

Lecture Notes

Statistical Signal Processing

Univ.-Prof. Dr.-Ing. Wolfgang Utschick

Technische Universität München

Summer 2019

© 2019 Univ.-Prof. Dr.-Ing. Wolfgang Utschick

A circulation of this document to other parties without a written consent of the author is forbidden.

Email: utschick@tum.de

Layout by \LaTeX _{2 ϵ}

Part I

Parameter Estimation

1. Statistic Modeling

STATISTIC ESTIMATION treats the problem of inferring underlying characteristics of unknown random variables on the basis of observations of outputs of those random variables.

The basic problem of statistical estimation is to infer the PROBABILITY MEASURE P , which the realizations of the respective random variables $X : \Omega \rightarrow \mathbb{X}$ are subject to.

The most critical part of any parameter estimation problem is the choice of a proper STATISTIC MODEL $(\Omega, \mathbb{F}, P_\theta)$, with the metric space (\mathbb{X}, \mathbb{B}) and

$$\text{OBSERVATION SPACE : } \mathbb{X}, \quad (1.1)$$

$$\text{SIGMA ALGEBRA : } \mathbb{F}, \quad (1.2)$$

$$\text{PROBABILITY MEASURE : } P_\theta, \theta \in \Theta. \quad (1.3)$$

The stochastic model is a set of PROBABILITY SPACES and the task of statistic estimation is to select the most appropriate candidate.

1.1 Standard Model

Definition. We call the introduced statistical model STANDARD MODEL and the inference problem PARAMETER ESTIMATION, if

$$\Theta \subset \mathbb{R}^D, \tag{1.4}$$

and the random variable X is either DISCRETE or CONTINUOUS.

Commonly used terminology:

- Random variables $X_i : \Omega \rightarrow \mathbb{X}$ are STATISTICS, and
- its realizations x_i of X_i are OBSERVATIONS, SAMPLES, MEASUREMENTS, etc.

Definition. A special statistic is the random variable $T : \mathbb{X} \rightarrow \Theta$, which maps one or multiple samples to $\theta \in \Theta$ or another parameter depending thereof. The random variable or statistic $T : \mathbb{X} \rightarrow \Theta$ is called ESTIMATOR.

1.2 Introductory Example

Given observations X_1, \dots, X_N of a uniquely distributed random variable

$$X : \Omega \rightarrow [0, \theta], \quad (1.5)$$

with $[0, \theta] \subset \mathbb{R}$ such that $F_X(\xi) = \frac{\xi}{\theta}$ and $f_X(\xi) = \frac{1}{\theta}$, if $0 \leq \xi \leq \theta$.

The unknown parameter θ , which describes the random variable X is DETERMINISTIC and UNKNOWN.

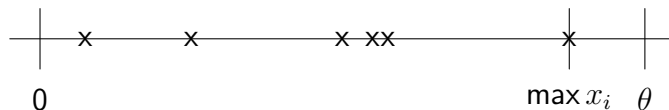


Fig. 1.1: Estimating the upper bound of an interval.

How to estimate the upper bound? Any guesses?

How to estimate the upper bound?

1. Attempt: Given $E[X] = \theta/2$, we conclude for the statistics X_i

$$T_1 = 2 \cdot \underbrace{\frac{1}{N} \sum_{i=1}^N X_i}_{\text{Average}} : \quad x_1, \dots, x_N \mapsto \hat{\theta}_1. \quad (1.6)$$

2. Attempt: Since for large N the maximum observed value will be close to the upper bound, we conclude

$$T_2 = \max_{i=1, \dots, N} \{X_i\} : \quad x_1, \dots, x_N \mapsto \hat{\theta}_2. \quad (1.7)$$

How reliable are these attempts?

Estimator T_1 : According to the LAW OF LARGE NUMBERS (CHEBYSHEV INEQUALITY):

$$P_{\theta} \left(\left| \frac{T_1}{2} - \frac{\theta}{2} \right| \geq \epsilon \right) \leq \frac{\text{Var}[X]}{N\epsilon^2} \xrightarrow{N \rightarrow \infty} 0. \quad (1.8)$$

Estimator T_2 : Again with an asymptotic approach:

$$\begin{aligned} P_{\theta} (|T_2 - \theta| \geq \epsilon) &= P_{\theta} (X_1 \leq \theta - \epsilon, \dots, X_N \leq \theta - \epsilon) \\ &= \prod_{i=1}^N P_{\theta} (X_i \leq \theta - \epsilon) \\ &= \prod_{i=1}^N \frac{\theta - \epsilon}{\theta} \\ &= \left(\frac{\theta - \epsilon}{\theta} \right)^N \\ &= \underbrace{\left(1 - \frac{\epsilon}{\theta} \right)}_{<1}^N \xrightarrow{N \rightarrow \infty} 0. \end{aligned} \quad (1.9)$$

2. Quality Criteria

2.1 Consistency and Unbiasedness

Definition. CONSISTENT ESTIMATORS are characterized by

$$\lim_{N \rightarrow \infty} T(x_1, \dots, x_N) = \theta. \quad (2.1)$$

Obviously, T_1 and T_2 are consistent.

Definition. UNBIASED ESTIMATORS are characterized by

$$\mathbb{E} [T(X_1, \dots, X_N)] = \theta. \quad (2.2)$$

T_1 and T_2 are random variables which depend on the randomly drawn observations X_1, \dots, X_N . UNBIASEDNESS means that averaging over the entire sample set gets the true parameter:

- T_1 is UNBIASED, since

$$\mathbb{E}[T_1(X_1, \dots, X_N)] = \mathbb{E}\left[\frac{2}{N} \sum_{i=1}^N X_i\right] = \frac{2}{N} \sum_{i=1}^N \mathbb{E}[X_i] = \frac{2}{N} \cdot N \cdot \frac{\theta}{2} = \theta. \quad (2.3)$$

- T_2 is ASYMPTOTICALLY UNBIASED, since

$$\mathbb{E}[T_2(X_1, \dots, X_N)] = \int_0^\theta \xi \cdot f_{T_2}(\xi) \, d\xi = \frac{N}{N+1} \theta, \quad (2.4)$$

where

$$f_{T_2}(\xi) = \frac{d}{d\xi} F_{T_2}(\xi) = \frac{d}{d\xi} \left(\frac{\xi}{\theta}\right)^N = N \cdot \frac{\xi^{N-1}}{\theta^N}, \quad \xi \in [0, \theta]. \quad (2.5)$$

- Consequently, we introduce an unbiased version of estimator T_2' by

$$T_2' = \frac{N+1}{N} T_2. \quad (2.6)$$

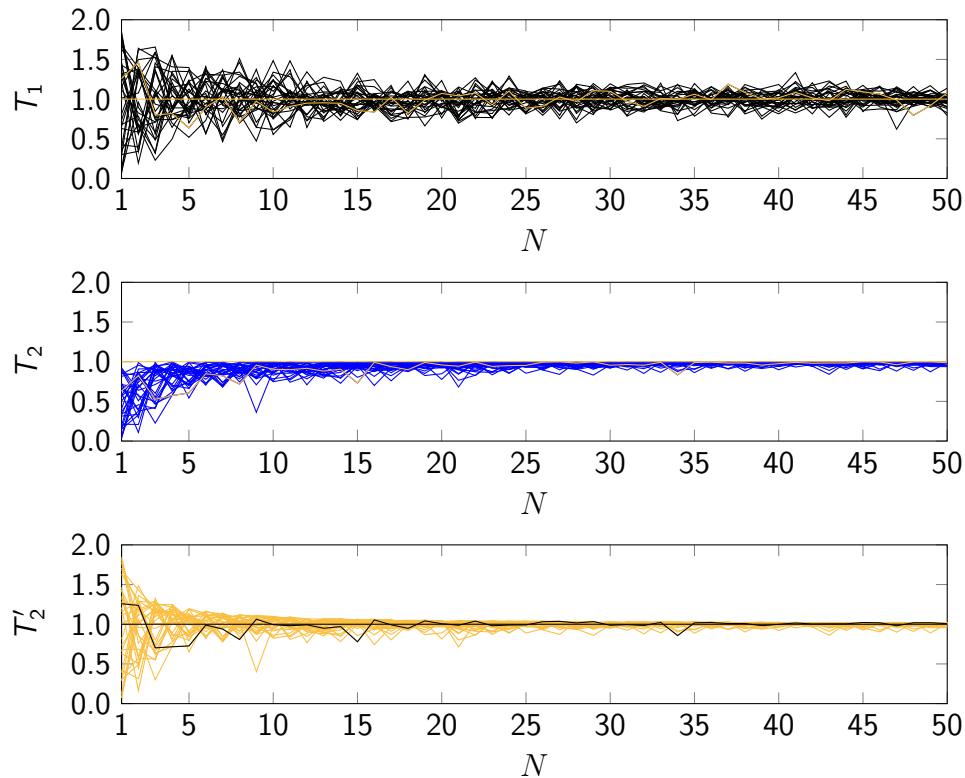


Fig. 2.1: T_1 , T_2 , and T'_2 for 30 uniformly chosen (independent and identically distributed) samples of size $N = 1 \dots 50$ with $\theta = 1$.

2.2 Variance

Definition. A further quality measure for an estimator is its VARIANCE

$$\text{Var} [T(X_1, \dots, X_N)] = E \left[(T - E [T])^2 \right]. \quad (2.7)$$

For the estimators T_1 , T_2 and T'_2 of the introductory example we obtain

$$\text{Var} [T_1] = \frac{\theta^2}{3N}, \quad (2.8)$$

$$\text{Var} [T_2] = \frac{N\theta^2}{(N+1)^2(N+2)}, \quad (2.9)$$

$$\text{Var} [T'_2] = \left(\frac{N+1}{N} \right)^2 \text{Var} [T_2] \quad (2.10)$$

$$= \frac{\theta^2}{N(N+2)}. \quad (2.11)$$

Table 2.1: Variances of the estimators T_1 , T_2 and T'_2 for $N = 10$ observations.

Estimator T	Variance σ_T^2
T_1	$\frac{5}{150}\theta^2$
T_2	$\frac{5}{726}\theta^2$
T'_2	$\frac{5}{600}\theta^2$

2.3 Mean Square Error

Definition. An extension of the VARIANCE is the MEAN SQUARE ERROR (MSE)

$$\varepsilon[T] = E \left[(T - \theta)^2 \right]. \quad (2.12)$$

For the estimators T_1 and T_2' of the introductory example the MSE is already known, since these estimators are unbiased and thus $E[T_1] = E[T_2'] = \theta$ and $\varepsilon[T] = \text{Var}[T]$.

The MSE of the 2nd estimator can be obtained by

$$\begin{aligned} \varepsilon[T_2] &= E \left[(T_2 - \theta)^2 \right] \\ &= E \left[(T_2)^2 \right] - 2 E[T_2] \theta + \theta^2 \\ &= \dots \\ &= \frac{2\theta^2}{(N+2)(N+1)}. \end{aligned} \quad (2.13)$$

Table 2.2: MSEs of the estimators T_1 , T_2 and T'_2 for $N = 10$ observations.

Estimator T	MSE $\varepsilon[T]$
T_1	$\frac{1}{30}\theta^2$
T_2	$\frac{1}{66}\theta^2$
T'_2	$\frac{1}{120}\theta^2$

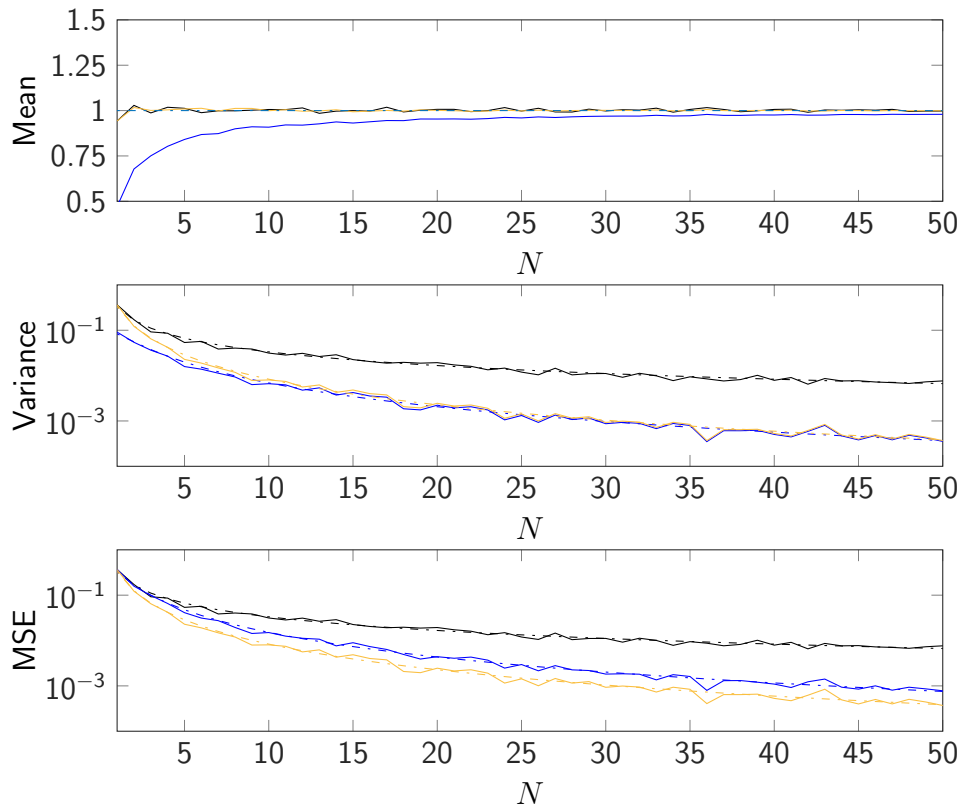


Fig. 2.2: Estimated (solid) and analytical (dashed-dotted) MEAN, VARIANCE, and MSE for T_1 (black), T_2 (blue), and T_2' (yellow) for 30 uniformly chosen (i.i.d.) samples of size $N = 1 \dots 50$ with $\theta = 1$.

2.4 Bias/Variance Trade–Off

The MSE of an estimator T ,

$$\varepsilon[T] = \mathbb{E} \left[(T(X_1, \dots, X_N) - \theta)^2 \right], \quad (2.14)$$

can be decomposed in its BIAS and its VARIANCE, i.e.,¹

$$\varepsilon[T] = \underbrace{\mathbb{E} \left[(T - \mathbb{E}[T])^2 \right]}_{\text{Var}[T]} + \left(\underbrace{\mathbb{E}[T] - \theta}_{\text{Bias}[T]} \right)^2 \quad (2.15)$$

$$= \text{Var}[T] + (\text{Bias}[T])^2. \quad (2.16)$$

Bias and Variance of an estimator cannot be minimized independently.

¹The decomposition can easily be shown by taking the expectation of the expansion of $(T - \mathbb{E}[T] + \mathbb{E}[T] - \theta)^2$.

2.5 Introductory Example (cont'd)

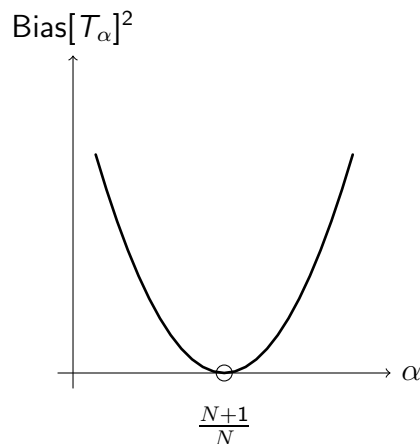
Another estimator T_α can be constructed by

$$T_\alpha = \alpha \cdot T_2 : \quad x_1, \dots, x_N \rightarrow \hat{\theta}_\alpha, \quad (2.17)$$

with

$$\text{Bias}[T_\alpha] = E[T_\alpha] - \theta = \alpha E[T_2] - \theta = \left(\frac{\alpha N}{N+1} - 1 \right) \theta. \quad (2.18)$$

The estimator T_α is obviously unbiased if $\alpha = \frac{N+1}{N}$.



3. Maximum Likelihood Estimation

3.1 Maximum Likelihood Principle

Given the statistical model $\{\mathbb{X}, \mathbb{F}, P_\theta; \theta \in \Theta\}$, the MAXIMUM-LIKELIHOOD (ML) principle suggest to select a candidate probability measure P_θ , such the observed outcomes of the experiment become most probable.

A maximum likelihood estimator T_{ML} picks the $\hat{\theta} \in \Theta$ which maximizes the LIKELIHOOD FUNCTION, i.e.

$$T_{\text{ML}} : x_1, \dots, x_N \mapsto \operatorname{argmax}_{\theta \in \Theta} \{L(x_1, \dots, x_N; \theta)\}. \quad (3.1)$$

Definition. The LIKELIHOOD FUNCTION depends on the statistical model $\{\mathbb{X}, \mathbb{F}, P_\theta; \theta \in \Theta\}$, i.e.,

$$\text{DISCRETE R.V. : } L(x_1, \dots, x_N; \theta) = P_\theta(x_1, \dots, x_N), \quad (3.2)$$

$$\text{CONTINUOUS R.V. : } L(x_1, \dots, x_N; \theta) = f_{X_1, \dots, X_N}(x_1, \dots, x_N; \theta). \quad (3.3)$$

Assuming that all observations are drawn from IDENTICALLY INDEPENDENTLY DISTRIBUTED (i.i.d.) random variables X_1, \dots, X_N ,

$$P_\theta(x_1, \dots, x_N) = \prod_{i=1}^N P_\theta(x_i), \quad (3.4)$$

$$f_{X_1, \dots, X_N}(x_1, \dots, x_N; \theta) = \prod_{i=1}^N f_{X_i}(x_i; \theta), \quad (3.5)$$

with $f_{X_1} = \dots = f_{X_N}$.

Log-Likelihood Function.

Due to the monotonicity of the log-function and assuming i.i.d. random variables, we obtain the following expressions:

$$T_{\text{ML}} : x_1, \dots, x_N \mapsto \operatorname{argmax}_{\theta \in \Theta} \{ \log(L(x_1, \dots, x_N; \theta)) \} \quad (3.6)$$

$$= \operatorname{argmax}_{\theta \in \Theta} \left\{ \log \left(\prod_{i=1}^N L(x_i; \theta) \right) \right\} \quad (3.7)$$

$$= \operatorname{argmax}_{\theta \in \Theta} \left\{ \sum_{i=1}^N \log(L(x_i; \theta)) \right\}. \quad (3.8)$$

Note. The computation of the maximum likelihood estimate obviously involves the solution of an optimization problem.

3.2 Channel Estimation Example

We consider the estimation of the attenuation coefficient h of a Single-Input Single-Output transmission channel with ADDITIVE WHITE GAUSSIAN NOISE (AWGN) $N_i \sim \mathcal{N}(0, \sigma^2)$ at the receiver.

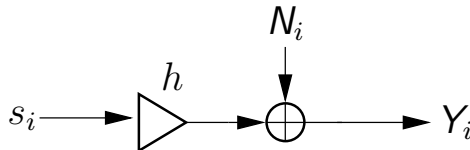


Fig. 3.1: Channel model.

The model of the receiver signal is

$$Y_i = hs_i + N_i, \quad Y_i \sim \mathcal{N}(hs_i, \sigma^2), \quad (3.9)$$

where s_i denotes the i -th training signal.

Assuming N observations y_1, \dots, y_N according to the training signals s_1, \dots, s_N , the likelihood function is given by $L(y_1, \dots, y_N; \theta) = \prod_{i=1}^N f_{Y_i}(y_i; \theta)$, with $\theta = h$ and

$$f_{Y_i}(y_i) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2} (y_i - \theta s_i)^2\right). \quad (3.10)$$

The maximum-likelihood estimation is derived by

$$\begin{aligned}
\hat{h}_{\text{ML}} &= \underset{\theta \in \mathbb{R}}{\operatorname{argmax}} \left\{ \sum_{i=1}^N \log \left(\frac{1}{\sqrt{2\pi}\sigma} \exp \left(-\frac{1}{2\sigma^2} \cdot (y_i - \theta s_i)^2 \right) \right) \right\} \\
&= \underset{\theta \in \mathbb{R}}{\operatorname{argmax}} \left\{ \sum_{i=1}^N -\frac{1}{2\sigma^2} \cdot (y_i - \theta s_i)^2 \right\} \\
&= \underset{\theta \in \mathbb{R}}{\operatorname{argmin}} \left\{ \sum_{i=1}^N \frac{1}{2\sigma^2} \cdot (y_i - \theta s_i)^2 \right\} \\
&= \underset{\theta \in \mathbb{R}}{\operatorname{argmin}} \left\{ \frac{1}{2\sigma^2} \cdot \|\mathbf{y} - \theta \mathbf{s}\|^2 \right\} \\
&= \underset{\theta \in \mathbb{R}}{\operatorname{argmin}} \{ \|\mathbf{y} - \theta \mathbf{s}\|^2 \}. \tag{3.11}
\end{aligned}$$

where $\mathbf{s} = [s_1, \dots, s_N]^\top$. The respective estimator \hat{h}_{ML} is given by the PSEUDO INVERSE $(\mathbf{s}^\top \mathbf{s})^{-1} \mathbf{s}^\top$ of the vector \mathbf{s} of training signals,¹ i.e.,

$$\hat{h}_{\text{ML}} = (\mathbf{s}^\top \mathbf{s})^{-1} \mathbf{s}^\top \mathbf{y}. \tag{3.12}$$

The ML estimator is obviously identical with the LEAST SQUARES estimator. This changes considerably when the statistics Y_i are correlated or when N is non-Gaussian distributed.

¹Using the optimality condition $\frac{d}{d\theta} \|\mathbf{y} - \theta \mathbf{s}\|^2 = -2\mathbf{s}^\top \mathbf{y} + 2\theta \mathbf{s}^\top \mathbf{s} = 0$, we obtain $\theta_{\text{ML}} = \frac{\mathbf{s}^\top \mathbf{y}}{\mathbf{s}^\top \mathbf{s}}$.

3.3 Introductory Example (cont'd)

The likelihood function for a single observation ($N = 1$) reads

$$L(x; \theta) = \begin{cases} \frac{1}{\theta}, & x \leq \theta \\ 0 & \text{otherwise.} \end{cases}$$

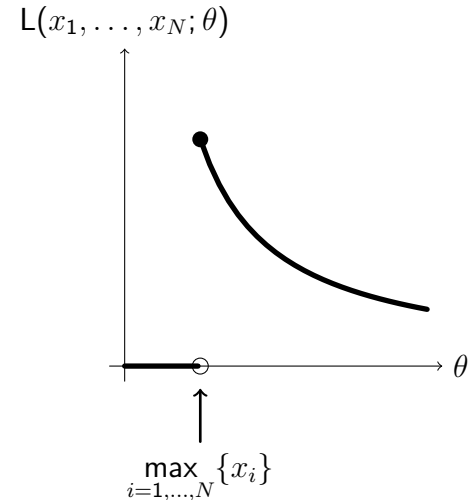
If $N > 1$:

$$L(x_1, \dots, x_N; \theta) = \begin{cases} \frac{1}{\theta^N}, & x_i \leq \theta \quad \forall i = 1, \dots, N \\ 0 & \text{otherwise.} \end{cases}$$

The ML Estimator is given by

$$\begin{aligned} \hat{\theta} &= \operatorname{argmax}_{\theta \in \mathbb{R}} \{L(x_1, \dots, x_N; \theta)\} \\ &= \max_{i=1, \dots, N} \{x_i\}, \end{aligned}$$

which corresponds to the 2nd intuitive attempt.



3.4 Bernoulli Experiment

We now consider a BERNOULLI EXPERIMENT with success probability θ , e.g. the transmission of a data packet over a link with probabilities

$$P(\text{no erasure}) = \theta, \quad (3.13)$$

$$P(\text{erasure}) = 1 - \theta. \quad (3.14)$$

We now study the maximum-likelihood framework for estimating the unknown parameter θ based on N independent observations:

The statistical model $\{\mathbb{X}, \mathbb{F}, P_\theta; \theta \in \Theta\}$ is

$$\left\{ \underbrace{\{0, 1, \dots, N\}}_{\substack{\text{number of successes} \\ \text{over } N \text{ observations}}}, \{0, 1, \dots, N\}^N, B_{N,\theta}; \theta \in [0, 1] \right\},$$

and the respective random variable X , which counts the number of successful Bernoulli trials within N attempts. X is BINOMIALLY DISTRIBUTED according to

$$B_{N,\theta}(x) = \binom{N}{x} \theta^x (1 - \theta)^{N-x}, \quad x \in \{0, 1, \dots, N\}.$$

Given the likelihood function

$$L(x; \theta) = \binom{N}{x} \theta^x (1 - \theta)^{N-x},$$

the log-likelihood function reads

$$\log L(x; \theta) = \log \binom{N}{x} + x \log \theta + (N - x) \log (1 - \theta).$$

The ML-Estimator can be obtained by

$$\hat{\theta} = \operatorname{argmax}_{\theta \in \Theta} \{\log L(x; \theta)\},$$

i.e.,

$$\begin{aligned} & \left. \frac{d}{d\theta} \log L(x; \theta) \right|_{\theta=\hat{\theta}} = 0 \\ \Leftrightarrow & \frac{x}{\hat{\theta}} + (N - x) \frac{1}{1 - \hat{\theta}} (-1) = 0 \\ \Leftrightarrow & T_{\text{ML}} : x \mapsto \hat{\theta} = \frac{x}{N}. \end{aligned}$$

Quality of $T_{\text{ML}} : x \mapsto \frac{x}{N}$.

Since X is binomially distributed

$$\mathbb{E}[X] = N\theta, \quad (3.15)$$

$$\text{Var}[X] = N\theta(1 - \theta), \quad (3.16)$$

we obtain

$$\mathbb{E}[T_{\text{ML}}] = \mathbb{E}\left[\frac{X}{N}\right] = \frac{1}{N} \mathbb{E}[X] = \frac{1}{N} \theta N = \theta, \quad (3.17)$$

$$\text{Var}[T_{\text{ML}}] = \text{Var}\left[\frac{X}{N}\right] = \frac{1}{N^2} \text{Var}[X] = \frac{1}{N^2} N\theta(1 - \theta) = \frac{\theta(1 - \theta)}{N}. \quad (3.18)$$

Obviously, the ML Estimator of the success probability θ is unbiased, and the MSE is equal to the variance of the estimator,

$$\text{Bias}[T_{\text{ML}}] = 0, \quad (3.19)$$

$$\varepsilon[T_{\text{ML}}] = \text{Var}[T_{\text{ML}}] = \frac{\theta(1 - \theta)}{N}. \quad (3.20)$$

Can we improve the estimator?

Alternative Solution: $T' : x \mapsto \frac{x+1}{N+2}$.

$$E[T'] = E\left[\frac{X+1}{N+2}\right] = \frac{N\theta + 1}{N+2} \Rightarrow \text{Bias}[T'] = \frac{1-2\theta}{N+2}, \quad (3.21)$$

$$\begin{aligned} \varepsilon[T'] &= \text{Var}[T'] + \text{Bias}[T']^2 \\ &= \text{Var}\left[\frac{X+1}{N+2}\right] + \left(\frac{1-2\theta}{N+2}\right)^2 \\ &= \frac{\text{Var}[X+1]}{(N+2)^2} + \left(\frac{1-2\theta}{N+2}\right)^2 \\ &= \frac{\text{Var}[X]}{(N+2)^2} + \left(\frac{1-2\theta}{N+2}\right)^2 \\ &= \frac{N\theta(1-\theta) + (1-2\theta)^2}{(N+2)^2}. \end{aligned} \quad (3.22)$$

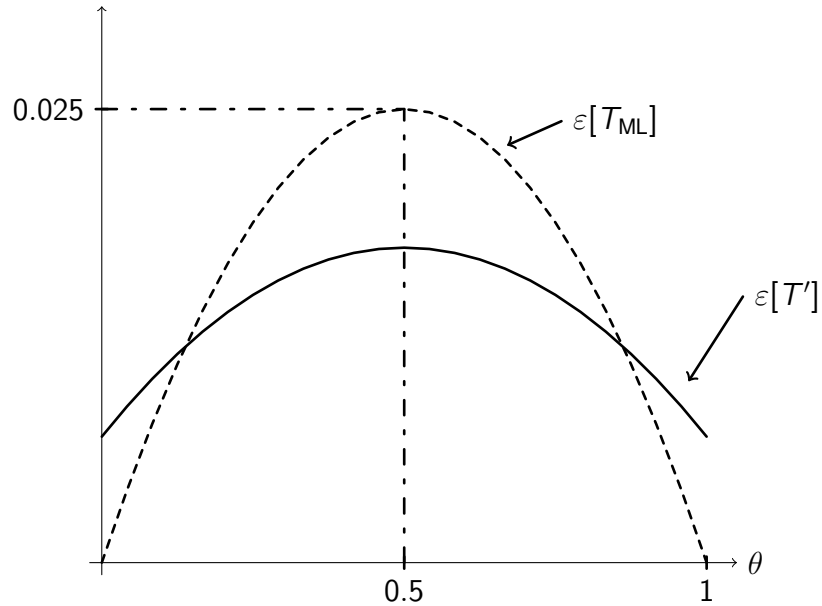


Fig. 3.2: Comparison of the achievable MSEs for $N = 10$ observations.

Note. Biased estimators can provide better estimates than unbiased estimators.

3.5 Best Unbiased Estimator

ML Estimators are not necessarily the best estimators. A wide class of estimators is defined by minimizing the MSE under an unbiasedness constraint.

Definition. Given the statistical model $\{\mathbb{X}, \mathbb{F}, P_\theta; \theta \in \Theta\}$, we call an estimator **BEST UNBIASED ESTIMATOR** if $E[T(X_1, \dots, X_N)] = \theta$ and

$$\text{Var}[T(X_1, \dots, X_N)] \leq \text{Var}[S(X_1, \dots, X_N)], \quad \forall \theta \in \Theta, \quad (3.23)$$

for any alternative unbiased estimator S .

Best unbiased estimators are also referred to as **UNIFORMLY MINIMUM VARIANCE UNBIASED (UMVU) Estimators**.

The estimator $T'_2 = \frac{N+1}{N} T_2$ of the introductory example is a UMVU Estimator.

4. Fisher's Information Inequality

An universal lower bound for the variance of an estimator can be introduced, if the following conditions are fulfilled:

$$L(x; \theta) > 0, \quad \forall x \in \mathbb{X}, \theta \in \Theta, \quad (4.1)$$

$$L(x; \theta) \text{ differentiable with respect to } \theta, \quad (4.2)$$

$$\int_{\mathbb{X}} \frac{\partial}{\partial \theta} L(x; \theta) \, dx = \frac{\partial}{\partial \theta} \int_{\mathbb{X}} L(x; \theta) \, dx. \quad (4.3)$$

Definition. Consequently, we define the SCORE FUNCTION as

$$g(x; \theta) = \frac{\partial}{\partial \theta} \log L(x; \theta) = \frac{\frac{\partial}{\partial \theta} L(x; \theta)}{L(x; \theta)}, \quad (4.4)$$

with $E[g(X; \theta)] = 0$.

4.1 Cramér-Rao Lower Bound

Given the score function $g(x; \theta)$ of an estimation problem, and the FISHER INFORMATION term $I_F(\theta) = \text{Var}[g(X, \theta)]$, the variance of any possible estimator can be lower bounded by

$$\text{Var}[T(X)] \geq \left(\frac{\partial \mathbb{E}[T(X)]}{\partial \theta} \right)^2 \frac{1}{I_F(\theta)}, \quad (4.5)$$

and

$$\text{Var}[T(X)] \geq \frac{1}{I_F(\theta)}, \quad (4.6)$$

if T is unbiased.

Proof.

$$\begin{aligned}
\frac{\partial \mathbb{E}[T(X)]}{\partial \theta} &= \frac{\partial}{\partial \theta} \int_{\mathbb{X}} T(x) f_X(x; \theta) dx = \int_{\mathbb{X}} T(x) \frac{\partial}{\partial \theta} f_X(x; \theta) dx \\
&= \int_{\mathbb{X}} (T(x) - \mathbb{E}[T(X)]) \frac{\partial}{\partial \theta} f_X(x; \theta) dx + \underbrace{\int_{\mathbb{X}} \mathbb{E}[T(X)] \frac{\partial}{\partial \theta} f_X(x; \theta) dx}_{=0} \\
&= \int_{\mathbb{X}} (T(x) - \mathbb{E}[T(X)]) \sqrt{f_X(x; \theta)} \frac{1}{\sqrt{f_X(x; \theta)}} \frac{\partial}{\partial \theta} f_X(x; \theta) dx \\
&= \left\langle (T(x) - \mathbb{E}[T(X)]) \sqrt{f_X(x; \theta)}, \frac{1}{\sqrt{f_X(x; \theta)}} \frac{\partial}{\partial \theta} f_X(x; \theta) \right\rangle. \tag{4.7}
\end{aligned}$$

Applying the CAUCHY-SCHWARZ INEQUALITY, $|\langle a, b \rangle|^2 \leq \langle a, a \rangle \cdot \langle b, b \rangle$, yields:

$$\begin{aligned}
\left(\frac{\partial \mathbb{E}[T(X)]}{\partial \theta} \right)^2 &\leq \int_{\mathbb{X}} (T(x) - \mathbb{E}[T(X)])^2 f_X(x; \theta) dx \cdot \int_{\mathbb{X}} \left(\frac{\partial}{\partial \theta} f_X(x; \theta) \right)^2 \frac{f_X(x; \theta)}{(f_X(x; \theta))^2} dx \\
&= \text{Var}[T(X)] \cdot \int_{\mathbb{X}} (g(x, \theta))^2 f_X(x; \theta) dx. \tag{4.8}
\end{aligned}$$

$$\Rightarrow \text{Var}[T(X)] \geq \left(\frac{\partial \mathbb{E}[T(X)]}{\partial \theta} \right)^2 \frac{1}{\text{Var}[g(X, \theta)]}. \tag{4.9}$$

4.2 2nd Version of Fisher Information

An alternative expression for the FISHER INFORMATION term $I_F(\theta)$ can be derived,

$$\begin{aligned} I_F(\theta) &= \text{Var} [g(X, \theta)] \\ &= \text{E} \left[g(X, \theta)^2 \right] - \underbrace{\text{E} [g(X, \theta)]^2}_{=0} \\ &= \text{E} \left[g(X, \theta)^2 \right] \end{aligned} \tag{4.10}$$

$$= -\text{E} \left[\frac{\partial^2}{\partial \theta^2} \log L(X; \theta) \right], \tag{4.11}$$

which can be interpreted as the negative MEAN CURVATURE of the LOG-LIKELIHOOD FUNCTION.

The last step from (4.10) to (4.11) is not obvious.

Proof.

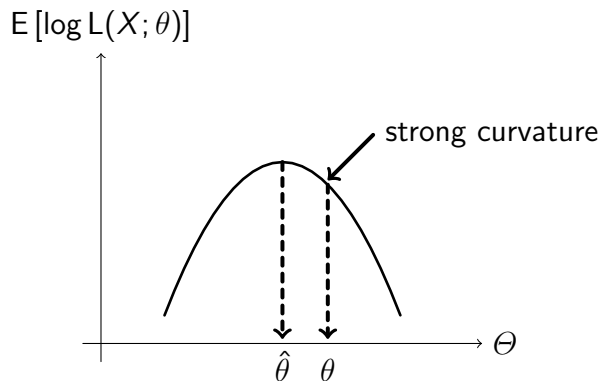
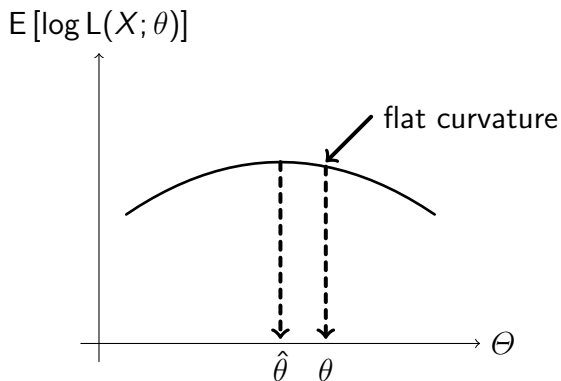
$$\begin{aligned}
\frac{\partial^2}{\partial \theta^2} \log L(x; \theta) &= \frac{\partial}{\partial \theta} \left(\frac{\partial}{\partial \theta} \log L(x; \theta) \right) \\
&= \frac{\partial}{\partial \theta} \left(\frac{1}{L(x; \theta)} \frac{\partial}{\partial \theta} L(x; \theta) \right) \\
&= \frac{\frac{\partial^2}{\partial \theta^2} L(x; \theta)}{L(x; \theta)} - \frac{\partial}{\partial \theta} L(x; \theta) \frac{\frac{\partial}{\partial \theta} L(x; \theta)}{L(x; \theta)^2} \\
&= \frac{\frac{\partial^2}{\partial \theta^2} L(x; \theta)}{L(x; \theta)} - \frac{\left(\frac{\partial}{\partial \theta} L(x; \theta) \right)^2}{L(x; \theta)^2} \\
&= \frac{\frac{\partial^2}{\partial \theta^2} L(x; \theta)}{L(x; \theta)} - g(x, \theta)^2.
\end{aligned} \tag{4.12}$$

$$\begin{aligned}
\Rightarrow E [g(X, \theta)^2] &= E \left[\frac{\frac{\partial^2}{\partial \theta^2} L(X; \theta)}{L(X; \theta)} - \frac{\partial^2}{\partial \theta^2} \log L(X; \theta) \right] \\
&= \underbrace{\int_{\mathbb{X}} \frac{\partial^2}{\partial \theta^2} L(x; \theta) dx}_{\frac{\partial^2}{\partial \theta^2} \int_{\mathbb{X}} L(x; \theta) dx = \frac{\partial^2}{\partial \theta^2} 1 = 0} - E \left[\frac{\partial^2}{\partial \theta^2} \log L(X; \theta) \right].
\end{aligned} \tag{4.13}$$

Interpretation.

The 2nd version of the FISHER INFORMATION $I_F(\theta)$ as the negative MEAN CURVATURE of the LOG-LIKELIHOOD FUNCTION allows an interpretation of the CRAMÉR-RAO LOWER BOUND.

$$I_F(\theta) = -E \left[\underbrace{\frac{\partial^2}{\partial \theta^2} \log L(X; \theta)}_{\text{curvature}} \right].$$



Properties of the Fisher Information:

- The log-likelihood function $\log(L(x_1, \dots, x_N))$ depends on given observations x_1, \dots, x_N .
- A weak curvature of $\log(L(x_1, \dots, x_N))$ corresponds to little information in x_1, \dots, x_N .
- A strong curvature of $\log(L(x_1, \dots, x_N))$ indicates more information in x_1, \dots, x_N .
- The Fisher information $I_F(\theta)$ refers to the negative mean curvature at the true parameter θ .
- $I_F(\theta)$ is based on the expectation of the negative mean curvature of $\log(L(X_1, \dots, X_N))$ with respect to the observation statistics X_1, \dots, X_N .
- $I_F(\theta)$ is a function of θ .
- $I_F(\theta)$ is monotonically increasing with the number of independent observation statistics, i.e., given the Fisher Information $I_F^{(1)}(\theta)$ for a single observation statistic ($N = 1$), then

$$I_F^{(N)}(\theta) = N \cdot I_F^{(1)}(\theta). \quad (4.14)$$

4.3 Exponential Models

Definition. A EXPONENTIAL MODEL is a statistical model $\{\mathbb{X}, \mathbb{F}, P_\theta; \theta \in \Theta\}$ with random variable X and

$$f_X(x; \theta) = \frac{h(x) \exp(a(\theta)t(x))}{\exp(b(\theta))}, \quad (4.15)$$

with $b(\theta) = \log \int_{\mathbb{X}} h(x) \exp(a(\theta)t(x)) dx$ for normalization.

Then the respective Fisher Information can be directly obtained by

$$I_F(\theta) = \frac{\partial a(\theta)}{\partial \theta} \frac{\partial E[t(X)]}{\partial \theta}. \quad (4.16)$$

For the special case that $f_X(x; \theta)$ can be arranged such that $E[t(X)] = \theta$, the unbiased estimator

$$T : x \mapsto t(x). \quad (4.17)$$

can directly be obtained. Otherwise, $t(x)$ at least provides a SUFFICIENT STATISTIC for estimating θ .

Given an EXPONENTIAL MODEL it can be shown that the

- UNBIASED ESTIMATOR $T : x \mapsto t(x)$
- is a UMVU ESTIMATOR,
- and achieves the CRAMÉR-RAO LOWER BOUND.

4.4 Mean Estimation Example

We consider the estimation of the unknown mean value $\theta = \mu$ of a random variable $X \sim \mathcal{N}(\mu, \sigma^2)$ based on a single observation.

The respective PDF can be arranged as follows:

$$\begin{aligned} f_X(x; \theta) &= L(x; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2} (x - \theta)^2\right) \\ &= \exp\left[\underbrace{\frac{\theta x}{\sigma^2}}_{a(\theta)t(x)} - \underbrace{\left(\frac{\theta^2}{2\sigma^2} + \frac{1}{2} \log(2\pi\sigma^2)\right)}_{b(\theta)}\right] \underbrace{\exp\left[-\frac{1}{2\sigma^2} x^2\right]}_{h(x)}, \end{aligned} \quad (4.18)$$

with $t(x) = x$ and $a(\theta) = \theta/\sigma^2$.

