Z^2 \mathbf{I}$$

and

$$\text{rank}(\mathbf{H}) = p.$$

Then the best linear unbiased estimator (BLUE) is the least squares estimator

$$\hat{\boldsymbol{\theta}} = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{X}.$$

Exercise 3.3-4:
Proof of the Gauß-Markov Theorem

Consequently, the minimum of the sum of squares is

$$\begin{aligned} q(\hat{\boldsymbol{\theta}}) &= \left(\mathbf{x} - \mathbf{H}(\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{x} \right)^T \left(\mathbf{x} - \mathbf{H}(\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{x} \right) \\ &= (\mathbf{x} - \mathbf{P}\mathbf{x})^T (\mathbf{x} - \mathbf{P}\mathbf{x}) = \mathbf{x}^T (\mathbf{I} - \mathbf{P})^T (\mathbf{I} - \mathbf{P}) \mathbf{x} \\ &= \mathbf{x}^T \mathbf{P}^\perp \mathbf{P}^\perp \mathbf{x} = \mathbf{x}^T \mathbf{P}^\perp \mathbf{x}, \end{aligned}$$

where

$$\mathbf{P} = \mathbf{H}(\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \quad \text{and} \quad \mathbf{P}^\perp = \mathbf{I} - \mathbf{P} = \mathbf{I} - \mathbf{H}(\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T$$

are projection matrices, which project a vector $\mathbf{a} \in \mathbb{R}^n$ by $\mathbf{P}\mathbf{a}$ and $\mathbf{P}^\perp \mathbf{a}$ into $R(\mathbf{H})$ and $N(\mathbf{H}^T)$, respectively.

$R(\mathbf{H})$: range of \mathbf{H} , i.e. the space spanned by the columns of \mathbf{H}

$N(\mathbf{H}^T)$: nullspace of \mathbf{H}^T , i.e. the space that is orthogonal to $R(\mathbf{H})$

Thus, $\mathbf{H}^T \mathbf{P}^\perp = \mathbf{0}$, $\mathbf{P}^\perp \mathbf{H} = \mathbf{0}$ and $\mathbf{P} \mathbf{P}^\perp = \mathbf{P}^\perp \mathbf{P} = \mathbf{0}$. The projection matrices are also symmetric and idempotent, i.e.

$$\mathbf{P} = \mathbf{P}^T, \quad \mathbf{P} = \mathbf{P} \mathbf{P} \quad \text{and} \quad \mathbf{P}^\perp = \mathbf{P}^{\perp T}, \quad \mathbf{P}^\perp = \mathbf{P}^\perp \mathbf{P}^\perp.$$

Moreover, by employing the trace of a square matrix

$$\text{tr}(\mathbf{A}) = \sum_{i=1}^n a_{ii} \quad \text{with} \quad \mathbf{A} = (a_{ij})_{i,j=1,\dots,n}$$

together with its property

$$\text{tr}(\mathbf{ABC}) = \text{tr}(\mathbf{BCA}) = \text{tr}(\mathbf{CAB}),$$

where \mathbf{A} , \mathbf{B} and \mathbf{C} are $n \times k$, $k \times l$ and $l \times n$ matrices the minimum of the sum of squares can be expressed by

$$q(\hat{\boldsymbol{\theta}}) = \mathbf{x}^T \mathbf{P}^\perp \mathbf{x} = \text{tr}(\mathbf{P}^\perp \mathbf{x} \mathbf{x}^T).$$

Now we want to consider how the variance of the noise model σ_Z^2 can be estimated. The expected value of the minimum of the sum of squares provides

$$\begin{aligned} E(q(\hat{\Theta})) &= E(\text{tr}(\mathbf{P}^\perp \mathbf{X} \mathbf{X}^T)) = \text{tr}(\mathbf{P}^\perp E((\mathbf{H}\boldsymbol{\theta} + \mathbf{Z})(\mathbf{H}\boldsymbol{\theta} + \mathbf{Z})^T)) \\ &= \text{tr}(\mathbf{P}^\perp (\mathbf{H}\boldsymbol{\theta}\boldsymbol{\theta}^T \mathbf{H}^T + E(\mathbf{Z})\boldsymbol{\theta}^T \mathbf{H}^T + \mathbf{H}\boldsymbol{\theta} E(\mathbf{Z}^T) + E(\mathbf{Z}\mathbf{Z}^T))) \\ &= \text{tr}(\mathbf{P}^\perp (\mathbf{H}\boldsymbol{\theta}\boldsymbol{\theta}^T \mathbf{H}^T + \sigma_Z^2 \mathbf{I})) = \sigma_Z^2 \text{tr}(\mathbf{P}^\perp) = (n - p)\sigma_Z^2, \end{aligned}$$

where

$$\begin{aligned} \text{tr}(\mathbf{P}^\perp) &= \text{tr}(\mathbf{I}_n - \mathbf{H}(\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T) \\ &= n - \text{tr}((\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{H}) = n - \text{tr}(\mathbf{I}_p) = n - p \end{aligned}$$

has been exploited.

This result motivates

$$S^2 = q(\hat{\Theta}) / (n - p)$$

as an unbiased estimator for σ_Z^2 .

Theorem:

If in addition to the Gauß-Markov theorem \mathbf{Z} obeys a normal distribution, i.e. $\mathbf{Z} \sim \mathcal{N}_n(\mathbf{0}, \sigma_Z^2 \mathbf{I})$, the following holds.

- a) $\hat{\Theta} = (\mathbf{H}^T \mathbf{H})^{-1} \mathbf{H}^T \mathbf{X} \sim \mathcal{N}_p(\boldsymbol{\theta}, \sigma_Z^2 (\mathbf{H}^T \mathbf{H})^{-1})$
- b) $\hat{\Theta}$ and S^2 are stochastically independent
- c) $(n - p) / \sigma_Z^2 \cdot S^2 \sim \chi_{n-p}^2$

Exercise 3.3-5:
Proof of the Theorem

Exercise 3.3-6:
Amplitude and phase estimation of sinusoids

Exercise 3.3-7:
System identification (FIR filter)

In the following we generalize the linear model to

$$\mathbf{X} = \mathbf{H}\boldsymbol{\theta} + \mathbf{U} \quad \text{with} \quad X_i = h_{i1}\theta_1 + \dots + h_{ip}\theta_p + U_i, \quad i = 1, \dots, n,$$

where the measurement noise possesses still the mean $E(\mathbf{U}) = \mathbf{0}$ but now the covariance matrix $\text{Cov}(\mathbf{U}) = \sigma_U^2 \mathbf{C}_{UU}$.

Prewhitening

If \mathbf{C}_{UU} is decomposed, e.g. by the Cholesky decomposition $\mathbf{C}_{UU} = \mathbf{C}\mathbf{C}^T$, and

$$\mathbf{Y} = \mathbf{C}^{-1} \mathbf{X}, \quad \mathbf{K} = \mathbf{C}^{-1} \mathbf{H}, \quad \mathbf{W} = \mathbf{C}^{-1} \mathbf{U}$$

be introduced, the linear model can be reformulated to

$$\mathbf{Y} = \mathbf{C}^{-1} \mathbf{X} = \mathbf{C}^{-1} \mathbf{H}\boldsymbol{\theta} + \mathbf{C}^{-1} \mathbf{U} = \mathbf{K}\boldsymbol{\theta} + \mathbf{W}$$

with

$$E(\mathbf{Y}) = \mathbf{K}\boldsymbol{\theta} = \mathbf{C}^{-1} \mathbf{H}\boldsymbol{\theta} \quad \text{and} \quad \text{Cov}(\mathbf{Y}) = \text{Cov}(\mathbf{W}) = \sigma_U^2 \mathbf{I}.$$

Hence, the weighted least squares criterion for the generalized linear model can be derived as follows.

$$\begin{aligned}
 q(\boldsymbol{\theta}) &= (\mathbf{y} - \mathbf{K}\boldsymbol{\theta})^T (\mathbf{y} - \mathbf{K}\boldsymbol{\theta}) \\
 &= (\mathbf{C}^{-1}(\mathbf{x} - \mathbf{H}\boldsymbol{\theta}))^T (\mathbf{C}^{-1}(\mathbf{x} - \mathbf{H}\boldsymbol{\theta})) \\
 &= (\mathbf{x} - \mathbf{H}\boldsymbol{\theta})^T (\mathbf{C}^{-1})^T \mathbf{C}^{-1}(\mathbf{x} - \mathbf{H}\boldsymbol{\theta}) = (\mathbf{x} - \mathbf{H}\boldsymbol{\theta})^T \mathbf{C}_{UU}^{-1}(\mathbf{x} - \mathbf{H}\boldsymbol{\theta}).
 \end{aligned}$$

The normal equation system

$$\mathbf{K}^T \mathbf{K} \hat{\boldsymbol{\theta}} = \mathbf{K}^T \mathbf{y} \quad \text{resp.} \quad \mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{H} \hat{\boldsymbol{\theta}} = \mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{x}$$

and its solution

$$\hat{\boldsymbol{\theta}} = (\mathbf{K}^T \mathbf{K})^{-1} \mathbf{K}^T \mathbf{y} = (\mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{x}$$

are obtained analogous to the white noise case.

Properties of the generalized least squares estimator

The mean vector is given by

$$\begin{aligned}
 E(\hat{\boldsymbol{\Theta}}) &= (\mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}_{UU}^{-1} E(\mathbf{X}) = (\mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{H} \boldsymbol{\theta} \\
 &= \begin{cases} \boldsymbol{\theta} & \text{rank}(\mathbf{H}) = p \\ \mathbf{V}_1 \mathbf{V}_1^T \boldsymbol{\theta} = (\mathbf{I} - \mathbf{V}_2 \mathbf{V}_2^T) \boldsymbol{\theta} & \text{rank}(\mathbf{H}) = M < p \end{cases}
 \end{aligned}$$

where the matrices \mathbf{V}_1 and \mathbf{V}_2 can be derived from the singular value decomposition

$$\mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{H} = \mathbf{V} \text{diag}(\lambda_1, \dots, \lambda_M, 0, \dots, 0) \mathbf{V}^T$$

with

$$\mathbf{V} = (\mathbf{V}_1, \mathbf{V}_2), \quad \mathbf{V}_1 = (\mathbf{v}_1, \dots, \mathbf{v}_M) \quad \mathbf{V}_2 = (\mathbf{v}_{M+1}, \dots, \mathbf{v}_p).$$

The covariance matrix of the generalized least squares estimator

$$\mathbf{C}_{\hat{\Theta}\hat{\Theta}} = \text{Cov}(\hat{\Theta}) = E\left(\left(\hat{\Theta} - E(\hat{\Theta})\right)\left(\hat{\Theta} - E(\hat{\Theta})\right)^T\right)$$

results after exploiting

$$\hat{\Theta} = (\mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{X} \quad \text{and} \quad E(\hat{\Theta}) = (\mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{H} \boldsymbol{\theta}$$

in the expression

$$\begin{aligned} \mathbf{C}_{\hat{\Theta}\hat{\Theta}} &= E\left(\left(\mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{H}\right)^{-1} \mathbf{H}^T \mathbf{C}_{UU}^{-1} (\mathbf{X} - \mathbf{H}\boldsymbol{\theta})(\mathbf{X} - \mathbf{H}\boldsymbol{\theta})^T \mathbf{C}_{UU}^{-1} \mathbf{H} \left(\mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{H}\right)^{-1}\right) \\ &= \left(\mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{H}\right)^{-1} \mathbf{H}^T \mathbf{C}_{UU}^{-1} E\left((\mathbf{X} - \mathbf{H}\boldsymbol{\theta})(\mathbf{X} - \mathbf{H}\boldsymbol{\theta})^T\right) \mathbf{C}_{UU}^{-1} \mathbf{H} \left(\mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{H}\right)^{-1} \\ &= \left(\mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{H}\right)^{-1} \mathbf{H}^T \mathbf{C}_{UU}^{-1} (\sigma_U^2 \mathbf{C}_{UU}) \mathbf{C}_{UU}^{-1} \mathbf{H} \left(\mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{H}\right)^{-1} \\ &= \sigma_U^2 \left(\mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{H}\right)^{-1}. \end{aligned}$$

Theorem: (Gauß-Markov)

Given the model

$$\mathbf{X} = \mathbf{H}\boldsymbol{\theta} + \mathbf{U} \quad \text{with} \quad X_n = h_{i1}\theta_1 + \dots + h_{ip}\theta_p + U_i, \quad i=1, \dots, n,$$

where

$$\mathbf{E}(\mathbf{X}) = \mathbf{H}\boldsymbol{\theta}, \quad \mathbf{Cov}(\mathbf{X}) = \mathbf{Cov}(\mathbf{U}) = \sigma_U^2 \mathbf{C}_{UU}$$

and

$$\text{rank}(\mathbf{H}) = p.$$

Then the best linear unbiased estimator (BLUE) is the least squares estimator

$$\hat{\boldsymbol{\theta}} = (\mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{X}.$$

Hence, the minimum of the weighted sum of squares is

$$\begin{aligned}
 q(\hat{\boldsymbol{\theta}}) &= \text{tr}(\mathbf{P}^\perp \mathbf{y} \mathbf{y}^T) = \text{tr}\left(\left(\mathbf{I} - \mathbf{K}(\mathbf{K}^T \mathbf{K})^{-1} \mathbf{K}^T\right) \mathbf{y} \mathbf{y}^T\right) \\
 &= \text{tr}\left(\left(\mathbf{I} - \mathbf{C}^{-1} \mathbf{H}(\mathbf{H}^T (\mathbf{C}^{-1})^T \mathbf{C}^{-1} \mathbf{H})^{-1} \mathbf{H}^T (\mathbf{C}^{-1})^T\right) \mathbf{C}^{-1} \mathbf{x} \mathbf{x}^T (\mathbf{C}^{-1})^T\right) \\
 &= \text{tr}\left(\mathbf{C}_{\text{UU}}^{-1} \left(\mathbf{I} - \mathbf{H}(\mathbf{H}^T \mathbf{C}_{\text{UU}}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}_{\text{UU}}^{-1}\right) \mathbf{x} \mathbf{x}^T\right) \\
 &= \text{tr}\left(\mathbf{C}_{\text{UU}}^{-1} \left(\mathbf{I} - \mathbf{P}_{\mathbf{C}_{\text{UU}}^{-1}}\right) \mathbf{x} \mathbf{x}^T\right) = \text{tr}\left(\mathbf{C}_{\text{UU}}^{-1} \mathbf{P}_{\mathbf{C}_{\text{UU}}^{-1}}^\perp \mathbf{x} \mathbf{x}^T\right),
 \end{aligned}$$

where

$$\mathbf{P}_{\mathbf{C}_{\text{UU}}^{-1}} = \mathbf{H}(\mathbf{H}^T \mathbf{C}_{\text{UU}}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}_{\text{UU}}^{-1} \quad \text{and} \quad \mathbf{P}_{\mathbf{C}_{\text{UU}}^{-1}}^\perp = \mathbf{I} - \mathbf{P}_{\mathbf{C}_{\text{UU}}^{-1}}$$

represent again projection matrices.

Again, we are interested in estimating the variance σ_Z^2 of the noise. The expected value of the minimum of the weighted sum of squares results in

$$\begin{aligned} E(q(\hat{\Theta})) &= E\left(\text{tr}(\mathbf{C}_{UU}^{-1} \mathbf{P}_{\mathbf{C}_{UU}^{-1}}^{\perp} \mathbf{X} \mathbf{X}^T)\right) = \text{tr}\left(\mathbf{C}_{UU}^{-1} \mathbf{P}_{\mathbf{C}_{UU}^{-1}}^{\perp} E\left((\mathbf{H}\boldsymbol{\theta} + \mathbf{U})(\mathbf{H}\boldsymbol{\theta} + \mathbf{U})^T\right)\right) \\ &= \text{tr}\left(\mathbf{C}_{UU}^{-1} \mathbf{P}_{\mathbf{C}_{UU}^{-1}}^{\perp} \left(\mathbf{H}\boldsymbol{\theta}\boldsymbol{\theta}^T \mathbf{H}^T + E(\mathbf{U})\boldsymbol{\theta}^T \mathbf{H}^T + \mathbf{H}\boldsymbol{\theta}E(\mathbf{U}^T) + E(\mathbf{U}\mathbf{U}^T)\right)\right) \\ &= \text{tr}\left(\mathbf{C}_{UU}^{-1} \mathbf{P}_{\mathbf{C}_{UU}^{-1}}^{\perp} \left(\mathbf{H}\boldsymbol{\theta}\boldsymbol{\theta}^T \mathbf{H}^T + \sigma_U^2 \mathbf{C}_{UU}\right)\right) = \sigma_U^2 \text{tr}(\mathbf{P}_{\mathbf{C}_{UU}^{-1}}^{\perp}) = (n-p)\sigma_U^2, \end{aligned}$$

where

$$\begin{aligned} \text{tr}(\mathbf{P}_{\mathbf{C}_{UU}^{-1}}^{\perp}) &= \text{tr}\left(\mathbf{I}_n - \mathbf{H}(\mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}_{UU}^{-1}\right) \\ &= n - \text{tr}\left((\mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{H}\right) = n - \text{tr}(\mathbf{I}_p) = n - p \end{aligned}$$

has been utilized.

This result motivates

$$S^2 = q(\hat{\Theta}) / (n - p)$$

as an unbiased estimator for σ_Z^2 .