Thus the single observation UMVU Estimator $T^{(1)}$ of $\theta = \mu$ can be obtained as

$$T^{(1)} : \hat{\theta} = x, \quad \text{with} \quad \text{Bias}[T^{(1)}(X)] = 0, \quad (4.19)$$

which minimizes both the variance and the MSE criterion to

$$\text{Var}[T^{(1)}] = \frac{1}{I_F} = \left(\frac{\partial a(\theta)}{\partial \theta}\right)^{-1} = \sigma^2. \quad (4.20)$$

4.5 Asymptotically Efficient Estimators

Definition. An estimator $T(X_1, \dots, X_N)$ is ASYMPTOTICALLY EFFICIENT if¹

$$\sqrt{N} (T(X_1, \dots, X_N) - \theta) \stackrel{N \rightarrow \infty}{\sim} \mathcal{N}(0, I_F^{(1)}(\theta)^{-1}). \quad (4.21)$$

Given the following requirements are fulfilled,

- Θ is an open set,
- $f_X(x; \theta)$ and $\int_{\mathbb{X}} f_X(x; \theta) dx$ is twice differentiable with respect to θ ,
- $I_F(\theta) < \infty$,
- ...

then a MAXIMUM-LIKELIHOOD ESTIMATOR is ASYMPTOTICALLY EFFICIENT.

¹Convergence in distribution.

Part II

Examples

5. ML Principle for Direction of Arrival Estimation

We consider the estimation of the DIRECTION OF ARRIVAL (DoA) θ of an impinging planar wavefront by means of an antenna array with M antennas.

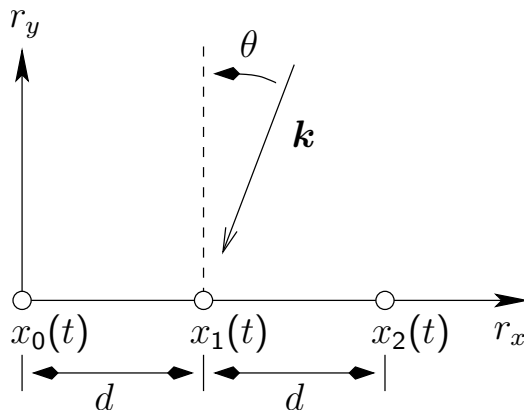


Fig. 5.1: Spatial sampling of a planar wavefront.

The received signal vector at the antenna array at time instant $t \in \mathbb{R}$ is equal to

$$\mathbf{x}(t) = \xi \mathbf{a}s(t) + \boldsymbol{\eta}(t) \in \mathbb{C}^M, \quad (5.1)$$

with the signal at the m th antenna element as

$$x_m(t) = \xi e^{j2\pi \sin(\theta)md/\lambda} s(t) + \eta_m(t), \quad m = 0, \dots, M-1. \quad (5.2)$$

The d and λ denote the distance between two adjacent antenna elements and the wavelength of the assumed NARROWBAND SIGNAL. The received signal is assumed to be corrupted by the Gaussian noise vector $\boldsymbol{\eta}(t) \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_\eta)$.

The $\xi \in \mathbb{R}$ represents the attenuation which the transmitted signal $s(t) \in \mathbb{R}$ experiences over the transmission path.

Without loss of generality we thus assume a UNIFORM LINEAR ARRAY (ULA) with M antenna elements, i.e.,

$$\mathbf{a} = \begin{bmatrix} \alpha^0 \\ \alpha^1 \\ \vdots \\ \alpha^{M-1} \end{bmatrix}, \quad \text{with} \quad \alpha = e^{j2\pi \sin(\theta)d/\lambda}. \quad (5.3)$$

In order to find the ML estimate of the DoA parameter θ , we consider the respective LIKELIHOOD FUNCTION for observations of the received vector $\mathbf{x}(t)$ at t_1, \dots, t_N ,¹ implicitly assuming a stationary scenario over the respective time intervall,

$$L(\mathbf{x}_1, \dots, \mathbf{x}_N; \theta) = \frac{1}{\pi^{MN} \det \mathbf{C}_\eta} \exp \left(- \sum_{i=1}^N (\mathbf{x}_i - \xi \mathbf{a}_{s_i})^H \mathbf{C}_\eta^{-1} (\mathbf{x}_i - \xi \mathbf{a}_{s_i}) \right). \quad (5.4)$$

For simplicity reasons we restrict the further analysis to the case of a single observation $N=1$ and ADDITIVE WHITE GAUSSIAN NOISE $\boldsymbol{\eta}(t) \sim \mathcal{N}(\mathbf{0}, \sigma_\eta^2 \mathbf{I})$, i.e.,

$$L(\mathbf{x}; \theta) = \frac{1}{(\pi \sigma_\eta^2)^M} \exp \left(- \frac{(\mathbf{x} - \xi \mathbf{a}_s)^H (\mathbf{x} - \xi \mathbf{a}_s)}{\sigma_\eta^2} \right). \quad (5.5)$$

¹For simplicity of notation we use \mathbf{x}_i instead of $\mathbf{x}(t_i)$.

After some simple reformulation steps the ML optimization problem is equal to

$$\min_{\theta} (\mathbf{x} - \xi \mathbf{a}s)^H (\mathbf{x} - \xi \mathbf{a}s) \quad (5.6)$$

Since the cost function can be expanded in

$$(\mathbf{x} - \xi \mathbf{a}s)^H (\mathbf{x} - \xi \mathbf{a}s) = \mathbf{x}^H \mathbf{x} - 2 \operatorname{Re} \{ \xi \mathbf{x}^H \mathbf{a}s \} + \xi^2 \mathbf{a}^H \mathbf{a} s^2, \quad (5.7)$$

and $\mathbf{a}^H \mathbf{a} = M$, it can be further reduced to

$$\max_{\theta} \operatorname{Re} \{ \mathbf{x}^H \mathbf{a}(\theta) \}. \quad (5.8)$$

The ML Estimator $\hat{\theta} = \operatorname{argmax}_{\theta} \{ \operatorname{Re} \{ \mathbf{x}^H \mathbf{a}(\theta) \} \}$ is obviously a NONLINEAR ESTIMATOR.

Note. In order to estimate the attenuation ξ , we need further information about the TRAINING SIGNAL (PILOT SIGNAL) $s(t)$.

The ML estimator for scenarios with a single impinging wavefront is equal to the so-called CONVENTIONAL BEAMFORMER.

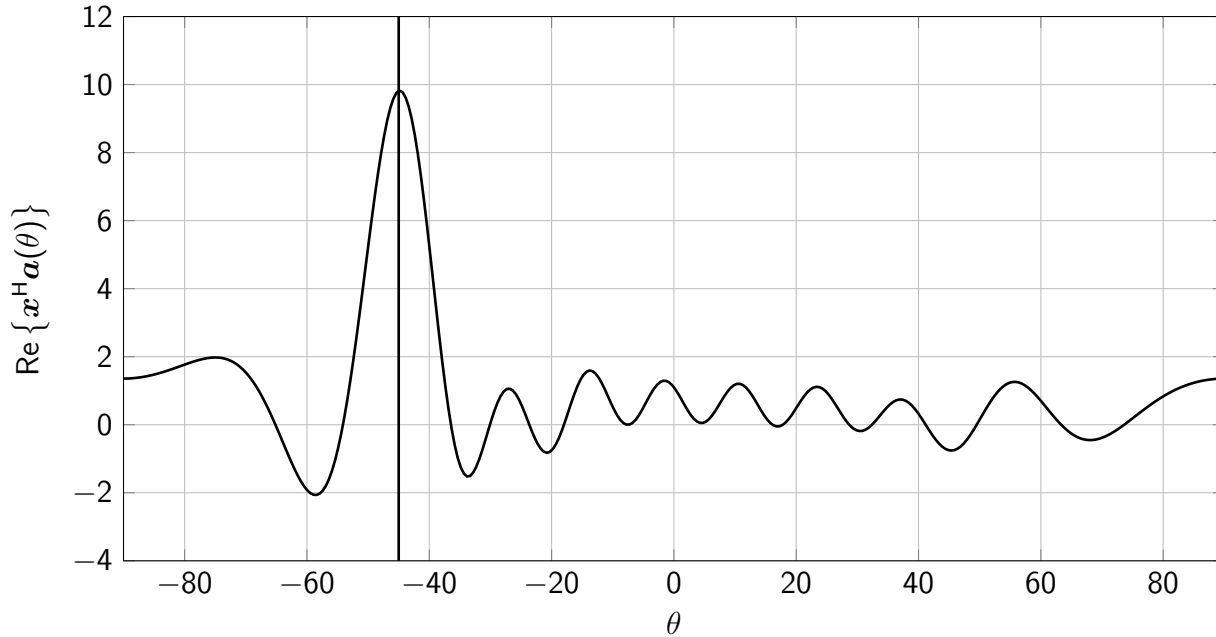


Fig. 5.2: $\text{Re}\{x^H a(\theta)\}$ for $M = 10$, $\xi = s = 1$, and $\sigma_\eta^2 = 0.1$ with DoA $\theta_1 = -45$ deg.

In scenarios with more than one impinging wavefront, the beamformer, i.e., the ML estimator for the case of a single wavefront, is only a heuristic method for DOA estimation.

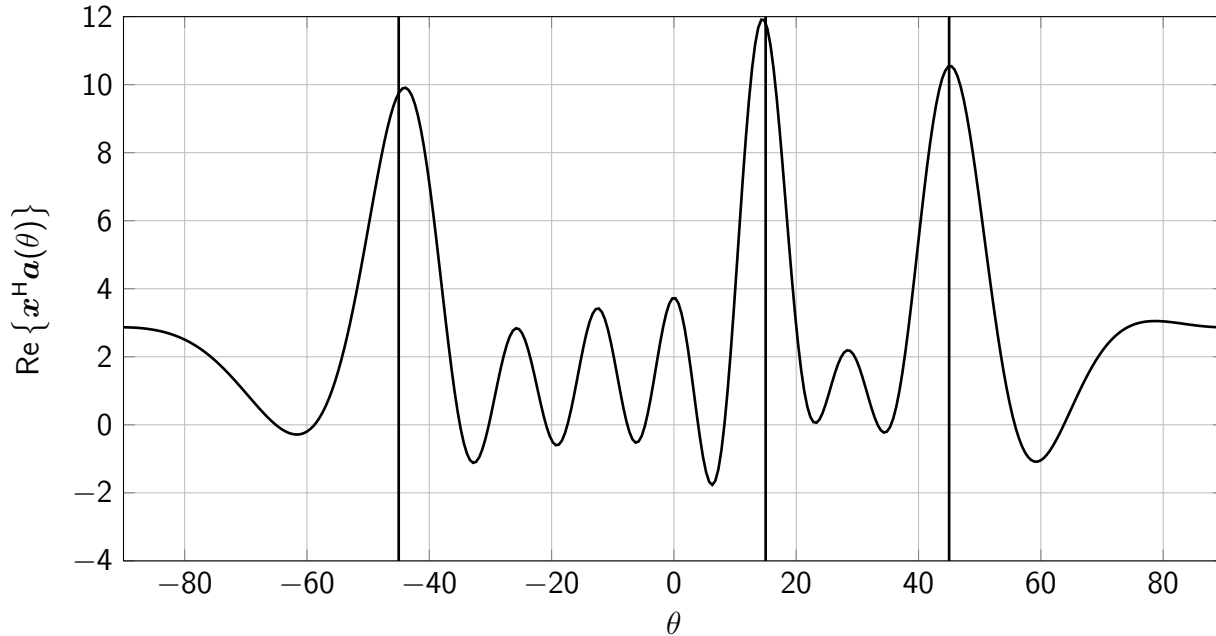


Fig. 5.3: $\text{Re}\{x^H a(\theta)\}$ for $M = 10$, $\xi = s = 1$, and $\sigma_\eta^2 = 0.1$ with DoAs $\theta_1 = -45^\circ$, $\theta_2 = 15^\circ$, and $\theta_3 = 45^\circ$.

5.1 Cramér-Rao Bound for DOA Estimation

The Likelihood function of the given estimation problem obviously belongs to the family of EXPONENTIAL DISTRIBUTIONS. We consider θ as parameter to be estimated. After some reformulation steps $L(\mathbf{x}; \theta)$ can be expressed as

$$L(\mathbf{x}; \theta) = \frac{1}{(\pi\sigma_\eta^2)^M} \exp\left(-\frac{(\mathbf{x} - \xi\mathbf{a}s)^H(\mathbf{x} - \xi\mathbf{a}s)}{\sigma_\eta^2}\right) = \frac{h(\mathbf{x}) \exp(\mathbf{c}^H(\theta)\mathbf{t}(\mathbf{x}))}{\exp(b(\theta))}, \quad (5.9)$$

with

$$h(\mathbf{x}) = \frac{1}{(\pi\sigma_\eta^2)^M} \exp\left(-\frac{\mathbf{x}^H\mathbf{x}}{\sigma_\eta^2}\right) \quad (5.10)$$

$$\mathbf{c}(\theta) = \frac{1}{\sigma_\eta^2} \begin{bmatrix} \xi\mathbf{a}s \\ (\xi\mathbf{a}s)^* \end{bmatrix} \quad (5.11)$$

$$\mathbf{t}(\mathbf{x}) = \begin{bmatrix} \mathbf{x} \\ \mathbf{x}^* \end{bmatrix} \quad (5.12)$$

$$b(\theta) = \frac{\xi^2 s^2 \mathbf{a}^H \mathbf{a}}{\sigma_\eta^2} = \frac{M\xi^2 s^2}{\sigma_\eta^2}. \quad (5.13)$$

Note, that $E[\mathbf{t}(\mathbf{x})] = \sigma_\eta^2 \mathbf{c}(\theta)$.

Consequently, the FISHER INFORMATION of the ML estimation problem can be obtained as²

$$I_F(\theta) = \left(\frac{\partial \mathbf{c}(\mathbf{x})}{\partial \theta} \right)^H \frac{\partial \mathbb{E}[\mathbf{t}(\mathbf{x})]}{\partial \theta} \quad (5.14)$$

$$= \frac{\xi^2 s^2}{\sigma_\eta^2} \frac{\partial}{\partial \theta} \left(\begin{bmatrix} \mathbf{a} \\ \mathbf{a}^* \end{bmatrix} \right)^H \frac{\partial}{\partial \theta} \left(\begin{bmatrix} \mathbf{a} \\ \mathbf{a}^* \end{bmatrix} \right). \quad (5.15)$$

Since θ is real,

$$I_F(\theta) = \frac{\xi^2 s^2}{\sigma_\eta^2} \left(\frac{\partial}{\partial \theta} \begin{bmatrix} \mathbf{a} \\ \mathbf{a}^* \end{bmatrix} \right)^H \left(\frac{\partial}{\partial \theta} \begin{bmatrix} \mathbf{a} \\ \mathbf{a}^* \end{bmatrix} \right) \quad (5.16)$$

$$= 2 \frac{\xi^2 s^2}{\sigma_\eta^2} \left\| \frac{\partial \mathbf{a}}{\partial \theta} \right\|^2. \quad (5.17)$$

Finally, with $\partial \alpha^m / \partial \theta = j m (2\pi d / \lambda) \cos \theta \alpha^m$ and $|\alpha^m| = 1$,

$$\begin{aligned} I_F(\theta) &= 2 \frac{\xi^2 s^2}{\sigma_\eta^2} \left(\frac{2\pi d}{\lambda} \cos \theta \right)^2 \sum_{m=1}^{M-1} m^2 \\ &= 2 \frac{\xi^2 s^2}{\sigma_\eta^2} \left(\frac{2\pi d}{\lambda} \cos \theta \right)^2 \frac{M(M-1)(2M-1)}{6} \propto \left(\frac{d}{\lambda} \right)^2 \frac{\xi^2 s^2}{\sigma_\eta^2} (\cos \theta)^2 M^3. \end{aligned} \quad (5.18)$$

² $\mathbb{E}[\mathbf{x}] = \xi \mathbf{a} s$.

References

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

Part III

Estimation of Random Variables

6. Bayes Estimation

In contrast to the Maximum-Likelihood principle, an A PRIORI information about the unknown parameter $\theta \in \Theta$ is now taken into account. We further assume that the a priori information of the unknown parameter is given by means of the

$$\text{PDF : } f_{\Theta}(\theta; \sigma) \quad (6.1)$$

$$\text{CONDITIONAL PDF : } f_{X|\Theta}(x|\theta), \quad (6.2)$$

where σ parameterizes the statistical model of the (now!) random variable Θ . We now consider the MEAN MSE with respect to the mean of the parameter Θ , i.e.,

$$\mathbb{E} [\mathbb{E} [(T(X) - \Theta)^2 | \Theta]] = \int_{\Theta} \underbrace{\int_{\mathbb{X}} (T(x) - \theta)^2 f_{X|\Theta}(x|\theta) dx}_{\mathbb{E} [(T(X) - \Theta)^2 | \Theta = \theta]} f_{\Theta}(\theta; \sigma) d\theta = \mathbb{E} [(T(X) - \Theta)^2]. \quad (6.3)$$

6.1 Conditional Mean Estimator – Minimizing the MSE

Theorem. The CONDITIONAL MEAN ESTIMATOR (BAYES ESTIMATOR) T_{CM} is MSE optimal,

$$T_{\text{CM}} : x \mapsto \mathbb{E} [\Theta | X = x] = \int_{\Theta} \theta f_{\Theta|X}(\theta|x) d\theta, \quad (6.4)$$

i.e., T_{CM} minimizes the MEAN MSE cost criterion

$$\mathbb{E} [\mathbb{E} [(T(X) - \Theta)^2 | \Theta]]. \quad (6.5)$$

Note. We distinguish

- $\mathbb{E} [(T(X) - \Theta)^2 | \Theta = \theta]$, which corresponds to the MSE subject to the condition of the (unknown) deterministic outcome θ of the random variable Θ ,
- and $\mathbb{E} [(T(X) - \Theta)^2] = \mathbb{E} [\mathbb{E} [(T(X) - \Theta)^2 | \Theta]]$, which assumes Θ to be a random variable, such that the expectation has to be taken over both random variables.

Proof.

$$\begin{aligned}
 \mathbb{E}[\varepsilon(T(X))] &= \int_{\Theta} \int_{\mathbb{X}} (T(x) - \theta)^2 \underbrace{f_{X|\Theta}(x|\theta) f_{\Theta}(\theta; \sigma)}_{=f_{X,\Theta}(x,\theta) \quad (\text{Bayes})} dx d\theta \\
 &= \int_{\mathbb{X}} \int_{\Theta} (T(x) - \theta)^2 f_{X,\Theta}(x, \theta) d\theta dx = \int_{\mathbb{X}} \underbrace{\int_{\Theta} (T(x) - \theta)^2 f_{\Theta|X}(\theta|x) d\theta}_{\text{min!}} f_X(x) dx
 \end{aligned}$$

$$\Rightarrow \boxed{\min_{T(x)} \int_{\Theta} (T(x) - \theta)^2 f_{\Theta|X}(\theta|x) d\theta,} \quad \text{i.e.,}$$

$$\Rightarrow \boxed{\forall x : T_x = T(x)}$$

$$\frac{\partial}{\partial T_x} \int_{\Theta} (T_x - \theta)^2 f_{\Theta|X}(\theta|x) d\theta \stackrel{!}{=} 0$$

$$2 \int_{\Theta} (T_x - \theta) f_{\Theta|X}(\theta|x) d\theta = 0$$

$$T_x \int_{\Theta} f_{\Theta|X}(\theta|x) d\theta = \int_{\Theta} \theta f_{\Theta|X}(\theta|x) d\theta, \text{ i.e.,}$$

$$\Rightarrow \boxed{T(x) = \mathbb{E}[\Theta|X = x].}$$

6.2 Bernoulli Experiment (Cont'd)

In Section 3.4, we found the MAXIMUM LIKELIHOOD ESTIMATE of the success probability of a BERNOULLI EXPERIMENT as

$$\hat{\theta}_{\text{ML}} = \operatorname{argmax}_{\theta \in \Theta} \{\log L(x; \theta)\} = \frac{x}{N}, \quad (6.6)$$

given the statistical model by means of the PROBABILITY MASS FUNCTION

$$B_{N,\theta}(x) = \binom{N}{x} \theta^x (1 - \theta)^{N-x}, \quad x \in \{0, 1, \dots, N\}.$$

We can expect a better estimator, if we design the estimator with respect to the MEAN SQUARED ERROR by exploiting the trivial side information

$$f_{\Theta}(\theta) = \begin{cases} 1 & ; \quad \theta \in [0, 1] \\ 0 & ; \quad \text{otherwise} \end{cases}.$$

The optimal estimator in terms of the MEAN SQUARED ERROR can be obtained by means of the CONDITIONAL MEAN ESTIMATOR

$$\hat{\theta}_{\text{CM}} = E[\Theta \mid X = x] = \int_0^1 \theta f_{\Theta|X}(\theta \mid x) d\theta. \quad (6.7)$$

The required CONDITIONAL PDF is

$$f_{\Theta|X}(\theta) = \frac{p_{X|\Theta}(x \mid \theta)}{p_X(x)} f_{\Theta}(\theta), \quad \theta \in [0, 1] \quad (6.8)$$

with¹

$$p_X(x) = \int_{\Theta} p_{X|\Theta}(x | \theta) f_{\Theta}(\theta) d\theta, \quad x \in \{0, 1, \dots, N\} \quad (6.9)$$

$$= \int_0^1 B_{N,\theta}(x) \cdot 1 d\theta = \int_0^1 \binom{N}{x} \theta^x (1 - \theta)^{N-x} d\theta \quad (6.10)$$

$$= \binom{N}{x} \frac{x!}{(N - x + 1)(N - x + 2) \cdots (N + 1)} = \frac{1}{N + 1}. \quad (6.11)$$

Then the CONDITIONAL MEAN ESTIMATOR equals (cf. alternative solution in Section 3.4)

$$\hat{\theta}_{\text{CM}} = p_X(x)^{-1} \int_0^1 \theta p_{X|\Theta}(x | \theta) d\theta \quad (6.12)$$

$$= p_X(x)^{-1} \int_0^1 \binom{N}{x} \theta^{x+1} (1 - \theta)^{N-x} d\theta \quad (6.13)$$

$$= p_X(x)^{-1} \binom{N}{x} \frac{(x + 1)!}{(N - x + 1)(N - x + 2) \cdots (N + 2)} = \frac{p_X(x)^{-1}}{N + 1} \frac{x + 1}{N + 2} \quad (6.14)$$

$$= \frac{x + 1}{N + 2} = \frac{N}{N + 2} \hat{\theta}_{\text{ML}} + \frac{1}{N + 2}. \quad (6.15)$$

¹ m -times applying the partial integration technique ($m, n \in \mathbb{N}$) we obtain

$$\int_0^1 x^m (1 - x)^n dx = \frac{m}{n + 1} \int_0^1 x^{m-1} (1 - x)^{n+1} dx = \cdots = \frac{m!}{(n + 1) \cdots (n + m)} \int_0^1 (1 - x)^{n+m} dx = \frac{m!}{(n + 1) \cdots (n + m + 1)}.$$

6.3 Mean Estimation Example

We consider the estimation of the unknown mean value θ of a random variable

$$X \sim \mathcal{N}(\theta, \sigma_{X|\Theta=\theta}^2), \quad (6.16)$$

now based on N i.i.d. observations

$$X_1, \dots, X_N, \quad (6.17)$$

drawn with respect to the CONDITIONAL PDF²

$$f_{X_1, \dots, X_N|\Theta}(x_1, \dots, x_N|\theta) = \frac{1}{(\sqrt{2\pi}\sigma_{X|\Theta=\theta})^N} \exp\left(-\frac{1}{2\sigma_{X|\Theta=\theta}^2} \sum_{i=1}^N (x_i - \theta)^2\right). \quad (6.18)$$

A priori knowledge about the unknown parameter θ is given by the PDF $\Theta \sim \mathcal{N}(m, \sigma_\Theta^2)$, i.e.,

$$\Theta \sim \frac{1}{\sqrt{2\pi}\sigma_\Theta} \exp\left(-\frac{1}{2\sigma_\Theta^2}(\theta - m)^2\right) = f_\Theta(\theta). \quad (6.19)$$

²The difference between statistics drawn from a CONDITIONAL DISTRIBUTION and statistics drawn from a UNCONDITIONAL DISTRIBUTION is essential for the understanding of the following results.

Note. Consider the difference between

$$\text{VARIANCE OF } \Theta : \sigma_{\Theta}^2, \quad (6.20)$$

$$\text{CONDITIONAL VARIANCE OF } X \text{ GIVEN } \theta : \sigma_{X|\Theta=\theta}^2, \quad (6.21)$$

$$\text{VARIANCE OF } X : \sigma_X^2, \quad (6.22)$$

with

$$\sigma_X^2 = \sigma_{X|\Theta=\theta}^2 + \sigma_{\Theta}^2. \quad (6.23)$$

Interpretation. Since X is the sum of independent random variables

$$X = \Theta + \eta, \quad (6.24)$$

with $\Theta \sim \mathcal{N}(m, \sigma_{\Theta}^2)$ and $\eta \sim \mathcal{N}(0, \sigma_{X|\Theta=\theta}^2)$, the variance of the sum is equal to the sum of variances and the expectation of the sum is equal to sum of expectations, i.e., $\mu_X = E[X] = E[\Theta + \eta] = E[\Theta] = m$. The covariance between Θ and X equals to the variance of Θ , i.e., $c_{\Theta,X} = \text{Cov}[\Theta, X] = \text{Cov}[\Theta, \Theta + \eta] = \text{Cov}[\Theta, \Theta] + \text{Cov}[\Theta, \eta] = \text{Cov}[\Theta, \Theta] = \text{Var}[\Theta]$. In essence,

$$\mu_X = \mu_{\Theta} \quad (6.25)$$

$$c_{\Theta,X} = \sigma_{\Theta}^2. \quad (6.26)$$

Note. Make sure to identify the difference between the random variable X given $\Theta = \theta$ which is drawn from the PDF $f_{X|\Theta}(x|\theta)$ and the random variable X which is drawn from the PDF $f_X(x)$.

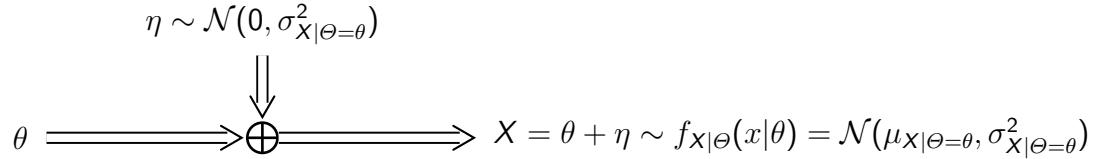


Fig. 6.1: Illustration of the random variable X given θ which is the sum of a specific **realization** θ and **random** additive noise η . Therefore, we have $\mu_{X|\Theta=\theta} = \theta$ and the **conditional variance** $\sigma_{X|\Theta=\theta}^2$.

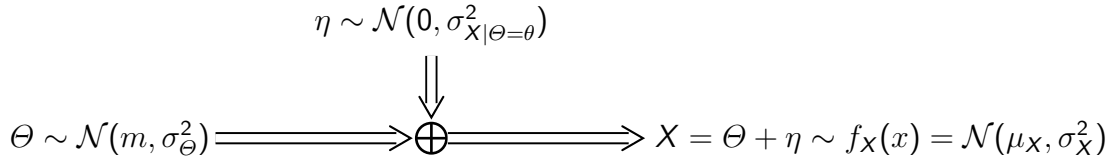


Fig. 6.2: Illustration of the random variable X which is the sum of the **random variable** Θ and **random** additive noise η . Therefore, we have $\mu_X = m$ and the **sum of variances** $\sigma_X^2 = \sigma_{X|\Theta=\theta}^2 + \sigma_{\Theta}^2$.

For the computation of the CONDITIONAL MEAN ESTIMATOR we need two steps:

a) Conditional PDF $f_{\Theta|X_1,\dots,X_N}(\theta|x_1,\dots,x_N)$:

$$f_{\Theta|X_1,\dots,X_N}(\theta|x_1,\dots,x_N) = \frac{f_{X_1,\dots,X_N|\Theta}(x_1,\dots,x_N|\theta)f_{\Theta}(\theta)}{f_{X_1,\dots,X_N}(x_1,\dots,x_N)} \quad (6.27)$$

$$= \frac{f_{X_1,\dots,X_N|\Theta}(x_1,\dots,x_N|\theta)f_{\Theta}(\theta)}{\underbrace{\int_{\Theta} f_{X_1,\dots,X_N|\Theta}(x_1,\dots,x_N|\theta)f_{\Theta}(\theta)d\theta}_{=f_{X_1,\dots,X_N}(x_1,\dots,x_N) \text{ (Marginalization)}}} \quad (6.28)$$

$$= \gamma \exp\left(-\frac{1}{2\sigma_{X|\Theta=\theta}^2} \sum_{i=1}^N (x_i - \theta)^2\right) \exp\left(-\frac{1}{2\sigma_{\Theta}^2}(\theta - m)^2\right), \quad (6.29)$$

with γ such that $\int_{\Theta} f_{\Theta|X_1,\dots,X_N}(\theta|x_1,\dots,x_N)d\theta = 1$.

b) Conditional Mean $E[\Theta|x_1, \dots, x_N]$.³

$$\hat{\theta}_{\text{CM}} = E[\Theta|x_1, \dots, x_N] \quad (6.30)$$

$$= \int_{\Theta} \theta f_{\Theta|X_1, \dots, X_N}(\theta|x_1, \dots, x_N) d\theta \quad \text{with CONDITIONAL PDF} \quad (6.31)$$

$$= \frac{\sum_{i=1}^N x_i \sigma_{\Theta}^2 + m \sigma_{X|\Theta=\theta}^2}{N \sigma_{\Theta}^2 + \sigma_{X|\Theta=\theta}^2},$$

which is LINEAR IN OBSERVATIONS x_1, \dots, x_N (!) and after some reformulation steps:

$$\hat{\theta}_{\text{CM}} = \frac{\sigma_{\Theta}^2}{\sigma_{\Theta}^2 + \frac{\sigma_{X|\Theta=\theta}^2}{N}} \hat{\theta}_{\text{ML}} + \frac{\frac{\sigma_{X|\Theta=\theta}^2}{N}}{\sigma_{\Theta}^2 + \frac{\sigma_{X|\Theta=\theta}^2}{N}} m, \quad \text{with} \quad \hat{\theta}_{\text{ML}} = \frac{1}{N} \sum_{i=1}^N x_i. \quad (6.32)$$

Note. For $N = 1$, we obtain $\hat{\theta}_{\text{CM}} = m + \frac{\sigma_{\Theta}^2}{\sigma_{\Theta}^2 + \sigma_{X|\Theta=\theta}^2} (x - m) = \mu_{\Theta} + \frac{c_{\Theta, X}}{\sigma_X^2} (x - \mu_X)$.

³Eq. (6.32) can be read out from Eq. (6.29) after some reformulation steps.

Derivation.

$$\begin{aligned}
f_{\theta|X_1, \dots, X_N}(\theta|x_1, \dots, x_N) &\propto \exp\left(-\frac{1}{2\sigma_{X|\theta=\theta}^2} \sum_{i=1}^N (x_i - \theta)^2\right) \exp\left(-\frac{1}{2\sigma_{\theta}^2}(\theta - m)^2\right) \\
&\propto \exp\left(-\frac{1}{2\sigma_{X|\theta=\theta}^2\sigma_{\theta}^2} \left(\sum_{i=1}^N (x_i^2 - 2x_i\theta + \theta^2) \sigma_{\theta}^2 + (\theta^2 - 2m\theta + m^2) \sigma_{X|\theta=\theta}^2\right)\right) \\
&\propto \exp\left(-\frac{1}{2\sigma_{X|\theta=\theta}^2\sigma_{\theta}^2} \left(-2\theta\sigma_{\theta}^2 \sum_{i=1}^N x_i + N\theta^2\sigma_{\theta}^2 + \theta^2\sigma_{X|\theta=\theta}^2 - 2m\theta\sigma_{X|\theta=\theta}^2 + \dots\right)\right) \\
&\propto \exp\left(-\frac{1}{2\sigma_{X|\theta=\theta}^2\sigma_{\theta}^2} \left(\theta^2 (N\sigma_{\theta}^2 + \sigma_{X|\theta=\theta}^2) - 2\theta \left(\sigma_{\theta}^2 \sum_{i=1}^N x_i + m\sigma_{X|\theta=\theta}^2\right) + \dots\right)\right) \\
&\propto \exp\left(-\frac{(N\sigma_{\theta}^2 + \sigma_{X|\theta=\theta}^2)}{2\sigma_{X|\theta=\theta}^2\sigma_{\theta}^2} \left(\theta - \frac{\sigma_{\theta}^2 \sum_{i=1}^N x_i + m\sigma_{X|\theta=\theta}^2}{(N\sigma_{\theta}^2 + \sigma_{X|\theta=\theta}^2)}\right)^2\right).
\end{aligned}$$

A comparison with the Gaussian standard model results into

$$\mu_{\theta|X_1, \dots, X_N} = \frac{\sigma_{\theta}^2 \sum_{i=1}^N x_i + m\sigma_{X|\theta=\theta}^2}{N\sigma_{\theta}^2 + \sigma_{X|\theta=\theta}^2} \quad \text{and} \quad \sigma_{\theta|X_1, \dots, X_N}^2 = \frac{\sigma_{X|\theta=\theta}^2\sigma_{\theta}^2}{N\sigma_{\theta}^2 + \sigma_{X|\theta=\theta}^2}. \quad (6.33)$$

Discussion.

- Given (a) a large number N of observations or (b) a small conditional variance $\sigma_{X|\Theta=\theta}^2$ of observations or (c) a large variance σ_Θ^2 of the a priori distribution of the unknown parameter θ , the derived solution obviously recommends to rely on the MAXIMUM-LIKELIHOOD ESTIMATOR, since

$$\lim_{N \rightarrow \infty} \hat{\theta}_{\text{CM}} = \lim_{\sigma_{X|\Theta=\theta}^2 \rightarrow 0} \hat{\theta}_{\text{CM}} = \lim_{\sigma_\Theta^2 \rightarrow \infty} \hat{\theta}_{\text{CM}} = \hat{\theta}_{\text{ML}}, \quad (6.34)$$

with

$$\hat{\theta}_{\text{ML}} = \frac{1}{N} \sum_{i=1}^N x_i. \quad (6.35)$$

- However, given (a) a small variance σ_Θ^2 of the a priori distribution of the unknown parameter θ or (b) a large conditional variance $\sigma_{X|\Theta=\theta}^2$ of observations, the derived solution suggests to ignore the given observations and instead to rely solely on the MEAN VALUE of Θ , since

$$\lim_{\sigma_\Theta^2 \rightarrow 0} \hat{\theta}_{\text{CM}} = \lim_{\sigma_{X|\Theta=\theta}^2 \rightarrow \infty} \hat{\theta}_{\text{CM}} = m. \quad (6.36)$$

- The respective MINIMUM MEAN SQUARE ERROR (MMSE) is given by⁴

$$\mathbb{E} [\mathbb{E} [(T_{\text{CM}} - \theta)^2 | \theta]] = \mathbb{E} [\mathbb{E} [(\mathbb{E} [\theta|X] - \theta)^2 | X]] \quad (6.37)$$

$$= \mathbb{E} [\mathbb{E} [\theta|X]^2 - 2 \mathbb{E} [\theta|X] \mathbb{E} [\theta|X] + \mathbb{E} [\theta^2|X]] \quad (6.38)$$

$$= \mathbb{E} [\mathbb{E} [\theta^2|X] - \mathbb{E} [\theta|X]^2] \quad (6.39)$$

$$= \mathbb{E} [\text{Var} [\theta|X]] \quad (6.40)$$

$$= \frac{\sigma_{X|\theta=\theta}^2 \sigma_{\theta}^2}{\underbrace{N \sigma_{\theta}^2 + \sigma_{X|\theta=\theta}^2}_{\text{refer to (6.33)}}} = \sigma_{\theta}^2 - \frac{(\sigma_{\theta}^2)^2}{\sigma_{\theta}^2 + \frac{\sigma_{X|\theta=\theta}^2}{N}}. \quad (6.41)$$

- Any requisites of the CM Estimator can be found from the mean vector and covariance matrix of the joint distribution of $(X, \theta)^T$, i.e.,

$$\mathbb{E} \left[\begin{bmatrix} X \\ \theta \end{bmatrix} \right] = \begin{bmatrix} \mu_X \\ \mu_{\theta} \end{bmatrix} = \begin{bmatrix} m \\ m \end{bmatrix}, \quad (6.42)$$

$$\text{Var} \left[\begin{bmatrix} X \\ \theta \end{bmatrix} \right] = \begin{bmatrix} c_{X,X} & c_{X,\theta} \\ c_{\theta,X} & c_{\theta,\theta} \end{bmatrix} = \begin{bmatrix} \sigma_{\theta}^2 + \sigma_{X|\theta=\theta}^2 & \sigma_{\theta}^2 \\ \sigma_{\theta}^2 & \sigma_{\theta}^2 \end{bmatrix}. \quad (6.43)$$

⁴Eq. (6.41) can be read out from Eq. (6.29) after some reformulation steps.

6.4 Jointly Gaussian Random Variables

Given two random vectors $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}_X, \mathbf{C}_X)$, $\boldsymbol{\Theta} \sim \mathcal{N}(\boldsymbol{\mu}_\Theta, \mathbf{C}_\Theta)$, and the covariance matrix \mathbf{C}_Z from the joint distribution of $\mathbf{Z} = (\mathbf{X}, \boldsymbol{\Theta})^\top \sim \mathcal{N}(\boldsymbol{\mu}_Z, \mathbf{C}_Z)$, i.e.,

$$\mathbf{C}_Z = \begin{bmatrix} \mathbf{C}_X & \mathbf{C}_{X,\Theta} \\ \mathbf{C}_{\Theta,X} & \mathbf{C}_\Theta \end{bmatrix}, \quad (6.44)$$

the MULTIVARIATE CONDITIONAL MEAN ESTIMATOR is obtained as

$$\mathbf{T}_{\text{CM}} : \quad x \mapsto \mathbb{E}[\boldsymbol{\Theta} | \mathbf{X} = x] = \boldsymbol{\mu}_\Theta + \mathbf{C}_{\Theta,X} \mathbf{C}_X^{-1} (x - \boldsymbol{\mu}_X). \quad (6.45)$$

The respective MMSE is equal to the TRACE of the conditional covariance matrix $\mathbf{C}_{\Theta|\mathbf{X}}$, this is

$$\mathbb{E} [\| \mathbf{T}_{\text{CM}} - \boldsymbol{\Theta} \|^2_2] = \text{tr} [\mathbf{C}_{\Theta|\mathbf{X}=x}] = \text{tr} [\mathbf{C}_\Theta - \mathbf{C}_{\Theta,X} \mathbf{C}_X^{-1} \mathbf{C}_{X,\Theta}]. \quad (6.46)$$

Note. Given JOINTLY GAUSSIAN RANDOM VARIABLES \mathbf{X} and \mathbf{Y} , the CONDITIONAL MEAN ESTIMATOR $\mathbb{E}[\mathbf{Y}|\mathbf{X}]$ is a LINEAR FUNCTION in \mathbf{X} . This does not hold for arbitrarily jointly distributed random variables!

6.5 Mean Estimation Example (Cont'd)

From a multivariate r.v. perspective, we obtain

$$\mathbf{X} = \begin{bmatrix} X_1 \\ \vdots \\ X_N \end{bmatrix} \sim \mathcal{N}(\mu_{\mathbf{X}}, \mathbf{C}_{\mathbf{X}}) \quad \text{and} \quad \Theta \sim \mathcal{N}(\mu_{\Theta}, \sigma_{\Theta}^2), \quad (6.47)$$

and thus

$$\mu_{\mathbf{X}} = m\mathbf{1}, \quad \mu_{\Theta} = m, \quad \mathbf{C}_{\Theta, \mathbf{X}} = \sigma_{\Theta}^2 \mathbf{1}\mathbf{1}^{\top}, \quad \mathbf{C}_{\mathbf{X}} = \sigma_{\Theta}^2 \mathbf{1}\mathbf{1}^{\top} + \sigma_{X|\Theta=\theta}^2 \mathbf{I}, \quad (6.48)$$

where $\mathbf{1} = [1, \dots, 1]^{\top}$ and \mathbf{I} is the unity matrix. With a closed form solution for the inverse $\mathbf{C}_{\mathbf{X}}^{-1}$, we obtain the CM Estimator

$$T_{\text{CM}} = m + \sigma_{\Theta}^2 \mathbf{1}^{\top} (\sigma_{\Theta}^2 \mathbf{1}\mathbf{1}^{\top} + \sigma_{X|\Theta=\theta}^2 \mathbf{I})^{-1} (\mathbf{x} - m\mathbf{1}) = \frac{\frac{\sum_{i=1}^N x_i}{\sigma_{X|\Theta=\theta}^2} + \frac{m}{\sigma_{\Theta}^2}}{\frac{N}{\sigma_{X|\Theta=\theta}^2} + \frac{1}{\sigma_{\Theta}^2}}. \quad (6.49)$$

Note. The computation of $C_{\mathbf{x}}^{-1}$ can be obtained in closed form by applying the SHERMAN-MORRISON FORMULA, a special version of the MATRIX-INVERSION LEMMA,⁵ i.e.,

$$C_{\mathbf{x}}^{-1} = (\sigma_{\Theta}^2 \mathbf{1}\mathbf{1}^T + \sigma_{X|\Theta=\theta}^2 \mathbf{I})^{-1} \quad (6.50)$$

$$= \frac{1}{\sigma_{X|\Theta=\theta}^2} \left[\mathbf{I} - \mathbf{1}\mathbf{1}^T \frac{\frac{1}{\sigma_{X|\Theta=\theta}^2}}{\frac{1}{\sigma_{\Theta}^2} + \frac{1}{\sigma_{X|\Theta=\theta}^2}} \right], \quad (6.51)$$

where

$$\mathbf{1}\mathbf{1}^T = \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & \cdots & 1 \end{bmatrix} = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ 1 & 1 & \cdots & 1 \\ \vdots & \vdots & & \vdots \\ 1 & 1 & \cdots & 1 \end{bmatrix}.$$

⁵Given the compatibility of all matrices and the existence of \mathbf{A}^{-1} , the following equality holds: $(\mathbf{A} + \mathbf{b}\mathbf{c}^T)^{-1} = \mathbf{A}^{-1} - \frac{\mathbf{A}^{-1}\mathbf{b}\mathbf{c}^T\mathbf{A}^{-1}}{(1 + \mathbf{c}^T\mathbf{A}^{-1}\mathbf{b})}$.

Note. The random variables X_1, \dots, X_N given θ are drawn from the conditional joint PDF with respect to the **same realization** θ of Θ , i.e., we have that X_1, \dots, X_N are drawn from the conditional joint PDF $f_{X_1, \dots, X_N | \Theta}(x_1, \dots, x_N | \theta)$.

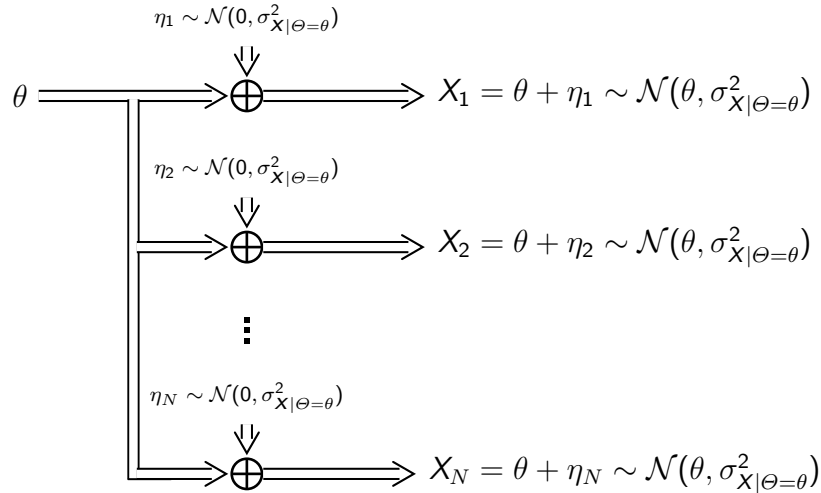


Fig. 6.3: Illustration of the random variables X_1, \dots, X_N given θ which each are the sum of the **same realization** θ and **independent random** additive noise η_i , i.e., with $E[\eta_i \eta_j] = 0$ for all $i \neq j$.

Note. Consequently, the random variables X_1, \dots, X_N and the random variable Θ are drawn from the joint PDF $f_{X_1, \dots, X_N, \Theta}(x_1, \dots, x_N, \theta)$.

6.6 Orthogonality Principle

The ORTHOGONALITY PRINCIPLE is an inherent property of the CONDITIONAL MEAN ESTIMATOR. It describes the inherent STOCHASTIC ORTHOGONALITY between the CM estimation error and any observations statistics or functions thereof, i.e.,

$$E[(T_{\text{CM}}(X_1, \dots, X_N) - \Theta)h(X_1, \dots, X_N)] = 0, \quad (6.52)$$

$$\Updownarrow$$

$$T_{\text{CM}}(X_1, \dots, X_N) - \Theta \perp h(X_1, \dots, X_N) \quad (6.53)$$

where $h : \mathbb{R}^N \rightarrow \mathbb{R}, x_1, \dots, x_N \mapsto h(x_1, \dots, x_N)$.

Proof.

$$E[(T_{\text{CM}}(X) - \Theta)h(X)] = E[(T_{\text{CM}}(X)h(X)] - E[\Theta h(X)] = E[E[\Theta|X]h(X)] - E[\Theta h(X)] \quad (6.54)$$

$$= E[E[\Theta h(X)|X]] - E[\Theta h(X)] = E[\Theta h(X)] - E[\Theta h(X)] = 0. \quad (6.55)$$

6.7 Mean Estimation Example (Cont'd)

The parameter θ to be estimated and the corresponding observation statistic X are JOINTLY GAUSSIAN DISTRIBUTED according to

$$\begin{bmatrix} X \\ \Theta \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mu_X \\ \mu_\Theta \end{bmatrix}, \begin{bmatrix} c_{X,X} & c_{X,\Theta} \\ c_{\Theta,X} & c_{\Theta,\Theta} \end{bmatrix} \right) = \mathcal{N} \left(\begin{bmatrix} m \\ m \end{bmatrix}, \begin{bmatrix} \sigma_\Theta^2 + \sigma_{X|\Theta=\theta}^2 & \sigma_\Theta^2 \\ \sigma_\Theta^2 & \sigma_\Theta^2 \end{bmatrix} \right). \quad (6.56)$$

As an SUFFICIENT STATISTIC⁶ for estimating the unknown realization of the random variable Θ we use⁷

$$S = \frac{1}{N} \sum_{i=1}^N X_i. \quad (6.57)$$

Applying the ORTHOGONALITY PRINCIPLE we obtain

$$\mathbb{E}[(T_{\text{CM}} - \Theta)S] = \mathbb{E}[T_{\text{CM}}S] - \mathbb{E}[\Theta S] = 0. \quad (6.58)$$

Knowing that the MSE optimal estimator is linear in S , we substitute T_{CM} by a linear model

$$T_{\text{CM}}(X) = aS + b. \quad (6.59)$$

For $h(X_1, \dots, X_N) = S$, the ORTHOGONALITY PRINCIPLE is now equal to

$$\mathbb{E}[(T_{\text{CM}} - \Theta)S] = \mathbb{E}[(aS + b - \Theta)S] = 0. \quad (6.60)$$

⁶The concept of SUFFICIENT STATISTIC is later introduced in this course.

⁷Note, that the realizations of the statistics X_1, \dots, X_N are conditioned on the unknown realization of Θ .

The ORTHOGONALITY PRINCIPLE leaves one DEGREE OF FREEDOM which is satisfied by a second orthogonality proposition with $h(X_1, \dots, X_N) = 1$,

$$\mathbb{E}[(aS + b - \Theta)1] = \mathbb{E}[aS + b] - \mathbb{E}[\Theta] = 0. \quad (6.61)$$

Then we obtain

$$a = \frac{c_{\Theta, S}}{\sigma_S^2} \quad \text{and} \quad b = \mu_\Theta - a\mu_S, \quad (6.62)$$

and the CONDITIONAL MEAN ESTIMATOR is given by⁸

$$T_{\text{CM}} : \quad s \mapsto \mu_\Theta + \frac{c_{\Theta, S}}{\sigma_S^2}(s - \mu_S) = m + \frac{\sigma_\Theta^2}{\frac{\sigma_{X|\Theta=\theta}^2}{N} + \sigma_\Theta^2}(s - m) = \frac{\frac{\sum_{i=1}^N x_i}{N} + \frac{m}{\sigma_\Theta^2}}{\frac{\sigma_{X|\Theta=\theta}^2}{N} + \frac{1}{\sigma_\Theta^2}}. \quad (6.63)$$

⁸Since the realizations of the statistics X_1, \dots, X_N are conditioned on the unknown realization of Θ , we obtain $c_{\Theta, S} = \mathbb{E}[\Theta S] - \mathbb{E}[\Theta] \mathbb{E}[S] = \mathbb{E}[\mathbb{E}[\Theta S | \Theta]] - \mathbb{E}[\Theta] \mathbb{E}[\mathbb{E}[S | \Theta]] = \mathbb{E}[\Theta^2] - (\mathbb{E}[\Theta])^2 = \sigma_\Theta^2$ and $\sigma_S^2 = \mathbb{E}[S^2] - (\mathbb{E}[S])^2 = \mathbb{E}[\mathbb{E}[S^2 | \Theta]] - (\mathbb{E}[\mathbb{E}[S | \Theta]])^2 = \mathbb{E}[\text{Var}[S | \Theta] + \mathbb{E}[S | \Theta]^2] - (\mathbb{E}[\Theta])^2 = \mathbb{E}\left[\frac{\sigma_{X|\Theta=\theta}^2}{N} + \Theta^2\right] - m^2 = \frac{\sigma_{X|\Theta=\theta}^2}{N} + \sigma_\Theta^2.$

Part IV

Linear Estimation

7. Linear Estimation

In this chapter, we directly focus on linear models, i.e., given an observation $\mathbf{x} \in \mathbb{R}^d$ corresponding to the respective random variable, we consider the estimation of the realization of the random variable y by

$$T_{\text{Lin}} : \mathbb{R}^d \rightarrow \mathbb{R}, \quad \mathbf{x} \mapsto \hat{y}, \quad \hat{y} = \mathbf{x}^\top \mathbf{t} + m. \quad (7.1)$$

In contrast to the MAXIMUM-LIKELIHOOD ESTIMATORS and the CONDITIONAL MEAN ESTIMATORS, where the inference of the estimator is based on statistical parameters of a given PROBABILITY DISTRIBUTION FUNCTION, here the estimation of y , by inference of \mathbf{t} and m , is based on N pairs of jointly drawn observations (\mathbf{x}_i, y_i) , where

$$\mathbf{y} = \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix} \in \mathbb{R}^N \quad \text{and} \quad \mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_N^\top \end{bmatrix} \in \mathbb{R}^{N \times d}. \quad (7.2)$$

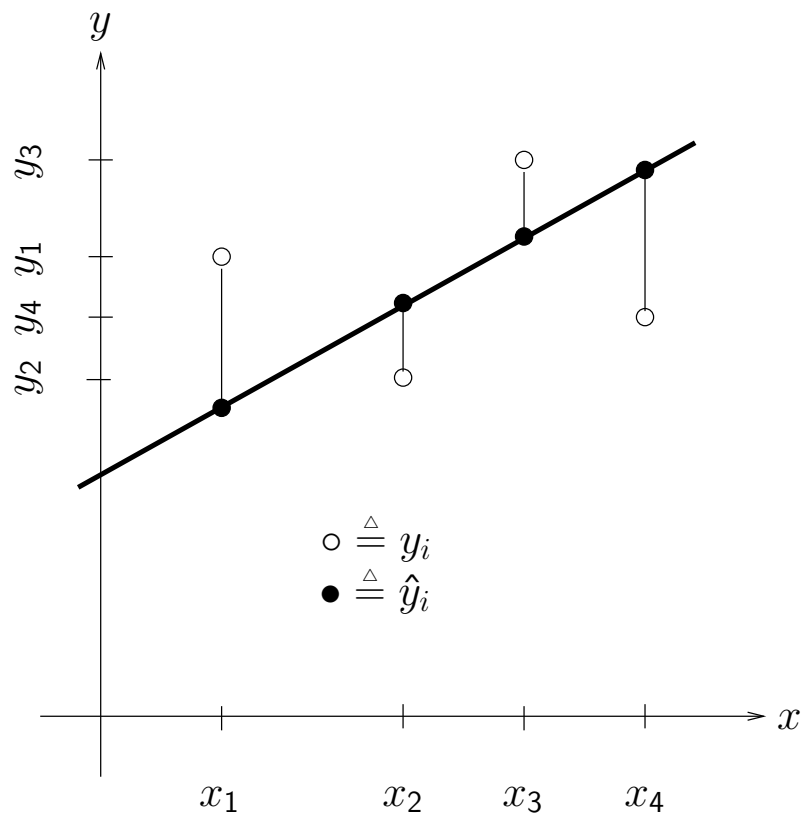


Fig. 7.1: Principle of linear estimation: $\hat{y}_{\text{Lin}} = xt_1 + t_2$

7.1 Least Squares Estimation

In LEAST SQUARES ESTIMATION (LS), the inference of \mathbf{t} is based on the minimization of the sum of squared errors between N pairs of observations y_i and outcomes \hat{y}_i of the LINEAR MODEL¹ $\hat{y} = \mathbf{x}^\top \mathbf{t}$ with \mathbf{x} and $\mathbf{t} \in \mathbb{R}^d$.

The related OPTIMIZATION PROBLEM is equal to

$$\min_{\mathbf{t} \in \mathbb{R}^d} \left\{ \sum_{i=1}^N (y_i - \mathbf{x}_i^\top \mathbf{t})^2 \right\} \Leftrightarrow \min_{\mathbf{t} \in \mathbb{R}^d} \left\| \begin{bmatrix} y_1 - \mathbf{x}_1^\top \mathbf{t} \\ \vdots \\ y_N - \mathbf{x}_N^\top \mathbf{t} \end{bmatrix} \right\|_2^2 \Leftrightarrow \min_{\mathbf{t} \in \mathbb{R}^d} \|\mathbf{y} - \mathbf{X}\mathbf{t}\|_2^2. \quad (7.3)$$

From a SUBSPACE PERSPECTIVE, we search for a vector $\hat{\mathbf{y}} = \mathbf{X}\mathbf{t}_{\text{LS}} \in \text{range}[\mathbf{X}] \in \mathbb{R}^N$, which is the BEST APPROXIMATION of $\mathbf{y} \in \mathbb{R}^N$, i.e., $\min_{\hat{\mathbf{y}} \in \text{range}[\mathbf{X}]} \|\mathbf{y} - \hat{\mathbf{y}}\|^2$ with

$$\mathbf{y} - \hat{\mathbf{y}} \perp \text{range}[\mathbf{X}] \Leftrightarrow \mathbf{y} - \hat{\mathbf{y}} \in \text{null}[\mathbf{X}^\top] \Leftrightarrow \mathbf{X}^\top(\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{0}. \quad (7.4)$$

The resulting vector solves $\mathbf{X}^\text{H}\mathbf{y} - \mathbf{X}^\top\mathbf{X}\mathbf{t}_{\text{LS}} = \mathbf{0}$ and thus (under favorable conditions)

$$\mathbf{t}_{\text{LS}} = (\mathbf{X}^\top\mathbf{X})^{-1} \mathbf{X}^\top\mathbf{y}. \quad (7.5)$$

¹The affine case $y = \mathbf{x}^\top \mathbf{t} + t_0$ can similarly be treated by introducing $y = \mathbf{x}'^\top \mathbf{t}'$, with $\mathbf{t}' = \begin{bmatrix} \mathbf{t} \\ t_0 \end{bmatrix}$ and $\mathbf{x}' = \begin{bmatrix} \mathbf{x} \\ 1 \end{bmatrix}$ and therefore $d' = d + 1$.

In other words, the BEST APPROXIMATION $\hat{\mathbf{y}}$ is found by the ORTHOGONAL PROJECTION onto $\text{range}[\mathbf{X}]$.

The appropriate LEFT INVERSE of $\mathbf{X} \in \mathbb{R}^{N \times d}$ in order to determine the ORTHOGONAL PROJECTOR on $\text{range}[\mathbf{X}]$ is given by the PSEUDO-INVERSE

$$\mathbf{X}^+ = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \quad (\text{LEFT INVERSE of } \mathbf{X}), \quad (7.6)$$

i.e.,

$$\hat{\mathbf{y}} = \mathbf{X} \mathbf{X}^+ \mathbf{y} = \underbrace{\mathbf{X} (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top}_{\text{ORTHOGONAL PROJECTOR}} \mathbf{y}. \quad (7.7)$$

Taking into account the estimation model $\hat{\mathbf{y}} = \mathbf{x}^\top \mathbf{t}$ and $\hat{\mathbf{y}} = \mathbf{X} \mathbf{t}$, we consequently again obtain

$$\mathbf{t}_{\text{LS}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}. \quad (7.8)$$

Note. In order to obtain a feasible solution, FAVORABLE CONDITIONS are required. To this end, the existence of the INVERSE MATRIX of $\mathbf{X}^\top \mathbf{X}$ is required, i.e., $N \geq d$ linear independent measurement vectors \mathbf{x}_i must be available.

7.2 SVD Perspective

By means of the SINGULAR VALUE DECOMPOSITION² of the OBSERVATION MATRIX $\mathbf{X} \in \mathbb{R}^{N \times d}$ with $\text{rank}[\mathbf{X}] = d$, we obtain

$$\|\mathbf{y} - \mathbf{X}\mathbf{t}\|_2^2, \quad \text{with} \quad \mathbf{X} = [\mathbf{U}, \mathbf{U}^\perp] \begin{bmatrix} \boldsymbol{\Sigma} \\ \mathbf{0} \end{bmatrix} \mathbf{V}^\top \quad (7.9)$$

where $\mathbf{U} \in \mathbb{R}^{N \times d}$, $\mathbf{U}^\perp \in \mathbb{R}^{N \times (N-d)}$, $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$, and $\mathbf{V} \in \mathbb{R}^{d \times d}$. After some reformulations,³ the cost function is equal to

$$\begin{aligned} \|\mathbf{y} - \mathbf{X}\mathbf{t}\|_2^2 &= \left\| \mathbf{y} - [\mathbf{U}, \mathbf{U}^\perp] \begin{bmatrix} \boldsymbol{\Sigma} \\ \mathbf{0} \end{bmatrix} \mathbf{V}^\top \mathbf{t} \right\|_2^2 = \left\| [\mathbf{U}, \mathbf{U}^\perp] \left[\begin{bmatrix} \mathbf{U}^\top \\ \mathbf{U}^{\perp, \top} \end{bmatrix} \mathbf{y} - \begin{bmatrix} \boldsymbol{\Sigma} \\ \mathbf{0} \end{bmatrix} \mathbf{V}^\top \mathbf{t} \right] \right\|_2^2 \\ &= \left\| \begin{bmatrix} \mathbf{U}^\top \\ \mathbf{U}^{\perp, \top} \end{bmatrix} \mathbf{y} - \begin{bmatrix} \boldsymbol{\Sigma} \\ \mathbf{0} \end{bmatrix} \mathbf{V}^\top \mathbf{t} \right\|_2^2 = \left\| \begin{bmatrix} \mathbf{U}^\top \mathbf{y} - \boldsymbol{\Sigma} \mathbf{V}^\top \mathbf{t} \\ \mathbf{U}^{\perp, \top} \mathbf{y} \end{bmatrix} \right\|_2^2. \end{aligned}$$

Since the lower part of the vector does not depend on \mathbf{t} , the inequality $\|\mathbf{y} - \mathbf{X}\mathbf{t}\|_2^2 \geq \|\mathbf{U}^{\perp, \top} \mathbf{y}\|_2^2$ holds and the minimum value is achieved if

$$\mathbf{U}^\top \mathbf{y} - \boldsymbol{\Sigma} \mathbf{V}^\top \mathbf{t}_{\text{LS}} = \mathbf{0}_d \quad \Leftrightarrow \quad \mathbf{t}_{\text{LS}} = \mathbf{V} \boldsymbol{\Sigma}^{-1} \mathbf{U}^\top \mathbf{y} \in \mathbb{R}^d. \quad (7.10)$$

Note: The residual error $\|\mathbf{U}^{\perp, \top} \mathbf{y}\|_2^2$ corresponds to the orthogonal projection of \mathbf{y} onto $\text{span}[\mathbf{U}^\perp]$.

²For real-valued matrices the singular vectors can be chosen real-valued, too.

³If \mathbf{Q} is unitary, $\|\mathbf{Q}\mathbf{x}\|_2^2 = \|\mathbf{x}\|_2^2 = (\mathbf{Q}\mathbf{x})^\top (\mathbf{Q}\mathbf{x}) = \mathbf{x} \mathbf{Q}^\top \mathbf{Q} \mathbf{x} = \|\mathbf{x}\|_2^2$.

7.3 Examples

Mean Estimation

In order to estimate the mean μ_X from an unknown distribution $f_X(x)$ based on N observations x_i by means of LEAST SQUARES ESTIMATION, we introduce the linear model

$$T_{\text{Lin}} : \mathbb{R} \rightarrow \mathbb{R}, \quad 1 \mapsto \hat{x}, \quad \hat{x}_i = t1, \quad \text{with} \quad t = \hat{\mu}_X, \quad (7.11)$$

and thus the LS optimization problem

$$\min_t \left\{ \sum_{i=1}^N (x_i - 1t)^2 \right\} \Leftrightarrow \min_t \left\{ \|\mathbf{x} - \mathbf{1}t\|^2 \right\}, \quad (7.12)$$

with

$$\mathbf{x} = \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} \quad \text{and} \quad \mathbf{1} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}. \quad (7.13)$$

The resulting LS Estimator is determined by (7.5) and reads

$$\hat{\mu}_X = t_{\text{LS}} = (\mathbf{1}^\top \mathbf{1})^{-1} \mathbf{1}^\top \mathbf{x} = \frac{\mathbf{1}^\top \mathbf{x}}{\mathbf{1}^\top \mathbf{1}} = \frac{\sum_{i=1}^N x_i}{\mathbf{1}^\top \mathbf{1}} = \frac{1}{N} \sum_{i=1}^N x_i. \quad (7.14)$$

LS Estimation (scalar and linear)

Given a sample set $\{(x_1, y_1), \dots, (x_N, y_N)\}$, i.e., the vector of observations $x_i \in \mathbb{R}$ and the vector of elements $y_i \in \mathbb{R}$,

$$\mathbf{x} \triangleq \begin{bmatrix} x_1 \\ \vdots \\ x_N \end{bmatrix} \quad \text{and} \quad \mathbf{y} \triangleq \begin{bmatrix} y_1 \\ \vdots \\ y_N \end{bmatrix},$$

a linear estimator is denoted by

$$T_{\text{Lin}} : \mathbb{R} \rightarrow \mathbb{R}, \quad x \mapsto \hat{y}, \quad \hat{y} = xt. \quad (7.15)$$

Minimizing the EUKLIDEAN NORM of $\mathbf{y} - \hat{\mathbf{y}}$ with $\hat{\mathbf{y}} = \mathbf{x}t$ by applying the BEST APPROXIMATION PRINCIPLE yields

$$\mathbf{y} - \mathbf{x}t_{\text{LS}} \perp \text{range}[\mathbf{x}] \quad \Leftrightarrow \quad \mathbf{y} - \mathbf{x}t_{\text{LS}} \in \text{null}[\mathbf{x}^{\text{T}}],$$

and thus

$$\mathbf{x}^{\text{T}}(\mathbf{y} - \mathbf{x}t_{\text{LS}}) = 0. \quad (7.16)$$

It follows for the LEAST SQUARES optimal weight:

$$t_{\text{LS}} = \frac{\mathbf{x}^{\text{T}}\mathbf{y}}{\mathbf{x}^{\text{T}}\mathbf{x}}. \quad (7.17)$$

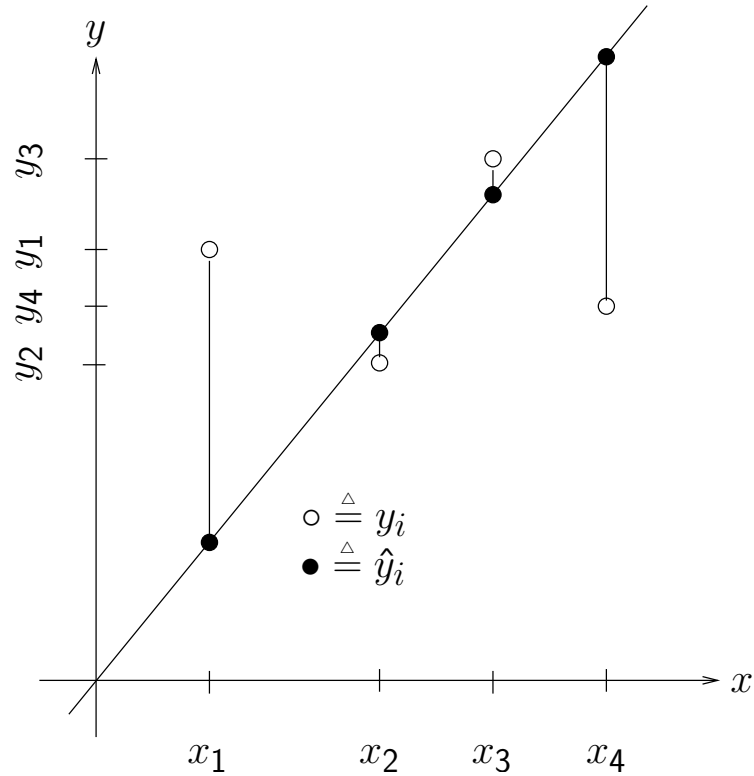


Fig. 7.2: Principle of scalar linear estimation: $\hat{y}_{\text{LS}} = xt_{\text{LS}}$

LS Estimation (scalar and affine)

Assuming the linear estimator as

$$T_{\text{Lin}} : \mathbb{R} \rightarrow \mathbb{R}, \quad x \mapsto \hat{y}, \quad \hat{y} = xt_1 + t_2 = [x \ 1] \begin{bmatrix} t_1 \\ t_2 \end{bmatrix} = \mathbf{x}^\top \mathbf{t}, \quad (7.18)$$

the vector of observations x_i becomes a matrix of observation vectors $\mathbf{x}_i^\top = [x_i \ 1]$, i.e.,

$$\mathbf{X} = \begin{bmatrix} \mathbf{x}_1^\top \\ \mathbf{x}_2^\top \\ \vdots \\ \mathbf{x}_N^\top \end{bmatrix} = \begin{bmatrix} [x_1 \ 1] \\ [x_2 \ 1] \\ \vdots \\ [x_N \ 1] \end{bmatrix}.$$

Since $\hat{\mathbf{y}} = \mathbf{X}\mathbf{t}$, or in other words $\hat{\mathbf{y}} \in \text{range}[\mathbf{X}]$, the BEST APPROXIMATION PRINCIPLE yields:

$$\mathbf{y} - \mathbf{X}\mathbf{t}_{\text{LS}} \perp \text{range}[\mathbf{X}] \quad \Leftrightarrow \quad \mathbf{y} - \mathbf{X}\mathbf{t}_{\text{LS}} \in \text{null}[\mathbf{X}^\top],$$

and

$$\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\mathbf{t}_{\text{LS}}) = \mathbf{0}. \quad (7.19)$$

Given $N \geq 2$ and at least 2 observation vectors form a linear independent basis of \mathbb{R}^2 ,

$$\mathbf{t}_{\text{LS}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}. \quad (7.20)$$

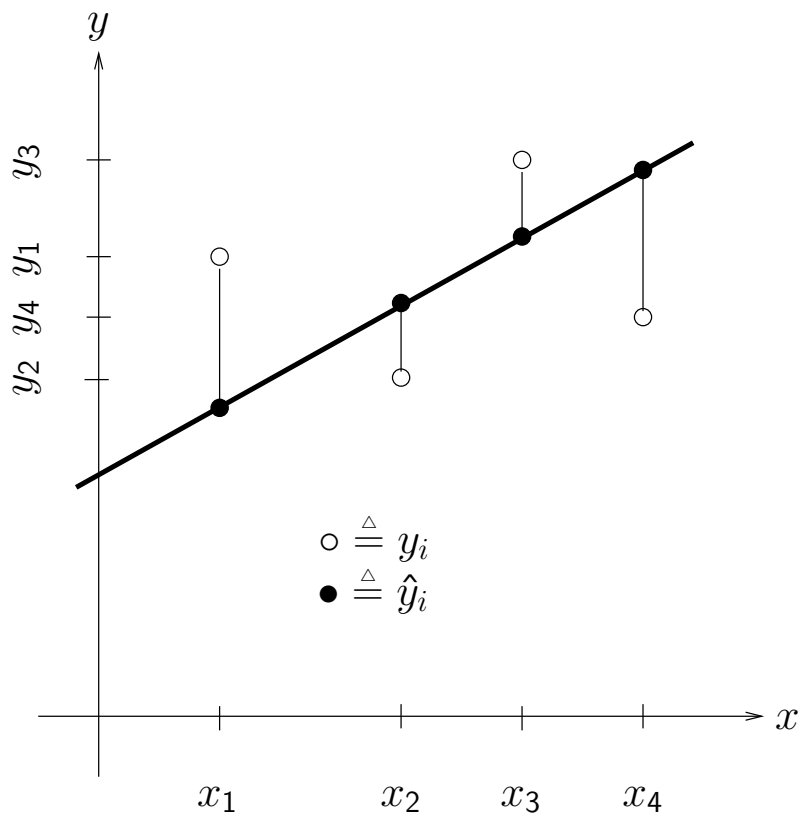


Fig. 7.3: Principle of scalar affine estimation: $\hat{y}_{\text{LS}} = xt_{\text{LS},1} + t_{\text{LS},2}$

LS Estimation (multi-dimensional and linear)

We now consider the most general case with $\{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)\}$, i.e., the matrix of observation vectors \mathbf{x}_i^\top and \mathbf{y}_i^\top ,

$$\mathbf{X} \triangleq \begin{bmatrix} \mathbf{x}_1^\top \\ \vdots \\ \mathbf{x}_N^\top \end{bmatrix} \quad \text{and} \quad \mathbf{Y} \triangleq \begin{bmatrix} \mathbf{y}_1^\top \\ \vdots \\ \mathbf{y}_N^\top \end{bmatrix}. \quad (7.21)$$

The optimal linear estimator with respect to the LS criterion, i.e.,

$$\hat{\mathbf{y}}_{\text{LS}}^\top = \mathbf{x}^\top \mathbf{T}_{\text{LS}}, \quad (7.22)$$

is found by applying the BEST APPROXIMATION PRINCIPLE columnwise:

$$\text{columnwise : } \mathbf{Y} - \mathbf{X}\mathbf{T}_{\text{LS}} \perp \text{range}[\mathbf{X}] \quad \Leftrightarrow \quad \mathbf{Y} - \mathbf{X}\mathbf{T}_{\text{LS}} \in \text{null}[\mathbf{X}^\top],$$

and

$$\mathbf{X}^\top (\mathbf{Y} - \mathbf{X}\mathbf{T}_{\text{LS}}) = \mathbf{0}. \quad (7.23)$$

If N is larger than the dimension of observation vectors \mathbf{x}_i , i.e., if $N > d$ and the columns of \mathbf{X} are linear independent such that $\text{rank}[\mathbf{X}] = d$, it follows that

$$\mathbf{T}_{\text{LS}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}. \quad (7.24)$$

7.4 Linear Minimum Mean Square Error Estimation

The LINEAR MINIMUM MEAN SQUARE ERROR (LMMSE) estimator is the estimator which minimizes the MSE based on a LINEAR MODEL for the estimator,

$$T_{\text{Lin}} : \mathbb{R}^d \rightarrow \mathbb{R}, \mathbf{x} \mapsto \hat{y}, \quad \hat{y} = \mathbf{x}^T \mathbf{t} + m \quad \text{or} \quad \hat{y} = \mathbf{t}^T \mathbf{x} + m, \quad (7.25)$$

i.e., the LMMSE Estimator is minimizer of the optimization problem⁴ $\min_{\mathbf{t}, m} \mathbb{E} \left[\|y - \mathbf{t}^T \mathbf{x} - m\|^2 \right]$.

Given the joint PDF of the random variables $\mathbf{z} = (x, y)^T$, with mean values and covariances

$$\boldsymbol{\mu}_{\mathbf{z}} = \begin{bmatrix} \mu_x \\ \mu_y \end{bmatrix} \quad \text{and} \quad \mathbf{C}_{\mathbf{z}} = \begin{bmatrix} \mathbf{C}_x & \mathbf{c}_{x,y} \\ \mathbf{c}_{y,x} & c_y \end{bmatrix}, \quad (7.26)$$

the LMMSE Estimator is given by

$$\hat{y} = \mu_y + \mathbf{c}_{y,x} \mathbf{C}_x^{-1} (\mathbf{x} - \boldsymbol{\mu}_x) = \underbrace{\mathbf{c}_{y,x} \mathbf{C}_x^{-1}}_{=\mathbf{t}^T} \mathbf{x} + \underbrace{\mu_y - \mathbf{c}_{y,x} \mathbf{C}_x^{-1} \boldsymbol{\mu}_x}_{=m}. \quad (7.27)$$

⁴In contrast to the standard notation, we denote random vectors by \mathbf{x} in order to avoid any confusions with matrix notation \mathbf{X} .

Proof.

Applying the ORTHOGONALITY PRINCIPLE

$$\begin{aligned} \mathbb{E}[(y - \mathbf{t}^\top \mathbf{x} - m)\mathbf{x}^\top] &= \mathbf{c}_{y,\mathbf{x}} + \mu_y \boldsymbol{\mu}_{\mathbf{x}}^\top - \mathbf{t}^\top (\mathbf{C}_{\mathbf{x}} + \boldsymbol{\mu}_{\mathbf{x}} \boldsymbol{\mu}_{\mathbf{x}}^\top) - m \boldsymbol{\mu}_{\mathbf{x}}^\top = 0 \\ \mathbb{E}[(y - \mathbf{t}^\top \mathbf{x} - m)1] &= \mu_y - \mathbf{t}^\top \boldsymbol{\mu}_{\mathbf{x}} - m = 0, \end{aligned}$$

we obtain $m = \mu_y - \mathbf{t}^\top \boldsymbol{\mu}_{\mathbf{x}}$ and

$$\begin{aligned} \mathbf{c}_{y,\mathbf{x}} + \mu_y \boldsymbol{\mu}_{\mathbf{x}}^\top - \mathbf{t}^\top (\mathbf{C}_{\mathbf{x}} + \boldsymbol{\mu}_{\mathbf{x}} \boldsymbol{\mu}_{\mathbf{x}}^\top) - (\mu_y - \mathbf{t}^\top \boldsymbol{\mu}_{\mathbf{x}}) \boldsymbol{\mu}_{\mathbf{x}}^\top &= 0 \\ \mathbf{c}_{y,\mathbf{x}} - \mathbf{t}^\top \mathbf{C}_{\mathbf{x}} &= 0, \end{aligned}$$

and thus $\mathbf{t}^\top = \mathbf{c}_{y,\mathbf{x}} \mathbf{C}_{\mathbf{x}}^{-1}$. The resulting estimator reads

$$\begin{aligned} \hat{y} &= \mathbf{t}^\top \mathbf{x} + \mu_y - \mathbf{t}^\top \boldsymbol{\mu}_{\mathbf{x}} \\ &= \mu_y + \mathbf{t}^\top (\mathbf{x} - \boldsymbol{\mu}_{\mathbf{x}}) \\ &= \mu_y + \mathbf{c}_{y,\mathbf{x}} \mathbf{C}_{\mathbf{x}}^{-1} (\mathbf{x} - \boldsymbol{\mu}_{\mathbf{x}}). \end{aligned}$$

Again applying the ORTHOGONALITY PRINCIPLE, the achievable MMSE is equal to

$$\begin{aligned} \mathbb{E}[\|y - \mathbf{t}^\top \mathbf{x} - m\|^2] &= \mathbb{E}[(y - \mathbf{t}^\top \mathbf{x} - m)y] \\ &= c_y + \mu_y^2 - \mathbf{t}^\top (\mathbf{c}_{\mathbf{x},y} + \boldsymbol{\mu}_{\mathbf{x}} \mu_y) - m \mu_y \\ &= c_y + \mu_y^2 - \mathbf{c}_{y,\mathbf{x}} \mathbf{C}_{\mathbf{x}}^{-1} \mathbf{c}_{\mathbf{x},y} + (m - \mu_y) \mu_y - m \mu_y \\ &= c_y - \mathbf{c}_{y,\mathbf{x}} \mathbf{C}_{\mathbf{x}}^{-1} \mathbf{c}_{\mathbf{x},y}. \end{aligned} \tag{7.28}$$

Special Case.

In the case of ZERO-MEAN random variables,

$$\mathbb{E} \left[\begin{bmatrix} \mathbf{x} \\ y \end{bmatrix} \right] = \begin{bmatrix} \boldsymbol{\mu}_x \\ \mu_y \end{bmatrix} = \mathbf{0}, \quad (7.29)$$

the LMMSE Estimator and its minimum MSE is

$$\text{LMMSE ESTIMATOR: } \hat{y} = \mathbf{c}_{y,\mathbf{x}} \mathbf{C}_x^{-1} \mathbf{x}, \quad (7.30)$$

$$\text{MINIMUM MSE: } \mathbb{E}[c_{y,\mathbf{x}}] = c_y - \mathbf{t}^\top \mathbf{c}_{\mathbf{x},y} = c_y - \mathbf{c}_{y,\mathbf{x}} \mathbf{C}_x^{-1} \mathbf{c}_{\mathbf{x},y}. \quad (7.31)$$

Part V

Examples

8. Estimation of a Matrix Channel

The estimation of a MULTIPLE-INPUT MULTIPLE OUTPUT (MIMO) channel is considered. In particular we assume a MIMO channel with K antenna elements at the transmitter and M antenna elements at the receiver, which means KM transmission channels to be estimated.

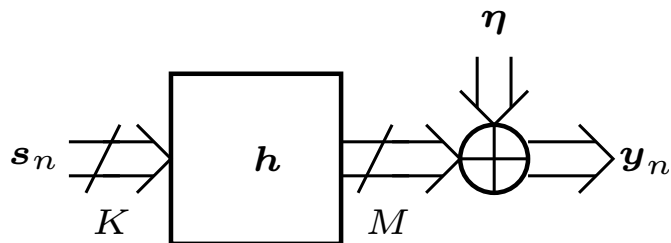


Fig. 8.1: Estimation of a MIMO channel.

Three LINEAR ESTIMATORS,

$$\hat{\mathbf{h}} = \mathbf{T}\mathbf{y} \in \mathbb{C}^{KM}, \quad (8.1)$$

are introduced and compared:

- The MINIMUM MEAN SQUARE ERROR ESTIMATOR (MMSE),
- the MAXIMUM LIKELIHOOD ESTIMATOR (ML),
- and the "MATCHED FILTER" ESTIMATOR (MF).

8.1 Model for Training Channel

The addressed model allows to cover a variety of scenarios, including TIME-VARIANT and DISPERSIVE channels. For the sake of simplicity, we assume the simplest case of an TIME-INVARIANT NON-DISPERSIVE MIMO channel. The task is to find good estimates of the channel coefficients

$$h_{m,k}, \quad \begin{aligned} m &= 1, \dots, M \\ k &= 1, \dots, K, \end{aligned} \quad (8.2)$$

where $h_{m,k}$ denotes the channel coefficient from the k th transmitter to the m th receiver.

To this end, we assume N vectors of TRAINING SIGNALS $\mathbf{s}_n \in \mathbb{C}^K$, $n = 1, \dots, N$. The estimation of the channel coefficients $h_{m,k}$ is based on the received signal vectors

$$\mathbf{y}_n = \mathbf{H} \mathbf{s}_n + \boldsymbol{\eta}_n, \quad n = 1, \dots, N, \quad (8.3)$$

where \mathbf{H} , \mathbf{y}_n and $\boldsymbol{\eta}_n$ denotes the matrix of channel coefficients, the received signal vector and the noise corruption at the receiver for the n th training vector, respectively.

The model for the training channel is thus given by

$$\underbrace{[\mathbf{y}_1, \dots, \mathbf{y}_N]}_{\mathbf{Y}} = \mathbf{H} \underbrace{[\mathbf{s}_1, \dots, \mathbf{s}_N]}_{\bar{\mathbf{S}}} + \underbrace{[\boldsymbol{\eta}_1, \dots, \boldsymbol{\eta}_N]}_{\mathbf{N}}, \quad (8.4)$$

i.e., $\mathbf{Y} = \mathbf{H}\bar{\mathbf{S}} + \mathbf{N}$.

By stacking the column vectors of the matrices, $\mathbf{y} = \text{vec}(\mathbf{Y})$, $\mathbf{h} = \text{vec}(\mathbf{H})$ and $\mathbf{n} = \text{vec}(\mathbf{N})$, we obtain¹

$$\mathbf{y} = \mathbf{S}\mathbf{h} + \mathbf{n} \quad (8.5)$$

$$= (\bar{\mathbf{S}}^T \otimes \mathbf{I}_M)\mathbf{h} + \mathbf{n}, \quad \mathbf{I}_M \in \mathbb{C}^{M \times M}. \quad (8.6)$$

The so called KRONECKER PRODUKT $\bar{\mathbf{S}}^T \otimes \mathbf{I}_M$ means

$$\bar{\mathbf{S}}^T \otimes \mathbf{I}_M = \begin{bmatrix} s_{1,1}\mathbf{I}_M & s_{1,2}\mathbf{I}_M & \cdots & s_{1,K}\mathbf{I}_M \\ s_{2,1}\mathbf{I}_M & s_{2,2}\mathbf{I}_M & \cdots & s_{2,K}\mathbf{I}_M \\ \vdots & \vdots & \ddots & \vdots \\ s_{N,1}\mathbf{I}_M & s_{N,2}\mathbf{I}_M & \cdots & s_{N,K}\mathbf{I}_M \end{bmatrix} \in \mathbb{C}^{NM \times KM}. \quad (8.7)$$

¹Here we use the following identity for matrix equations: $\mathbf{A}\mathbf{X}\mathbf{B} = \mathbf{C} \Leftrightarrow (\mathbf{B}^T \otimes \mathbf{A})\text{vec}(\mathbf{X}) = \text{vec}(\mathbf{C})$.