Theorem:

If in addition to the Gauß-Markov theorem \mathbf{U} obeys a normal distribution, i.e. $\mathbf{U} \sim \mathcal{N}_n(\mathbf{0}, \sigma_U^2 \mathbf{C}_{UU})$, the following holds.

a) $\hat{\Theta} = (\mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{H})^{-1} \mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{X} \sim \mathcal{N}_p(\boldsymbol{\theta}, \sigma_U^2 (\mathbf{H}^T \mathbf{C}_{UU}^{-1} \mathbf{H})^{-1})$

b) $\hat{\Theta}$ and S^2 are stochastically independent

c) $(n - p) / \sigma_U^2 \cdot S^2 \sim \chi_{n-p}^2$

3.4 Confidence Intervals

Now, an unknown parameter θ is considered, where the density $f_{\mathbf{X}}(\mathbf{x}|\theta)$ of \mathbf{X} and an estimator $\hat{\Theta} = \hat{\theta}(\mathbf{X})$ for θ possessing the density $f_{\hat{\Theta}}(\hat{\theta}|\theta)$ are given.

With the knowledge of $f_{\hat{\Theta}}(\hat{\theta}|\theta)$ and a given α , e.g. $\alpha = 0.05$, we can derive from

$$1 - \alpha = P(\theta - a_1 < \hat{\Theta} \leq \theta + a_2) = F_{\hat{\Theta}}(\theta + a_2 | \theta) - F_{\hat{\Theta}}(\theta - a_1 | \theta)$$

the probability equation

$$\begin{aligned} 1 - \alpha &= P(\theta - a_1 < \hat{\Theta} \leq \theta + a_2) = P(-a_1 < \hat{\Theta} - \theta \leq a_2) \\ &= P(-\hat{\Theta} - a_1 < -\theta \leq a_2 - \hat{\Theta}) = P(\hat{\Theta} - a_2 \leq \theta < \hat{\Theta} + a_1) \end{aligned}$$

which states that random interval $[\hat{\theta} - a_2, \hat{\theta} + a_1)$ covers the parameter θ with probability $1 - \alpha$.

For given observations \mathbf{x} the set

$$C(\mathbf{x}) = \left\{ \theta : \hat{\theta}(\mathbf{x}) - a_2 \leq \theta < \hat{\theta}(\mathbf{x}) + a_1 \right\}$$

is called

- confidence interval for θ with the confidence coefficient $1 - \alpha$

or alternatively

- interval estimate to the confidence level $1 - \alpha$.

In case of an unknown parameter vector $\boldsymbol{\theta}$ an interval estimate can be constructed if a function $g(\mathbf{x} | \boldsymbol{\theta})$ exists such that the corresponding random variable

$$Y = g(\mathbf{x} | \boldsymbol{\theta})$$

obeys a known distribution which is independent of $\boldsymbol{\theta}$.

Frequently one can find such a function, which depends on \mathbf{x} only over the estimate $\hat{\boldsymbol{\theta}}(\mathbf{x})$, i.e.

$$g(\mathbf{x} | \boldsymbol{\theta}) = h(\hat{\boldsymbol{\theta}} | \boldsymbol{\theta}).$$

Examples are given in the subsequent exercises.

Exercise 3.4-1:

Signal + white noise, confidence interval for the

- *signal amplitude estimate (known noise variance)*
- *signal amplitude estimate (unknown noise variance)*
- *noise variance estimate*

Exercise 3.4-2:

Sinusoids + white noise, confidence interval for the amplitude estimates (unknown noise variance)

3.5 Cramer-Rao Lower Bound

In the following only unbiased estimators, i.e.

$$E(\hat{\Theta}) = \int_{\mathbb{R}^n} \hat{\Theta}(\mathbf{x}) f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}) d\mathbf{x} = \boldsymbol{\theta},$$

are considered, whose covariance matrix

$$\mathbf{C}_{\hat{\Theta}\hat{\Theta}} = \text{Cov}(\hat{\Theta}) = E\left((\hat{\Theta} - \boldsymbol{\theta})(\hat{\Theta} - \boldsymbol{\theta})^T\right)$$

exist. The Cramer-Rao inequality provides now a lower bound to the covariance matrix of any unbiased estimator of $\boldsymbol{\theta}$. Due to the following regularity conditions certain expected values can be determined, which are of importance for deriving the Cramer-Rao inequality.

Regularity conditions:

- a) The parameter set Ω is an interval in \mathbb{R}^p .
- b) The gradient $\nabla_{\boldsymbol{\theta}} \ln(f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}))$ exists for any \mathbf{x} and $\boldsymbol{\theta}$.
- c) The gradient of $\int_{\mathbb{R}^n} f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}) d\mathbf{x}$ with respect to $\boldsymbol{\theta}$ can be obtained by taking the gradient under the integral sign, i.e. $\nabla_{\boldsymbol{\theta}} \int_{\mathbb{R}^n} f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}) d\mathbf{x} = \int_{\mathbb{R}^n} \nabla_{\boldsymbol{\theta}} f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}) d\mathbf{x}$.
- d) The gradient of $E(\hat{\boldsymbol{\theta}}(\mathbf{X}))$ with respect to $\boldsymbol{\theta}$ can be obtained by taking the gradient under the integral sign, i.e. $\nabla_{\boldsymbol{\theta}} \int_{\mathbb{R}^n} \hat{\boldsymbol{\theta}}(\mathbf{x}) f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}) d\mathbf{x} = \int_{\mathbb{R}^n} \hat{\boldsymbol{\theta}}(\mathbf{x}) \nabla_{\boldsymbol{\theta}} f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}) d\mathbf{x}$.

Corollary:

If the former regularity conditions hold, then

$$1) \ E(\nabla_{\boldsymbol{\theta}} \ln(f_{\mathbf{X}}(\mathbf{X} | \boldsymbol{\theta}))) = \mathbf{0}$$

$$2) \ E\left(\left(\nabla_{\boldsymbol{\theta}} \ln(f_{\mathbf{X}}(\mathbf{X} | \boldsymbol{\theta}))\right)(\hat{\boldsymbol{\Theta}} - \boldsymbol{\theta})^T\right) = \mathbf{I}$$

Proof of 1):

$$\begin{aligned} E(\nabla_{\boldsymbol{\theta}} \ln(f_{\mathbf{X}}(\mathbf{X} | \boldsymbol{\theta}))) &= \int_{\mathbb{R}^n} \nabla_{\boldsymbol{\theta}} \ln(f_{\mathbf{X}}(\mathbf{x} | \boldsymbol{\theta})) f_{\mathbf{X}}(\mathbf{x} | \boldsymbol{\theta}) d\mathbf{x} = \\ &= \int_{\mathbb{R}^n} \frac{\nabla_{\boldsymbol{\theta}} f_{\mathbf{X}}(\mathbf{x} | \boldsymbol{\theta})}{f_{\mathbf{X}}(\mathbf{x} | \boldsymbol{\theta})} f_{\mathbf{X}}(\mathbf{x} | \boldsymbol{\theta}) d\mathbf{x} = \int_{\mathbb{R}^n} \nabla_{\boldsymbol{\theta}} f_{\mathbf{X}}(\mathbf{x} | \boldsymbol{\theta}) d\mathbf{x} \\ &= \nabla_{\boldsymbol{\theta}} \int_{\mathbb{R}^n} f_{\mathbf{X}}(\mathbf{x} | \boldsymbol{\theta}) d\mathbf{x} = \nabla_{\boldsymbol{\theta}} 1 = \mathbf{0} \end{aligned}$$

Proof of 2):

$$\begin{aligned}
 E\left(\left(\nabla_{\boldsymbol{\theta}} \ln(f_{\mathbf{x}}(\mathbf{X} | \boldsymbol{\theta}))\right)\left(\hat{\boldsymbol{\Theta}} - \boldsymbol{\theta}\right)^T\right) &= \\
 &= \int_{\mathbb{R}^n} \nabla_{\boldsymbol{\theta}} \ln(f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta})) \left(\hat{\boldsymbol{\Theta}}(\mathbf{x}) - \boldsymbol{\theta}\right)^T f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}) d\mathbf{x} \\
 &= \int_{\mathbb{R}^n} \frac{\nabla_{\boldsymbol{\theta}} f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta})}{f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta})} \left(\hat{\boldsymbol{\Theta}}(\mathbf{x}) - \boldsymbol{\theta}\right)^T f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}) d\mathbf{x} \\
 &= \int_{\mathbb{R}^n} \nabla_{\boldsymbol{\theta}} f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}) \left(\hat{\boldsymbol{\Theta}}(\mathbf{x}) - \boldsymbol{\theta}\right)^T d\mathbf{x} \\
 &= \nabla_{\boldsymbol{\theta}} \int_{\mathbb{R}^n} f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}) \hat{\boldsymbol{\Theta}}(\mathbf{x})^T d\mathbf{x} - \nabla_{\boldsymbol{\theta}} \int_{\mathbb{R}^n} f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}) d\mathbf{x} \boldsymbol{\theta}^T \\
 &= \nabla_{\boldsymbol{\theta}} \boldsymbol{\theta}^T - (\nabla_{\boldsymbol{\theta}} 1) \boldsymbol{\theta}^T = \mathbf{I}
 \end{aligned}$$