Example. $K=M=N=2$.

$$\begin{aligned}
\begin{bmatrix} y_{11} & y_{12} \\ y_{21} & y_{22} \end{bmatrix} &= \begin{bmatrix} h_{11} & h_{12} \\ h_{21} & h_{22} \end{bmatrix} \begin{bmatrix} s_{11} & s_{12} \\ s_{21} & s_{22} \end{bmatrix} + \begin{bmatrix} \eta_{11} & \eta_{12} \\ \eta_{21} & \eta_{22} \end{bmatrix} \\
&= \begin{bmatrix} h_{11}s_{11} + h_{12}s_{21} & h_{11}s_{12} + h_{12}s_{22} \\ h_{21}s_{11} + h_{22}s_{21} & h_{21}s_{12} + h_{22}s_{22} \end{bmatrix} + \begin{bmatrix} \eta_{11} & \eta_{12} \\ \eta_{21} & \eta_{22} \end{bmatrix} \\
\begin{bmatrix} y_{11} \\ y_{21} \\ y_{12} \\ y_{22} \end{bmatrix} &= \begin{bmatrix} h_{11}s_{11} + h_{12}s_{21} \\ h_{21}s_{11} + h_{22}s_{21} \\ h_{11}s_{12} + h_{12}s_{22} \\ h_{21}s_{12} + h_{22}s_{22} \end{bmatrix} + \begin{bmatrix} \eta_{11} \\ \eta_{21} \\ \eta_{12} \\ \eta_{22} \end{bmatrix} \\
&= \begin{bmatrix} s_{11} & 0 & s_{21} & 0 \\ 0 & s_{11} & 0 & s_{21} \\ s_{12} & 0 & s_{22} & 0 \\ 0 & s_{12} & 0 & s_{22} \end{bmatrix} \begin{bmatrix} h_{11} \\ h_{21} \\ h_{12} \\ h_{22} \end{bmatrix} + \begin{bmatrix} \eta_{11} \\ \eta_{21} \\ \eta_{12} \\ \eta_{22} \end{bmatrix} \\
&= \begin{bmatrix} s_{11} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} & s_{21} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \\ s_{12} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} & s_{22} \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \end{bmatrix} \begin{bmatrix} h_{11} \\ h_{21} \\ h_{12} \\ h_{22} \end{bmatrix} + \begin{bmatrix} \eta_{11} \\ \eta_{21} \\ \eta_{12} \\ \eta_{22} \end{bmatrix} \\
&= \begin{bmatrix} s_{11} & s_{21} \\ s_{12} & s_{22} \end{bmatrix} \otimes \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} h_{11} \\ h_{21} \\ h_{12} \\ h_{22} \end{bmatrix} + \begin{bmatrix} \eta_{11} \\ \eta_{21} \\ \eta_{12} \\ \eta_{22} \end{bmatrix}.
\end{aligned}$$

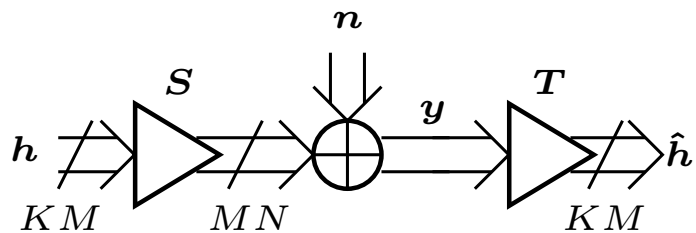


Fig. 8.2: Model for the estimation of a MIMO channel.

Further Assumptions.

We further assume that the stacked vector of channel coefficients $\mathbf{h} \in \mathbb{C}^{KM}$ and the stacked vector of distortions $\mathbf{n} \in \mathbb{C}^{NM}$ are GAUSSIAN DISTRIBUTED as²

$$\mathbf{h} \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_h) \quad \text{and} \quad \mathbf{n} \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_n). \quad (8.8)$$

In the following, we assume no correlations between noise vectors at different time instances, thus leading to a covariance $\mathbf{C}_n = \mathbf{I}_N \otimes \mathbf{C}_\eta \in \mathbb{C}^{NM \times NM}$.

Additionally, assuming no correlations between distortions of adjacent antenna elements leads to $\mathbf{C}_n = \sigma_\eta^2 \mathbf{I}_{NM}$.

The covariance matrix of the channel vector in general also shows additional structural properties, which are not taken into account in the following.

The channel vectors \mathbf{h} and the noise distortions \mathbf{n} are assumed to be STOCHASTICALLY INDEPENDENT, and thus UNCORRELATED ($\text{Cov} [\mathbf{h}, \mathbf{n}^H] = \mathbf{0}$).

Note. Not taking into account these structural properties does not change the channel estimates, however, might be very useful to design more efficient algorithms.

²In contrast to the standard notation, we again denote random vectors by \mathbf{x} in order to avoid any confusions with matrix notation \mathbf{X} .

Due to the LINEAR CHANNEL MODEL $\mathbf{y} = \mathbf{S}\mathbf{h} + \mathbf{n}$, we conclude that the stacked vector of signal vector at the receiver and the unknown channel vector $\mathbf{z} = [\mathbf{y}^T, \mathbf{h}^T]^T$ are JOINTLY GAUSSIAN DISTRIBUTED. The covariances of the JOINT GAUSSIAN PDF is equal to

$$\mathbf{C}_z = \begin{bmatrix} \mathbf{C}_y & \mathbf{C}_{y,h} \\ \mathbf{C}_{h,y} & \mathbf{C}_h \end{bmatrix}, \quad (8.9)$$

with

$$\mathbf{C}_y = \text{Var}[\mathbf{y}] = \text{E}[\mathbf{y}\mathbf{y}^H] = \text{E}[(\mathbf{S}\mathbf{h} + \mathbf{n})(\mathbf{S}\mathbf{h} + \mathbf{n})^H] \cdots = \mathbf{S}\mathbf{C}_h\mathbf{S}^H + \mathbf{C}_n, \quad (8.10)$$

$$\mathbf{C}_{y,h} = \text{Cov}[\mathbf{y}, \mathbf{h}] = \text{E}[\mathbf{y}\mathbf{h}^H] = \text{E}[(\mathbf{S}\mathbf{h} + \mathbf{n})\mathbf{h}^H] \cdots = \mathbf{S}\mathbf{C}_h, \quad (8.11)$$

$$\mathbf{C}_{h,y} = \mathbf{C}_{y,h}^H = \mathbf{C}_h\mathbf{S}^H, \quad (8.12)$$

where we intensively applied $\text{Cov}[\mathbf{h}, \mathbf{n}^H] = \text{E}[\mathbf{h}\mathbf{n}^H] - \boldsymbol{\mu}_h\boldsymbol{\mu}_n^T = \mathbf{0}$.

Note. For the following estimators, we assume full knowledge of the covariance matrices

$$\mathbf{C}_y \in \mathbb{C}^{NM \times NM} \quad \text{and} \quad \mathbf{C}_{y,h} \in \mathbb{C}^{NM \times KM}. \quad (8.13)$$

In practice the required covariance matrices must be estimated as well!

8.2 Minimum Mean Square Error Estimator

Since we assume JOINTLY GAUSSIAN DISTRIBUTED RANDOM VARIABLES, the CONDITIONAL MEAN ESTIMATOR is identical with the linear MMSE Estimator \mathbf{T}_{MMSE} , the minimizer of

$$\min_{\mathbf{T}} \left\{ \mathbb{E} \left[\|\mathbf{h} - \mathbf{T}(\mathbf{S}\mathbf{h} + \mathbf{n})\|^2 \right] \right\}, \quad (8.14)$$

which is given by

$$\hat{\mathbf{h}}_{\text{MMSE}} = \mathbf{T}_{\text{MMSE}} \mathbf{y}, \quad (8.15)$$

$$\mathbf{T}_{\text{MMSE}} = \mathbf{C}_{\mathbf{h},\mathbf{y}} \mathbf{C}_{\mathbf{y}}^{-1} \quad (8.16)$$

$$= \mathbf{C}_{\mathbf{h}} \mathbf{S}^{\text{H}} (\mathbf{S} \mathbf{C}_{\mathbf{h}} \mathbf{S}^{\text{H}} + \mathbf{C}_{\mathbf{n}})^{-1} \quad (8.17)$$

$$= (\mathbf{C}_{\mathbf{h}} \mathbf{S}^{\text{H}} \mathbf{C}_{\mathbf{n}}^{-1} \mathbf{S} + \mathbf{I}_{KM})^{-1} \mathbf{C}_{\mathbf{h}} \mathbf{S}^{\text{H}} \mathbf{C}_{\mathbf{n}}^{-1}. \quad (8.18)$$

Reformulation.

The reformulation of the MMSE ESTIMATOR in (8.18) is based on the MATRIX-INVERSION LEMMA,³

$$(A + BCD)^{-1} = A^{-1} - A^{-1}B(C^{-1} + DA^{-1}B)^{-1}DA^{-1}. \quad (8.19)$$

It follows

$$\begin{aligned} T_{\text{MMSE}} &= C_h S^H (C_n + SC_h S^H)^{-1} \\ &= C_h S^H \left(C_n^{-1} - C_n^{-1}S(C_h^{-1} + S^H C_n^{-1}S)^{-1}S^H C_n^{-1} \right) \\ &= C_h S^H C_n^{-1} - C_h S^H C_n^{-1}S(C_h^{-1} + S^H C_n^{-1}S)^{-1}S^H C_n^{-1} \\ &= C_h \left(I - S^H C_n^{-1}S(C_h^{-1} + S^H C_n^{-1}S)^{-1} \right) S^H C_n^{-1} \\ &= C_h \left((C_h^{-1} + S^H C_n^{-1}S) - S^H C_n^{-1}S \right) (C_h^{-1} + S^H C_n^{-1}S)^{-1} S^H C_n^{-1} \\ &= (C_h^{-1} + S^H C_n^{-1}S)^{-1} S^H C_n^{-1} \\ &= (I_{KM} + C_h S^H C_n^{-1}S)^{-1} C_h S^H C_n^{-1}. \end{aligned}$$

³Assuming the compatibility of the matrices and the existence of the corresponding inverse matrices, it can be shown that

$$\begin{aligned} (A + BCD)(A^{-1} - A^{-1}B(C^{-1} + DA^{-1}B)^{-1}DA^{-1}) &= I + BCDA^{-1} - B(C^{-1} + DA^{-1}B)^{-1}DA^{-1} - BCDA^{-1}B(C^{-1} + DA^{-1}B)^{-1}DA^{-1} \\ &= I + BCDA^{-1} - B(I + CDA^{-1}B)(C^{-1} + DA^{-1}B)^{-1}DA^{-1} \\ &= I + BCDA^{-1} - BC \underbrace{(C^{-1} + DA^{-1}B)(C^{-1} + DA^{-1}B)^{-1}}_I DA^{-1} = I. \end{aligned}$$

8.3 Maximum-Likelihood Estimator

If we consider the unknown channel vector as deterministic, the Jointly Distributed Random Variables \mathbf{y} and \mathbf{h} are replaced by $\mathbf{y} \sim \mathcal{N}(\mathbf{S}\mathbf{h}, \mathbf{C}_{\mathbf{y}|\mathbf{h}=\mathbf{h}})$, with $\mathbf{C}_{\mathbf{y}|\mathbf{h}=\mathbf{h}} = \mathbf{C}_n$, and the CONDITIONAL PDF

$$f_{\mathbf{y}|\mathbf{h}}(\mathbf{y}|\mathbf{h}) = (\pi^{NM} \det \mathbf{C}_n)^{-1} \exp(-(\mathbf{y} - \mathbf{S}\mathbf{h})^H \mathbf{C}_n^{-1} (\mathbf{y} - \mathbf{S}\mathbf{h})). \quad (8.20)$$

Consequently, the ML Estimator \mathbf{T}_{ML} is the minimizer of the optimization problem

$$\min_{\mathbf{T}} \{(\mathbf{y} - \mathbf{S}\mathbf{h})^H \mathbf{C}_n^{-1} (\mathbf{y} - \mathbf{S}\mathbf{h})\}, \quad (8.21)$$

which results into

$$\hat{\mathbf{h}}_{\text{ML}} = \mathbf{T}_{\text{ML}} \mathbf{y}, \quad (8.22)$$

$$\mathbf{T}_{\text{ML}} = (\mathbf{S}^H \mathbf{C}_n^{-1} \mathbf{S})^{-1} \mathbf{S}^H \mathbf{C}_n^{-1}. \quad (8.23)$$

Note. Given the ML Estimator \mathbf{T}_{ML} the MMSE Estimator, can be denoted as

$$\mathbf{T}_{\text{MMSE}} = \left(\mathbf{I} + (\mathbf{S}^H \mathbf{C}_n^{-1} \mathbf{S})^{-1} \mathbf{C}_h^{-1} \right)^{-1} \mathbf{T}_{\text{ML}} = \dots \quad (8.24)$$

$$= (\mathbf{S}^H \mathbf{C}_n \mathbf{S})^{-\frac{1}{2}} \left(\mathbf{I}_{KM} + (\mathbf{S}^H \mathbf{C}_n \mathbf{S})^{-\frac{1}{2}} \mathbf{C}_h^{-1} (\mathbf{S}^H \mathbf{C}_n \mathbf{S})^{-\frac{1}{2}} \right)^{-1} (\mathbf{S}^H \mathbf{C}_n \mathbf{S})^{\frac{1}{2}} \mathbf{T}_{\text{ML}}. \quad (8.25)$$

8.4 Correlation Estimator

Assuming training signals such that $\mathbf{S}^H \mathbf{S} \propto N \mathbf{I}_{KM \times KM}$ and further assuming AGWN, the ML Estimator is equal to the so called CORRELATION ESTIMATOR (C)

$$\hat{\mathbf{h}}_C = \mathbf{T}_C \mathbf{y}, \quad (8.26)$$

$$\mathbf{T}_C = (\mathbf{S}^H \mathbf{C}_n^{-1} \mathbf{S})^{-1} \mathbf{S}^H \mathbf{C}_n^{-1} \propto \mathbf{S}^H. \quad (8.27)$$

Note. The CORRELATION ESTIMATOR is the simplest estimator which is often applied beyond cases where it is identical with the ML Estimator.

8.5 Matched Filter Estimator

The MATCHED FILTER ESTIMATOR (MF) takes the principle of STRONG CORRELATIONS between the channel vector \mathbf{h} and its $\hat{\mathbf{h}}$ as a cost function. In other words, the MF Estimator \mathbf{T}_{MF} is the maximizer of the optimization problem

$$\max_{\mathbf{T}} \left\{ \frac{\left| \mathbb{E} [\hat{\mathbf{h}}^H \mathbf{h}] \right|^2}{\text{tr} [\text{Var} [\mathbf{T} \mathbf{n}]]} \right\}. \quad (8.28)$$

The corresponding MF estimator is equal to

$$\hat{\mathbf{h}}_{\text{MF}} = \mathbf{T}_{\text{MF}} \mathbf{y}, \quad (8.29)$$

$$\mathbf{T}_{\text{MF}} \propto \mathbf{C}_{\mathbf{h}} \mathbf{S}^H \mathbf{C}_{\mathbf{n}}^{-1}. \quad (8.30)$$

8.6 Special Case

In order to simplify the further analysis, we assume training signals such that

$$\mathbf{S}^H \mathbf{S} = N \sigma_s^2 \mathbf{I}_{KM} \quad (8.31)$$

$$\mathbf{C}_n = \sigma_\eta^2 \mathbf{I}_{NM}, \quad (8.32)$$

where σ_s^2 and σ_η^2 are the variances of single training and noise distortions signals.

Using these assumptions the proposed estimators can be denoted as follows:

$$\mathbf{T}_{\text{MF}} \propto \mathbf{C}_h \mathbf{S}^H, \quad (8.33)$$

$$\mathbf{T}_{\text{ML}} = (\mathbf{S}^H \mathbf{C}_n^{-1} \mathbf{S})^{-1} \mathbf{S}^H \mathbf{C}_n^{-1} = \frac{1}{N \sigma_s^2} \mathbf{S}^H, \quad (8.34)$$

$$\mathbf{T}_{\text{MMSE}} = (\mathbf{C}_h \mathbf{S}^H \mathbf{C}_n^{-1} \mathbf{S} + \mathbf{I}_{KM})^{-1} \mathbf{C}_h \mathbf{S}^H \mathbf{C}_n^{-1} = \left(\mathbf{C}_h + \frac{\sigma_\eta^2}{N \sigma_s^2} \mathbf{I}_{KM} \right)^{-1} \mathbf{C}_h \underbrace{\frac{1}{N \sigma_s^2} \mathbf{S}^H}_{\mathbf{T}_{\text{ML}}}. \quad (8.35)$$

The simplified case allows a convenient asymptotic analysis for the three estimators.

$$\text{LOW NOISE REGIME:} \quad \lim_{\sigma_\eta^2 \rightarrow 0} \mathbf{T}_{\text{MMSE}} = \mathbf{T}_{\text{ML}} \quad (8.36)$$

$$\text{HIGH NOISE REGIME :} \quad \lim_{\sigma_\eta^2 \rightarrow \infty} \sigma_\eta^2 \mathbf{T}_{\text{MMSE}} \propto \mathbf{T}_{\text{MF}}. \quad (8.37)$$

Using the EIGEN VALUE DECOMPOSITION of the channel covariance matrix

$$\mathbf{C}_h = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^H = \sum_{i=1}^{KM} \lambda_i \mathbf{u}_i \mathbf{u}_i^H, \quad (8.38)$$

the estimators are equal to

$$\mathbf{T}_{\text{MF}} \propto \mathbf{U} \mathbf{\Lambda} \mathbf{U}^H \mathbf{S}^H \quad (8.39)$$

$$\mathbf{T}_{\text{ML}} = \frac{1}{N\sigma_s^2} \mathbf{S}^H, \quad (8.40)$$

$$\mathbf{T}_{\text{MMSE}} = \mathbf{U} \left(\mathbf{\Lambda} + \frac{\sigma_\eta^2}{N\sigma_s^2} \mathbf{I}_{KM} \right)^{-1} \mathbf{\Lambda} \mathbf{U}^H \mathbf{T}_{\text{ML}}. \quad (8.41)$$

Interpretation. Denoting the MMSE estimate in a series of the eigenvectors \mathbf{u}_i , $i = 1, \dots, KM$,

$$\hat{\mathbf{h}}_{\text{MMSE}} = \sum_{i=1}^{KM} \frac{\lambda_i}{\lambda_i + \frac{\sigma_\eta^2}{N\sigma_s^2}} \mathbf{u}_i \mathbf{u}_i^H \hat{\mathbf{h}}_{\text{ML}}, \quad (8.42)$$

the improvement by the MMSE Estimator can be interpreted as a WEIGHTING OF THE ML ESTIMATE which is optimally adapted to the eigenspaces of the random channel vector.

8.7 Bias/Variance Trade–Off

For the analysis of all introduced channel estimators, we consider the decomposition of the MEAN SQUARE ERROR into the estimator's SQUARED BIAS and VARIANCE. Since we compare estimators of different paradigms with respect to A PRIORI INFORMATION about the unknown parameter, we study the mean average of the respective Bias/Variance Decomposition. To this end, the MSE is decomposed by the following steps⁴

$$\mathbb{E} \left[\| \mathbf{h} - \mathbf{T}(\mathbf{S}\mathbf{h} + \mathbf{n}) \|^2 \right] = \mathbb{E} \left[\mathbb{E} \left[\| \mathbf{h} - \mathbf{T}(\mathbf{S}\mathbf{h} + \mathbf{n}) \|^2 \mid \mathbf{h} \right] \right] \quad (8.43)$$

$$= \mathbb{E} \left[\mathbb{E} \left[\left\| \underbrace{(\mathbf{I}_{KM} - \mathbf{T}\mathbf{S})\mathbf{h}}_{\text{Bias}} + \mathbf{T}\mathbf{n} \right\|^2 \mid \mathbf{h} \right] \right] \quad (8.44)$$

$$= \underbrace{\mathbb{E} \left[\underbrace{\|(\mathbf{I}_{KM} - \mathbf{T}\mathbf{S})\mathbf{h}\|^2}_{\text{Squared Bias}} \mid \mathbf{h} \right]}_{\text{Average Squared Bias}} + \underbrace{\mathbb{E} \left[\underbrace{\|\mathbf{T}\mathbf{n}\|^2}_{\text{Variance}} \mid \mathbf{h} \right]}_{\text{Average Variance = Variance}} \quad (8.45)$$

$$= \text{tr} \left[(\mathbf{I}_{KM} - \mathbf{T}\mathbf{S})\mathbf{C}_h(\mathbf{I}_{KM} - \mathbf{T}\mathbf{S})^H \right] + \text{tr} \left[\mathbf{T}\mathbf{C}_n\mathbf{T}^H \right]. \quad (8.46)$$

⁴The last step is found by using the identity $\mathbb{E} \left[\|\mathbf{x}\|^2 \right] = \mathbb{E} \left[\mathbf{x}^H \mathbf{x} \right] = \mathbb{E} \left[\text{tr} \left[\mathbf{x}\mathbf{x}^H \right] \right] = \text{tr} \left[\mathbb{E} \left[\mathbf{x}\mathbf{x}^H \right] \right] = \text{tr} \left[\mathbf{R}_x \right]$ and $\mathbf{R}_x = \mathbf{C}_x$ if \mathbf{x} is zero mean.

For the simplified assumptions $\mathbf{S}^H \mathbf{S} = N\sigma_s^2 \mathbf{I}_{KM}$ and $\mathbf{C}_n = \sigma_\eta^2 \mathbf{I}_{NM}$.⁵

Estimator	Averaged Squared Bias	Variance
ML/Correlator	0	$KM \frac{\sigma_\eta^2}{N\sigma_s^2}$
Matched Filter	$\sum_{i=1}^{KM} \lambda_i \left(\frac{\lambda_i}{\lambda_1} - 1 \right)^2$	$\sum_{i=1}^{KM} \left(\frac{\lambda_i}{\lambda_1} \right)^2 \frac{\sigma_\eta^2}{N\sigma_s^2}$
MMSE	$\sum_{i=1}^{KM} \lambda_i \left(\frac{1}{1 + \frac{\sigma_\eta^2}{\lambda_i N \sigma_s^2}} - 1 \right)^2$	$\sum_{i=1}^{KM} \frac{1}{\left(1 + \frac{\sigma_\eta^2}{\lambda_i N \sigma_s^2} \right)^2} \frac{\sigma_\eta^2}{N\sigma_s^2}$

The minimum MSE achieved by the MMSE Estimator is given by

$$\mathbb{E} \left[\|\mathbf{h} - \mathbf{T}(\mathbf{S}\mathbf{h} + \mathbf{n})\|^2 \right] = \sum_{i=1}^{KM} \frac{\sigma_\eta^2}{N\sigma_s^2} \lambda_i \left(\lambda_i + \frac{\sigma_\eta^2}{N\sigma_s^2} \right)^{-1}. \quad (8.47)$$

⁵For this analysis the MF estimator is assumed to be scaled as $\mathbf{T}_{MF} = \frac{1}{\lambda_1} \mathbf{U} \mathbf{\Lambda} \mathbf{U}^H \mathbf{T}_{ML}$.

8.8 Numerical Experiments

We consider

- a NON-DISPERSIVE MIMO channel $\mathbf{h} = \text{vec}[\mathbf{H}]$,
- with $M = K = 8$ antenna elements,
- a UNIFORM LINEAR ARRAYS at the receiver,
- and $N = 16$ training signal vectors $\mathbf{s}_n \in \{-1, +1\}^K$ (BINARY).
- The covariance matrix \mathbf{C}_h corresponds to 8 impinging planar wavefronts at

$$-45^\circ, -32.1^\circ, -19.3^\circ, -6.4^\circ, 6.4^\circ, 19.3^\circ, 32.1^\circ, 45^\circ, \quad (8.48)$$

where each AZIMUTH ANGLE represents a cluster of wavefronts spread according to a LAPLACE ANGULAR POWER SPECTRUM with spread $\sigma = 5^\circ$

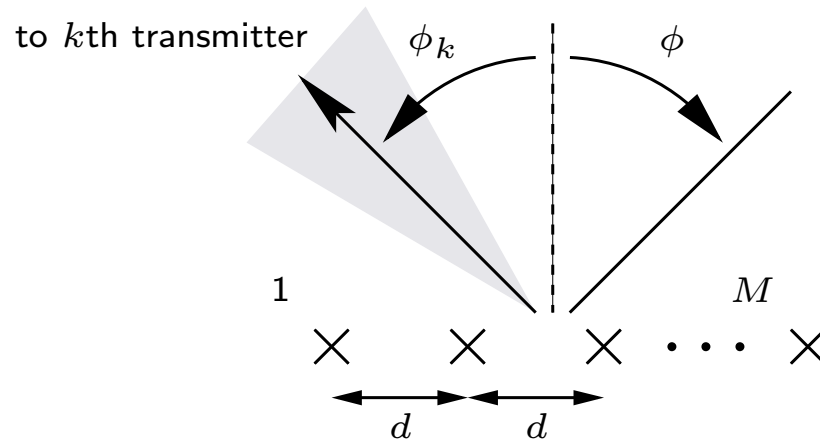


Fig. 8.3: Array geometry.

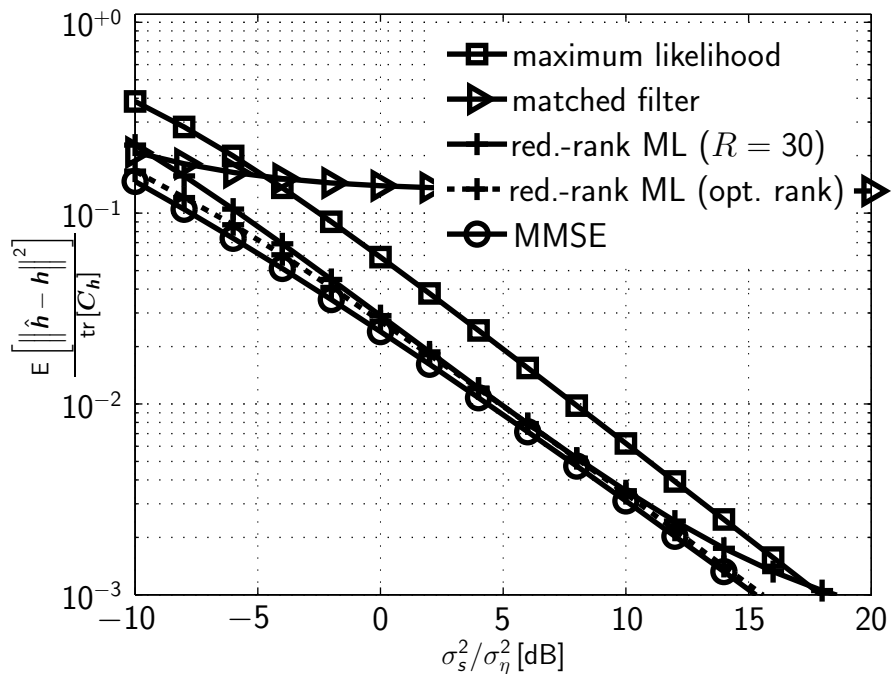


Fig. 8.4: Normalized MSE of channel estimators for $M = 8$, $K = 8$, and $N = 16$.

8.9 Numerical Experiments (cont'd)

We now consider

- a DISPERSIVE SINGLE-INPUT MULTIPLE-OUTPUT (SIMO) channel $\mathbf{h} = \text{vec}[\mathbf{H}]$,
- with $M = 8$ receiver antenna elements, $K = 1$ transmitter antenna elements,
- a UNIFORM LINEAR ARRAYS at the receiver,
- with dispersion order $L = 4$,
- and $N = 16$ training signal vectors $\mathbf{s}_n \in \{-1, +1\}^K$ (BINARY).
- The covariance matrix \mathbf{C}_h corresponds to 5 impinging planar wavefronts at

$$-45^\circ, -22.5^\circ, 0^\circ, 22.5^\circ, 45^\circ \quad (8.49)$$

per delay, where each AZIMUTH ANGLE represents a cluster of wavefronts spread according to a LAPLACE ANGULAR POWER SPECTRUM with spread $\sigma = 5^\circ$

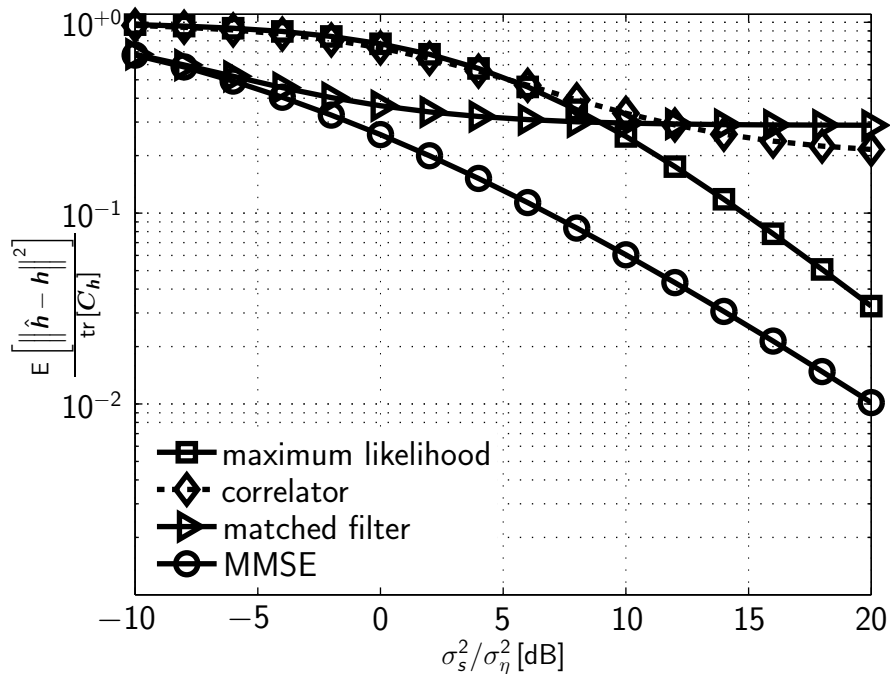


Fig. 8.5: Normalized MSE of channel estimators for $M = 8$, $K = 1$, $L = 4$ (DISPERSIVE), and $N = 16$.

References

- S. M. Kay. *Fundamentals of Statistical Signal Processing: Estimation Theory*, volume I, Prentice Hall Signal Processing Series, 1993.
- L. L. Scharf. *Statistical Signal Processing: Detection, Estimation, and Time Series Analysis*, 1st edition, Prentice Hall, 1991.
- F. A. Dietrich and W. Utschick. *Pilot assisted channel estimation based on second order statistics*. IEEE Transactions on Signal Processing, 53(3):1178–1193, 2005.

Part VI

Estimation of Random Sequences

9. Random Sequences

In the following, we study SEQUENCES OF RANDOM VARIABLES $X_n : \Omega \rightarrow \mathbb{X}_n, n \in \mathbb{N}$, with a probability space (Ω, \mathbb{F}, P) and a measurement space (STATE SPACE) $\mathbb{X} \supseteq \mathbb{X}_n$, which are uniquely defined by the distribution functions

$$F_{X_{i_1}, \dots, X_{i_n}}(x_{i_1}, \dots, x_{i_n}) = P(X_{i_1} \leq x_{i_1}, \dots, X_{i_n} \leq x_{i_n}), \quad (9.1)$$

for all $i_1, \dots, i_n \in \mathbb{N}$ and for all $n \in \mathbb{N}$.

Given the existence of probability density functions of the corresponding distribution functions, the random sequences can be defined alternatively by means of

$$f_{X_{i_1}, \dots, X_{i_n}}(x_{i_1}, \dots, x_{i_n}), \quad (9.2)$$

for all $i_1, \dots, i_n \in \mathbb{N}$ and for all $n \in \mathbb{N}$.

Note. The definitions of the mean, variance, auto-correlation, and cross-correlation function of random sequences directly follow from the definitions of single random variables.

9.1 Conditional Stochastical Independence

For the subsequent definition of MARKOV SEQUENCES, we introduce the CONDITIONAL STOCHASTICAL INDEPENDENCE of random variables.

Definition. Given Y , the random variables X and Z are conditionally stochastically independent, if

$$f_{Z|Y,X}(z | y, x) = f_{Z|Y}(z | y), \quad (9.3)$$

and $f_{Y,X}(y, x) > 0$.

Alternatively, following the definition of stochastical independence, the CONDITIONAL STOCHASTICAL INDEPENDENCE of Z and X given Y , we obtain

$$f_{Z,X|Y}(z, x | y) = f_{Z|Y}(z | y)f_{X|Y}(x | y), \quad (9.4)$$

respectively. As a short notation for conditional stochastical independence, we use

$$X \rightarrow Y \rightarrow Z. \quad (9.5)$$

9.2 Markov Sequences

Definition. A MARKOV SEQUENCE (PROCESS) is a random sequence $X_n : \Omega \rightarrow \mathbb{X}_n$, $n \in \mathbb{N}$, with a probability space (Ω, \mathbb{F}, P) and a state space $\mathbb{X} \supseteq \mathbb{X}_n$, where for any ASCENDING SEQUENCE of indices $n_i \in \mathbb{N}$ and $k > 2$ the following CONDITIONAL STOCHASTICAL INDEPENDENCE holds, namely

$$(X_{n_1}, \dots, X_{n_{k-2}}) \rightarrow X_{n_{k-1}} \rightarrow X_{n_k}, \quad (9.6)$$

which means

$$f_{X_{n_k} | X_{n_{k-1}}, \dots, X_{n_1}}(x_{n_k} | x_{n_{k-1}}, \dots, x_{n_1}) = f_{X_{n_k} | X_{n_{k-1}}}(x_{n_k} | x_{n_{k-1}}), \quad (9.7)$$

and $f_{X_{n_{k-1}}, \dots, X_{n_1}}(x_{n_{k-1}}, \dots, x_{n_1}) > 0$.¹

Definition. The conditional $f_{X_{n_k} | X_{n_{k-1}}}(x_{n_k} | x_{n_{k-1}})$ defines the STATE-TRANSITION DENSITY.

¹In this lecture, we consider so-called FIRST-ORDER MARKOV SEQUENCES with $f_{X_{n_k} | X_{n_{k-1}}, \dots, X_{n_1}} = f_{X_{n_k} | X_{n_{k-1}}}$. If the underlying conditional independence holds only if we condition on more than one previous state, we have HIGHER-ORDER MARKOV SEQUENCES. If, for example, $f_{X_{n_k} | X_{n_{k-1}}, \dots, X_{n_1}} = f_{X_{n_k} | X_{n_{k-1}}, X_{n_{k-2}}}$, we have an order of two. The order of a MARKOV SEQUENCE can be interpreted as the memory of the stochastic sequence.

As a consequence of this MARKOV PROPERTY we obtain for the joint PDF of the first n elements

$$\begin{aligned}
 f_{X_1, \dots, X_n}(x_1, \dots, x_n) &= f_{X_n|X_{n-1}, \dots, X_1}(x_n|x_{n-1}, \dots, x_1) f_{X_{n-1}|X_{n-2}, \dots, X_1}(x_{n-1}|x_{n-2}, \dots, x_1) \dots f_{X_1}(x_1) \\
 &= f_{X_n|X_{n-1}}(x_n|x_{n-1}) f_{X_{n-1}|X_{n-2}}(x_{n-1}|x_{n-2}) \dots f_{X_1}(x_1) \\
 &= f_{X_1}(x_1) \prod_{i=2}^n f_{X_i|X_{i-1}}(x_i|x_{i-1}),
 \end{aligned} \tag{9.8}$$

where

- the first line follows from the MULTIPLICATION RULE²
- and the second line follows from the STOCHASTICAL INDEPENDENCE between X_n and X_{n-2} given an event related to X_{n-1} .

² $P(A \cap B \cap C) = P(A|B \cap C) P(B|C) P(C)$.

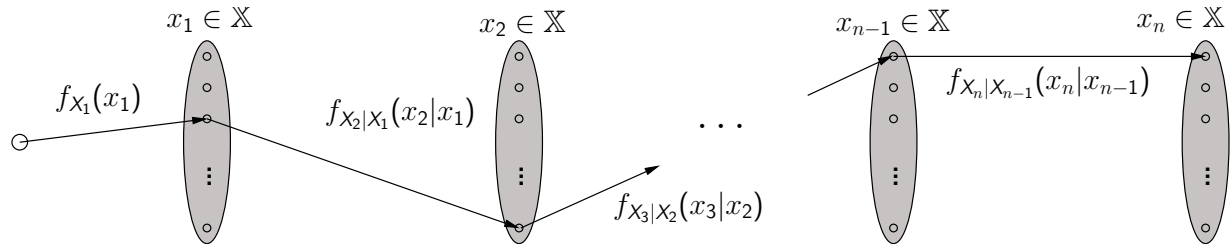


Fig. 9.1: The joint PDF $f_{X_1, \dots, X_n}(x_1, \dots, x_n)$ obviously results from a chain STATE-TRANSITION DENSITIES starting with the PDF of the random variable X_1 .

9.3 Chapman-Kolmogorov Equation

Given a Markov sequence $X_n : \Omega \rightarrow \mathbb{X}_n$ and $f_{X_n}(x_n) > 0$ for all $n \in \mathbb{N}$, the $(m + \ell)$ -step state-transition density starting from X_n depends on the m -step (from X_n) and ℓ -step (from X_{n+m}) state-transition density according to the CHAPMAN-KOLMOGOROV EQUATION

$$f_{X_{n+m+\ell}|X_n}(x_{n+m+\ell}|x_n) = \int_{x_{n+m} \in \mathbb{X}_n} f_{X_{n+m+\ell}|X_{n+m}}(x_{n+m+\ell}|x_{n+m}) f_{X_{n+m}|X_n}(x_{n+m}|x_n) dx_{n+m}. \quad (9.9)$$

The CHAPMAN-KOLMOGOROV EQUATIONS (9.9) can be derived

- by marginalization of $f_{X_{n+m+\ell}, X_{n+m}|X_n}(x_{n+m+\ell}, x_{n+m}|x_n)$
- and further taking into account the Markov condition of X_n , i.e.,

$$f_{X_{n+m+\ell}|X_{n+m}, X_n}(x_{n+m+\ell}|x_{n+m}, x_n) = f_{X_{n+m+\ell}|X_{n+m}}(x_{n+m+\ell}|x_{n+m}). \quad (9.10)$$

Example.

We consider the 2-step state-transition density based on two 1-step state-transition densities, i.e.,

$$\begin{aligned}
 f_{X_{n+2}|X_n}(x_{n+2}|x_n) &= \underbrace{\int_{-\infty}^{\infty} f_{X_{n+2}, X_{n+1}|X_n}(x_{n+2}, \xi|x_n) d\xi}_{\text{Marginalization}} \\
 &= \int_{-\infty}^{\infty} \underbrace{f_{X_{n+2}|X_{n+1}, X_n}(x_{n+2}|\xi, x_n) f_{X_{n+1}|X_n}(\xi|x_n)}_{\text{Conditional PDF} \times \text{Prior PDF}} d\xi \\
 &= \int_{-\infty}^{\infty} \underbrace{f_{X_{n+2}|X_{n+1}}(x_{n+2}|\xi)}_{\text{Markov Property}} f_{X_{n+1}|X_n}(\xi|x_n) d\xi.
 \end{aligned}$$

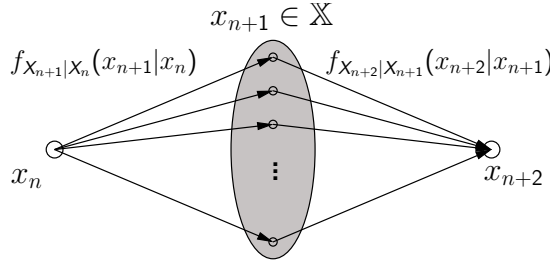


Fig. 9.2: State-Transitions from $\{X_n = x_n\}$ to $\{X_{n+2} = x_{n+2}\}$.

The General Case.

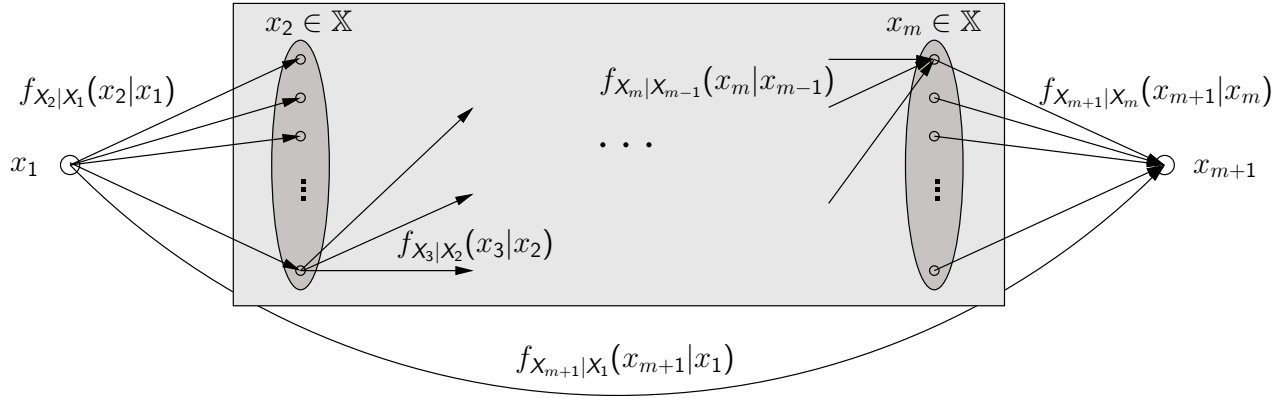


Fig. 9.3: Illustration of (9.9) (m -step state-transition density) for $n = 1$.

9.4 Estimation of Markov Sequences

A HIDDEN MARKOV PROCESS, is a Markov Process which states cannot be observed directly, i.e., the states are *hidden*. The state variables \mathbf{X}_n are observed indirectly via the random observations \mathbf{Y}_n drawn from the conditional PDF $f_{\mathbf{Y}_n|\mathbf{X}_n}(\mathbf{y}_n|\mathbf{x}_n)$, which is the likelihood of \mathbf{Y}_n given that $\mathbf{X}_n = \mathbf{x}_n$.

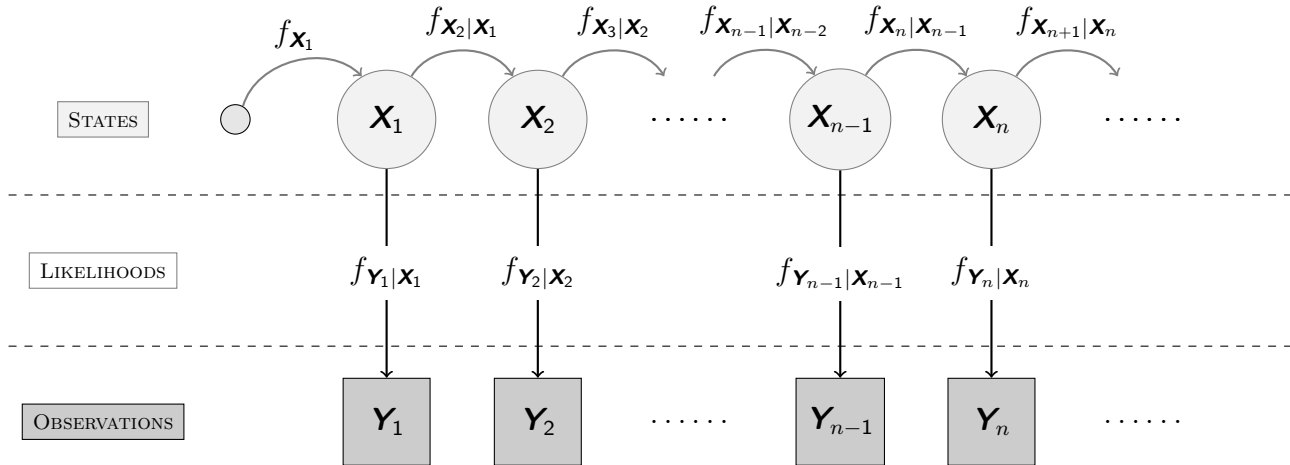


Fig. 9.4: Illustration of a HIDDEN MARKOV PROCESS.

The CHAPMAN-KOLMOGOROV EQUATION is the key element of

A recursive computation of the conditional PDFs of a Markov sequence \mathbf{X}_n conditioned on observations of $\mathbf{Y}_{(n)}$ thereof, this is \mathbf{Y}_n and $\mathbf{Y}_{n-1}, \dots, \mathbf{Y}_1$.^a

^aIn the following, the short notation $\mathbf{Y}_{(n)}$ will be used instead of $\mathbf{Y}_n, \dots, \mathbf{Y}_1$. The same holds accordingly for $\mathbf{Y}_{n-1}, \dots, \mathbf{Y}_1$, etc.

The recursive updates can be obtained in two steps:

a) CHAPMAN-KOLMOGOROV EQUATION and the MARKOV CONDITION,

$$\begin{aligned}
 f_{\mathbf{X}_n | \mathbf{Y}_{(n-1)}}(\mathbf{x}_n | \mathbf{y}_{(n-1)}) &= \\
 &= \int_{\mathbb{X}} f_{\mathbf{X}_n | \mathbf{X}_{n-1}, \mathbf{Y}_{(n-1)}}(\mathbf{x}_n | \mathbf{x}_{n-1}, \mathbf{y}_{(n-1)}) f_{\mathbf{X}_{n-1} | \mathbf{Y}_{(n-1)}}(\mathbf{x}_{n-1} | \mathbf{y}_{(n-1)}) d\mathbf{x}_{n-1} \\
 &= \int_{\mathbb{X}} f_{\mathbf{X}_n | \mathbf{X}_{n-1}}(\mathbf{x}_n | \mathbf{x}_{n-1}) f_{\mathbf{X}_{n-1} | \mathbf{Y}_{(n-1)}}(\mathbf{x}_{n-1} | \mathbf{y}_{(n-1)}) d\mathbf{x}_{n-1}.
 \end{aligned} \tag{9.11}$$

b) BAYES RULE

$$\begin{aligned}
 f_{\mathbf{X}_n | \mathbf{Y}_n, \mathbf{Y}_{(n-1)}}(\mathbf{x}_n | \mathbf{y}_n, \mathbf{y}_{(n-1)}) &= \\
 &= \frac{f_{\mathbf{Y}_n | \mathbf{X}_n, \mathbf{Y}_{(n-1)}}(\mathbf{y}_n | \mathbf{x}_n, \mathbf{y}_{(n-1)}) f_{\mathbf{X}_n | \mathbf{Y}_{(n-1)}}(\mathbf{x}_n | \mathbf{y}_{(n-1)})}{f_{\mathbf{Y}_n | \mathbf{Y}_{(n-1)}}(\mathbf{y}_n | \mathbf{y}_{(n-1)})} \\
 &= \frac{f_{\mathbf{Y}_n | \mathbf{X}_n}(\mathbf{y}_n | \mathbf{x}_n) f_{\mathbf{X}_n | \mathbf{Y}_{(n-1)}}(\mathbf{x}_n | \mathbf{y}_{(n-1)})}{f_{\mathbf{Y}_n | \mathbf{Y}_{(n-1)}}(\mathbf{y}_n | \mathbf{y}_{(n-1)})}. \tag{9.12}
 \end{aligned}$$

Both steps together are equivalent (except some scaling constant) to³

$$\begin{aligned}
 f_{\mathbf{X}_n | \mathbf{Y}_{(n)}}(\mathbf{x}_n | \mathbf{y}_{(n)}) &\propto \\
 f_{\mathbf{Y}_n | \mathbf{X}_n}(\mathbf{y}_n | \mathbf{x}_n) \int_{\mathbb{X}} f_{\mathbf{X}_n | \mathbf{X}_{n-1}}(\mathbf{x}_n | \mathbf{x}_{n-1}) f_{\mathbf{X}_{n-1} | \mathbf{Y}_{(n-1)}}(\mathbf{x}_{n-1} | \mathbf{y}_{(n-1)}) d\mathbf{x}_{n-1}. \tag{9.13}
 \end{aligned}$$

³The denominator of Eq. (9.12) does not depend on the random variable \mathbf{X}_n to be estimated and is implicitly given as the normalization constant $f_{\mathbf{Y}_n | \mathbf{Y}_{(n-1)}}(\mathbf{y}_n | \mathbf{y}_{(n-1)}) = \int_{\mathbb{X}} f_{\mathbf{Y}_n | \mathbf{X}_n}(\mathbf{y}_n | \mathbf{x}_n) f_{\mathbf{X}_n | \mathbf{Y}_{(n-1)}}(\mathbf{x}_n | \mathbf{y}_{(n-1)}) d\mathbf{x}_n$.

Note.

In essence, there are two conditional PDFs and one state-transition PDF which form the basis for the recursive computation of the conditional PDFs of the state variables.

desired conditional PDF in the n th step

$$\underbrace{f_{\mathbf{x}_n|\mathbf{y}_{(n)}}(\mathbf{x}_n|\mathbf{y}_{(n)})}_{\text{likelihood of the new observation}} \propto \int_{\mathbb{X}} \underbrace{f_{\mathbf{x}_n|\mathbf{x}_{n-1}}(\mathbf{x}_n|\mathbf{x}_{n-1})}_{\text{state-transition PDF}} \underbrace{f_{\mathbf{x}_{n-1}|\mathbf{y}_{(n-1)}}(\mathbf{x}_{n-1}|\mathbf{y}_{(n-1)})}_{\text{conditional PDF from the } (n-1)\text{th-step before}} d\mathbf{x}_{n-1} . \quad (9.14)$$

The principal idea of a variety of estimation methods for TRACKING OF STATE VARIABLES is to exploit the recursive computation of the conditional PDFs for the recursive estimation of the characteristic parameters of the distributions, cf. Chapter 10 and 12.

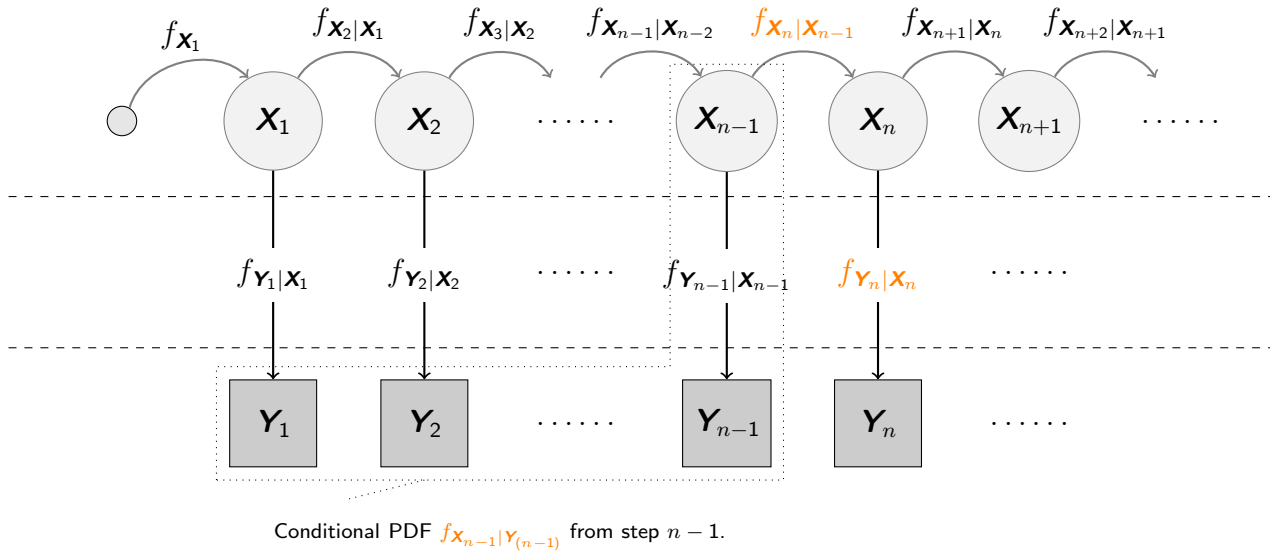


Fig. 9.5: In order to determine the cond. PDF $f_{\mathbf{x}_n|\mathbf{y}_{(n)}}$, we need the cond. PDF $f_{\mathbf{x}_{n-1}|\mathbf{y}_{(n-1)}}$ of the previous step, the likelihood $f_{\mathbf{y}_n|\mathbf{x}_n}$, and the state-transition PDF $f_{\mathbf{x}_n|\mathbf{x}_{n-1}}$, cf. Eq. (9.13).

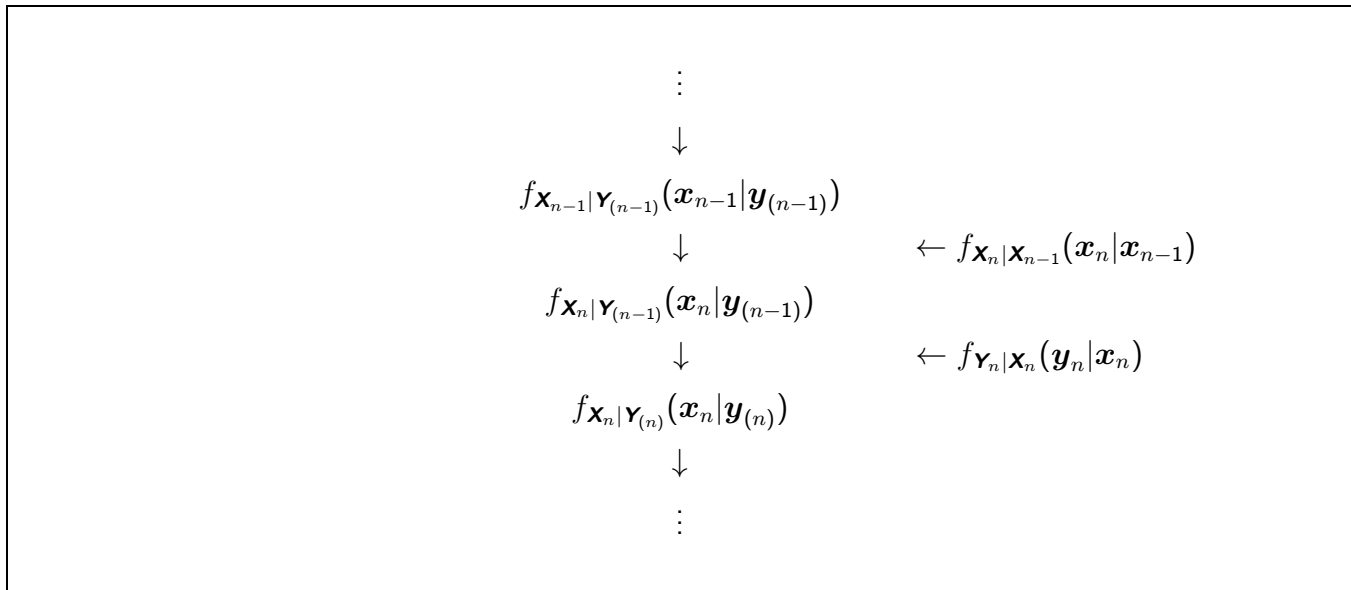


Fig. 9.6: The recursive structure of the computation of the conditional PDFs.

9.5 Gauß-Markov Sequences

Definition. GAUSS-MARKOV SEQUENCES are Gaussian sequences which fulfill the Markov condition, i.e., any ensemble of random variables of the sequence is jointly Gaussian distributed and the state-transition probability density functions⁴ are fully described by conditional means and covariances and the observation of the last recent random variable of the sequence. Since for jointly Gaussian variables the conditional mean of a random variable is a linear function of the others, $\mu_{n|n-1} = E[X_n|X_{n-1} = x_{n-1}] = \rho_n x_{n-1}$, we obtain

$$X_0 \sim \mathcal{N}(0, \sigma_0^2) \quad (9.15)$$

$$X_1 \sim \mathcal{N}(\mu_{X_1|X_0=x_0}, \sigma_{X_1|X_0=x_0}^2), \quad (9.16)$$

$$\vdots \quad (9.17)$$

$$X_n \sim \mathcal{N}(\mu_{X_n|X_{n-1}=x_{n-1}}, \sigma_{X_n|X_{n-1}=x_{n-1}}^2). \quad (9.18)$$

Alternatively,⁵ the GAUSS-MARKOV SEQUENCE $X_{n-2} \rightarrow X_{n-1} \rightarrow X_n$ is fully described by

$$X_n = \rho_n X_{n-1} + W_n, \quad (9.19)$$

with $\rho_n = c_{X_n, X_{n-1}} \sigma_{X_{n-1}}^{-2}$ and $W_n \sim \mathcal{N}(0, \sigma_{X_n|X_{n-1}=x_{n-1}}^2)$.

⁴W.l.o.g. the further considerations are presented for the case of continuous random variables.

⁵W.l.o.g. we assume zero mean random variables for the unconditioned state variables X_n , i.e., $E[X_n] = 0, \forall n \in \mathbb{N} \cup \{0\}$.

10. Kalman Filter

In the following, we study the problem of estimating the stochastic path of a Gauß-Markov sequence \mathbf{X}_n based on linear observations $\mathbf{Y}_n, \mathbf{Y}_{n-1}, \dots$ of the random state variables,¹

$$\mathbf{X}_n = \mathbf{G}_n \mathbf{X}_{n-1} + \mathbf{V}_n, \quad (10.1)$$

$$\mathbf{Y}_n = \mathbf{H}_n \mathbf{X}_n + \mathbf{W}_n, \quad (10.2)$$

with

$$\mathbf{V}_n \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_{\mathbf{V}_n}),$$

$$\mathbf{W}_n \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_{\mathbf{W}_n}).$$

In the context of KALMAN FILTERING, the random variables \mathbf{V}_n and \mathbf{W}_n are called PROCESS NOISE and MEASUREMENT NOISE. The random variables $\mathbf{W}_m, \mathbf{W}_n$, and $\mathbf{V}_k, \mathbf{V}_\ell$ are stochastically independent for different indices m, n and k, ℓ , and $\mathbf{V}_m, \mathbf{W}_n$ are stochastically independent for any m and n .

Note. Eqs. (10.1)-(10.2) are a direct consequence of Eq. (9.19). In the following, we assume that the states \mathbf{X}_n are OBSERVABLE.

¹For the sake of generality, we consider the case of a sequence of random vectors.

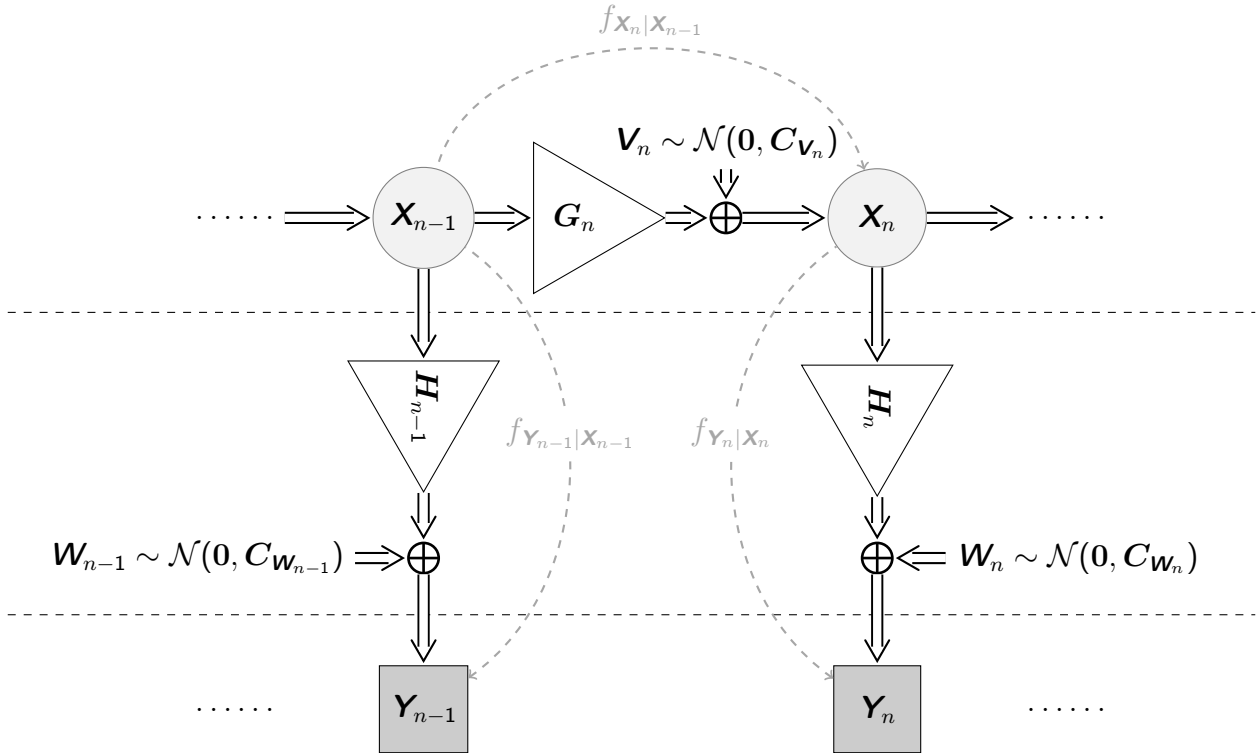


Fig. 10.1: In a HIDDEN GAUSS-MARKOV PROCESS, the states \mathbf{X}_n are given as $\mathbf{X}_n = \mathbf{G}_n \mathbf{X}_{n-1} + \mathbf{V}_n$ with $\mathbf{V}_n \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_{\mathbf{V}_n})$ and the observations \mathbf{Y}_n are given as $\mathbf{Y}_n = \mathbf{H}_n \mathbf{X}_n + \mathbf{W}_n$ with $\mathbf{W}_n \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_{\mathbf{W}_n})$. Thus, the conditional PDFs $f_{\mathbf{X}_n|\mathbf{Y}_{(n)}}$, the likelihoods $f_{\mathbf{Y}_n|\mathbf{X}_n}$, and the state-transition PDFs $f_{\mathbf{X}_n|\mathbf{X}_{n-1}}$ are Gaussian.

In general, the MMSE optimal estimate is obtained by the **CONDITIONAL MEAN**

$$\mathbb{E} [\mathbf{X}_n | \mathbf{Y}_{(n)} = \mathbf{y}_{(n)}] = \int_{\mathbb{X}} \mathbf{x}_n f_{\mathbf{X}_n | \mathbf{Y}_{(n)}}(\mathbf{x}_n | \mathbf{y}_{(n)}) d\mathbf{x}_n, \quad (10.3)$$

which can be derived from the **CONDITIONAL PDF**

$$f_{\mathbf{X}_n | \mathbf{Y}_{(n)}}(\mathbf{x}_n | \mathbf{y}_{(n)}). \quad (10.4)$$

Recall that the **CONDITIONAL PDF** of any the state of the Gauß-Markov sequence is uniquely characterized by the first and second order conditional moments

$$\boldsymbol{\mu}_{\mathbf{X}_n | \mathbf{Y}_{(n)} = \mathbf{y}_{(n)}} = \mathbb{E} [\mathbf{X}_n | \mathbf{Y}_{(n)} = \mathbf{y}_{(n)}], \quad (10.5)$$

$$\mathbf{C}_{\mathbf{X}_n | \mathbf{Y}_{(n)} = \mathbf{y}_{(n)}} = \mathbb{E} \left[\left(\mathbf{X}_n - \boldsymbol{\mu}_{\mathbf{X}_n | \mathbf{Y}_{(n)} = \mathbf{y}_{(n)}} \right) \left(\mathbf{X}_n - \boldsymbol{\mu}_{\mathbf{X}_n | \mathbf{Y}_{(n)} = \mathbf{y}_{(n)}} \right)^{\top} \middle| \mathbf{Y}_{(n)} = \mathbf{y}_{(n)} \right]. \quad (10.6)$$

Note.

$\mathbf{C}_{\mathbf{X}_n | \mathbf{Y}_{(n)} = \mathbf{y}_{(n)}} = \mathbb{E} \left[\left(\mathbf{X}_n - \boldsymbol{\mu}_{\mathbf{X}_n | \mathbf{Y}_{(n)} = \mathbf{y}_{(n)}} \right) \left(\mathbf{X}_n - \boldsymbol{\mu}_{\mathbf{X}_n | \mathbf{Y}_{(n)} = \mathbf{y}_{(n)}} \right)^{\top} \middle| \mathbf{Y}_{(n)} = \mathbf{y}_{(n)} \right]$ does not depend on the actual realization $\mathbf{y}_{(n)}$, since all random variables are Gaussian distributed.

The same holds for $\mathbf{C}_{\mathbf{X}_n | \mathbf{Y}_{(n-1)} = \mathbf{y}_{(n-1)}}$.

Short Notations.

- In the following, $\hat{\mathbf{x}}_{n|n}$ and $\hat{\mathbf{x}}_{n|n-1}$ are used as short notations for $\boldsymbol{\mu}_{\mathbf{X}_n|\mathbf{Y}_{(n)}=\mathbf{y}_{(n)}}$ and $\boldsymbol{\mu}_{\mathbf{X}_n|\mathbf{Y}_{(n-1)}=\mathbf{y}_{(n-1)}}$, i.e.,

$$\begin{aligned}\hat{\mathbf{x}}_{n|n} &= \mathbb{E} [\mathbf{X}_n | \mathbf{Y}_{(n)} = \mathbf{y}_{(n)}], \\ \hat{\mathbf{x}}_{n|n-1} &= \mathbb{E} [\mathbf{X}_n | \mathbf{Y}_{(n-1)} = \mathbf{y}_{(n-1)}].\end{aligned}$$

- The same holds for $\mathbf{C}_{\mathbf{X}_{n|n}}$ and $\mathbf{C}_{\mathbf{X}_{n|n-1}}$, which are short notations for

$$\begin{aligned}\mathbf{C}_{\mathbf{X}_{n|n}} &= \mathbb{E} \left[(\mathbf{X}_n - \hat{\mathbf{x}}_{n|n}) (\mathbf{X}_n - \hat{\mathbf{x}}_{n|n})^\top \mid \mathbf{Y}_{(n)} = \mathbf{y}_{(n)} \right], \\ \mathbf{C}_{\mathbf{X}_{n|n-1}} &= \mathbb{E} \left[(\mathbf{X}_n - \hat{\mathbf{x}}_{n|n-1}) (\mathbf{X}_n - \hat{\mathbf{x}}_{n|n-1})^\top \mid \mathbf{Y}_{(n-1)} = \mathbf{y}_{(n-1)} \right].\end{aligned}$$

10.1 Recursive Computation of Conditional Means and Covariances

- a) Exploiting the CHAPMAN-KOLMOGOROV EQUATION (9.11) and the LINEAR STATE SPACE MODEL in Eq. (10.1), we obtain

$$\begin{aligned}\mu_{\mathbf{x}_n|\mathbf{y}_{(n-1)}=\mathbf{y}_{(n-1)}} &= \int_{\mathbb{X}} \underbrace{\int_{\mathbb{X}} \mathbf{x}_n f_{\mathbf{x}_n|\mathbf{x}_{n-1}}(\mathbf{x}_n|\mathbf{x}_{n-1}) d\mathbf{x}_n}_{\text{cond. expectation of } \mathbf{x}_n \text{ (refer to 10.1)}} f_{\mathbf{x}_{n-1}|\mathbf{y}_{(n-1)}}(\mathbf{x}_{n-1}|\mathbf{y}_{(n-1)}) d\mathbf{x}_{n-1} \\ &= \mathbf{G}_n \mu_{\mathbf{x}_{n-1}|\mathbf{y}_{(n-1)}=\mathbf{y}_{(n-1)}},\end{aligned}$$

and thus the STATE PREDICTION is given by

$$\hat{\mathbf{x}}_{n|n-1} = \mathbf{G}_n \hat{\mathbf{x}}_{n-1|n-1}. \quad (10.7)$$

The respective CONDITIONAL STATE COVARIANCE MATRIX of \mathbf{X}_n is obtained as

$$\mathbf{C}_{\mathbf{x}_{n|n-1}} = \mathbf{G}_n \mathbf{C}_{\mathbf{x}_{n-1|n-1}} \mathbf{G}_n^\top + \mathbf{C}_{\mathbf{v}_n}. \quad (10.8)$$

b) Since both sides of the BAYES RULE in (9.12) are Gaussian PDFs, a comparison of the exponents,²

$$\log \left(f_{\mathbf{x}_n | \mathbf{y}_{(n)}}(\mathbf{x}_n | \mathbf{y}_{(n)}) \right) = \log \left(f_{\mathbf{y}_n | \mathbf{x}_n}(\mathbf{y}_n | \mathbf{x}_n) \right) + \log \left(f_{\mathbf{x}_n | \mathbf{y}_{(n-1)}}(\mathbf{x}_n | \mathbf{y}_{(n-1)}) \right) + \dots,$$

and taking into account the LINEAR OBSERVATION MODEL in Eq. (10.2) results in

$$\hat{\mathbf{x}}_{n|n} = \hat{\mathbf{x}}_{n|n-1} + \mathbf{K}_n (\mathbf{y}_n - \mathbf{H}_n \hat{\mathbf{x}}_{n|n-1}), \quad (10.9)$$

$$\mathbf{C} \mathbf{x}_{n|n} = \mathbf{C} \mathbf{x}_{n|n-1} - \mathbf{K}_n \mathbf{H}_n \mathbf{C} \mathbf{x}_{n|n-1}, \quad (10.10)$$

with the KALMAN GAIN MATRIX

$$\mathbf{K}_n = \mathbf{C} \mathbf{x}_{n|n-1} \mathbf{H}_n^\top \left(\mathbf{H}_n \mathbf{C} \mathbf{x}_{n|n-1} \mathbf{H}_n^\top + \mathbf{C} \mathbf{w}_n \right)^{-1}. \quad (10.11)$$

²Note, that we are using (9.12) instead of (9.13).

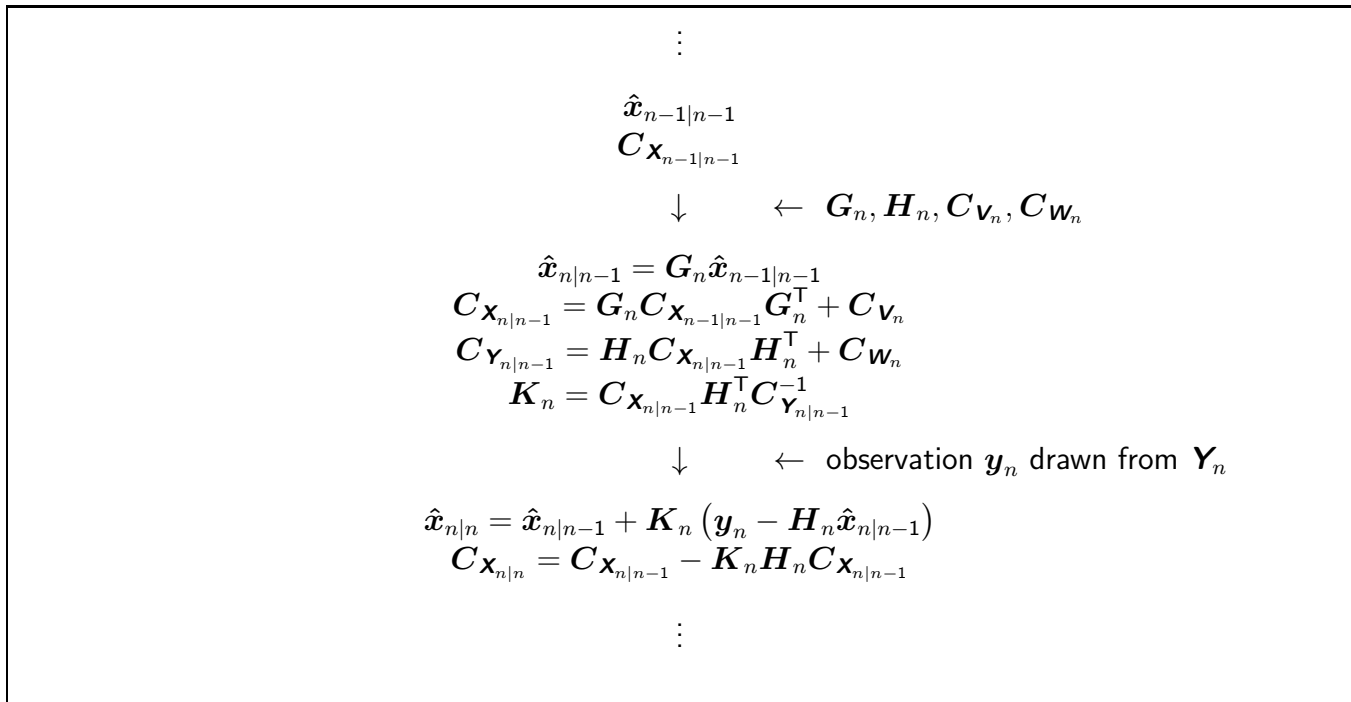


Fig. 10.2: The recursive STATE TRACKING procedure. In many applications, the system parameter matrices $\mathbf{G}_n, \mathbf{H}_n, \mathbf{C}_{\mathbf{v}_n}, \mathbf{C}_{\mathbf{w}_n}$ are assumed to be constant and we instead consider \mathbf{G}, \mathbf{H} and $\mathbf{C}_{\mathbf{v}}, \mathbf{C}_{\mathbf{w}}$.

10.2 Proof

A proof for a scalar version is presented. See the appendix for the multivariate version.

a) Prediction step:

$$\begin{aligned}
 \boxed{\hat{x}_{n|n-1}} &= \mathbb{E} [X_n \mid Y_{(n-1)} = y_{(n-1)}] \\
 &= \int_{\mathbb{X}} x_n f_{X_n|Y_{(n-1)}}(x_n \mid y_{(n-1)}) \mathrm{d}x_n \\
 &= \int_{\mathbb{X}} \int_{\mathbb{X}} x_n f_{X_n|X_{n-1}}(x_n|x_{n-1}) \mathrm{d}x_n f_{X_{n-1}|Y_{(n-1)}}(x_{n-1}|y_{(n-1)}) \mathrm{d}x_{n-1} \\
 &= \int_{\mathbb{X}} \mathbb{E} [X_n \mid X_{n-1} = x_{n-1}] f_{X_{n-1}|Y_{(n-1)}}(x_{n-1}|y_{(n-1)}) \mathrm{d}x_{n-1} \\
 &= \int_{\mathbb{X}} \mathbb{E} [g_n X_{n-1} + V_n \mid X_{n-1} = x_{n-1}] f_{X_{n-1}|Y_{(n-1)}}(x_{n-1}|y_{(n-1)}) \mathrm{d}x_{n-1} \\
 &= g_n \int_{\mathbb{X}} x_{n-1} f_{X_{n-1}|Y_{(n-1)}}(x_{n-1}|y_{(n-1)}) \mathrm{d}x_{n-1} = g_n \mathbb{E} [X_{n-1} \mid Y_{(n-1)} = y_{(n-1)}] = g_n \boxed{\hat{x}_{n-1|n-1}},
 \end{aligned}$$

$$\begin{aligned}
 \boxed{\sigma_{X_{n|n-1}}^2} &= \mathbb{E} [(X_n - \hat{x}_{n|n-1})^2 \mid Y_{(n-1)} = y_{(n-1)}] \\
 &= \mathbb{E} [(g_n X_{n-1} + V_n - g_n \hat{x}_{n-1|n-1})^2 \mid Y_{(n-1)} = y_{(n-1)}] \\
 &= g_n^2 \mathbb{E} [(X_{n-1} - \hat{x}_{n-1|n-1})^2 \mid Y_{(n-1)} = y_{(n-1)}] + \mathbb{E} [V_n^2 \mid Y_{(n-1)} = y_{(n-1)}] = g_n^2 \boxed{\sigma_{X_{n-1|n-1}}^2} + \sigma_{V_n}^2.
 \end{aligned}$$

b) Correction step:

$$\begin{aligned}
\log f_{X_n|Y_{(n)}}(x_n|y_{(n)}) &= \log f_{Y_n|X_n}(y_n | x_n) + \log f_{X_n|Y_{(n-1)}}(x_n|y_{(n-1)}) + \dots \\
&= -\frac{(y_n - \mathbb{E}[Y_n | X_n = x_n])^2}{2 \text{Var}[Y_n | X_n = x_n]} - \frac{(x_n - \hat{x}_{n|n-1})^2}{2\sigma_{X_{n|n-1}}^2} + \dots \\
&= -\frac{(y_n - h_n x_n)^2}{2\sigma_{W_n}^2} - \frac{(x_n - \hat{x}_{n|n-1})^2}{2\sigma_{X_{n|n-1}}^2} + \dots \\
&= -\frac{\sigma_{X_{n|n-1}}^2 (y_n^2 - 2y_n h_n x_n + h_n^2 x_n^2) + \sigma_{W_n}^2 (x_n^2 - 2x_n \hat{x}_{n|n-1} + \hat{x}_{n|n-1}^2)}{2\sigma_{W_n}^2 \sigma_{X_{n|n-1}}^2} + \dots \\
&= -\frac{(\sigma_{X_{n|n-1}}^2 h_n^2 + \sigma_{W_n}^2) x_n^2 - 2(\sigma_{X_{n|n-1}}^2 y_n h_n + \sigma_{W_n}^2 \hat{x}_{n|n-1}) x_n}{2\sigma_{W_n}^2 \sigma_{X_{n|n-1}}^2} + \dots \\
&= -\frac{\sigma_{X_{n|n-1}}^2 h_n^2 + \sigma_{W_n}^2}{2\sigma_{W_n}^2 \sigma_{X_{n|n-1}}^2} \left(x_n - \frac{\sigma_{X_{n|n-1}}^2 y_n h_n + \sigma_{W_n}^2 \hat{x}_{n|n-1}}{\sigma_{X_{n|n-1}}^2 h_n^2 + \sigma_{W_n}^2} \right)^2 + \dots \\
&= -\frac{(x_n - \boxed{\hat{x}_{n|n}})^2}{2\boxed{\sigma_{X_{n|n}}^2}} + \dots
\end{aligned}$$

Comparison:

$$\begin{aligned}
\boxed{\hat{x}_{n|n}} &= \frac{\sigma_{X_{n|n-1}}^2 y_n h_n + \sigma_{W_n}^2 \hat{x}_{n|n-1}}{\sigma_{X_{n|n-1}}^2 h_n^2 + \sigma_{W_n}^2} \\
&= \hat{x}_{n|n-1} + \frac{\sigma_{X_{n|n-1}}^2 y_n h_n + \sigma_{W_n}^2 \hat{x}_{n|n-1} - (\sigma_{X_{n|n-1}}^2 h_n^2 + \sigma_{W_n}^2) \hat{x}_{n|n-1}}{\sigma_{X_{n|n-1}}^2 h_n^2 + \sigma_{W_n}^2} \\
&= \hat{x}_{n|n-1} + \frac{\sigma_{X_{n|n-1}}^2 y_n h_n - \sigma_{X_{n|n-1}}^2 h_n^2 \hat{x}_{n|n-1}}{\sigma_{X_{n|n-1}}^2 h_n^2 + \sigma_{W_n}^2} \\
&= \hat{x}_{n|n-1} + \frac{\sigma_{X_{n|n-1}}^2 h_n}{\sigma_{X_{n|n-1}}^2 h_n^2 + \sigma_{W_n}^2} (y_n - h_n \hat{x}_{n|n-1}) \\
&= \boxed{\hat{x}_{n|n-1}} + \boxed{k_n} (y_n - h_n \boxed{\hat{x}_{n|n-1}}),
\end{aligned}$$

$$\boxed{\sigma_{X_{n|n}}^2} = \frac{\sigma_{W_n}^2 \boxed{\sigma_{X_{n|n-1}}^2}}{\boxed{\sigma_{X_{n|n-1}}^2} h_n^2 + \sigma_{W_n}^2} = \frac{\sigma_{W_n}^2 \boxed{\sigma_{X_{n|n-1}}^2}}{\boxed{\sigma_{Y_{n|n-1}}^2}},$$

$$\text{with } \boxed{k_n} = \frac{\boxed{\sigma_{X_{n|n-1}}^2} h_n}{\boxed{\sigma_{X_{n|n-1}}^2} h_n^2 + \sigma_{W_n}^2} = \frac{\boxed{\sigma_{X_{n|n-1}}^2} h_n}{\boxed{\sigma_{Y_{n|n-1}}^2}} \text{ and } \boxed{\sigma_{Y_{n|n-1}}^2} = \boxed{\sigma_{X_{n|n-1}}^2} h_n^2 + \sigma_{W_n}^2.$$

10.3 Discussion

Innovation Sequence.

The sequence of random variables $\Delta \mathbf{Y}_n = \mathbf{Y}_n - \hat{\mathbf{Y}}_{n|n-1}$ in Eq. (10.9) forms an i.i.d. random sequence which can be constructed by a bijective linear transform of the random observation sequence \mathbf{Y}_n ,

$$\begin{aligned}
 \Delta \mathbf{Y}_n &= \mathbf{Y}_n - \hat{\mathbf{Y}}_{n|n-1} \\
 &= \mathbf{Y}_n - \mathbb{E} [\mathbf{Y}_n | \mathbf{Y}_{(n-1)} = \mathbf{y}_{(n-1)}] \\
 &= \mathbf{Y}_n - \mathbb{E} [\mathbf{H}_n \mathbf{X}_n + \mathbf{W}_n | \mathbf{Y}_{(n-1)} = \mathbf{y}_{(n-1)}] \\
 &= \mathbf{Y}_n - \mathbf{H}_n \underbrace{\mathbb{E} [\mathbf{X}_n | \mathbf{Y}_{(n-1)} = \mathbf{y}_{(n-1)}]}_{\text{linear in } \mathbf{Y}_{(n-1)}} \\
 &= \mathbf{Y}_n - \mathbf{H}_n \hat{\mathbf{x}}_{n|n-1},
 \end{aligned} \tag{10.12}$$

and thus fulfills the requirements of an INNOVATION SEQUENCE according to the observations \mathbf{Y}_n .

An INNOVATION SEQUENCE consists of i.i.d. elements and is a bijective linear function of the observations.

Predictor-Corrector-Step.

Since the innovation variable is stochastically independent of previous observations of $\mathbf{Y}_{(n-1)}$, Eq. (10.9) can be interpreted as the sum of two stochastically independent contributions for the state estimate $\hat{\mathbf{x}}_{n|n}$.

Another interpretation is to consider the state estimate $\hat{\mathbf{x}}_{n|n}$ as the accumulative contribution of a PREDICTOR STEP and a CORRECTOR STEP.

$$\hat{\mathbf{x}}_{n|n} = \underbrace{\hat{\mathbf{x}}_{n|n-1}}_{\text{estimation step: } \mathbb{E}[\mathbf{X}_n | \mathbf{Y}_{(n-1)} = \mathbf{y}_{(n-1)}]} + \underbrace{\overbrace{\mathbf{K}_n (\mathbf{y}_n - \mathbf{H}_n \hat{\mathbf{x}}_{n|n-1})}^{\text{correction step: } \mathbb{E}[\mathbf{X}_n | \Delta \mathbf{Y}_n = \Delta \mathbf{y}_{(n)}]}}_{\text{innovation: } \Delta \mathbf{y}_n} . \quad (10.13)$$

Kalman Gain Matrix.