Definition: (Fisher information matrix)

The $p \times p$ Fisher information matrix is defined by

$$\begin{aligned} \mathcal{I}(\boldsymbol{\theta}) &= \mathbf{E} \left(\nabla_{\boldsymbol{\theta}} \ln(f_{\mathbf{X}}(\mathbf{X} | \boldsymbol{\theta})) \nabla_{\boldsymbol{\theta}}^T \ln(f_{\mathbf{X}}(\mathbf{X} | \boldsymbol{\theta})) \right) = \\ &= \int_{\mathbb{R}^n} \nabla_{\boldsymbol{\theta}} \ln(f_{\mathbf{X}}(\mathbf{x} | \boldsymbol{\theta})) \nabla_{\boldsymbol{\theta}}^T \ln(f_{\mathbf{X}}(\mathbf{X} | \boldsymbol{\theta})) f_{\mathbf{X}}(\mathbf{x} | \boldsymbol{\theta}) d\mathbf{x}. \end{aligned}$$

Hence, the entries of the Fisher information matrix are for $i, j = 1, \dots, p$ given by

$$\begin{aligned} \mathcal{I}_{ij}(\boldsymbol{\theta}) &= \mathbf{E} \left(\frac{\partial \ln(f_{\mathbf{X}}(\mathbf{x} | \boldsymbol{\theta}))}{\partial \theta_i} \cdot \frac{\partial \ln(f_{\mathbf{X}}(\mathbf{X} | \boldsymbol{\theta}))}{\partial \theta_j} \right) \\ &= \int_{\mathbb{R}^n} \frac{\partial \ln(f_{\mathbf{X}}(\mathbf{x} | \boldsymbol{\theta}))}{\partial \theta_i} \cdot \frac{\partial \ln(f_{\mathbf{X}}(\mathbf{X} | \boldsymbol{\theta}))}{\partial \theta_j} f_{\mathbf{X}}(\mathbf{x} | \boldsymbol{\theta}) d\mathbf{x}. \end{aligned}$$

Theorem: (Cramer-Rao inequality)

Suppose the former assumptions and regularity conditions hold and the Fisher information matrix $\mathcal{I}(\boldsymbol{\theta})$ is positive definite. Then the variance of any estimator $\hat{\Theta}_{\mathbf{a}} = \mathbf{a}^T \hat{\Theta}$ is bounded downwards by

$$\text{Var}(\hat{\Theta}_{\mathbf{a}}) = \mathbf{a}^T \mathbf{C}_{\hat{\Theta}\hat{\Theta}} \mathbf{a} \geq \mathbf{a}^T \mathcal{I}(\boldsymbol{\theta})^{-1} \mathbf{a} \quad \forall \mathbf{a} \in \mathbb{R}^p.$$

If particularly $\mathbf{a} = \mathbf{e}_i$, i.e. the i -th Cartesian unit vector, the inequality provides the lower limit for the variance of the i -th component of $\hat{\Theta}$.

$$\text{Var}(\hat{\Theta}_i) = \mathbf{e}_i^T \mathbf{C}_{\hat{\Theta}\hat{\Theta}} \mathbf{e}_i \geq \mathbf{e}_i^T \mathcal{I}(\boldsymbol{\theta})^{-1} \mathbf{e}_i = \mathcal{I}(\boldsymbol{\theta})_{ii}^{-1}$$

$$i = 1, \dots, p$$

Proof:

First, the random vector

$$\mathbf{Y} = \hat{\boldsymbol{\Theta}} - \boldsymbol{\theta} - \mathcal{I}(\boldsymbol{\theta})^{-1} \nabla_{\boldsymbol{\theta}} \ln(f_{\mathbf{X}}(\mathbf{X} | \boldsymbol{\theta}))$$

is introduced and its second order moment matrix

$$\begin{aligned} \mathbf{E}(\mathbf{Y}\mathbf{Y}^T) &= \mathbf{E}\left((\hat{\boldsymbol{\Theta}} - \boldsymbol{\theta})(\hat{\boldsymbol{\Theta}} - \boldsymbol{\theta})^T\right) \\ &\quad - \mathbf{E}\left((\hat{\boldsymbol{\Theta}} - \boldsymbol{\theta})\nabla_{\boldsymbol{\theta}}^T \ln(f_{\mathbf{X}}(\mathbf{X} | \boldsymbol{\theta}))\right) \mathcal{I}(\boldsymbol{\theta})^{-1} \\ &\quad - \mathcal{I}(\boldsymbol{\theta})^{-1} \mathbf{E}\left(\nabla_{\boldsymbol{\theta}} \ln(f_{\mathbf{X}}(\mathbf{X} | \boldsymbol{\theta}))(\hat{\boldsymbol{\Theta}} - \boldsymbol{\theta})^T\right) \\ &\quad + \mathcal{I}(\boldsymbol{\theta})^{-1} \mathbf{E}\left(\nabla_{\boldsymbol{\theta}} \ln(f_{\mathbf{X}}(\mathbf{X} | \boldsymbol{\theta}))\nabla_{\boldsymbol{\theta}}^T \ln(f_{\mathbf{X}}(\mathbf{X} | \boldsymbol{\theta}))\right) \mathcal{I}(\boldsymbol{\theta})^{-1}. \end{aligned}$$

is determined. After exploiting the former corollary the

second order moment matrix can be simplified to

$$E(\mathbf{Y}\mathbf{Y}^T) = E\left((\hat{\boldsymbol{\Theta}} - \boldsymbol{\theta})(\hat{\boldsymbol{\Theta}} - \boldsymbol{\theta})^T\right) - \mathcal{I}(\boldsymbol{\theta})^{-1}.$$

Since

$$\mathbf{a}^T E(\mathbf{Y}\mathbf{Y}^T) \mathbf{a} = E(\mathbf{a}^T \mathbf{Y}\mathbf{Y}^T \mathbf{a}) = E\left((\mathbf{a}^T \mathbf{Y})^2\right) \geq 0 \quad \forall \mathbf{a} \in \mathbb{R}^p$$

the second order moment matrix is positive semidefinite.

Thus, we can write

$$\mathbf{a}^T E(\mathbf{Y}\mathbf{Y}^T) \mathbf{a} = \mathbf{a}^T \mathbf{C}_{\hat{\boldsymbol{\Theta}}\hat{\boldsymbol{\Theta}}} \mathbf{a} - \mathbf{a}^T \mathcal{I}(\boldsymbol{\theta})^{-1} \mathbf{a} \geq 0 \quad \forall \mathbf{a} \in \mathbb{R}^p$$

and the assertion

$$\mathbf{a}^T \mathbf{C}_{\hat{\boldsymbol{\Theta}}\hat{\boldsymbol{\Theta}}} \mathbf{a} \geq \mathbf{a}^T \mathcal{I}(\boldsymbol{\theta})^{-1} \mathbf{a} \quad \forall \mathbf{a} \in \mathbb{R}^p$$

of the theorem follows.

Under the following regularity condition an alternative and often more convenient expression for calculating the information matrix can be presented.

Regularity condition:

e) The Hessian matrix of

$$\int_{\mathbb{R}^n} f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}) d\mathbf{x}$$

with respect to $\boldsymbol{\theta}$ can be obtained by taking all the required second partial derivatives under the integral sign, i.e.

$$\nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}}^T \int_{\mathbb{R}^n} f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}) d\mathbf{x} = \int_{\mathbb{R}^n} \nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}}^T f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}) d\mathbf{x}.$$

Corollary:

If the former regularity condition holds, the Fisher information matrix can be alternatively determined by

$$\begin{aligned}\mathcal{I}(\boldsymbol{\theta}) &= \mathbf{E}\left(\nabla_{\boldsymbol{\theta}} \ln(f_{\mathbf{X}}(\mathbf{X} | \boldsymbol{\theta})) \nabla_{\boldsymbol{\theta}}^T \ln(f_{\mathbf{X}}(\mathbf{X} | \boldsymbol{\theta}))\right) \\ &= -\mathbf{E}\left(\nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}}^T \ln(f_{\mathbf{X}}(\mathbf{X} | \boldsymbol{\theta}))\right).\end{aligned}$$

Consequently, the elements of the Fisher information matrix can be expressed by

$$\begin{aligned}\mathcal{I}_{ij}(\boldsymbol{\theta}) &= \mathbf{E}\left(\frac{\partial \ln(f_{\mathbf{X}}(\mathbf{X} | \boldsymbol{\theta}))}{\partial \theta_i} \cdot \frac{\partial \ln(f_{\mathbf{X}}(\mathbf{X} | \boldsymbol{\theta}))}{\partial \theta_j}\right) = -\mathbf{E}\left(\frac{\partial^2 \ln(f_{\mathbf{X}}(\mathbf{X} | \boldsymbol{\theta}))}{\partial \theta_i \partial \theta_j}\right) \\ & \quad i, j = 1, \dots, p.\end{aligned}$$

Proof:

$$\begin{aligned}
 E\left(\nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}}^T \ln(f_{\mathbf{X}}(\mathbf{X}|\boldsymbol{\theta}))\right) &= E\left(\nabla_{\boldsymbol{\theta}} \frac{\nabla_{\boldsymbol{\theta}}^T f_{\mathbf{X}}(\mathbf{X}|\boldsymbol{\theta})}{f_{\mathbf{X}}(\mathbf{X}|\boldsymbol{\theta})}\right) = \\
 &= E\left(\frac{\nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}}^T f_{\mathbf{X}}(\mathbf{X}|\boldsymbol{\theta}) \cdot f_{\mathbf{X}}(\mathbf{X}|\boldsymbol{\theta}) - \nabla_{\boldsymbol{\theta}} f_{\mathbf{X}}(\mathbf{X}|\boldsymbol{\theta}) \cdot \nabla_{\boldsymbol{\theta}}^T f_{\mathbf{X}}(\mathbf{X}|\boldsymbol{\theta})}{f_{\mathbf{X}}^2(\mathbf{X}|\boldsymbol{\theta})}\right) \\
 &= E\left(\frac{\nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}}^T f_{\mathbf{X}}(\mathbf{X}|\boldsymbol{\theta})}{f_{\mathbf{X}}(\mathbf{X}|\boldsymbol{\theta})}\right) - E\left(\frac{\nabla_{\boldsymbol{\theta}} f_{\mathbf{X}}(\mathbf{X}|\boldsymbol{\theta}) \cdot \nabla_{\boldsymbol{\theta}}^T f_{\mathbf{X}}(\mathbf{X}|\boldsymbol{\theta})}{f_{\mathbf{X}}^2(\mathbf{X}|\boldsymbol{\theta})}\right) \\
 &= \int_{\mathbb{R}^n} \nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}}^T f_{\mathbf{X}}(\mathbf{x}|\boldsymbol{\theta}) d\mathbf{x} - E\left(\nabla_{\boldsymbol{\theta}} \ln(f_{\mathbf{X}}(\mathbf{X}|\boldsymbol{\theta})) \nabla_{\boldsymbol{\theta}}^T \ln(f_{\mathbf{X}}(\mathbf{X}|\boldsymbol{\theta}))\right) \\
 &= \nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}}^T \int_{\mathbb{R}^n} f_{\mathbf{X}}(\mathbf{x}|\boldsymbol{\theta}) d\mathbf{x} - \mathcal{I}(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}}^T 1 - \mathcal{I}(\boldsymbol{\theta}) = -\mathcal{I}(\boldsymbol{\theta})
 \end{aligned}$$

Theorem:

Let $f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\xi})$ belong to the exponential family in canonical form and \mathbf{T} denote the corresponding sufficient statistic.

- 1) The regularity conditions a) – e) are satisfied and the information matrix is positive definite

$$\begin{aligned}\mathcal{I}(\boldsymbol{\xi}) &= \mathbf{E}\left(\nabla_{\boldsymbol{\xi}} \ln(f_{\mathbf{x}}(\mathbf{X} | \boldsymbol{\xi})) \nabla_{\boldsymbol{\xi}}^T \ln(f_{\mathbf{x}}(\mathbf{X} | \boldsymbol{\xi}))\right) \\ &= -\mathbf{E}\left(\nabla_{\boldsymbol{\xi}} \nabla_{\boldsymbol{\xi}}^T \ln(f_{\mathbf{x}}(\mathbf{X} | \boldsymbol{\xi}))\right) \\ &= \nabla_{\boldsymbol{\xi}} \nabla_{\boldsymbol{\xi}}^T A(\boldsymbol{\xi}) = \text{Cov}(\mathbf{T}) = \mathbf{C}_{\mathbf{TT}}.\end{aligned}$$

Thus, for unbiased estimators the Cramer-Rao inequality holds.

2) Let the vector valued function

$$\boldsymbol{\xi} = \mathbf{g}(\boldsymbol{\theta}), \quad \boldsymbol{\xi} \in \mathbb{R}^k, \quad \boldsymbol{\theta} \in \mathbb{R}^p \quad \text{with } p \leq k$$

be one-to-one and differentiable with

$$\mathbf{G}(\boldsymbol{\theta}) = \nabla_{\boldsymbol{\theta}} \mathbf{g}(\boldsymbol{\theta})^T = \left(\frac{\partial g_j(\boldsymbol{\theta})}{\partial \theta_i} \right)_{i=1, \dots, p; j=1, \dots, k} .$$

Then the regularity conditions a) - e) are satisfied with respect to $\boldsymbol{\theta}$ and the information matrix can be expressed by

$$\begin{aligned} \mathcal{I}(\boldsymbol{\theta}) &= -E\left(\nabla_{\boldsymbol{\theta}} \nabla_{\boldsymbol{\theta}}^T \ln(f_{\mathbf{X}}(\mathbf{X} | \boldsymbol{\theta}))\right) = \mathbf{G}(\boldsymbol{\theta}) \mathcal{I}(\boldsymbol{\xi}) \Big|_{\boldsymbol{\xi}=\mathbf{g}(\boldsymbol{\theta})} \mathbf{G}(\boldsymbol{\theta})^T \\ &= \mathbf{G}(\boldsymbol{\theta}) \left(\nabla_{\boldsymbol{\xi}} \nabla_{\boldsymbol{\xi}}^T A(\boldsymbol{\xi}) \right) \Big|_{\boldsymbol{\xi}=\mathbf{g}(\boldsymbol{\theta})} \mathbf{G}(\boldsymbol{\theta})^T . \end{aligned}$$

Now, one asks oneself immediately, does there exist an unbiased estimator for which the information inequality takes the equality sign, i.e. whose covariance matrix is equal to the inverse Fisher information matrix.

An estimator for which the information inequality takes the equality sign, i.e. whose covariance matrix coincides with the inverse Fisher information matrix (Cramer-Rao lower bound), is called efficient.

Efficient estimators can be achieved only in special cases, since the random vector \mathbf{Y} , introduced for proving the information inequality, must be the zero-vector with probability 1.

Exercise 3.5-1:

Non/efficiency of linear least squares estimators when the noise is i.i.d. and normally distributed

The previous exercise showed, that the deviation of the covariance matrix from the inverse Fisher information matrix can be neglected for large sample sizes.

This motivates the definition of the limit

$$\Gamma(\boldsymbol{\theta}) = \lim_{n \rightarrow \infty} \frac{1}{n} \mathcal{I}_n(\boldsymbol{\theta}),$$

where n indicates the growing sample size. If the matrix $\Gamma(\boldsymbol{\theta})$ exists and is not singular, then unbiased and consistent estimators $\hat{\boldsymbol{\Theta}}_n$ are of interest that satisfy the limit

$$\lim_{n \rightarrow \infty} n \mathbf{C}_{\hat{\boldsymbol{\Theta}}_n \hat{\boldsymbol{\Theta}}_n} = \Gamma(\boldsymbol{\theta})^{-1}.$$

Simultaneously, such estimators are then even consistent in the mean square sense.

Estimators used in practice are often neither unbiased nor mean square consistent. In these cases the limiting distributions of such estimators are considered. If the limiting distribution satisfies the property

$$\lim_{n \rightarrow \infty} \sqrt{n} (\hat{\Theta}_n - \boldsymbol{\theta}) \sim \mathcal{N}_p(\mathbf{0}, \mathbf{D}(\boldsymbol{\theta})) \quad \text{and} \quad \mathbf{D}(\boldsymbol{\theta}) = \boldsymbol{\Gamma}(\boldsymbol{\theta})^{-1},$$

the estimator is said to be asymptotically efficient.

Furthermore, under certain regularity conditions (similar to the former) one can show that in the class of asymptotically normally distributed estimators the following lower bound can be proved.

$$\mathbf{a}^T \mathbf{D}(\boldsymbol{\theta}) \mathbf{a} \geq \mathbf{a}^T \boldsymbol{\Gamma}(\boldsymbol{\theta})^{-1} \mathbf{a} \quad \mathbf{a} \in \mathbb{R}^p$$

Exercise 3.5-2:

*Asymptotic efficiency of linear least squares estimators
when the noise is i.i.d. and normally distributed*

3.6 Maximum Likelihood Estimation

The maximum Likelihood method is one of the most important procedures for the construction of estimators.