The KALMAN GAIN MATRIX K_n in Eq. (10.9) represents the LINEAR PREDICTOR for the state variable X_n based on the single innovation variable $Y_n - H_n \hat{x}_{n|n-1}$.

$$\begin{aligned}
 K_n &= E [X_n \Delta Y_n^T] E [\Delta Y_n \Delta Y_n^T]^{-1} \\
 &= E [X_n (Y_n - H_n \hat{x}_{n|n-1})^T] E [(Y_n - H_n \hat{x}_{n|n-1})(Y_n - H_n \hat{x}_{n|n-1})^T]^{-1} \\
 &= E [X_n (H_n X_n - H_n \hat{x}_{n|n-1})^T] E [(Y_n - H_n \hat{x}_{n|n-1})(Y_n - H_n \hat{x}_{n|n-1})^T]^{-1} \\
 &= E [X_n (X_n - \hat{x}_{n|n-1})^T] H_n^T C_{Y_{n|n-1}}^{-1} \\
 &= E [(X_n - \hat{x}_{n|n-1})(X_n - \hat{x}_{n|n-1})^T] H_n^T C_{Y_{n|n-1}}^{-1} = C_{X_{n|n-1}} H_n^T C_{Y_{n|n-1}}^{-1}. \tag{10.14}
 \end{aligned}$$

The single steps of the derivation:

- from the 1st to 2nd line by the definition of the innovation sequence,
- from the 2nd to 3rd line by the stochastic independence between X_n and W_n ,
- from the 3rd to 4th line by the definition of the conditional covariance of the observation variable,
- and from the 4th to 5th line by the orthogonality principle and the definition of the conditional covariance of the state variable.

Innovation Covariance Matrix.

The INNOVATION COVARIANCE MATRIX is derived according to the definition of covariance matrices as

$$\begin{aligned} C_{\Delta \mathbf{Y}_{n|n-1}} &= E \left[(\mathbf{Y}_n - \mathbf{H}_n \hat{\mathbf{x}}_{n|n-1}) (\mathbf{Y}_n - \mathbf{H}_n \hat{\mathbf{x}}_{n|n-1})^\top \mid \mathbf{Y}_{(n-1)} = \mathbf{y}_{(n-1)} \right] \\ &= E \left[(\mathbf{H}_n \mathbf{X}_n + \mathbf{W}_n - \mathbf{H}_n \hat{\mathbf{x}}_{n|n-1}) (\mathbf{H}_n \mathbf{X}_n + \mathbf{W}_n - \mathbf{H}_n \hat{\mathbf{x}}_{n|n-1})^\top \mid \mathbf{Y}_{(n-1)} = \mathbf{y}_{(n-1)} \right] \\ &= E \left[(\mathbf{H}_n (\mathbf{X}_n - \hat{\mathbf{x}}_{n|n-1}) + \mathbf{W}_n) (\mathbf{H}_n (\mathbf{X}_n - \hat{\mathbf{x}}_{n|n-1}) + \mathbf{W}_n)^\top \mid \mathbf{Y}_{(n-1)} = \mathbf{y}_{(n-1)} \right] \\ &= \mathbf{H}_n C_{\mathbf{X}_{n|n-1}} \mathbf{H}_n^\top + C_{\mathbf{W}_n}. \end{aligned} \tag{10.15}$$

Note. $C_{\mathbf{Y}_{n|n-1}} = C_{\Delta \mathbf{Y}_{n|n-1}}$, since the random variable $\Delta \mathbf{Y}_{n|n-1}$ is zero mean:

$$\begin{aligned} E [\Delta \mathbf{Y}_n \mid \mathbf{Y}_{(n-1)} = \mathbf{y}_{(n-1)}] &= E [\mathbf{Y}_n - \hat{\mathbf{Y}}_{n|n-1} \mid \mathbf{Y}_{(n-1)} = \mathbf{y}_{(n-1)}] \\ &= E [\mathbf{Y}_n \mid \mathbf{Y}_{(n-1)} = \mathbf{y}_{(n-1)}] - \hat{\mathbf{Y}}_{n|n-1} \\ &= \mathbf{0}. \end{aligned}$$

Predicted State Covariance Matrix.

The covariance matrix of the predicted state $\mathbf{X}_{n|n-1}$ or the PREDICTION ERROR COVARIANCE MATRIX can be derived by

$$\begin{aligned} C_{\mathbf{X}_{n|n-1}} &= \mathbb{E} \left[(\mathbf{X}_n - \hat{\mathbf{x}}_{n|n-1}) (\mathbf{X}_n - \hat{\mathbf{x}}_{n|n-1})^\top \middle| \mathbf{Y}_{(n-1)} = \mathbf{y}_{(n-1)} \right] \\ &= \mathbb{E} \left[(\mathbf{G}_n \mathbf{X}_{n-1} + \mathbf{V}_n - \mathbf{G}_n \hat{\mathbf{x}}_{n-1|n-1}) (\mathbf{G}_n \mathbf{X}_{n-1} + \mathbf{V}_n - \mathbf{G}_n \hat{\mathbf{x}}_{n-1|n-1})^\top \middle| \mathbf{Y}_{(n-1)} = \mathbf{y}_{(n-1)} \right] \\ &= \mathbb{E} \left[(\mathbf{G}_n (\mathbf{X}_{n-1} - \hat{\mathbf{x}}_{n-1|n-1}) + \mathbf{V}_n) (\mathbf{G}_n (\mathbf{X}_{n-1} - \hat{\mathbf{x}}_{n-1|n-1}) + \mathbf{V}_n)^\top \middle| \mathbf{Y}_{(n-1)} = \mathbf{y}_{(n-1)} \right] \\ &= \mathbf{G}_n C_{\mathbf{X}_{n-1|n-1}} \mathbf{G}_n^\top + C_{\mathbf{V}_n}. \end{aligned} \tag{10.16}$$

Filtered State Covariance Matrix.

The resulting FILTERED STATE COVARIANCE MATRIX is analogously given by

$$\begin{aligned} C_{\mathbf{x}_{n|n}} &= E \left[(\mathbf{X}_n - \hat{\mathbf{x}}_{n|n}) (\mathbf{X}_n - \hat{\mathbf{x}}_{n|n})^\top \middle| \mathbf{Y}_{(n)} = \mathbf{y}_{(n)} \right] \\ &= E \left[(\mathbf{X}_n - \hat{\mathbf{x}}_{n|n}) (\mathbf{X}_n - \hat{\mathbf{x}}_{n|n})^\top \right] \\ &= E \left[(\mathbf{X}_n - \hat{\mathbf{x}}_{n|n-1} - \mathbf{K}_n \Delta \mathbf{Y}_n) (\mathbf{X}_n - \hat{\mathbf{x}}_{n|n-1} - \mathbf{K}_n \Delta \mathbf{Y}_n)^\top \right] \\ &= C_{\mathbf{x}_{n|n-1}} + \mathbf{K}_n C_{\mathbf{y}_{n|n-1}} \mathbf{K}_n^\top - 2\mathbf{K}_n E \left[\Delta \mathbf{Y}_n (\mathbf{X}_n - \hat{\mathbf{x}}_{n|n-1})^\top \right] \\ &= C_{\mathbf{x}_{n|n-1}} + \mathbf{K}_n C_{\mathbf{y}_{n|n-1}} \mathbf{K}_n^\top - 2\mathbf{K}_n \mathbf{H}_n E \left[(\mathbf{X}_n - \hat{\mathbf{x}}_{n|n-1}) (\mathbf{X}_n - \hat{\mathbf{x}}_{n|n-1})^\top \right] \\ &= C_{\mathbf{x}_{n|n-1}} + \mathbf{K}_n C_{\mathbf{y}_{n|n-1}} \mathbf{K}_n^\top - 2\mathbf{K}_n \mathbf{H}_n C_{\mathbf{x}_{n|n-1}} \\ &= C_{\mathbf{x}_{n|n-1}} + \mathbf{K}_n C_{\mathbf{y}_{n|n-1}} C_{\mathbf{y}_{n|n-1}}^{-1} \mathbf{H}_n C_{\mathbf{x}_{n|n-1}} - 2\mathbf{K}_n \mathbf{H}_n C_{\mathbf{x}_{n|n-1}} \\ &= C_{\mathbf{x}_{n|n-1}} - \mathbf{K}_n \mathbf{H}_n C_{\mathbf{x}_{n|n-1}}. \end{aligned} \tag{10.17}$$

11. Two Examples

11.1 Unknown Parameter Estimation

The Kalman Filter approach can be applied for the recursive estimation of a deterministic unknown parameter x by viewing the parameter as a non-changing state variable $\mathbf{X}_n = \mathbf{X}_{n-1}$, where only the state observations $\mathbf{Y}_n = \mathbf{H}_n \mathbf{X}_n + \mathbf{W}_n$ appear as a random sequence.

For the sake of simplicity, we further assume univariate random variables, which leads to a state space formulation as

$$X_n = X_{n-1} \tag{11.1}$$

$$Y_n = h_n X_n + W_n, \quad \text{i.i.d. } W_n \sim \mathcal{N}(0, \sigma_{W_n}^2). \tag{11.2}$$

The respective update rules of the Kalman Filter are

$$\hat{x}_{n|n-1} = \hat{x}_{n-1|n-1} \quad (11.3)$$

$$\sigma_{X_{n|n-1}}^2 = \sigma_{X_{n-1|n-1}}^2 \quad (11.4)$$

$$\sigma_{Y_{n|n-1}}^2 = h_n^2 \sigma_{X_{n|n-1}}^2 + \sigma_{W_n}^2 \quad (11.5)$$

$$k_n = h_n \frac{\sigma_{X_{n|n-1}}^2}{\sigma_{Y_{n|n-1}}^2} \quad (11.6)$$

$$\hat{x}_{n|n} = \hat{x}_{n|n-1} + k_n(y_n - h_n \hat{x}_{n|n-1}) \quad (11.7)$$

$$\sigma_{X_{n|n}}^2 = \sigma_{X_{n|n-1}}^2 - k_n h_n \sigma_{X_{n|n-1}}^2. \quad (11.8)$$

which can be derived from the general recursive state tracking procedure of the Kalman Filter.

For the numerical analysis, we assume

- $x = 1$ for the parameter to be estimated,
- the initial conditions $x_1 = 0$ and $\sigma_{X_{1|1}}^2 = 1$,
- and the parameters $\sigma_{W_n}^2 \equiv 0.01$, $h_n = \sin(\pi n/100)$ for all $n = 1, \dots, 100$.

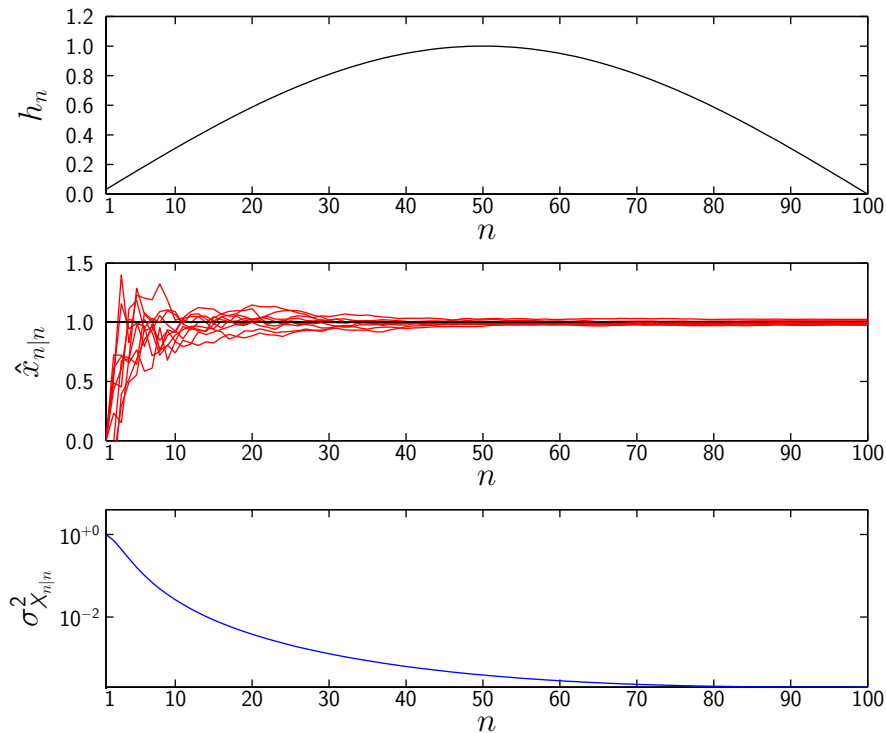


Fig. 11.1: Stochastic paths $\hat{x}_{n|n}$ of the estimation of a constant parameter $x = 1$ for 10 uniformly chosen noise sequences W_n and $Y_n = \sin(\pi n/100) + W_n$, $\forall n = 1, \dots, 100$, based on the Kalman Filter.

11.2 Tracking Example

The discretization of a stochastically excited dynamical system¹

$$\ddot{X}_t = V_t, \quad \text{i.i.d. } V_n \sim \mathcal{N}(0, \sigma_{V_n}^2), \quad (11.9)$$

with an equidistant sampling interval of T , yields the discrete linear system²

$$X_n = X_{n-1} + T\dot{X}_{n-1} \quad (11.10)$$

$$\dot{X}_n = \dot{X}_{n-1} + TV_n. \quad (11.11)$$

When applying a Kalman Filter for tracking the state space variable X_n , based on observations

$$Y_n = X_n + W_n, \quad \text{i.i.d. } W_n \sim \mathcal{N}(0, \sigma_{W_n}^2), \quad (11.12)$$

and taking into account Eqs. (10.1)-(10.2), we obtain

$$\mathbf{X}_n = \mathbf{G}_n \mathbf{X}_{n-1} + \begin{bmatrix} 0 \\ T \end{bmatrix} V_n, \quad (11.13)$$

$$Y_n = \mathbf{h}_n^\top \mathbf{X}_n + W_n, \quad (11.14)$$

with $\mathbf{X}_n = [X_n \ \dot{X}_n]^\top$ and

$$\mathbf{G}_n = \begin{bmatrix} 1 & T \\ 0 & 1 \end{bmatrix}, \quad \mathbf{h}_n = \begin{bmatrix} 1 \\ 0 \end{bmatrix}. \quad (11.15)$$

¹Eq. (11.9) denotes a STOCHASTICAL DIFFERENTIAL EQUATION (SDE).

²The discretization is based on $\dot{X}_n - \dot{X}_{n-1} = T\ddot{X}_{n-1}$ and $X_n - X_{n-1} = T\dot{X}_{n-1}$.

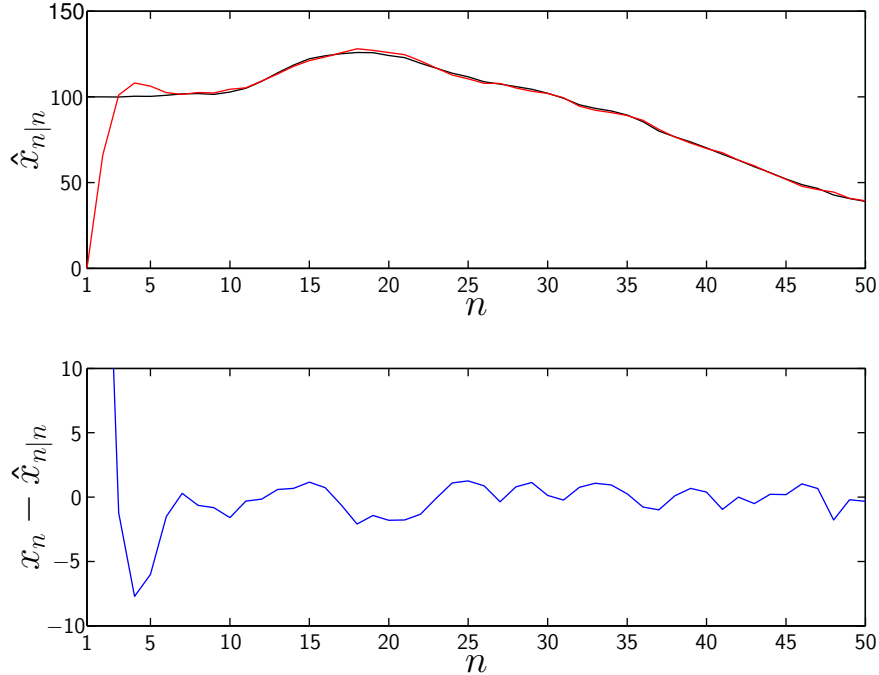


Fig. 11.2: Estimation of the trajectory $\hat{x}_{n|n}$, $n = 1, \dots, 50$, of one single stochastic path from the SDE $X_n = X_{n-1} + \dot{X}_{n-1}$ and $\dot{X}_n = \dot{X}_{n-1} + V_{n-1}$ for the initial conditions $x_0 = 100$ and $\dot{x}_0 = 0$. The Kalman Filter is based on observations $Y_n = X_n + W_n$ and initial conditions $\hat{x}_{0|0} = \mathbf{0}$ and $\mathbf{C}_{0|0} = \mathbf{I}$. The parameter $T = 1$.

12. Particle Filter

The KALMAN FILTER is a powerful linear approach for estimating the states of a random sequence, otherwise it only shows optimum performance for GAUSS-MARKOV SEQUENCES, this is if there is a linear relation between states of the sequence and if there is an additional Gaussian distortion at most.

For the general state-space model,

$$\mathbf{X}_n = g_n(\mathbf{X}_{n-1}, \mathbf{V}_n), \quad (12.1)$$

$$\mathbf{Y}_n = h_n(\mathbf{X}_n, \mathbf{W}_n), \quad (12.2)$$

where g_n and h_n are nonlinear functions and \mathbf{V}_n and \mathbf{W}_n are non-Gaussian distributed random variables, alternative techniques beyond the Kalman Filter are required.

The general approach is to substitute the linear Kalman Filter by SUBOPTIMUM NONLINEAR FILTERS.

The three most well known nonlinear alternatives to the KALMAN FILTER are

- the EXTENDED KALMAN FILTER,
- the UNSCENTED KALMAN FILTER,
- and the PARTICLE FILTER.

a) The EXTENDED KALMAN FILTER is typically applied to state-space models like

$$\mathbf{X}_n = g_n(\mathbf{X}_{n-1}) + \mathbf{V}_n, \quad (12.3)$$

$$\mathbf{Y}_n = h_n(\mathbf{X}_n) + \mathbf{W}_n, \quad (12.4)$$

where \mathbf{V}_n and \mathbf{W}_n are or are assumed to be Gaussian distributed random variables.

The recursive estimation of the state-space variables in a standard KALMAN FILTER type fashion is based on a linear approximation of the nonlinear functions g_{n-1} and h_n in (12.3) and (12.4). The approximation must be performed in each step of the recursive estimation process. Obviously, the EXTENDED KALMAN FILTER in general cannot achieve optimal performance.

- b) The so-called **UNSCENTED KALMAN FILTER** approximates the desired conditional PDF $f_{\mathbf{x}_n|\mathbf{y}_{(n)}}(\mathbf{x}_n|\mathbf{y}_{(n)})$ by a Gaussian PDF of the original recursive computation process.

The recursive estimation process again resembles the standard **KALMAN FILTER** process, but replacing the estimation of the conditional state covariance matrix in each step by a weighted sample covariance matrix of the non-linear transformed samples of the state-space variables.

- c) Whereas the **EXTENDED KALMAN FILTER** and the **UNSCENTED KALMAN FILTER** strongly depend on the validity of the Gaussian assumption, the **PARTICLE FILTER** approach completely breaks with these requirements. It is purely based on the fundamental update rule of all tracking approaches, namely the recursive computation of the posterior conditional PDF.

$$\begin{array}{c}
 \text{desired conditional PDF in the } n\text{th step} \\
 \overbrace{f_{\mathbf{x}_n|\mathbf{y}_{(n)}}(\mathbf{x}_n|\mathbf{y}_{(n)})} \quad \propto \\
 \underbrace{f_{\mathbf{y}_n|\mathbf{x}_n}(\mathbf{y}_n|\mathbf{x}_n)}_{\text{likelihood of the new observation}} \times \int_{\mathbb{X}} \underbrace{f_{\mathbf{x}_n|\mathbf{x}_{n-1}}(\mathbf{x}_n|\mathbf{x}_{n-1})}_{\text{state-transition PDF}} \underbrace{f_{\mathbf{x}_{n-1}|\mathbf{y}_{(n-1)}}(\mathbf{x}_{n-1}|\mathbf{y}_{(n-1)})}_{\text{conditional PDF from the } (n-1)\text{-th-step before}} d\mathbf{x}_{n-1}.
 \end{array} \quad (12.5)$$

It remains the task to compute the integral of the recursive update of the posterior conditional PDF!

12.1 Monte Carlo Integration

Since the integration of (9.13) is in general impossible, the respective integral is solved via MONTE CARLO INTEGRATION. MONTE CARLO INTEGRATION means to replace an integral by its sample mean

$$I = \int_{\mathbb{X}} g(\xi) f_X(\xi) d\xi \approx I_N = \frac{1}{N} \sum_{i=1}^N g(x^i), \quad (12.6)$$

where the samples x^i , $i = 1, \dots, N$, are drawn according to the PDF $f_X(x)$, i.e., $x^i \sim f_X(x)$.

Theorem.

If the samples x^i are stochastically independent, the approximation I_N is an unbiased estimate of I with

$$\lim_{N \rightarrow \infty} I_N \rightarrow I, \quad (12.7)$$

and

$$\lim_{N \rightarrow \infty} \sqrt{N}(I_N - I) \sim \mathcal{N}(0, \sigma^2), \quad (12.8)$$

with $\sigma^2 = \text{Var}[g(X)]$.

Proof.

The proof is based on the law of large numbers and the fact that $I = \mathbb{E}[g(X)]$.

Importance Sampling.

Assume – as usual in the PARTICLE FILTER framework – instead of using the PDF $f_X(x)$ for drawing samples, we use an IMPORTANCE DENSITY $q_X(x)$ with $q_X(x) > 0$ if $f_X(x) > 0$ for all x .¹ In this case, the correct MONTE CARLO INTEGRATION can be still guaranteed by appropriate weights for the samples,

$$I = \int_{\mathbb{X}} g(\xi) f_X(\xi) d\xi = \int_{\mathbb{X}} g(\xi) \frac{f_X(\xi)}{q_X(\xi)} q_X(\xi) d\xi \approx I_N = \frac{1}{N} \sum_{i=1}^N g(x^i) \frac{f_X(x^i)}{q_X(x^i)}, \quad (12.9)$$

where the samples x^i , $i = 1, \dots, N$, are drawn according to the PDF $q_X(x)$, i.e.,

$$I \approx I_N = \frac{1}{N} \sum_{i=1}^N \tilde{w}^i g(x^i), \quad (12.10)$$

with $\tilde{w}^i \triangleq \frac{f_X(x^i)}{q_X(x^i)}$ and $x^i \sim q_X(x)$.

¹This might be motivated by several arguments, e.g., drawing samples from $q_X(x)$ might be simpler than drawing samples from $f_X(x)$.

Importance Sampling if $\int_{\mathbb{X}} f_X(\xi) \, d\xi \neq 1$.

If $f_X(x)$ fulfills the properties of a PDF except a normalization constant, the IMPORTANCE DENSITY results in

$$I \approx I_N = \sum_{i=1}^N w^i g(x^i), \quad (12.11)$$

with normalized weights

$$w^i \triangleq \frac{\tilde{w}^i}{\sum_{i=1}^N \tilde{w}^i}, \quad \tilde{w}^i \triangleq \frac{f_X(x^i)}{q_X(x^i)} \quad \text{and} \quad x^i \sim q_X(x).$$

Note.

The MONTE CARLO INTEGRAL I_N can be interpreted as the true expected value of $g(X)$ with respect to the approximate PDF²

$$\sum_{i=1}^N w^i \delta(x - x^i) \approx f_X(x). \quad (12.12)$$

² $\mathbb{E}[g(X)] = \int_{\mathbb{X}} g(\xi) \sum_{i=1}^N w^i \delta(\xi - x^i) \, d\xi = \sum_{i=1}^N w^i \int_{\mathbb{X}} g(\xi) \delta(\xi - x^i) \, d\xi = \sum_{i=1}^N w^i g(x^i).$

12.2 Sequential Importance Sampling

We now apply the IMPORTANCE SAMPLING to the POSTERIOR CONDITIONAL PDF $f_{\mathbf{X}_n|\mathbf{Y}_{(n)}}(\mathbf{x}_n|\mathbf{y}_{(n)})$. Due to the recursive computation, we obtain a SEQUENTIAL IMPORTANCE SAMPLING rule.

First, we again consider the IMPORTANCE DENSITY and the POSTERIOR CONDITIONAL PDF for the construction of the IMPORTANCE WEIGHTS in a slightly different version, viz.³

$$\tilde{w}_n^i = \frac{f_{\mathbf{X}_{(n)}|\mathbf{Y}_{(n)}}(\mathbf{x}_{(n)}^i|\mathbf{y}_{(n)})}{q_{\mathbf{X}_{(n)}|\mathbf{Y}_{(n)}}(\mathbf{x}_{(n)}^i|\mathbf{y}_{(n)})} = \frac{f_{\mathbf{X}_n,\mathbf{X}_{(n-1)}|\mathbf{Y}_{(n)}}(\mathbf{x}_n^i, \mathbf{x}_{(n-1)}^i|\mathbf{y}_{(n)})}{q_{\mathbf{X}_n,\mathbf{X}_{(n-1)}|\mathbf{Y}_{(n)}}(\mathbf{x}_n^i, \mathbf{x}_{(n-1)}^i|\mathbf{y}_{(n)})}. \quad (12.13)$$

We then expand the nominator by

$$\begin{aligned} & f_{\mathbf{X}_n,\mathbf{X}_{(n-1)}|\mathbf{Y}_{(n)}}(\mathbf{x}_n, \mathbf{x}_{(n-1)} | \mathbf{y}_{(n)}) \\ \propto & f_{\mathbf{Y}_n|\mathbf{X}_n,\mathbf{X}_{(n-1)},\mathbf{Y}_{(n-1)}}(\dots) f_{\mathbf{X}_n,\mathbf{X}_{(n-1)}|\mathbf{Y}_{(n-1)}}(\dots) \\ \propto & \underbrace{f_{\mathbf{Y}_n|\mathbf{X}_n,\mathbf{X}_{(n-1)},\mathbf{Y}_{(n-1)}}(\dots)}_{\text{Markov property!}} \underbrace{f_{\mathbf{X}_n|\mathbf{X}_{(n-1)},\mathbf{Y}_{(n-1)}}(\dots)}_{\text{Markov property!}} f_{\mathbf{X}_{(n-1)}|\mathbf{Y}_{(n-1)}}(\dots) \\ \propto & f_{\mathbf{Y}_n|\mathbf{X}_n}(\mathbf{y}_n | \mathbf{x}_n) f_{\mathbf{X}_n|\mathbf{X}_{n-1}}(\mathbf{x}_n | \mathbf{x}_{n-1}) f_{\mathbf{X}_{(n-1)}|\mathbf{Y}_{(n-1)}}(\mathbf{x}_{(n-1)} | \mathbf{y}_{(n-1)}), \end{aligned} \quad (12.14)$$

³Note, that this version is equivalent with the original definition of IMPORTANCE WEIGHTS up to scaling factors, since at the n -th time instant the realizations of the $(n-1)$ -th state vectors $\mathbf{X}_{(n-1)}$ by means of its particles $\mathbf{x}_{(n-1)}^i$ are no longer random. The scaling factors are without any impact, due to the subsequent normalization steps.

and approximate the denominator as⁴

$$\begin{aligned}
& q_{\mathbf{x}_n, \mathbf{x}_{(n-1)} | \mathbf{y}_{(n)}}(\mathbf{x}_n, \mathbf{x}_{(n-1)} | \mathbf{y}_{(n)}) \\
& \approx \underbrace{q_{\mathbf{x}_n | \mathbf{x}_{(n-1)}, \mathbf{y}_{(n)}}(\mathbf{x}_n | \mathbf{x}_{(n-1)}, \mathbf{y}_{(n)})}_{\text{Markov property!}} q_{\mathbf{x}_{(n-1)} | \mathbf{y}_{(n-1)}}(\mathbf{x}_{(n-1)} | \mathbf{y}_{(n-1)}) \\
& = q_{\mathbf{x}_n | \mathbf{x}_{n-1}, \mathbf{y}_n}(\mathbf{x}_n | \mathbf{x}_{n-1}, \mathbf{y}_n) q_{\mathbf{x}_{(n-1)} | \mathbf{y}_{(n-1)}}(\mathbf{x}_{(n-1)} | \mathbf{y}_{(n-1)}).
\end{aligned} \tag{12.15}$$

Consequently, the weights are obtained by

$$\tilde{w}_n^i = \frac{f_{\mathbf{y}_n | \mathbf{x}_n}(\mathbf{y}_n | \mathbf{x}_n^i) f_{\mathbf{x}_n | \mathbf{x}_{n-1}}(\mathbf{x}_n^i | \mathbf{x}_{n-1}^i) f_{\mathbf{x}_{(n-1)} | \mathbf{y}_{(n-1)}}(\mathbf{x}_{(n-1)}^i | \mathbf{y}_{(n-1)})}{q_{\mathbf{x}_n | \mathbf{x}_{n-1}, \mathbf{y}_n}(\mathbf{x}_n^i | \mathbf{x}_{n-1}^i, \mathbf{y}_n) q_{\mathbf{x}_{(n-1)} | \mathbf{y}_{(n-1)}}(\mathbf{x}_{(n-1)}^i | \mathbf{y}_{(n-1)})}. \tag{12.16}$$

⁴It is an approximation, since it does not equal with the correct expansion of the IMPORTANCE DENSITY. The correct expansion following the BAYES RULE would read $q_{\mathbf{x}_n, \mathbf{x}_{(n-1)} | \mathbf{y}_{(n)}}(\mathbf{x}_n, \mathbf{x}_{(n-1)} | \mathbf{y}_{(n)}) = q_{\mathbf{x}_n | \mathbf{x}_{(n-1)}, \mathbf{y}_{(n)}}(\mathbf{x}_n | \mathbf{x}_{(n-1)}, \mathbf{y}_{(n)}) q_{\mathbf{x}_{(n-1)} | \mathbf{y}_{(n)}}(\mathbf{x}_{(n-1)} | \mathbf{y}_{(n)})$. In other words, the defined version neglects the information from the latest observation \mathbf{y}_n .

Finally, taking into account the definition in (12.13) the weights of the PARTICLES (samples) \mathbf{x}_n^i in the n th step are given by

$$\tilde{w}_n^i = \tilde{w}_{n-1}^i \times \frac{f_{\mathbf{Y}_n|\mathbf{X}_n}(\mathbf{y}_n | \mathbf{x}_n^i) f_{\mathbf{X}_n|\mathbf{X}_{n-1}}(\mathbf{x}_n^i | \mathbf{x}_{n-1}^i)}{q_{\mathbf{X}_n|\mathbf{X}_{n-1}, \mathbf{Y}_n}(\mathbf{x}_n^i | \mathbf{x}_{n-1}^i, \mathbf{y}_n)}, \quad \forall i = 1, \dots, N, \quad (12.17)$$

where the particles \mathbf{x}_n^i are drawn from the conditional PDF $q_{\mathbf{X}_n|\mathbf{X}_{n-1}, \mathbf{Y}_n}(\mathbf{x}_n | \mathbf{x}_{n-1}^i, \mathbf{y}_n)$.

Given the particles \mathbf{x}_{n-1}^i with respect to $q_{\mathbf{X}_n|\mathbf{X}_{n-1}, \mathbf{Y}_n}(\mathbf{x}_n | \mathbf{x}_{n-1}^i, \mathbf{y}_n)$ in the latest step of the recursion process, and given an observation \mathbf{y}_n with respect to the state-space model, e.g.,

$$\mathbf{X}_n = g_n(\mathbf{X}_{n-1}) + \mathbf{V}_n, \quad (12.18)$$

$$\mathbf{Y}_n = h_n(\mathbf{X}_n) + \mathbf{W}_n, \quad (12.19)$$

where \mathbf{V}_n and \mathbf{W}_n are random variables with any PDF $f_{\mathbf{V}_n}(\mathbf{v})$ and $f_{\mathbf{W}_n}(\mathbf{w})$, the weights of the IMPORTANCE SAMPLING are updated according to (12.17) with

$$f_{\mathbf{Y}_n|\mathbf{X}_n}(\mathbf{y}_n | \mathbf{x}_n^i) = f_{\mathbf{W}_n}(\mathbf{y}_n - h_n(\mathbf{x}_n^i)) \quad (12.20)$$

$$f_{\mathbf{X}_n|\mathbf{X}_{n-1}}(\mathbf{x}_n^i | \mathbf{x}_{n-1}^i) = f_{\mathbf{V}_n}(\mathbf{x}_n^i - g_n(\mathbf{x}_{n-1}^i)). \quad (12.21)$$

12.3 Degeneracy Problem and Resampling

Degeneracy.

It has been shown that using an IMPORTANCE DENSITY with

$$q_{\mathbf{x}_n|\mathbf{x}_{n-1},\mathbf{y}_n}(\dots)q_{\mathbf{x}_{n-1}|\mathbf{y}_{(n-1)}}(\dots) \neq f_{\mathbf{x}_n,\mathbf{x}_{n-1}|\mathbf{y}_{(n)}}(\dots)$$

leads to a monotonic increase of the variance of the weights over the number of iterations n . This in turn leads to the DEGENERACY EFFECT, which means that after a number of recursive steps all but one particle have vanishing normalized IMPORTANCE WEIGHTS.

DEGENERACY is typically detected by a threshold test. If

$$\frac{1}{\sum_{i=1}^N (w_n^i)^2} \leq w_{\text{thr}} \quad (12.22)$$

is smaller than a certain threshold w_{thr} , a resampling process must be performed.

It is the purpose of the RESAMPLING STEP to eliminate particles with small weights and to replicate particles with strong weights. REPLICATION means that when drawing new PARTICLES those \mathbf{x}_n^i with strong weights might have more than one follower in the next recursive step.

Resampling Step.

Particles are randomly drawn (with replacement) from the finite set of particles $\mathbb{X}_n^N \triangleq \{\mathbf{x}_n^1, \dots, \mathbf{x}_n^N\}$, where the probability of each particle is equal to its current IMPORTANCE WEIGHT, i.e., $P(\{\mathbf{x}_n^i\}) = w_n^i$.

Using an alternative formulation of the RESAMPLING STEP by means of the respective PMF⁵, we obtain

$$\forall j = 1, \dots, N : \quad \mathbf{x}_n^j \sim p_{\mathbf{x}_n}(\mathbf{x}_n^j) \quad \text{with} \quad p_{\mathbf{x}_n}(\mathbf{x}_n) = \sum_{i=1}^N w_n^i \delta(\mathbf{x}_n - \mathbf{x}_n^i). \quad (12.23)$$

The PMF $p_{\mathbf{x}_n}(\mathbf{x}_n)$ obviously corresponds to the approximated posterior conditional PDF

$$f_{\mathbf{x}_n|\mathbf{y}_{(n)}}(\mathbf{x}_n \mid \mathbf{y}_{(n)}) \approx p_{\mathbf{x}_n}(\mathbf{x}_n). \quad (12.24)$$

See (12.11) and (12.12).⁶

Note.

After the RESAMPLING STEP the IMPORTANCE WEIGHTS are uniformly set to $w_n^i = N^{-1}$. Thereby the phenomenon of degeneracy is broken.

⁵PROBABILITY MASS FUNCTION.

⁶Directly following (12.11) and (12.12) leads to $f_{\mathbf{x}_{(n)}|\mathbf{y}_{(n)}}(\mathbf{x}_{(n)} \mid \mathbf{y}_{(n)}) \approx \sum_{i=1}^N w_n^i \delta(\mathbf{x}_{(n)} - \mathbf{x}_{(n)}^i) = \sum_{i=1}^N w_n^i \delta(\mathbf{x}_n - \mathbf{x}_n^i) \delta(\mathbf{x}_{(n-1)} - \mathbf{x}_{(n-1)}^i)$ as we started with $f_{\mathbf{x}_{(n)}|\mathbf{y}_{(n)}}(\mathbf{x}_{(n)} \mid \mathbf{y}_{(n)})$ in (12.13). A subsequent marginalization step with respect to $\mathbf{x}_{(n-1)}$ finally yields (12.24).

12.4 Importance Density

Optimal Importance Density.

So far we did not discuss the choice of the IMPORTANCE DENSITY $q_{\mathbf{x}_n|\mathbf{x}_{n-1}, \mathbf{y}_n}(\dots)$, which is of course essential for the performance of the PARTICLE FILTER. It has been shown that the optimal choice which minimizes the variance of the IMPORTANCE WEIGHTS is equal to the true posterior conditional PDF, i.e.,

$$\begin{aligned} q_{\mathbf{x}_n|\mathbf{x}_{n-1}, \mathbf{y}_n}^{\text{opt}}(\mathbf{x}_n | \mathbf{x}_{n-1}^i, \mathbf{y}_n) &= f_{\mathbf{x}_n|\mathbf{x}_{n-1}, \mathbf{y}_n}(\mathbf{x}_n | \mathbf{x}_{n-1}^i, \mathbf{y}_n) \\ &= \frac{f_{\mathbf{y}_n|\mathbf{x}_n, \mathbf{x}_{n-1}}(\mathbf{y}_n | \mathbf{x}_n, \mathbf{x}_{n-1}^i) f_{\mathbf{x}_n|\mathbf{x}_{n-1}}(\mathbf{x}_n | \mathbf{x}_{n-1}^i)}{f_{\mathbf{y}_n|\mathbf{x}_{n-1}}(\mathbf{y}_n | \mathbf{x}_{n-1}^i)}, \end{aligned} \quad (12.25)$$

with

$$f_{\mathbf{y}_n|\mathbf{x}_{n-1}}(\mathbf{y}_n | \mathbf{x}_{n-1}^i) = \int_{\mathbb{X}} f_{\mathbf{y}_n|\mathbf{x}_n}(\mathbf{y}_n | \boldsymbol{\xi}) f_{\mathbf{x}_n|\mathbf{x}_{n-1}}(\boldsymbol{\xi} | \mathbf{x}_{n-1}^i) d\boldsymbol{\xi}. \quad (12.26)$$

Applying the OPTIMAL IMPORTANCE DENSITY in (12.17) results in the OPTIMAL SEQUENTIAL IMPORTANCE SAMPLING given by

$$\tilde{w}_n^{\text{opt}, i} = \tilde{w}_{n-1}^{\text{opt}, i} \times f_{\mathbf{y}_n|\mathbf{x}_{n-1}}(\mathbf{y}_n | \mathbf{x}_{n-1}^i), \quad \forall i = 1, \dots, N. \quad (12.27)$$

Suboptimal Importance Density.

Since in general (12.25) and (12.26) are hard to compute in real-world application, the OPTIMAL IMPORTANCE DENSITY can be only applied in special cases.⁷

Otherwise, the most popular SUBOPTIMAL IMPORTANCE DENSITY is given by

$$q_{\mathbf{x}_n|\mathbf{x}_{n-1}, \mathbf{y}_n}(\mathbf{x}_n | \mathbf{x}_{n-1}^i, \mathbf{y}_n) \triangleq f_{\mathbf{x}_n|\mathbf{x}_{n-1}}(\mathbf{x}_n | \mathbf{x}_{n-1}^i). \quad (12.28)$$

Applying this suboptimal version of the IMPORTANCE DENSITY in (12.17) results in the SUBOPTIMAL SEQUENTIAL IMPORTANCE SAMPLING given by

$$\tilde{w}_n^i = \tilde{w}_{n-1}^i \times f_{\mathbf{y}_n|\mathbf{x}_n}(\mathbf{y}_n | \mathbf{x}_n^i), \quad \forall i = 1, \dots, N. \quad (12.29)$$

It has been observed that the SUBOPTIMAL SEQUENTIAL IMPORTANCE SAMPLING works quite well, except for the case where the TRANSITIONAL PDF $f_{\mathbf{x}_n|\mathbf{x}_{n-1}}(\mathbf{x}_n | \mathbf{x}_{n-1}^i)$ has a broader spread than the likelihood function $f_{\mathbf{y}_n|\mathbf{x}_n}(\mathbf{y}_n | \mathbf{x}_n)$. For the latter a rapid degeneracy of the weights has been observed.

There is huge room for further suboptimal techniques...

⁷One of this special cases is again the ADDITIVE GAUSSIAN MODEL. The linearity assumption is not required for the PARTICLE FILTERING.

13. Tracking the Trajectory of a Moving Object

In the following, we consider the STATE-SPACE MODELS given as

$$\mathbf{X}_{n+1} = g(\mathbf{X}_n) + \mathbf{V}_n, \quad (13.1)$$

$$\mathbf{Y}_n = h(\mathbf{X}_n) + \mathbf{W}_n, \quad (13.2)$$

where for simplicity additive noise is assumed. \mathbf{X}_n and \mathbf{Y}_n are state space and measurements vectors. The functions $g(\cdot)$ and $h(\cdot)$ are in general non-linear. \mathbf{v}_k and \mathbf{w}_k are not necessarily Gaussian distributed.

In the following, we will estimate the results of this state-space model using a KALMAN FILTER (KF) and using a simple PARTICLE FILTER (PF) with SEQUENTIAL IMPORTANCE RE-SAMPLING (SIR).

13.1 Linear Model

We consider a two-dimensional version of the example in Section 11.2, i.e.,

$$\ddot{\mathbf{X}}_t = \mathbf{V}_t, \quad \text{i.i.d. } \mathbf{V}_n \sim \mathcal{N}(0, \mathbf{C}_V). \quad (13.3)$$

The discretization of the STOCHASTICAL DIFFERENTIAL EQUATION (13.3) yields

$$\mathbf{X}_{n+1} = \mathbf{G}\mathbf{X}_n + \mathbf{V}_n, \quad (13.4)$$

$$\mathbf{Y}_n = \mathbf{H}\mathbf{X}_n + \mathbf{W}_n, \quad (13.5)$$

with

$$\mathbf{G} = \begin{bmatrix} 1 & 0 & T & 0 \\ 0 & 1 & 0 & T \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad \text{and} \quad \mathbf{H} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}. \quad (13.6)$$

The state space vector and the observation vector are

$$\mathbf{X}_n = \begin{bmatrix} X_{1,n} \\ X_{2,n} \\ \dot{X}_{1,n} \\ \dot{X}_{2,n} \end{bmatrix} \in \mathbb{R}^4 \quad \text{and} \quad \mathbf{Y}_n \in \mathbb{R}^2. \quad (13.7)$$

The $X_{1,n}$ and $X_{2,n}$ represent the random variables for a position on the (x_1, x_2) -plane and $\dot{X}_{1,n}$ and $\dot{X}_{2,n}$ represent the random variables for the respective velocities along the x_1 and x_2 direction. The $Y_{1,n}$ and $Y_{2,n}$ denote noisy observations of $X_{1,n}$ and $X_{2,n}$.

The \mathbf{V}_n and \mathbf{W}_n are assumed to be Gaussian distributed with

$$\mathbf{V}_n \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_V) \quad \text{and} \quad \mathbf{W}_n \sim \mathcal{N}(\mathbf{0}, \mathbf{C}_W), \quad (13.8)$$

and

$$\mathbf{C}_V = T \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 50 & 0 \\ 0 & 0 & 0 & 50 \end{bmatrix} \quad \text{and} \quad \mathbf{C}_W = \begin{bmatrix} 0.1 & 0 \\ 0 & 0.1 \end{bmatrix}. \quad (13.9)$$

Note.

For the PARTICLE FILTER the number of particles is set to $N = 100$ and resampling is applied if

$$N = \frac{1}{\sum_{i=1}^N (w_n^i)^2} \leq 50. \quad (13.10)$$

The w_n^i are the particle weights at time-step n .

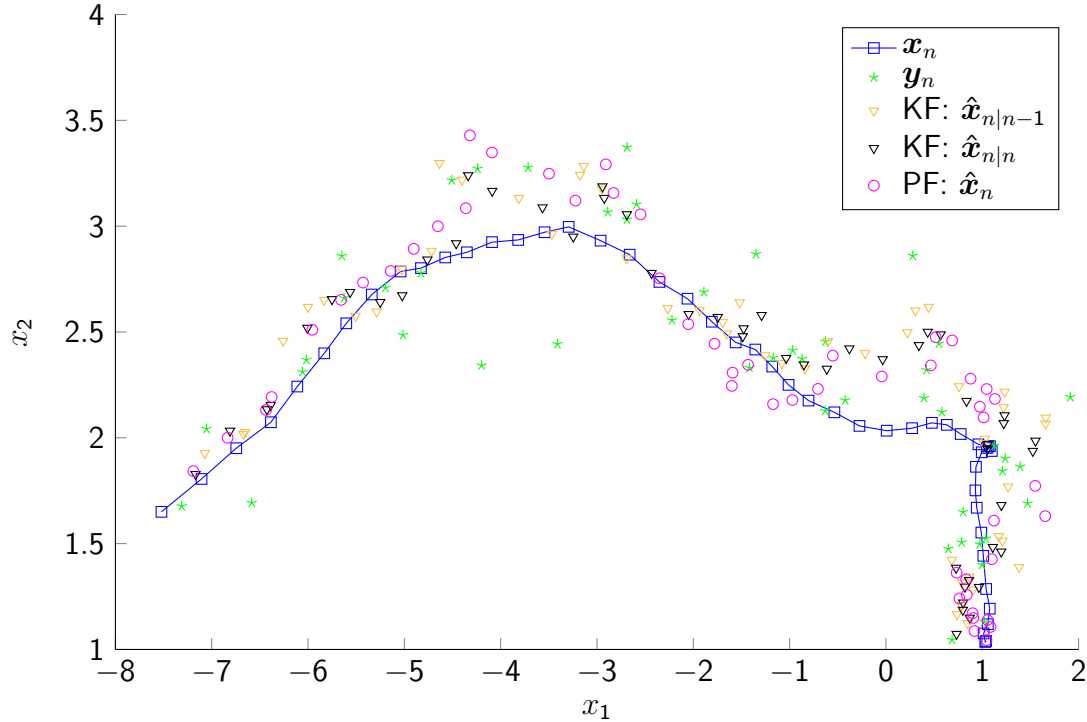


Fig. 13.1: Tracking of an exemplary random trajectory in the (x_1, x_2) -plane based on KALMAN (KF) and PARTICLE FILTERING (PF) applied to observations \mathbf{y}_n from a LINEAR SYSTEM MODEL with initial condition $\mathbf{x}_0 = [1, 0, 0, 0]^T$, $T = 0.02$ and $n = 1 \dots 50$.

13.2 Nonlinear Model

The nonlinear model differs from the model in Section 13.1 by a modification of the first row of Eq. (13.4), i.e., we replace $\dot{X}_{1,n+1} = X_{1,n} + T\dot{X}_{1,n}$ by

$$\dot{X}_{1,n+1} = \sin(X_{1,n}) + T\dot{X}_{1,n}. \quad (13.11)$$

Additionally, the measurement noise covariance matrix is changed to

$$C_W = \begin{bmatrix} 0.025 & 0 \\ 0 & 0.025 \end{bmatrix}. \quad (13.12)$$

Note.

For the `PARTICLE FILTER` the number of particles is again set to $N = 100$ and resampling is applied if

$$N = \frac{1}{\sum_{i=1}^N (w_n^i)^2} \leq 50. \quad (13.13)$$

For comparison we use the Kalman Filter with the linearization (Jacobi matrix) of $g(\mathbf{x}_n)$ at the initial condition \mathbf{x}_0 as constant matrix \mathbf{G} , i.e.,

$$[\mathbf{G}]_{k\ell} = \left. \frac{\partial [g(\mathbf{x})]_k}{\partial [\mathbf{x}]_\ell} \right|_{\mathbf{x}=\mathbf{x}_0}. \quad (13.14)$$

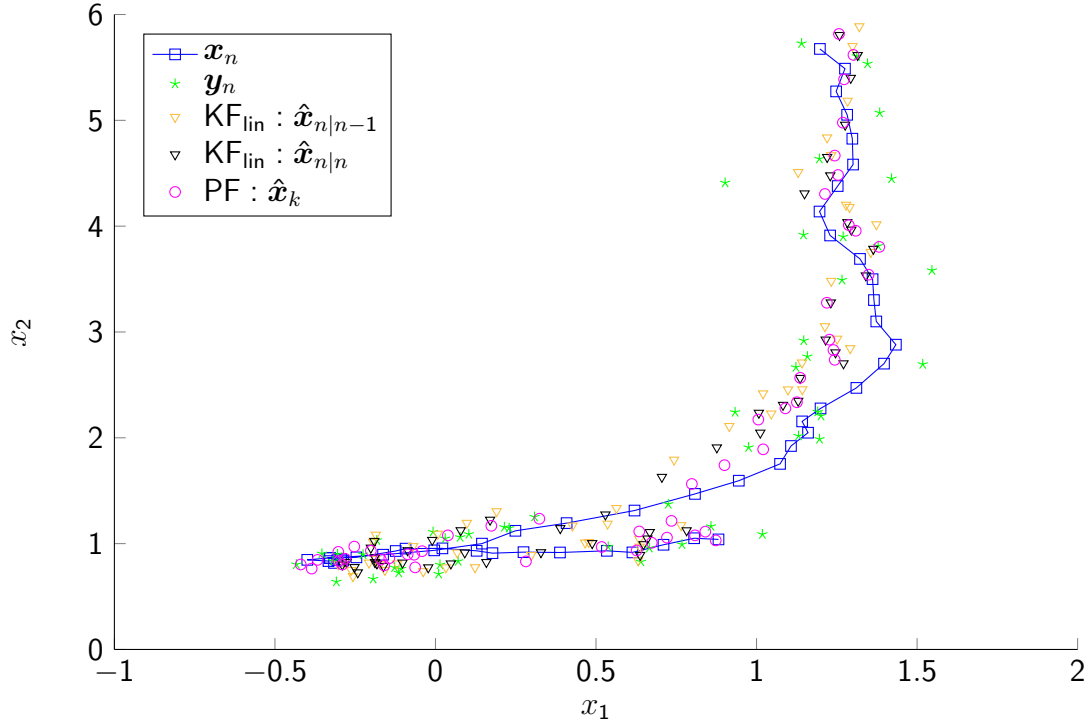


Fig. 13.2: Tracking of an exemplary random trajectory in the (x_1, x_2) -plane based on KALMAN (KF) and PARTICLE FILTERING (PF) applied to observations y_n from a NONLINEAR SYSTEM MODEL with initial condition $x_0 = [1, 0, 0, 0]^T$, $T = 0.02$ and $n = 1 \dots 50$.

14. Demo

<https://syncandshare.lrz.de/dl/fiDKhkL7MKJpUZ9R7YxnUQa9>

14.1 Gauss-Markov Model with Standard Measurement Noise

We illustrate the PARTICLE FILTER using a very simple GAUSS-MARKOV model

$$X_n = X_{n-1} + V_n, \quad (14.1)$$

$$Y_n = X_n + W_n \quad (14.2)$$

with $x_0 = 0$, $V_n \sim \mathcal{N}(0, 1)$, and $W_n \sim \mathcal{N}(0, 1)$.

- We use the sub-optimal IMPORTANCE DENSITY $f_{X_n|X_{n-1}}(x_n|x_{n-1})$, i.e.,

$$q_{X_n|X_{n-1}, Y_n}(x_n|x_{n-1}, y_n) = f_{X_n|X_{n-1}}(x_n|x_{n-1}). \quad (14.3)$$

- We use $N_{\text{part}} = 500$ particles with $w_{\text{thr}} = 25$.

- The PARTICLE FILTER approximates the optimal KALMAN FILTER which in turn estimates the true trajectory is $X_n = x_n$ with $n = 1, \dots, 51$.
- After each new measurement $Y_n = y_n$, particles are propagated by using the IMPORTANCE DENSITY $f_{X_n|X_{n-1}}(x_n|x_{n-1}^i)$.
- In this example, $f_{X_n|X_{n-1}}(x_n|x_{n-1}^i) = \mathcal{N}(x_{n-1}^i, 1)$.
- Then, the weights are updated by using

$$\tilde{w}_n^i = w_{n-1}^i \times f_{Y_n|X_n}(y_n|x_n^i) \quad (14.4)$$

and scaled by using

$$w_n^i = \frac{\tilde{w}_n^i}{\sum w_n^i}. \quad (14.5)$$

- In this example, $f_{Y_n|X_n}(y_n|x_n^i) = \mathcal{N}(x_n^i, 1)$ and $w_0^i = \frac{1}{N}$.
- Eventually, $E[X_n|Y_{(n)} = y_{(n)}]$ is approximated by using $\hat{x}_{n|n} \approx \sum x_n^i w_n^i$.
- If WEIGHT DEGENERACY occurs, i.e.,

$$\frac{1}{\sum (w_n^i)^2} \leq w_{\text{thr}} = 25, \quad (14.6)$$

particles are resampled by means of $P(x_n^i) = w_n^i$ and the weights are set to $w_n^i = \frac{1}{N}$.

14.2 Gauss-Markov Model with Non-Standard Measurement Noise

We again illustrate the PARTICLE FILTER using a very simple GAUSS-MARKOV model

$$X_n = X_{n-1} + V_n, \quad (14.7)$$

$$Y_n = X_n + W_n \quad (14.8)$$

but now with $x_0 = 0$, $V_n \sim \mathcal{N}(0, 2)$, and $W_n \sim \frac{1}{2}\mathcal{N}(0, 0.1) + \frac{1}{2}\mathcal{N}(10, 0.1)$.

In this case of a BIMODAL GAUSSIAN DISTRIBUTION, the PARTICLE FILTER will clearly outperform the KALMAN FILTER, which is not able to interpret the likelihood of the measurements in an appropriate way.

In this example,

$$- f_{Y_n|X_n}(y_n|x_n^i) = \frac{1}{2}\mathcal{N}(x_n^i, 0.1) + \frac{1}{2}\mathcal{N}(x_n^i + 10, 0.1) \text{ with } w_0^i = \frac{1}{N}.$$

References

- D. J. C. MacKay. *Information Theory, Inference, and Learning Algorithms*, Cambridge University Press.
- B. Hajek. *An Exploration of Random Processes for Engineers*. Lecture Notes at University of Illinois, Urbana-Champaign. URL: <http://www.freotechbooks.com>.
- G. R. Grimmett, D. R. Stirzaker. *Probability and Random Processes*, Oxford University Press. (↑advanced reader!)
- P. M. Djurić et al. "Particle Filtering: A review of the theory and how it can be used for solving problems in wireless communications", *IEEE Signal Processing Magazine*, p. 19–38, September 2003.
- B. Ristic, S. Arulampalam, and N. Gordon. *Beyond the Kalman Filter*, Artech House, 2004.
- D. Simon. *Optimal State Estimation*, Wiley, 2006. (↑advanced reader!)

Part VII

Hypothesis Testing

15. Hypothesis Testing – Decision Making

Contrary to parameter estimation, hypothesis testing is about making a decision based on the observations, so-called decision making problems. Again the statistical model is in the main focus.

15.1 Statistical Model

Given the statistical model $\{\mathbb{X}, \mathbb{F}, P_\theta; \theta \in \Theta\}$, we can decompose the model parameter $\Theta = \Theta_0 \cup \Theta_1$ into a so-called null hypothesis Θ_0 and alternative Θ_1 . In this lecture we restrict ourselves to the singleton sets $\Theta_0 = \{\theta_0\}$ and $\Theta_1 = \{\theta_1\}$.

Definition. A test is a statistic

$$d : \mathbb{X} \rightarrow [0, 1].$$

The statistic d corresponds to a decision rule and the set $\{x \in \mathbb{X} | d(x) = 1\}$ is called rejection region.

Definition. The cost criterion for a test d is

$$G_d : \Theta \rightarrow [0, 1], \theta \mapsto E [d(X)|\theta] .$$

Definition. A test has an error level α if

$$E [d(X)|\theta_0] \leq \alpha$$

Two Error Types

false alarm: Although $\theta \in \Theta_0$, the null hypothesis is canceled, that is $d(x) = 1$. The false alarm error probability is

$$G_d(\theta_0) = E [d(X)|\theta_0] .$$

detection error: Although $\theta \in \Theta_1$, the alternative is not detected, that is $d(x) = 0$. The detection error probability is

$$1 - G_d(\theta_1) = 1 - E [d(X)|\theta_1] .$$

15.2 Design of a Test

A test should be designed under the following requirements

- the false alarm error probability should not exceed a given value α :

$$G_d(\theta_0) \leq \alpha.$$

- the detection error probability should be as small as possible:

$$\max \{G_d(\theta_1)\}.$$

Thus, the design of a test can be formulated as an optimization problem:

$$\begin{array}{ll} \max_d & G_d(\theta_1) \\ \text{s. t.} & G_d(\theta_0) \leq \alpha \end{array}$$

Definition. A test d of θ_0 against θ_1 is called uniformly most powerful (UMP) test in case

$$G_d(\theta_1) \geq G_{d'}(\theta_1)$$

for all tests $d' \in \mathbb{D}$, which satisfy $G_d(\theta_0) = G_{d'}(\theta_0)$.

Definition. A test is called a true test with error level α if

$$G_d(\theta_0) \leq \alpha \leq G_d(\theta_1)$$

A true test decides with higher probability for the alternative, if it is correct, than incorrectly for the null hypothesis. This means, a correct detection is more probable than a false alarm.

15.3 Alternative Test

Here we regard the two probability measures $P(\theta_0)$, $P(\theta_1)$ respectively P_0, P_1 . The statistical model is $\{\mathbb{X}, \mathbb{F}, \{P_0, P_1\}\}$ with probability density functions $f_X(x; \theta_0)$, $f_X(x; \theta_1)$, such that $f_X(x; \theta_0) + f_X(x; \theta_1) > 0$, $\forall x \in \mathbb{X}$.

What is the best test d of P_0 against P_1 for a given error probability α ?

Given an observation x , according to the ML-principle one could decide for the alternative in case $f_X(x; \theta_1) > f_X(x; \theta_0)$. The decision threshold is given by the ML-ratio:

$$R(x) = \begin{cases} \frac{f_X(x; \theta_1)}{f_X(x; \theta_0)} & ; \quad f_X(x; \theta_0) > 0 \\ \infty & ; \quad f_X(x; \theta_0) = 0 \text{ and } f_X(x; \theta_1) > 0. \end{cases}$$

Neyman-Pearson-Test

Given a statistical model $\{\mathbb{X}, \mathbb{F}, \{P_0, P_1\}\}$ with discrete hypothesis and alternative, for any error level $0 < \alpha < 1$:

a) The best test of P_0 against P_1 is

$$d(x) = \begin{cases} 1 & ; R > c \\ \gamma & ; R = c \\ 0 & ; R < c, \end{cases}$$

where $\gamma = \frac{\alpha - P_0(\{R > c\})}{P_0(\{R = c\})}$ or $\gamma = 0$ if $P_0(\{R = c\}) = 0$.¹

b) There exists a test that fully utilizes the error level α .

c) Any Neyman-Pearson (NP) test with error level α is a best test.

¹

$$\begin{aligned} E[d(x)|\theta_0] &= 1 \cdot P_0(\{R > c\}) + \gamma \cdot P_0(\{R = c\}) \\ &= \begin{cases} \alpha & ; P_0(\{R = c\}) = 0 \\ P_0(\{R > c\}) & ; P_0(\{R = c\}) > 0 \end{cases} + P_0(\{R = c\}) \cdot \begin{cases} 0 & ; P_0(\{R = c\}) = 0 \\ \frac{\alpha - P_0(\{R > c\})}{P_0(\{R = c\})} & ; P_0(\{R = c\}) > 0 \end{cases} \\ &= \alpha. \end{aligned}$$

Example 1

$$f_X(x; \theta_0) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}x^2\right)$$
$$f_X(x; \theta_1) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x-1)^2\right).$$

The Neyman-Pearson test is

$$d(x) = \begin{cases} 1 & ; R(x) > c \\ 0 & ; R(x) < c. \end{cases}$$

with decision rule

$$R(x) = \frac{f_X(x; \theta_1)}{f_X(x; \theta_0)} = \frac{\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x-1)^2\right)}{\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}x^2\right)} = \exp\left(\frac{1}{\sigma^2}\left(x - \frac{1}{2}\right)\right) \geq c.$$

Further, we have

$$x \geq x_\alpha \Rightarrow R(x) \geq R(x_\alpha),$$

where $R(x_\alpha) = c$. The α -fractil of $\mathcal{N}(0, \sigma^2)$ is x_α and therefore

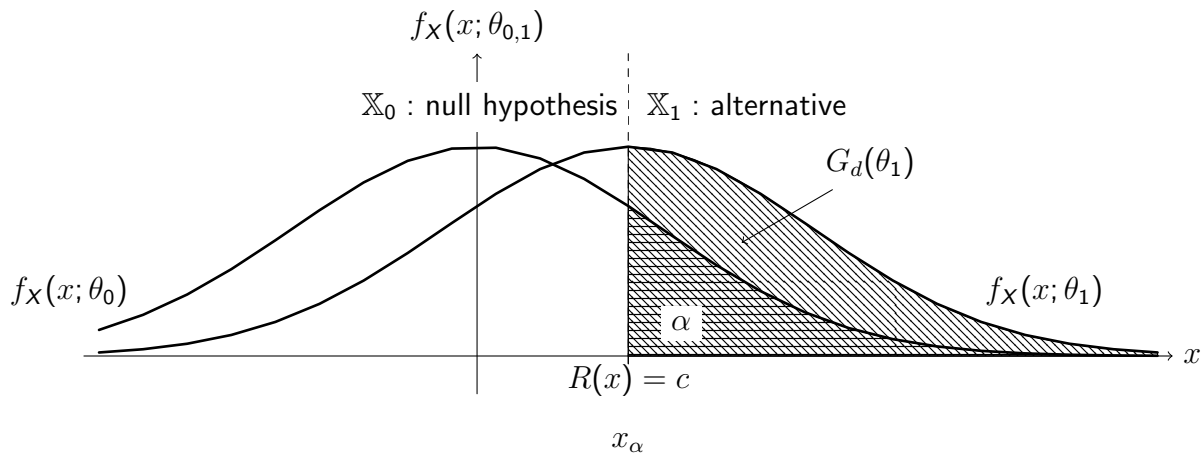
$$\alpha = \int_{x_\alpha}^{\infty} \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}x^2\right) dx = Q\left(\frac{x_\alpha}{\sigma}\right)$$
$$\Rightarrow x_\alpha = \sigma Q^{-1}(\alpha).$$

As $R(x_\alpha) = c$, we obtain

$$c = \exp \left(\frac{1}{\sigma^2} \left(x_\alpha - \frac{1}{2} \right) \right),$$

and the Neyman-Pearson test can be written as

$$d(x) = \begin{cases} 1 & ; R(x) > \exp \left(\frac{1}{\sigma^2} \left(\sigma Q^{-1}(\alpha) - \frac{1}{2} \right) \right) \\ 0 & ; \text{otherwise.} \end{cases}$$



Example 2

$$f_X(x; \theta_0) = \frac{1}{\sqrt{2\pi}\sigma_0} \exp\left(-\frac{1}{2\sigma_0^2}x^2\right)$$
$$f_X(x; \theta_1) = \frac{1}{\sqrt{2\pi}\sigma_1} \exp\left(-\frac{1}{2\sigma_1^2}x^2\right).$$

The decision rule is

$$R(x) = \frac{\frac{1}{\sigma_1} \exp\left(-\frac{1}{2\sigma_1^2}x^2\right)}{\frac{1}{\sigma_0} \exp\left(-\frac{1}{2\sigma_0^2}x^2\right)} \geq c$$
$$-\frac{1}{2\sigma_1^2}x^2 + \frac{1}{2\sigma_0^2}x^2 \geq \ln c + \ln \frac{\sigma_1}{\sigma_0}$$
$$-\frac{1}{2} \left(\frac{1}{\sigma_1^2} - \frac{1}{\sigma_0^2} \right) x^2 \geq \ln c + \frac{1}{2} \ln \frac{\sigma_1^2}{\sigma_0^2}.$$

In case $\sigma_1^2 > \sigma_0^2$, we have

$$x^2 \geq \frac{2 \ln c + \ln \frac{\sigma_1^2}{\sigma_0^2}}{\frac{1}{\sigma_0^2} - \frac{1}{\sigma_1^2}}.$$

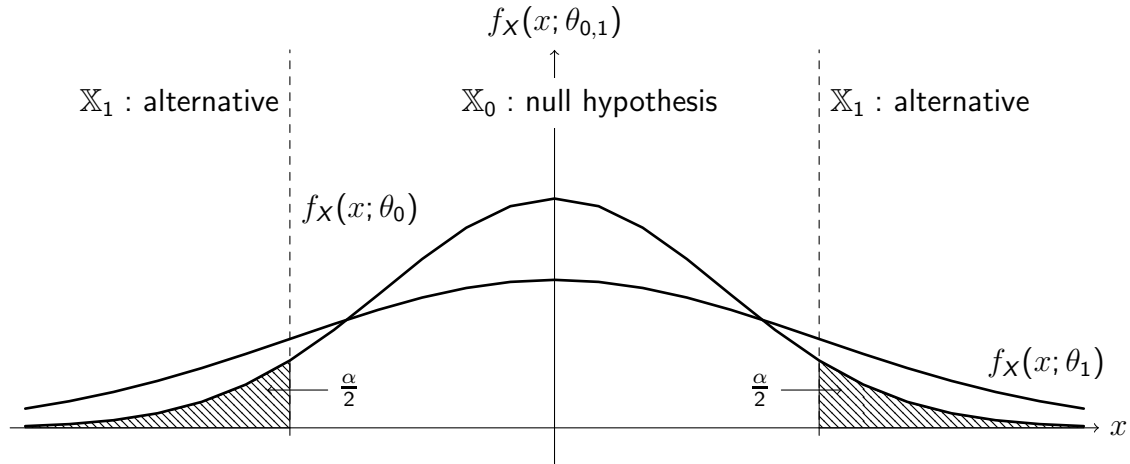
From $R(x_\alpha) = c$, we conclude

$$\ln c = -\frac{1}{2} \left(\frac{1}{\sigma_1^2} - \frac{1}{\sigma_0^2} \right) x_\alpha^2 - \frac{1}{2} \ln \frac{\sigma_1^2}{\sigma_0^2},$$

where x_α is the α -fractil of $\mathcal{N}(0, \sigma_0^2)$. Therefore, the Neyman-Pearson test can be formulated as

$$d(x) = \begin{cases} 1 & ; \quad x^2 > \frac{2 \ln c + \ln \frac{\sigma_1^2}{\sigma_0^2}}{\frac{1}{\sigma_0^2} - \frac{1}{\sigma_1^2}} = x_\alpha^2, \\ 0 & ; \quad \text{otherwise.} \end{cases}$$

Note that $R(x)$ is monotone increasing/decreasing in $\mathbb{R}_+/\mathbb{R}_-$.



Example 3

$$\begin{aligned} B_{N,\theta_0}(x) &= \binom{N}{x} \theta_0^x (1 - \theta_0)^{N-x}, \quad x \in \{0, 1, \dots, N\} \\ B_{N,\theta_1}(x) &= \binom{N}{x} \theta_1^x (1 - \theta_1)^{N-x}, \quad x \in \{0, 1, \dots, N\}. \end{aligned}$$

The Neyman-Pearson test subject to an error level α is given as²

$$d(x) = \begin{cases} 1 & ; \quad R(x) > c(\alpha) \\ \gamma & ; \quad R(x) = c(\alpha) \\ 0 & ; \quad R(x) < c(\alpha), \end{cases}$$

with the likelihood ratio given as

$$R(x) = \frac{B_{N,\theta_1}(x)}{B_{N,\theta_0}(x)} = \frac{\binom{N}{x} \theta_1^x (1 - \theta_1)^{N-x}}{\binom{N}{x} \theta_0^x (1 - \theta_0)^{N-x}} = \left(\frac{\theta_1}{\theta_0} \right)^x \left(\frac{1 - \theta_1}{1 - \theta_0} \right)^{N-x} \geq c(\alpha), \quad x \in \{0, 1, \dots, N\}.$$

We observe that if $\theta_1 > \theta_0$, $R(x)$ is monotonically increasing in x as $\frac{\theta_1}{\theta_0} > 1$ and $\frac{1-\theta_1}{1-\theta_0} < 1$. Similarly, if $\theta_1 < \theta_0$, $R(x)$ is monotonically decreasing in x . In the following, let $\theta_1 > \theta_0$.

²Note that for the Binomial distribution it may occur that $P(\{R(x) = c(\alpha)\}) > 0$ if $\alpha, \theta_0, \theta_1 \in (0, 1)$.

Therefore, we can determine the decision threshold x_T of a Neyman-Pearson test subject to an error level of α by looking for an α -fractil of $B_{N,\theta_0}(x)$, i.e., by considering

$$P_{\theta_0}(\{X \geq x_T\}) = \sum_{x \geq x_T}^N \binom{N}{x} \theta_0^x (1 - \theta_0)^{N-x} \geq \alpha.$$

For instance, let $N = 10$ and $\theta_0 = \frac{1}{2}$. It follows that

$$P_{\theta_0}(\{X = 10\}) = \binom{10}{10} \left(\frac{1}{2}\right)^{10} = \frac{1}{1024},$$

$$P_{\theta_0}(\{X = 9\}) = \binom{10}{9} \left(\frac{1}{2}\right)^{10} = \frac{10}{1024},$$

$$P_{\theta_0}(\{X = 8\}) = \binom{10}{8} \left(\frac{1}{2}\right)^{10} = \frac{45}{1024}.$$

Using simple decision thresholds, we can only construct tests that achieve certain error levels. For example, using $x_T = 8$ and $x_T = 9$ we achieve error levels of $P_{\theta_0}(\{X \geq 8\}) = \frac{56}{1024}$ and $P_{\theta_0}(\{X \geq 9\}) = \frac{11}{1024}$. The Neyman-Pearson test subject to an error level of, for example, $\alpha = \frac{50}{1024}$ is given as

$$d(x) = \begin{cases} 1 & ; \quad x > 8 \\ \gamma & ; \quad x = 8 \\ 0 & ; \quad x < 8, \end{cases}$$

with

$$\gamma = \frac{\alpha - P_{\theta_0}(\{X > 8\})}{P_{\theta_0}(\{X = 8\})} = \frac{\frac{50}{1024} - \frac{11}{1024}}{\frac{45}{1024}} = \frac{39}{45}.$$

The decision at the threshold $x_T = 8$ has to be randomized in order to fully utilize the error level $\alpha = \frac{50}{1024}$ of the Neyman-Pearson test. If we observe $X = 8$, we reject the null hypothesis with probability $\gamma = \frac{39}{45}$ in order to achieve that

$$\begin{aligned} \alpha &= P_{\theta_0}(\{X > 8\}) + \gamma P_{\theta_0}(\{X = 8\}) \\ &= \frac{11}{1024} + \frac{39}{45} \frac{45}{1024} = \frac{50}{1024}. \end{aligned}$$

Special Case

The test

$$d(x) = \begin{cases} 1 & ; R(x) > 1 \\ 0 & ; \text{otherwise.} \end{cases}$$

is called ML detector, for instance set $c = 1$ in example 1.

Receiver-Operating-Characteristic (ROC) Graphs

The performance of an NP test can be illustrated by a so-called Reciever-Operating-Characteristics, which plots $G_d(\theta_1)$ as a function of $G_d(\theta_0)$.

ROC for Example 1 ($\sigma^2 = 1$)

$$\begin{aligned} P_{\text{FA}} = G_d(\theta_0) &= E[d(X)|\theta_0] = \int_{x_\alpha}^{\infty} 1 \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}x^2\right) dx = Q(x_\alpha) = \alpha \\ \Rightarrow x_\alpha &= Q^{-1}(P_{\text{FA}}). \end{aligned}$$

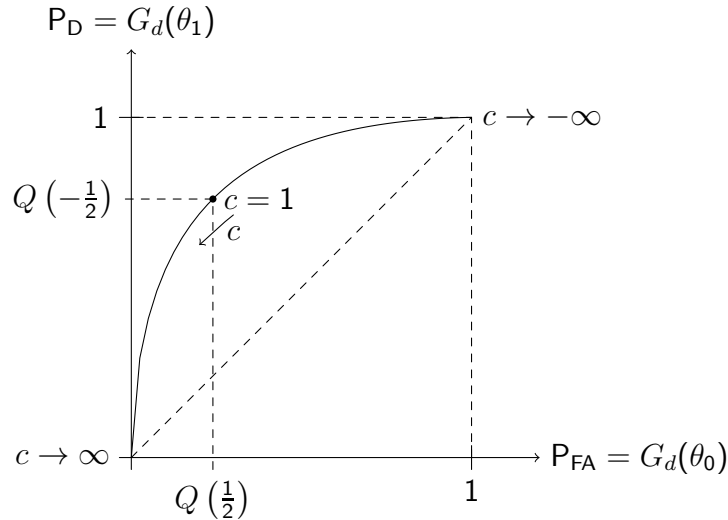
$$\begin{aligned} P_{\text{D}} = G_d(\theta_1) &= E[d(X)|\theta_1] = \int_{x_\alpha}^{\infty} 1 \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x-1)^2\right) dx \\ &= \int_{x_\alpha-1}^{\infty} 1 \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}x^2\right) dx = Q(x_\alpha - 1) \\ \Rightarrow P_{\text{D}} &= Q(Q^{-1}(P_{\text{FA}}) - 1). \end{aligned}$$

For $c = 1$, we have (see page 181):

$$1 = \exp\left(\frac{1}{\sigma^2}\left(x_\alpha - \frac{1}{2}\right)\right) \Rightarrow x_\alpha = \frac{1}{2}$$

$$\Rightarrow Q(x_\alpha) = Q\left(\frac{1}{2}\right)$$

$$\Rightarrow P_D = Q(x_\alpha - 1) = Q\left(-\frac{1}{2}\right).$$



15.4 Bayes Test

Similar to Bayes estimation, we regard the case where prior knowledge on possible hypotheses are available, given by so-called *a priori probabilities*. These are $P(\theta_0)$ and $P(\theta_1)$, with $P(\theta_0) + P(\theta_1) = 1$

Definition. A test $d : \mathbb{X} \rightarrow \{0, 1\}$ that, given $P(\theta_0)$ and $P(\theta_1)$, minimizes the probability of a wrong decision, is called Bayes test or Bayes detector:

$$d_{\text{Bayes}} = \underset{d}{\operatorname{argmin}} \{P_\varepsilon\},$$

with $P_\varepsilon = E[G_d(\theta)(1 - \Pi(\theta)) + (1 - G_d(\theta))\Pi(\theta)]$ and

$$\Pi(\theta) = \begin{cases} 0 & ; \theta = \theta_0 \\ 1 & ; \theta = \theta_1. \end{cases}$$

For P_ε we can write:

$$\begin{aligned} P_\varepsilon &= P(\theta_0)[G_d(\theta_0)(1 - \Pi(\theta_0)) + (1 - G_d(\theta_0))\Pi(\theta_0)] + P(\theta_1)[G_d(\theta_1)(1 - \Pi(\theta_1)) + (1 - G_d(\theta_1))\Pi(\theta_1)] \\ &= P(\theta_0)G_d(\theta_0) + P(\theta_1)(1 - G_d(\theta_1)). \end{aligned}$$

The decision rule d partitions the observation space $\mathbb{X} = \mathbb{X}_{d,0} \cup \mathbb{X}_{d,1}$, such that $\mathbb{X}_{d,0} \cap \mathbb{X}_{d,1} = \emptyset$.

As $G_d(\theta) = E[d(X)|\theta] = \int_{d(\mathbb{X})} d f_d(d|\theta) dd = \int_{\mathbb{X}} d(x) f_X(x|\theta) dx$, we have:

$$\begin{aligned} P_\varepsilon &= P(\theta_0) E[d(X)|\theta_0] + P(\theta_1) (1 - E[d(X)|\theta_1]) \\ &= P(\theta_0) \int_{\mathbb{X}_{d,1}} 1 f_X(x|\theta_0) dx + P(\theta_1) - P(\theta_1) \int_{\mathbb{X}_{d,1}} 1 f_X(x|\theta_1) dx \\ &= \int_{\mathbb{X}_{d,1}} (P(\theta_0) f_X(x|\theta_0) - P(\theta_1) f_X(x|\theta_1)) dx + P(\theta_1). \end{aligned}$$

P_ε is minimized, by choosing $\mathbb{X}_{d,1}$ such, that the integrand is negative for every x . We therefore get

$$\mathbb{X}_{d,1} = \{x | P(\theta_0) f_X(x|\theta_0) - P(\theta_1) f_X(x|\theta_1) < 0\}.$$

The resulting Bayes detector is

$$d_{\text{Bayes}} = \begin{cases} 1 & ; \quad \frac{f_X(x|\theta_1)}{f_X(x|\theta_0)} > \frac{P(\theta_0)}{P(\theta_1)} \\ 0 & ; \quad \text{otherwise.} \end{cases}$$

For the special case, where $P(\theta_0) = P(\theta_1)$ the Bayes test is equivalent to the ML test.

By applying the Bayes rule,

$$P(\theta_i|x) = \frac{f_X(x|\theta_i) P(\theta_i)}{f_X(x)} \text{ (with } f_X(x) \neq 0),$$

on the test condition, we obtain

$$d_{\text{Bayes}} = \begin{cases} 1 & ; \quad P(\theta_1|x) > P(\theta_0|x), \\ 0 & ; \quad \text{otherwise.} \end{cases}$$

Therefore the Bayes test is also known as *Maximum A Posteriori Detector* or *MAP Detector*.

Bayes Risk Function

In the following we intend to minimize the risk of a wrong decision, instead of the error probability. Therefore the two error types are assigned different weights C_0 and C_1 and the risk function is

$$\begin{aligned}\text{risk}(d) &= E [C_0 G_d(\theta)(1 - \Pi(\theta)) + C_1(1 - G_d(\theta))\Pi(\theta)] \\ &= C_0 G_d(\theta_0) P(\theta_0) + C_1(1 - G_d(\theta_1)) P(\theta_1).\end{aligned}$$

The resulting decision rule is

$$d_{\text{Bayes}} = \begin{cases} 1 & ; \quad \frac{f_X(x|\theta_1)}{f_X(x|\theta_0)} > \frac{C_0 P(\theta_0)}{C_1 P(\theta_1)}, \\ 0 & ; \quad \text{otherwise.} \end{cases}$$

While for a transmitted bit, the risk of a wrong decision will have the same weight, this is different for scenarios where the alternative corresponds to a hazardous incident, in this case $\Rightarrow C_1 \gg C_0$!

15.5 Bayes Test for Multiple Hypothesis

As a preparation we start with a reformulation of the alternative test:

$$\begin{aligned}
 \text{Risk}(d) &= C_0 E [d(X)|\theta_0] P(\theta_0) + C_1 E [1 - d(X)|\theta_1] P(\theta_1) \\
 &= E [C_0 d(X)|\theta_0] P(\theta_0) + E [C_1 (1 - d(X))|\theta_1] P(\theta_1) \\
 &= E [L(d(X), \theta_0)] P(\theta_0) + E [L(d(X), \theta_1)] P(\theta_1), \\
 &= E [E [L(d(X), \theta)|\theta]] \\
 &= E [L(d(X), \theta)],
 \end{aligned}$$

where $L(d(x), \theta)$ is a loss function defined as

$$L : \{0, 1\} \times \{\theta_0, \theta_1\} \rightarrow \mathbb{R}_+ : L(d(x), \theta) \mapsto \begin{cases} C_0 & ; \text{ for a false alarm, } d(x) = 1, \text{ but } \theta = \theta_0, \\ C_1 & ; \text{ for a detection error, } d(x) = 0, \text{ but } \theta = \theta_1, \\ 0 & ; \text{ otherwise.} \end{cases}$$

In the following, we generalize this result to multiple hypothesis, where $\Theta = \{\theta_1, \dots, \theta_K\}$, and the loss function is

$$L : \{1, \dots, K\} \times \Theta \rightarrow \mathbb{R}_+ : L(d(x), \theta_k) \mapsto \begin{cases} C_k & ; d(x) \neq \theta_k, \\ 0 & ; \text{ otherwise.} \end{cases}$$

The expected risk is:

$$\text{risk}(d) = E [E [L(d(X), \theta)|\theta]] = E [L(d(X), \theta)] = E [E [L(d(x), \theta)|x = X]].$$

The risk is minimized, if by the choice of d the term $E[L(d(x), \theta)|x]$ is minimized for all $x \in \mathbb{X}$, that is $\forall x \in \mathbb{X}$ and $d(x) = k' \in \{1, \dots, K\}$:

$$d : x \mapsto \operatorname{argmin}_{k' \in \{1, \dots, K\}} \{E[L(k', \theta)|x]\}.$$

As we have,

$$\begin{aligned} E[L(k', \theta)|x] &= \sum_{k \in \{1, \dots, K\}} P(\theta_k|x) \underbrace{L(k', \theta_k)}_{\in \{C_k, 0\}} \\ &= \sum_{k \in \{1, \dots, K\} \setminus k'} P(\theta_k|x) C_k \\ &= \sum_{k \in \{1, \dots, K\} \setminus k'} P(\theta_k|x) C_k + P(\theta_{k'}|x) C_{k'} - P(\theta_{k'}|x) C_{k'} \\ &= \text{const} - P(\theta_{k'}|x) C_{k'}, \end{aligned}$$

the decision rule becomes

$$d : x \mapsto \operatorname{argmin}_{k' \in \{1, \dots, K\}} \{E[L(k', \theta)|x]\} = \operatorname{argmax}_{k' \in \{1, \dots, K\}} \{P(\theta_{k'}|x) C_{k'}\}.$$

For the special case $C_k = 1$, $\forall \theta_k \in \Theta$, the loss function corresponds to the error probability and the Bayes test is

$$d : x \mapsto \operatorname{argmax}_{k' \in \{1, \dots, K\}} \{P(\theta_{k'}|x)\}.$$

This test is known as the *maximum a posteriori (MAP) detector*.