For given observations \mathbf{x} the procedure selects that $\boldsymbol{\theta}$ as maximum likelihood estimate $\hat{\boldsymbol{\theta}}(\mathbf{x})$, for which the density function $f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta})$ takes its maximum.

Thus, the maximum likelihood estimate can be interpreted heuristically as that parameter vector that makes the occurrence of the data observed most likely.

Instead of determining the maximum of $f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta})$, one can due to the monotonicity of the logarithmic function alternatively maximize $\ln(f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}))$.

Definition:

Suppose \mathbf{X} possesses the density $f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta})$ with $\boldsymbol{\theta} \in \Omega$ and the vector of observations $\mathbf{x} = (x_1, \dots, x_n)^T$ is given. Then $f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta})$ and $\ln(f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}))$ are called likelihood function and log-likelihood function respectively, and

$$\hat{\boldsymbol{\theta}} = \underset{\boldsymbol{\theta} \in \Omega}{\operatorname{argmax}} f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}) = \underset{\boldsymbol{\theta} \in \Omega}{\operatorname{argmax}} \ln(f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}))$$

is called the maximum likelihood estimate (MLE) of $\boldsymbol{\theta}$.

If the gradient of $f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta})$ with respect to $\boldsymbol{\theta}$ exists and if $f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta})$ is positive for the given \mathbf{x} , one can try to find the MLE by solving the likelihood equation system

$$\nabla_{\boldsymbol{\theta}} \ln(f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta})) \Big|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}(\mathbf{x})} = \mathbf{0}.$$

Remarks:

- Generally the likelihood equation system is nonlinear.
- The likelihood equation system can possess several solutions.
- The solutions of the likelihood equation system do not necessarily correspond to relative maxima.
- If the absolute maximum of the likelihood function lies on the boundary of Ω it can not be described by a solution of the likelihood equation system.
- Consequently, one has generally to employ sophisticated numerical optimization techniques for determining maximum likelihood estimates.

Properties of Maximum Likelihood Estimation:

1) Let $\mathbf{g}: \mathbb{R}^q \rightarrow \mathbb{R}^p$, $p \leq q$ be a continuous mapping of the parameter vector $\boldsymbol{\eta}$ and $\hat{\boldsymbol{\eta}}$ the MLE of $\boldsymbol{\eta}$. Then the MLE of $\boldsymbol{\theta} = \mathbf{g}(\boldsymbol{\eta})$ is given by $\hat{\boldsymbol{\theta}} = \mathbf{g}(\hat{\boldsymbol{\eta}})$.

2) Suppose $\mathbf{t}(\mathbf{x})$ is a sufficient transformation for $\boldsymbol{\theta}$, i.e.

$$f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}) = g(\mathbf{t}(\mathbf{x}) | \boldsymbol{\theta}) \cdot h(\mathbf{x}).$$

Then the maximum likelihood estimate $\hat{\boldsymbol{\theta}}$ only depends over $\mathbf{t}(\mathbf{x})$ on \mathbf{x} .

3) If $f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\xi})$ belongs to the exponential family in canonical form the related likelihood equation system

$$\nabla_{\boldsymbol{\xi}} \ln(f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\xi})) = \mathbf{t}(\mathbf{x}) - \nabla_{\boldsymbol{\xi}} A(\boldsymbol{\xi}) = \mathbf{0}$$

possesses a unique solution.

Exercise 3.6-1:

Maximum likelihood estimation for linear models when the noise is i.i.d. and normally distributed

Exercise 3.6-2:

Sinusoids + white noise, maximum likelihood estimation of amplitude, phase and frequency when the noise is i.i.d. and normally distributed

Exercise 3.6-3:

Example for inconsistent maximum likelihood estimation

3.7 Bayesian Estimation

In the Bayesian theory of parameter estimation the unknown parameter vector $\boldsymbol{\theta}$ is treated itself as a realization of a random experiment that obeys its own so called prior distribution $f_{\boldsymbol{\theta}}(\boldsymbol{\theta})$.

The objective is to use $f_{\boldsymbol{\theta}}(\boldsymbol{\theta})$ together with the measurements \mathbf{x} drawn from the distribution $f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta})$ to turn the prior distribution $f_{\boldsymbol{\theta}}(\boldsymbol{\theta})$ into a posterior distribution

$$f_{\boldsymbol{\theta}}(\boldsymbol{\theta} | \mathbf{x}) = \frac{f_{\mathbf{x}\boldsymbol{\theta}}(\mathbf{x}, \boldsymbol{\theta})}{f_{\mathbf{x}}(\mathbf{x})} = \frac{f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta})f_{\boldsymbol{\theta}}(\boldsymbol{\theta})}{\int_{\mathbb{R}^p} f_{\mathbf{x}\boldsymbol{\theta}}(\mathbf{x}, \boldsymbol{\theta})d\boldsymbol{\theta}} = \frac{f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta})f_{\boldsymbol{\theta}}(\boldsymbol{\theta})}{\int_{\mathbb{R}^p} f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta})f_{\boldsymbol{\theta}}(\boldsymbol{\theta})d\boldsymbol{\theta}}$$

that is used for the statistical inference.

Loss function

The quality of the estimate $\hat{\boldsymbol{\theta}}(\mathbf{x})$ is measured by a real-valued loss function $L(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}(\mathbf{x}))$, e.g. by the Euclidean distance between $\boldsymbol{\theta}$ and $\hat{\boldsymbol{\theta}}(\mathbf{x})$, i.e.

$$L(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}(\mathbf{x})) = \left(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}(\mathbf{x}) \right)^T \left(\boldsymbol{\theta} - \hat{\boldsymbol{\theta}}(\mathbf{x}) \right).$$

Risk or loss average

The risk is the loss $L(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}(\mathbf{x}))$ averaged over the distribution of the measurements for any fixed parameter vector $\boldsymbol{\theta} \in \Omega$, i.e.

$$R_L(\hat{\boldsymbol{\theta}} | \boldsymbol{\theta}) = E_{\mathbf{x}|\boldsymbol{\theta}} \left(L(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}(\mathbf{X})) \right) = \int_{\mathbb{R}^n} L(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}(\mathbf{x})) f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}) d\mathbf{x}.$$

Bayes risk

The Bayes risk is the risk $R_L(\hat{\boldsymbol{\theta}} | \boldsymbol{\Theta})$ averaged over the prior distribution $f_{\boldsymbol{\theta}}(\boldsymbol{\theta})$, i.e.

$$\begin{aligned} R_L(\hat{\boldsymbol{\theta}}) &= E_{\boldsymbol{\theta}} \left(R_L(\hat{\boldsymbol{\theta}} | \boldsymbol{\Theta}) \right) = E_{\boldsymbol{\theta}} \left(E_{\mathbf{x}|\boldsymbol{\theta}} \left(L(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}(\mathbf{X})) \right) \right) \\ &= \int_{\mathbb{R}^p} \int_{\mathbb{R}^n} L(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}(\mathbf{x})) f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}) f_{\boldsymbol{\theta}}(\boldsymbol{\theta}) d\mathbf{x} d\boldsymbol{\theta} \\ &= \int_{\mathbb{R}^{n+p}} L(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}(\mathbf{x})) f_{\mathbf{x}\boldsymbol{\theta}}(\mathbf{x}, \boldsymbol{\theta}) d\mathbf{x} d\boldsymbol{\theta} = E_{\mathbf{x}\boldsymbol{\theta}} \left(L(\boldsymbol{\theta}, \hat{\boldsymbol{\theta}}(\mathbf{X})) \right) \end{aligned}$$

The Bayes risk measures the average risk one incurs if the estimator $\hat{\boldsymbol{\theta}} = \hat{\boldsymbol{\theta}}(\mathbf{X})$ is used for the random experiment at hand. Hence, $R_L(\hat{\boldsymbol{\theta}})$ is depending only on the estimating function $\hat{\boldsymbol{\theta}}(\mathbf{x})$.

Bayes estimator

Minimization of the Bayes risk $R_L(\hat{\Theta})$ over the class of all estimating functions $\hat{\Theta}(\mathbf{x})$ for which $R_L(\hat{\Theta})$ exists, provides the Bayes estimating function

$$\hat{\Theta}_L(\mathbf{x}) = \underset{\hat{\Theta}(\mathbf{x})}{\operatorname{argmin}} R_L(\hat{\Theta}) = \underset{\hat{\Theta}(\mathbf{x})}{\operatorname{argmin}} E_{\mathbf{x}, \Theta} (L(\Theta, \hat{\Theta}(\mathbf{X})))$$

and therefore the Bayes estimator $\hat{\Theta}_L = \hat{\Theta}_L(\mathbf{X})$.

In the following Bayes estimators are considered in more detail for the square error, the absolute error and uniform error loss functions.

3.7.1 Minimum Mean Square Error Estimation

Non-Linear Minimum Mean Square Error Estimation

To estimate the random parameter vector Θ by means of the realization of the random vector \mathbf{X} in the minimum mean square error sense we have to find an estimating function $\hat{\Theta}(\mathbf{x})$ such that

$$\begin{aligned} R_{MSE}(\hat{\Theta}) &= E\left(\left(\Theta - \hat{\Theta}(\mathbf{X})\right)^T \left(\Theta - \hat{\Theta}(\mathbf{X})\right)\right) \\ &= \int_{\mathbb{R}^{n+p}} \left(\theta - \hat{\Theta}(\mathbf{x})\right)^T \left(\theta - \hat{\Theta}(\mathbf{x})\right) f_{\mathbf{x}\Theta}(\mathbf{x}, \theta) d\mathbf{x} d\theta \end{aligned}$$

is minimum. Over the class of all functions $\hat{\Theta}(\mathbf{x})$ for which the expected value exists, $R_{MSE}(\hat{\Theta})$ is minimized by

$$\hat{\boldsymbol{\theta}}_{MMSE}(\mathbf{x}) = \mathbf{E}(\boldsymbol{\theta} | \mathbf{X} = \mathbf{x}),$$

cf. Exercise 1.9-18.

Linear Minimum Mean Square Error Estimation

Now, we wish to estimate the random parameter vector $\boldsymbol{\Theta}$ by the linear estimating function

$$\hat{\boldsymbol{\theta}}(\mathbf{x}) = \mathbf{A}\mathbf{x} + \mathbf{b}$$

such that the mean square error

$$\begin{aligned} R_{LMSE}(\mathbf{A}, \mathbf{b}) &= \mathbf{E}\left(\left(\boldsymbol{\Theta} - (\mathbf{A}\mathbf{X} + \mathbf{b})\right)^T \left(\boldsymbol{\Theta} - (\mathbf{A}\mathbf{X} + \mathbf{b})\right)\right) \\ &= \int_{\mathbb{R}^{n+p}} \left(\boldsymbol{\theta} - (\mathbf{A}\mathbf{x} + \mathbf{b})\right)^T \left(\boldsymbol{\theta} - (\mathbf{A}\mathbf{x} + \mathbf{b})\right) f_{\mathbf{x}\boldsymbol{\theta}}(\mathbf{x}, \boldsymbol{\theta}) d\mathbf{x} d\boldsymbol{\theta} \end{aligned}$$

is minimized by varying the parameter matrix \mathbf{A} and vector \mathbf{b} . Thus, we have to solve the minimization problem

$$(\hat{\mathbf{A}}, \hat{\mathbf{b}}) = \underset{\mathbf{A}, \mathbf{b}}{\operatorname{argmin}} (R_{LMSE}(\mathbf{A}, \mathbf{b})).$$

The mean square error is minimum if

$$\frac{\partial}{\partial \mathbf{A}} R_{LMSE}(\mathbf{A}, \mathbf{b}) = -2E\left(\left(\boldsymbol{\Theta} - (\hat{\mathbf{A}}\mathbf{X} + \hat{\mathbf{b}})\right)\mathbf{X}^T\right) = \mathbf{0}$$

$$\nabla_{\mathbf{b}} R_{LMSE}(\mathbf{A}, \mathbf{b}) = -2E\left(\boldsymbol{\Theta} - (\hat{\mathbf{A}}\mathbf{X} + \hat{\mathbf{b}})\right) = \mathbf{0}$$

are satisfied. Applying the expectation operator we obtain

$$\mathbf{R}_{\Theta\mathbf{X}} - \hat{\mathbf{A}}\mathbf{R}_{\mathbf{X}\mathbf{X}} - \hat{\mathbf{b}}\boldsymbol{\mu}_{\mathbf{X}}^T = \mathbf{0} \quad \text{and} \quad \boldsymbol{\mu}_{\Theta} - \hat{\mathbf{A}}\boldsymbol{\mu}_{\mathbf{X}} - \hat{\mathbf{b}} = \mathbf{0},$$

where

$$\mathbf{R}_{\mathbf{X}\mathbf{X}} = E(\mathbf{X}\mathbf{X}^T), \quad \mathbf{R}_{\Theta\mathbf{X}} = E(\boldsymbol{\Theta}\mathbf{X}^T), \quad \boldsymbol{\mu}_{\mathbf{X}} = E(\mathbf{X}) \quad \text{and} \quad \boldsymbol{\mu}_{\Theta} = E(\boldsymbol{\Theta}).$$

Resolving for $\hat{\mathbf{A}}$ and $\hat{\mathbf{b}}$ leads after some manipulations to

$$\hat{\mathbf{A}} = (\mathbf{R}_{\Theta\mathbf{X}} - \boldsymbol{\mu}_{\Theta}\boldsymbol{\mu}_{\mathbf{X}}^T)(\mathbf{R}_{\mathbf{X}\mathbf{X}} - \boldsymbol{\mu}_{\mathbf{X}}\boldsymbol{\mu}_{\mathbf{X}}^T)^{-1}$$

$$\hat{\mathbf{b}} = \boldsymbol{\mu}_{\Theta} - (\mathbf{R}_{\Theta\mathbf{X}} - \boldsymbol{\mu}_{\Theta}\boldsymbol{\mu}_{\mathbf{X}}^T)(\mathbf{R}_{\mathbf{X}\mathbf{X}} - \boldsymbol{\mu}_{\mathbf{X}}\boldsymbol{\mu}_{\mathbf{X}}^T)^{-1}\boldsymbol{\mu}_{\mathbf{X}}$$

and therefore to the estimating function

$$\hat{\boldsymbol{\theta}}_{LMMSE}(\mathbf{x}) = \hat{\mathbf{A}}\mathbf{x} + \hat{\mathbf{b}} = \boldsymbol{\mu}_{\Theta} + (\mathbf{R}_{\Theta\mathbf{X}} - \boldsymbol{\mu}_{\Theta}\boldsymbol{\mu}_{\mathbf{X}}^T)(\mathbf{R}_{\mathbf{X}\mathbf{X}} - \boldsymbol{\mu}_{\mathbf{X}}\boldsymbol{\mu}_{\mathbf{X}}^T)^{-1}(\mathbf{x} - \boldsymbol{\mu}_{\mathbf{X}}).$$

After exploiting the relationships between the covariance and the first and second order moments

$$\mathbf{C}_{\mathbf{X}\mathbf{X}} = \mathbf{R}_{\mathbf{X}\mathbf{X}} - \boldsymbol{\mu}_{\mathbf{X}}\boldsymbol{\mu}_{\mathbf{X}}^T \quad \text{and} \quad \mathbf{C}_{\Theta\mathbf{X}} = \mathbf{R}_{\Theta\mathbf{X}} - \boldsymbol{\mu}_{\Theta}\boldsymbol{\mu}_{\mathbf{X}}^T$$

the estimating function can also be expressed by

$$\hat{\boldsymbol{\theta}}_{LMMSE}(\mathbf{x}) = \boldsymbol{\mu}_{\Theta} + \mathbf{C}_{\Theta\mathbf{X}}\mathbf{C}_{\mathbf{X}\mathbf{X}}^{-1}(\mathbf{x} - \boldsymbol{\mu}_{\mathbf{X}}).$$

Theorem:

Let \mathbf{Y} be a random vector composed of the random vectors \mathbf{X} and Θ and obeying the distribution

$$\mathbf{Y} = \begin{pmatrix} \mathbf{X} \\ \Theta \end{pmatrix} \sim \mathcal{N}_{n+p} \left(\begin{pmatrix} \boldsymbol{\mu}_X \\ \boldsymbol{\mu}_\Theta \end{pmatrix}, \begin{pmatrix} \boldsymbol{\Sigma}_{XX} & \boldsymbol{\Sigma}_{X\Theta} \\ \boldsymbol{\Sigma}_{\Theta X} & \boldsymbol{\Sigma}_{\Theta\Theta} \end{pmatrix} \right).$$

Then \mathbf{X} and Θ possessing the marginal distributions

$$\mathbf{X} \sim \mathcal{N}_n(\boldsymbol{\mu}_X, \boldsymbol{\Sigma}_{XX}), \quad \Theta \sim \mathcal{N}_p(\boldsymbol{\mu}_\Theta, \boldsymbol{\Sigma}_{\Theta\Theta})$$

and the conditional distribution

$$\Theta | \mathbf{X} = \mathbf{x} \sim \mathcal{N}_p \left(\boldsymbol{\mu}_\Theta + \boldsymbol{\Sigma}_{\Theta X} \boldsymbol{\Sigma}_{XX}^{-1} (\mathbf{x} - \boldsymbol{\mu}_X), \boldsymbol{\Sigma}_{\Theta\Theta} - \boldsymbol{\Sigma}_{\Theta X} \boldsymbol{\Sigma}_{XX}^{-1} \boldsymbol{\Sigma}_{X\Theta} \right)$$

cf. exercise 1.10-7.

Theorem:

Given the linear model

$$\mathbf{X} = \mathbf{H}\Theta + \mathbf{Z},$$

where \mathbf{H} is a known matrix, Θ and \mathbf{Z} are statistically independent random vectors with $\Theta \sim \mathcal{N}_p(\boldsymbol{\mu}_\Theta, \boldsymbol{\Sigma}_{\Theta\Theta})$ and $\mathbf{Z} \sim \mathcal{N}_n(\mathbf{0}, \boldsymbol{\Sigma}_{\mathbf{Z}\mathbf{Z}})$. Then

$$\mathbf{Y} = \begin{pmatrix} \mathbf{X} \\ \Theta \end{pmatrix} \sim \mathcal{N}_{n+p} \left(\begin{pmatrix} \mathbf{H}\boldsymbol{\mu}_\Theta \\ \boldsymbol{\mu}_\Theta \end{pmatrix}, \begin{pmatrix} \mathbf{H}\boldsymbol{\Sigma}_{\Theta\Theta}\mathbf{H}^T + \boldsymbol{\Sigma}_{\mathbf{Z}\mathbf{Z}} & \mathbf{H}\boldsymbol{\Sigma}_{\Theta\Theta} \\ \boldsymbol{\Sigma}_{\Theta\Theta}\mathbf{H}^T & \boldsymbol{\Sigma}_{\Theta\Theta} \end{pmatrix} \right)$$

and the MMSE estimate can be expressed by

$$\hat{\boldsymbol{\theta}}_{MMSE}(\mathbf{x}) = \boldsymbol{\mu}_\Theta + \boldsymbol{\Sigma}_{\Theta\Theta}\mathbf{H}^T (\mathbf{H}\boldsymbol{\Sigma}_{\Theta\Theta}\mathbf{H}^T + \boldsymbol{\Sigma}_{\mathbf{Z}\mathbf{Z}})^{-1} (\mathbf{x} - \mathbf{H}\boldsymbol{\mu}_\Theta).$$

Exercise 3.7-1:
Proof of the Theorem

3.7.2 Minimum Mean Absolute Error Estimation

The mean absolute error estimation is considered for the single parameter case only. To estimate the random parameter Θ by means of the realization of the random vector \mathbf{X} in the minimum mean absolute error sense

$$R_{MAE}(\hat{\theta}) = E(|\Theta - \hat{\theta}(\mathbf{X})|) = \int_{\mathbb{R}^{n+1}} |\theta - \hat{\theta}(\mathbf{x})| f_{\mathbf{x}\Theta}(\mathbf{x}, \theta) d\mathbf{x} d\theta$$

has to be minimized. As shown in the subsequent exercise the minimization of $R_{MAE}(\hat{\theta})$ leads to

$$\int_{-\infty}^{\hat{\theta}_{MMAE}(\mathbf{x})} f_{\Theta}(\theta | \mathbf{x}) d\theta = \int_{\hat{\theta}_{MMAE}(\mathbf{x})}^{\infty} f_{\Theta}(\theta | \mathbf{x}) d\theta.$$

Hence, the $\hat{\theta}_{MMAE}(\mathbf{x})$ is the median of the posterior density function $f_{\Theta}(\theta | \mathbf{x})$.

Exercise 3.7-2:
Proof of median

3.7.3 Maximum A Posteriori Estimation

Another estimation technique that fits within the Bayesian framework is the maximum a posteriori estimation (MAP).

Single parameter case

To motivate the approach the uniform loss function

$$L(\theta, \hat{\theta}(\mathbf{x})) = \begin{cases} 0 & |\theta - \hat{\theta}(\mathbf{x})| \leq \delta \\ 1 & |\theta - \hat{\theta}(\mathbf{x})| > \delta \end{cases} \quad \text{where } \delta > 0$$

is introduced. The Bayes risk is then given by

$$\begin{aligned} R_L(\hat{\theta}) &= \mathbb{E}(L(\theta, \hat{\theta}(\mathbf{X}))) = \int_{\mathbb{R}^{n+1}} L(\theta, \hat{\theta}(\mathbf{x})) f_{\mathbf{x}\theta}(\mathbf{x}, \theta) d\mathbf{x} d\theta \\ &= \int_{\mathbb{R}^n} \left(\int_{-\infty}^{\infty} L(\theta, \hat{\theta}(\mathbf{x})) f_{\theta}(\theta | \mathbf{x}) d\theta \right) f_{\mathbf{x}}(\mathbf{x}) d\mathbf{x}. \end{aligned}$$

Since $f_{\mathbf{x}}(\mathbf{x})$ and the integral in the brackets are always positive it is sufficient to minimize

$$\int_{-\infty}^{\infty} L(\theta, \hat{\theta}(\mathbf{x})) f_{\theta}(\theta | \mathbf{x}) d\theta = \int_{-\infty}^{\hat{\theta}(\mathbf{x}) - \delta} f_{\theta}(\theta | \mathbf{x}) d\theta + \int_{\hat{\theta}(\mathbf{x}) + \delta}^{\infty} f_{\theta}(\theta | \mathbf{x}) d\theta$$

or equivalently to maximize

$$1 - \left(\int_{-\infty}^{\hat{\theta}(\mathbf{x}) - \delta} f_{\theta}(\theta | \mathbf{x}) d\theta + \int_{\hat{\theta}(\mathbf{x}) + \delta}^{\infty} f_{\theta}(\theta | \mathbf{x}) d\theta \right) = \int_{\hat{\theta}(\mathbf{x}) - \delta}^{\hat{\theta}(\mathbf{x}) + \delta} f_{\theta}(\theta | \mathbf{x}) d\theta$$

for every \mathbf{x} .

For δ arbitrary small $\hat{\theta}(\mathbf{x})$ determines the mode of $f_{\theta}(\theta | \mathbf{x})$ and is termed maximum a posteriori (MAP) estimate, i.e.

$$\begin{aligned} \hat{\theta}_{MAP}(\mathbf{x}) &= \operatorname{argmax}_{\theta} f_{\theta}(\theta | \mathbf{x}) = \operatorname{argmax}_{\theta} f_{\mathbf{x}}(\mathbf{x} | \theta) f_{\theta}(\theta) \\ &= \operatorname{argmax}_{\theta} \left(\ln(f_{\mathbf{x}}(\mathbf{x} | \theta)) + \ln(f_{\theta}(\theta)) \right) \end{aligned}$$

Multiple parameter case

Employing the marginal posterior density functions

$$f_{\theta_i}(\theta_i | \mathbf{x}) = \int_{\mathbb{R}^{p-1}} f_{\boldsymbol{\theta}}(\boldsymbol{\theta} | \mathbf{x}) d\theta_1 \dots d\theta_{i-1} d\theta_{i+1} \dots d\theta_p$$

the MAP estimates are given by

$$\hat{\theta}_{i,MAP}(\mathbf{x}) = \underset{\theta_i}{\operatorname{argmax}} f_{\theta_i}(\theta_i | \mathbf{x}) \quad i = 1, \dots, p.$$

However, due to the required integration we might propose the alternative vector valued MAP estimate

$$\hat{\boldsymbol{\theta}}_{MAP}(\mathbf{x}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} f_{\boldsymbol{\theta}}(\boldsymbol{\theta} | \mathbf{x}) = \underset{\boldsymbol{\theta}}{\operatorname{argmax}} f_{\mathbf{x}}(\mathbf{x} | \boldsymbol{\theta}) f_{\boldsymbol{\theta}}(\boldsymbol{\theta}).$$

Remark:

Generally, both MAP estimates are not equivalent.

Exercise 3.7-3:

MMSE, MMAE and MAP estimation of the parameter of an exponential distribution, where the parameter obeys also an exponential distribution

Exercise 3.7-4:

Signal + noise, MMSE signal amplitude estimation

References to Chapter 3

- [1] M. Fisz, *Probability Theory and Mathematical Statistics*, Krieger Publishing Company, 1980
- [2] C.W. Helstrom, *Elements of Signal Detection and Estimation*, Prentice Hall, 1994
- [3] S. Kay, *Fundamentals of Statistical Signal Processing, Vol. 1: Estimation Theory*, Prentice Hall, 1993
- [4] E.L. Lehmann and G. Casella, *Theory of Point Estimation*, Springer, 2003
- [5] H.V. Poor, *An Introduction to Signal Detection and Estimation*, Springer, 1994
- [6] C.R. Rao, *Linear Statistical Inference and Its Application*, John Wiley, 1973
- [7] L.L. Scherf, *Statistical Signal Processing*, Addison Wesley, 1990