Example

$$f_{X|\Theta}(x|\theta_k) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2}(x - \theta_k)^2\right),$$

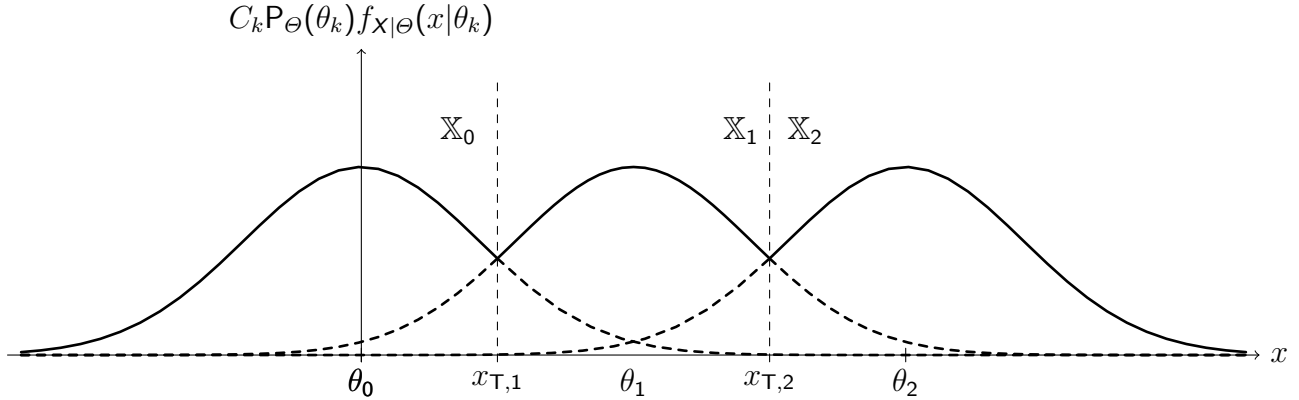
with $\theta_0 < \theta_1 < \theta_2$, a prior $P_\Theta(\theta_k)$ for Θ , and costs C_k . The Bayes test is given as

$$d(x) = \operatorname{argmax}_{k \in \{0,1,2\}} \{C_k P_{\Theta|X}(\theta_k|x)\} = \operatorname{argmax}_{k \in \{0,1,2\}} \{C_k P_\Theta(\theta_k) f_{X|\Theta}(x|\theta_k)\}.$$

The Bayes test $d(x)$ decomposes the observation space \mathbb{X} into three subsets \mathbb{X}_k with

$$d(x) = \begin{cases} 0 & ; & x \in \mathbb{X}_0 \\ 1 & ; & x \in \mathbb{X}_1 \\ 2 & ; & x \in \mathbb{X}_2. \end{cases}$$

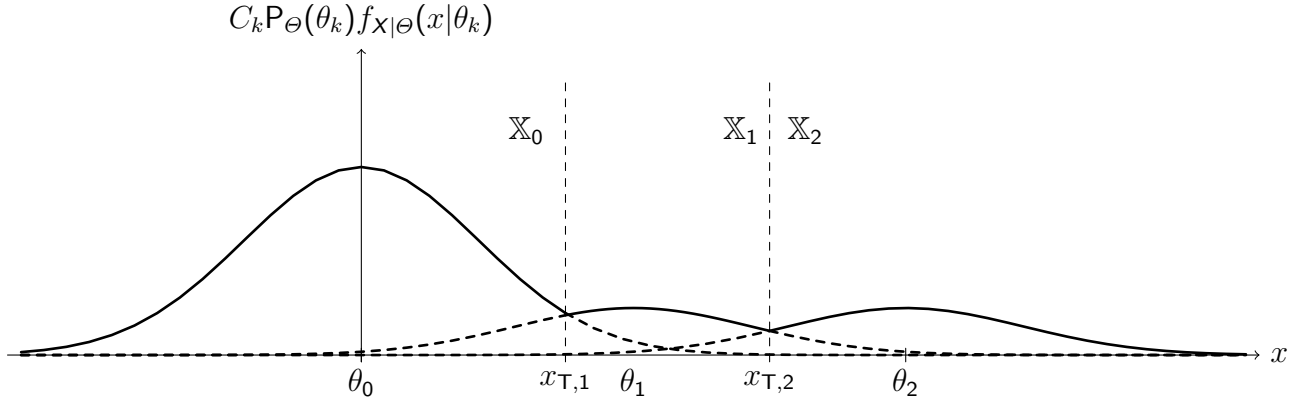
First, consider $P_{\Theta}(\theta_k) = \frac{1}{3}$ and $C_k = 1 \forall k$.



For the case of a uniform prior and equal costs, the Bayes test equals a Maximum-Likelihood test. The test statistic is given as

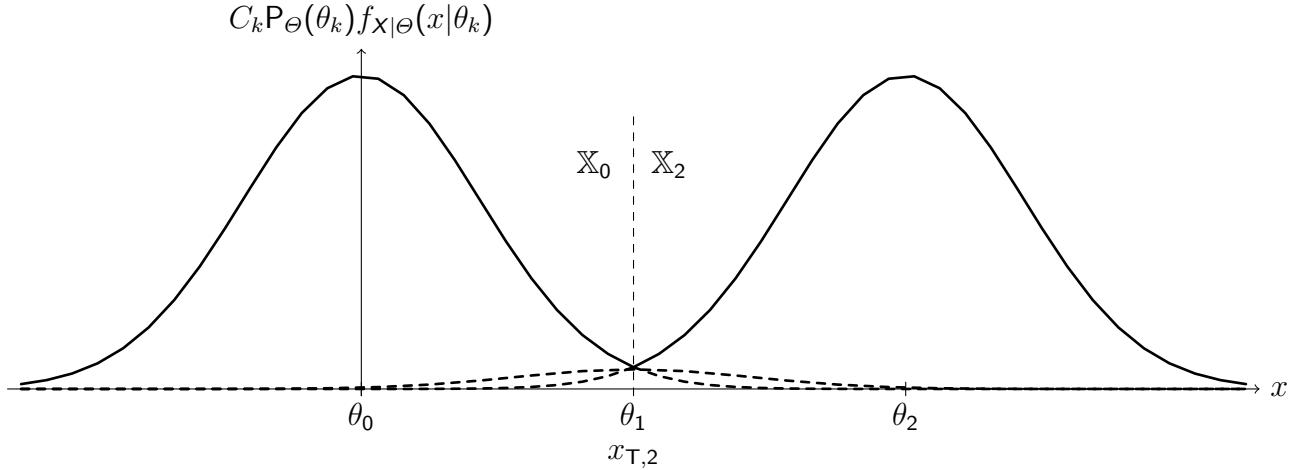
$$d(x) = \begin{cases} 0 & ; \quad x \leq x_{T,1} \\ 1 & ; \quad x_{T,1} < x < x_{T,2} \\ 2 & ; \quad x \geq x_{T,2}. \end{cases}$$

Now, consider $P_{\theta}(\theta_0) = \frac{2}{3}$, $P_{\theta}(\theta_1) = P_{\theta}(\theta_2) = \frac{1}{6}$, and $C_k = 1 \forall k$.



As $P_{\theta}(\theta_0) > P_{\theta}(\theta_1) = P_{\theta}(\theta_2)$, the threshold $x_{T,1}$ is shifted towards θ_1 .

Finally, assume that $P_{\theta}(\theta_0) = P_{\theta}(\theta_2) = \frac{48}{100}$, $P_{\theta}(\theta_1) = \frac{4}{100}$, and $C_k = 1; \forall k$.



Although $P_{\theta}(\theta_1) > 0$, the Bayes test is given as

$$d(x) = \begin{cases} 0 & ; \quad x \in \mathbb{X}_0 \\ 2 & ; \quad x \in \mathbb{X}_2 \end{cases}$$

as we have that $\mathbb{X}_1 = \emptyset$.

15.6 Linear Alternative Tests

In the following, we consider the special case of *linear alternative tests*, i.e., tests with

$$d : \mathbb{X} \rightarrow \mathbb{R}, \mathbf{x} \mapsto \begin{cases} 1 & ; \mathbf{w}^\top \mathbf{x} - w_0 > 0 \\ 0 & ; \text{otherwise.} \end{cases} \quad (15.1)$$

Motivation: Gaussian $f_{\mathbf{X}}(\mathbf{x}; \theta)$

Consider the ML test (Neyman-Pearson test with $c = 1$) for Gaussian PDFs, i.e., we have

$$f_{\mathbf{X}}(\mathbf{x}; \underbrace{\boldsymbol{\mu}_k, \mathbf{C}_k}_{\theta_k}) = \frac{1}{\sqrt{\det(2\pi\mathbf{C}_k)}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^\top \mathbf{C}_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)\right) \quad (15.2)$$

with $\theta_0 = \{\boldsymbol{\mu}_0, \mathbf{C}_0\}$ and $\theta_1 = \{\boldsymbol{\mu}_1, \mathbf{C}_1\}$. The likelihood ratio $R(\mathbf{x})$ is given as

$$R(\mathbf{x}) = \frac{\sqrt{\det(2\pi\mathbf{C}_0)}}{\sqrt{\det(2\pi\mathbf{C}_1)}} \cdot \frac{\exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^\top \mathbf{C}_1^{-1}(\mathbf{x} - \boldsymbol{\mu}_1)\right)}{\exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_0)^\top \mathbf{C}_0^{-1}(\mathbf{x} - \boldsymbol{\mu}_0)\right)} \geq 1. \quad (15.3)$$

Using the log-likelihood ratio $\log R(\mathbf{x})$, the separating surface in $\mathbb{X} = \mathbb{R}^N$ which decomposes \mathbb{X} into \mathbb{X}_0 and \mathbb{X}_1 can be described as the set of all $\mathbf{x} \in \mathbb{X}$ for which

$$\log R(\mathbf{x}) = \frac{1}{2} \ln \left(\frac{\det(\mathbf{C}_0)}{\det(\mathbf{C}_1)} \right) + \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_0)^\top \mathbf{C}_0^{-1}(\mathbf{x} - \boldsymbol{\mu}_0) - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^\top \mathbf{C}_1^{-1}(\mathbf{x} - \boldsymbol{\mu}_1) = 0. \quad (15.4)$$

Motivation: $\log R(x)$ as an affine function of x

In the special case of $C_0 = C_1 = C$, the quadratic terms in $\log R(x)$ cancel out, and we have

$$\log R(x) = \frac{1}{2} \ln \left(\frac{\det(C_0)}{\det(C_1)} \right) + \frac{1}{2}(x - \mu_0)^\top C_0^{-1}(x - \mu_0) - \frac{1}{2}(x - \mu_1)^\top C_1^{-1}(x - \mu_1) \quad (15.5)$$

$$= \mathbf{w}^\top x - w_0 \geq 0 \quad (15.6)$$

where the normal vector of the separating hyperplane is given as

$$\mathbf{w}^\top = (\mu_1 - \mu_0)^\top C^{-1} \quad (15.7)$$

and the constant translation w_0 with respect to the origin is given as

$$w_0 = \frac{(\mu_1 - \mu_0)^\top C^{-1}(\mu_1 + \mu_0)}{2}. \quad (15.8)$$

The ML estimator reduces to the affine decision rule given as

$$d : \mathbb{X} \rightarrow \mathbb{R}, x \mapsto \begin{cases} 1 & ; \mathbf{w}^\top x - w_0 > 0 \\ 0 & ; \text{otherwise.} \end{cases} \quad (15.9)$$

Motivation: $\log R(x)$ as a linear function of x

If in addition to $C_0 = C_1$, we have that $\mu_0 = -\mu_1$, the log-likelihood ratio simplifies further to

$$\log R(x) = w^\top x \geq 0 \quad (15.10)$$

as from $\mu_0 = -\mu_1$ it follows that $w_0 = 0$. Therefore, the ML estimator reduces to the linear decision rule

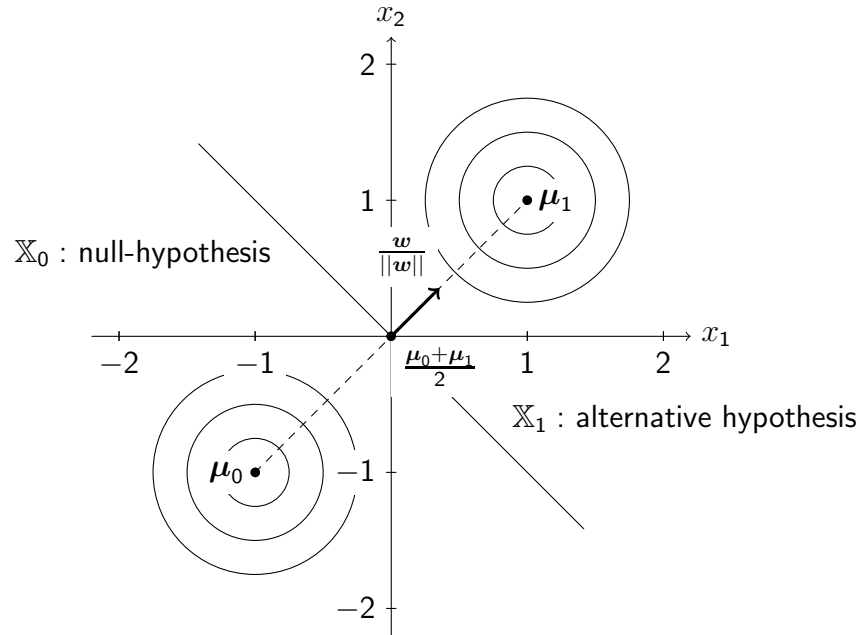
$$d : \mathbb{X} \rightarrow \mathbb{R}, x \mapsto \begin{cases} 1 & ; w^\top x > 0 \\ 0 & ; \text{otherwise.} \end{cases} \quad (15.11)$$

For Gaussian $f_X(x; \mu_k, C_k)$ with $\theta_0 = \{\mu_0, C_0\}$ and $\theta_1 = \{\mu_1, C_1\}$, it follows that

- if $C_0 \neq C_1$, $\log R(x) = 0$ is non-linear and the separating surfaces are surfaces of second order: parabolic, hyperbolic, or elliptic surfaces.
- if $C_0 = C_1$, $\log R(x) = 0$ is affine and thus defines a hyperplane in \mathbb{X} which decomposes \mathbb{X} into \mathbb{X}_0 and \mathbb{X}_1 .
- if $C_0 = C_1$ and $\mu_0 = -\mu_1$, $\log R(x) = 0$ is linear and defines a separating hyperplane in \mathbb{X} which contains the origin.

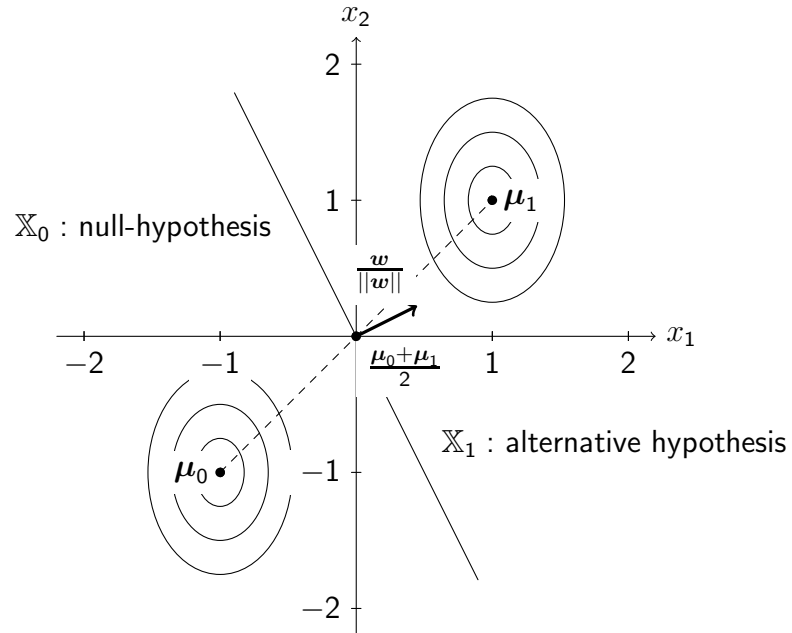
Example 1:

$$\mu_0 = (-1, -1)^\top, \mu_1 = (1, 1)^\top, C = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, w^\top = (2, 2)$$



Example 2:

$$\mu_0 = (-1, -1)^\top, \mu_1 = (1, 1)^\top, C = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}, w^\top = (2, 1)$$



Statistical Properties of $T(\mathbf{X}) = \mathbf{w}^\top \mathbf{X}$

Consider an alternative test with Gaussian $f_{\mathbf{X}}(\mathbf{x}; \theta_k)$, with $\theta_0 = \{\boldsymbol{\mu}_0, \mathbf{C}_0\}$, and $\theta_1 = \{\boldsymbol{\mu}_1, \mathbf{C}_1\}$ where $\mathbf{C} = \mathbf{C}_0 = \mathbf{C}_1$. In the following, we will discuss a test statistic given as

$$T(\mathbf{X}) = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^\top \mathbf{C}^{-1} \mathbf{X} = \mathbf{w}^\top \mathbf{X}. \quad (15.12)$$

As \mathbf{X} is Gaussian, $T(\mathbf{X})$ is Gaussian as well (linear transformation of Gaussian RV is Gaussian). Thus, the distribution of $T(\mathbf{X})$ is determined by

$$\mathbb{E}[T(\mathbf{X})] = \mathbb{E}[\mathbf{w}^\top \mathbf{X}] = \mathbf{w}^\top \boldsymbol{\mu}_k = \mu_{T,k}, \quad k = 0, 1, \quad (15.13)$$

$$\text{Var}[T(\mathbf{X})] = \mathbf{w}^\top \text{Var}[\mathbf{X}] \mathbf{w} = \mathbf{w}^\top \mathbf{C} \mathbf{w} = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^\top \mathbf{C}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) = \sigma_T^2. \quad (15.14)$$

The PDF of $T = T(\mathbf{X})$ is given as

$$f_{T,k}(t) = \frac{1}{\sqrt{2\pi}\sigma_T} \exp\left(-\frac{1}{2\sigma_T^2}(t - \mu_{T,k})^2\right). \quad (15.15)$$

The ML test can be expressed as

$$\begin{aligned}
\ln(R(T)) &= \frac{f_{T,1}(t)}{f_{T,0}(t)} = \ln \left(\frac{\exp \left(-\frac{1}{2\sigma_T^2} (T - \mu_{T,1})^2 \right)}{\exp \left(-\frac{1}{2\sigma_T^2} (T - \mu_{T,0})^2 \right)} \right) \geq 0 \\
&\Rightarrow (\mu_{T,1} - \mu_{T,0})T - \frac{(\mu_{T,1}^2 - \mu_{T,0}^2)}{2} \geq 0 \\
&\Rightarrow T - \frac{\mu_{T,1} + \mu_{T,0}}{2} \geq 0 \\
&\Rightarrow T - w_0 \geq 0 \\
&\Rightarrow T(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} \geq w_0
\end{aligned}$$

Note that that $w_0 = \frac{\mu_{T,1} + \mu_{T,0}}{2} = \frac{(\mu_1 - \mu_0)^\top C^{-1}(\mu_1 + \mu_0)}{2}$ as $\mu_{T,k} = \mathbf{w}^\top \mu_k$ and $\mathbf{w}^\top = (\mu_1 - \mu_0)^\top C^{-1}$.

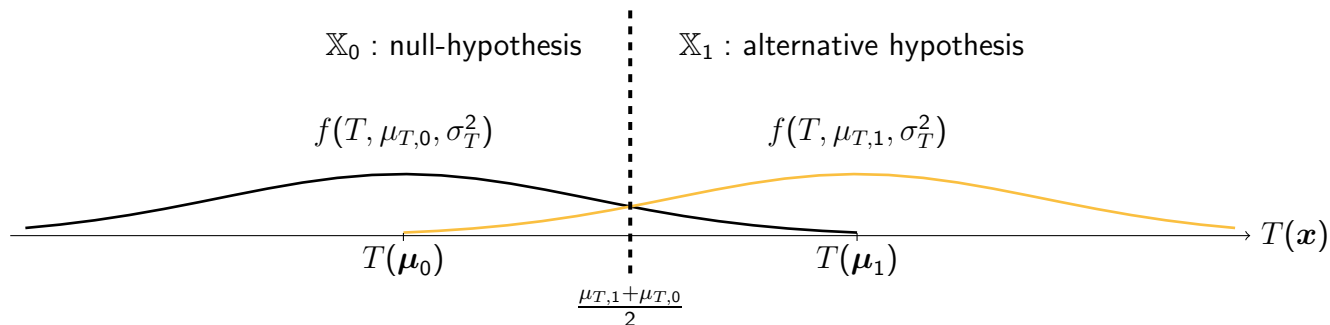


Fig. 15.1: Illustration of the decision rule $T(x) \geq \frac{\mu_{T,1} + \mu_{T,0}}{2}$ for Gaussian PDFs with $\mu_0 = -\mu_1$ and $C_0 = C_1$.

15.7 Sufficient Statistics

An important property of a test statistic is whether it is *sufficient* in order to decide for a parameter $\theta \in \Theta$ or not. In this context, sufficiency means that no other test statistic, i.e., function of the observations, contains additional information about the parameter to be estimated.

Definition: Sufficient Statistic

Consider a statistical model $\{\mathbb{X}, \mathbb{F}, P_\theta; \theta \in \Theta\}$. Let \mathbf{X} be a random sample drawn from $f_{\mathbf{X}}(\mathbf{x}; \theta)$. If for a test statistic $T(\mathbf{X})$, $f_{\mathbf{X}|T}(\mathbf{x}|T(\mathbf{x}) = t)$ is independent of θ for all $\theta \in \Theta$, i.e., if

$$f_{\mathbf{X}|T}(\mathbf{x}|T(\mathbf{x}) = t, \theta) = f_{\mathbf{X}|T}(\mathbf{x}|T(\mathbf{x}) = t), \quad (15.16)$$

the test statistic $T(\mathbf{X})$ is a SUFFICIENT STATISTIC for the underlying parameter θ , or more precisely, for the family of probability distributions P_θ which are parametrized by θ .

Example:

Transmission of $d = \pm 1$ from a source to a sink:

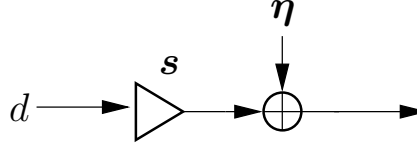


Fig. 15.2: System model.

The receive signal at the sink is given as

$$\mathbf{x} = d\mathbf{s} + \boldsymbol{\eta}. \quad (15.17)$$

As $d \in \{-1, +1\}$ and $\boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \mathbf{C})$, the ML test is given as:

$$R(\mathbf{x}) = \frac{\exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{s})^T \mathbf{C}^{-1}(\mathbf{x} - \mathbf{s})\right)}{\exp\left(-\frac{1}{2}(\mathbf{x} + \mathbf{s})^T \mathbf{C}^{-1}(\mathbf{x} + \mathbf{s})\right)} \gtrless 1 \quad (15.18)$$

$$\Leftrightarrow \log R(\mathbf{x}) = \mathbf{s}^T \mathbf{C}^{-1} \mathbf{x} \gtrless 0 \quad (15.19)$$

Using a *Whitening Matched Filter* $\mathbf{g}^T = \mathbf{s}^T \mathbf{C}^{-1}$, the ML test is given as³:

$$d : \mathbf{x} \mapsto \begin{cases} +1 & ; \mathbf{g}^T \mathbf{x} > 0 \\ -1 & ; \text{otherwise.} \end{cases} \quad (15.20)$$

³Here, we define the test as a function which maps to the transmit alphabet $\{-1, +1\}$ directly.

Remark: The test statistic $G(\mathbf{x}) = \mathbf{g}^\top \mathbf{x} = \mathbf{s}^\top \mathbf{C}^{-1} \mathbf{x}$ is a sufficient statistic, i.e., for $G(\mathbf{x}) = c$, it follows that $f_{\mathbf{X}}(\mathbf{x} | G(\mathbf{x}) = c)$ is independent of d .

Sketch of Proof: Assume that $G(\mathbf{x}) = 0$, i.e., that

$$\mathbf{x} \in \{\mathbf{x} | G(\mathbf{x}) = 0\} = \{\mathbf{x} | \mathbf{g}^\top \mathbf{x} = 0\} = \{\mathbf{x} | \mathbf{g}^\perp t, t \in \mathbb{R}\}, \quad (15.21)$$

where \mathbf{g}^\perp such that $\mathbf{g}^\top \mathbf{g}^\perp = 0$. It remains to show that $f(\mathbf{x} | G(\mathbf{x}) = 0; d)$ is independent of d :

$$\begin{aligned} f(\mathbf{x} | \mathbf{x} = \mathbf{g}^\perp t; d) &\propto \exp \left[-\frac{1}{2} (\mathbf{g}^\perp t - \mathbf{s}d)^\top \mathbf{C}^{-1} (\mathbf{g}^\perp t - \mathbf{s}d) \right] \\ &\propto \exp \left[-\frac{1}{2} \mathbf{g}^{\perp, \top} \mathbf{C}^{-1} \mathbf{g}^\perp t^2 + d \underbrace{\mathbf{s}^\top \mathbf{C}^{-1} \mathbf{g}^\perp}_{=0} t \right] \\ &\propto \exp \left[-\frac{1}{2} \mathbf{g}^{\perp, \top} \mathbf{C}^{-1} \mathbf{g}^\perp t^2 \right] = f(\mathbf{x} | \mathbf{x} = \mathbf{g}^\perp t) = f_t(t). \end{aligned}$$

Appendix

A. Probability Theory and Stochastics

A.1 Probability Space

A PROBABILITY SPACE,

$$(\Omega, \mathbb{F}, P), \tag{A.1}$$

represents a mathematical model consisting of an

- OBSERVATION SPACE Ω :
the nonempty set of potential outputs/observations of an experiment,
- SIGMA ALGEBRA \mathbb{F} :
a set of subsets of potential outputs/observations of an experiment,
- PROBABILITY MEASURE P :
a function which maps each element of $A \in \mathbb{F}$ to the interval $[0, 1]$.

Definition. A SIGMA ALGEBRA \mathbb{F} is a set of subsets (EVENTS) of Ω with

$$\Omega \in \mathbb{F}, \quad (\text{A.2})$$

$$A \in \mathbb{F} \Rightarrow A^c \in \mathbb{F}, \quad (\text{A.3})$$

$$A_1, \dots, A_k \in \mathbb{F} \Rightarrow \bigcup_{i=1}^k A_i \in \mathbb{F}. \quad (\text{A.4})$$

Definition. A PROBABILITY MEASURE P maps $A \in \mathbb{F}$ to $P(A)$, the PROBABILITY of the EVENT A , with

$$\text{Nonnegativity :} \quad P(A) \geq 0, \quad (\text{A.5})$$

$$\text{Normation :} \quad P(\Omega) = 1, \quad (\text{A.6})$$

$$\text{Additivity :} \quad P\left(\bigcup_{i=1}^{\infty} A_i\right) = \sum_{i=1}^{\infty} P(A_i), \quad \text{wenn } A_i \cap A_j = \emptyset, \forall i \neq j. \quad (\text{A.7})$$

Consequently, the following general properties of probability measures are obtained:

$$P(A^c) = 1 - P(A), \quad (\text{A.8})$$

$$P(\emptyset) = 0, \quad (\text{A.9})$$

$$P(A \setminus B) = P(A \cap B^c) = P(A) - P(A \cap B), \quad (\text{A.10})$$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B), \quad (\text{A.11})$$

$$A \subset B \Rightarrow P(A) \leq P(B), \quad (\text{A.12})$$

$$P\left(\bigcup_{i=1}^k A_i\right) \leq \sum_{i=1}^k P(A_i). \quad (\text{A.13})$$

A.2 Conditional Probability and Stochastical Independence

Conditional Probability

The probability of an event changes when we are informed about the realization of another event within the same probability space.

Definition. Given $P(A) > 0$ for $A \in \mathbb{F}$, the **CONDITIONAL PROBABILITY** of B conditioned by A is defined as

$$P(B | A) \triangleq P_A(B) = \frac{P(A \cap B)}{P(A)}. \quad (\text{A.14})$$

Bayes Theorem

Given a partition of Ω in disjoint sets $B_i, i \in I$, such that

$$\bigcup_{i \in I} B_i = \Omega, \quad (\text{A.15})$$

$$B_i \cap B_j = \emptyset, \forall i \neq j, \quad (\text{A.16})$$

then for any $A \in \mathbb{F}$ with $P(A) > 0$:

THEOREM OF TOTAL PROBABILITY :	$P(A) = \sum_{i \in I} P(A B_i) P(B_i), \quad (\text{A.17})$
BAYES THEOREM :	$P(B_j A) = \frac{P(A B_j) P(B_j)}{\sum_{i \in I} P(A B_i) P(B_i)}. \quad (\text{A.18})$

Special Case. Given $A, B \subset \Omega$ and $\Omega = B \cup B^c$, the BAYES THEOREM is given by

$$P(B | A) = \frac{P(A | B) P(B)}{P(A)} = \frac{P(A | B) P(B)}{P(A | B) P(B) + P(A | B^c) (1 - P(B))}. \quad (\text{A.19})$$

Illustration

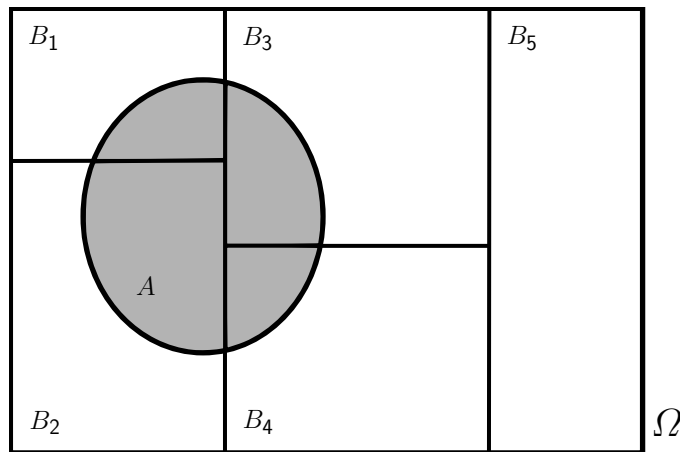


Fig. A.1: Partition of Ω in B_1, \dots, B_5

$$P(A) = \sum_{i=1}^4 P(A | B_i) P(B_i) + 0 \cdot P(B_5), \quad (\text{A.20})$$

$$P(B_2 | A) = \frac{P(A | B_2) P(B_2)}{\sum_{i=1}^4 P(A | B_i) P(B_i)}. \quad (\text{A.21})$$

Stochastical Independence

A und B are stochastically independent if the probability measure for B is not changed by any knowledge about the realization of A , and vice versa, i.e., given $P(A) > 0$,

$$P(B | A) = P(B). \quad (\text{A.22})$$

Otherwise, with

$$P(B | A) = \frac{P(A \cap B)}{P(A)} \quad (\text{A.23})$$

we obtain the following

Definition. Given (Ω, \mathbb{F}, P) is a probability space, then $A, B \in \mathbb{F}$ are STOCHASTICALLY INDEPENDENT with respect to P , if

$$P(A \cap B) = P(A) P(B). \quad (\text{A.24})$$

A.3 Random Variables and Distributions

Definition. Given a probability space (Ω, \mathbb{F}, P) and a measure space (Ω', \mathbb{F}') , we call the mapping $X : \Omega \rightarrow \Omega'$ a **RANDOM VARIABLE**, if and only if for each $A' \in \mathbb{F}'$ there exists an element $A \in \mathbb{F}$, such that $A = \{\omega \in \Omega \mid X(\omega) \in A'\} \in \mathbb{F}$ or — in shorter notation — $A = \{X \in A'\} \in \mathbb{F}$.

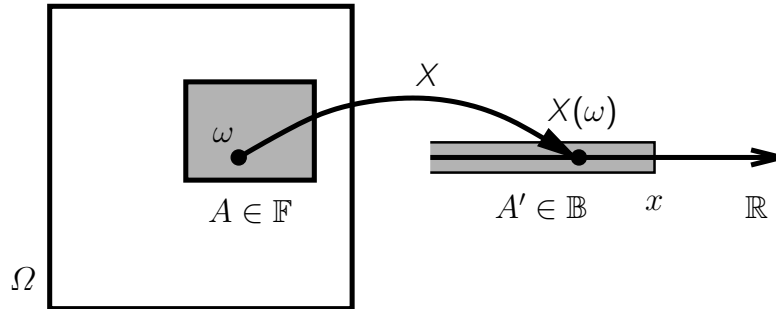


Fig. A.2: Real-valued random variable: $x : \Omega \rightarrow \mathbb{R}$

Distribution of Real-Valued Random Variables

Definition. Given a real-valued random variable X on (Ω, \mathbb{F}, P) , the function $F_X : \mathbb{R} \rightarrow [0, 1]$, with

$$F_X(x) = P(\{X \leq x\}), \quad (\text{A.25})$$

defines the CUMULATIVE DISTRIBUTION FUNCTION (CDF) of X .

Properties of a CDF $F_X : \mathbb{R} \rightarrow [0, 1]$:

- $F_X(x)$ is monotonically increasing.
- $\lim_{h \rightarrow 0} F_X(x + h) = F_X(x), \quad \forall x \in \mathbb{R}.$
- $\lim_{x \rightarrow -\infty} F_X(x) = 0, \quad \lim_{x \rightarrow \infty} F_X(x) = 1.$

Definition. We call a random variable X CONTINUOUS if its CDF F_X can be determined by means of

$$F_X(x) = \int_{-\infty}^x f_X(\xi) \, \mathrm{d}\xi, \quad (\text{A.26})$$

with $f_X : \mathbb{R} \rightarrow [0, \infty[$, which is the PROBABILITY DENSITY FUNCTION (PDF) of X .

If F_X is continous and can be differentiated almost everywhere, then the random variable X is continuous and

$$f_X(x) = \frac{\mathrm{d}F_X(x)}{\mathrm{d}x}. \quad (\text{A.27})$$

Examples

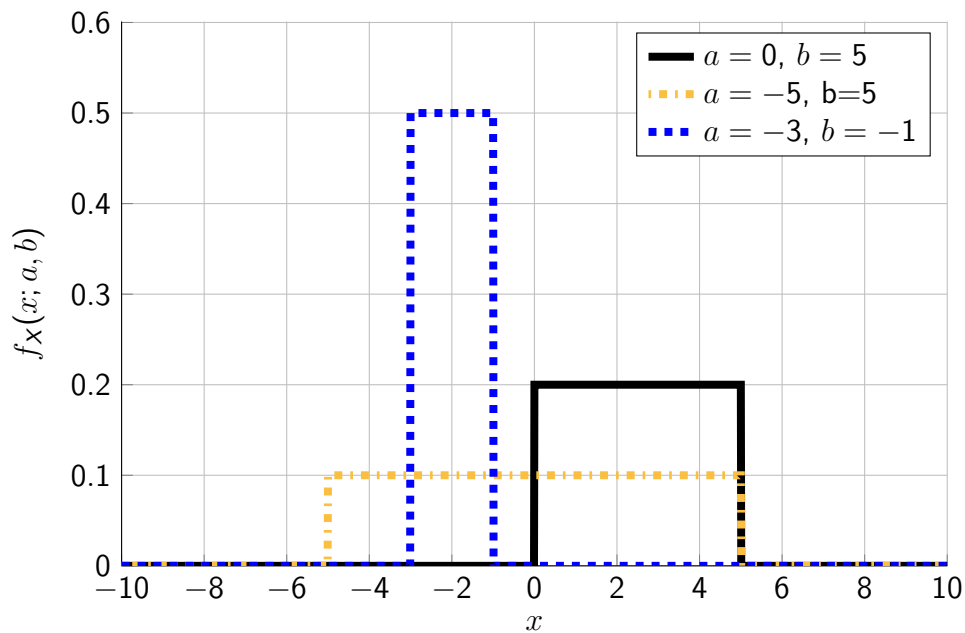


Fig. A.3: Uniform PDF: $f_X(x) = \frac{1}{b-a}, a \leq x \leq b$

Examples (Cont'd)

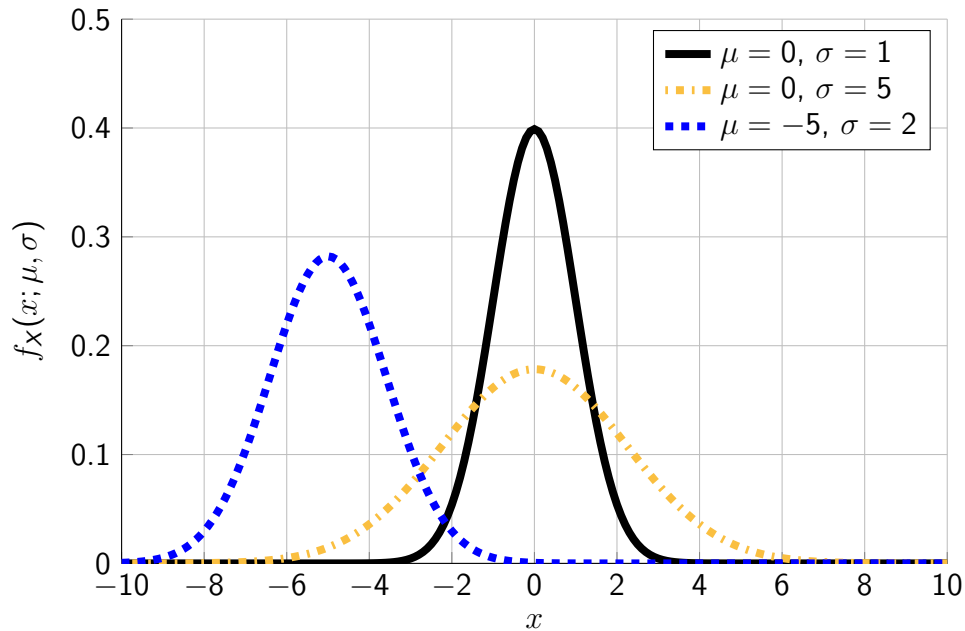


Fig. A.4: Gaussian PDF: $f_X(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, \quad x \in \mathbb{R}$

Examples (Cont'd)

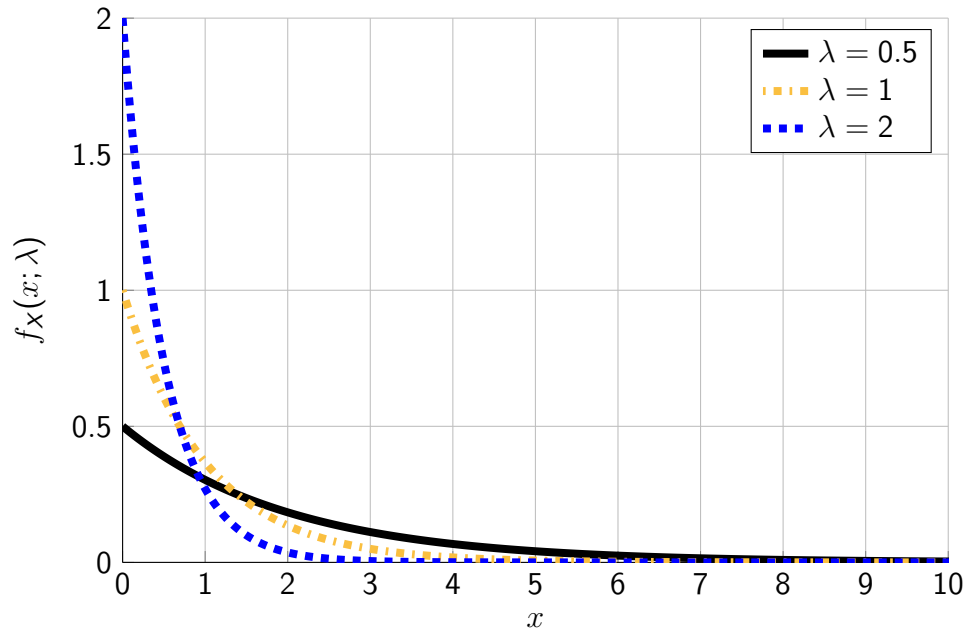


Fig. A.5: Exponential PDF: $f_X(x) = \lambda e^{-\lambda x}$, $x \geq 0$

Multidimensional Distributions

Definition. Given a n -dimensional real-valued random variable (random vector) $\mathbf{X} = [X_1, \dots, X_n]^T$ on the probability (Ω, \mathbb{F}, P) , then the JOINT CDF $F_{X_1, \dots, X_n} : \mathbb{R}^n \rightarrow [0, 1]$ is denoted as

$$F_{\mathbf{X}}(\mathbf{x}) = P(\{\mathbf{X} \leq \mathbf{x}\}) \triangleq P(\{X_1 \leq x_1, \dots, X_n \leq x_n\}). \quad (\text{A.28})$$

The respective JOINT PDF is defined as

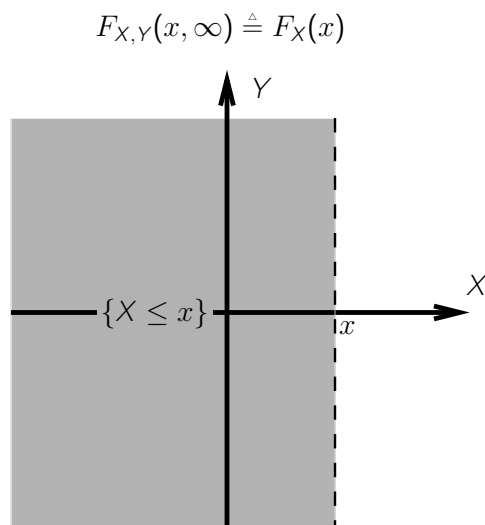
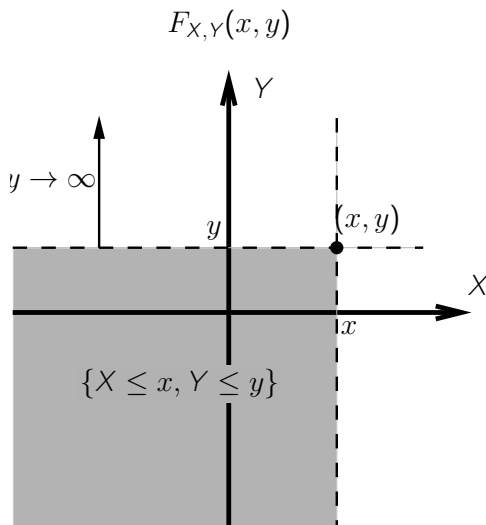
$$F_{X_1, \dots, X_n}(x_1, \dots, x_n) = \int_{-\infty}^{x_1} \dots \int_{-\infty}^{x_n} f_{X_1, \dots, X_n}(\xi_1, \dots, \xi_n) d\xi_n \dots d\xi_1 \quad (\text{A.29})$$

$$f_{X_1, \dots, X_n}(x_1, \dots, x_n) = \frac{\partial^n F_{\mathbf{X}}(x_1, \dots, x_n)}{\partial x_1 \dots \partial x_n}. \quad (\text{A.30})$$

Marginal CDF and PDF

Definition. Given a pair of r.v. $[X, Y]$, the MARGINAL CDF and PDF of X can be obtained by

$$F_X(x) = \lim_{y \rightarrow \infty} F_{X,Y}(x, y), \quad f_X(x) = \lim_{y \rightarrow \infty} \int_{-\infty}^y f_{X,Y}(x, \eta) d\eta = \int_{-\infty}^{\infty} f_{X,Y}(x, y) dy. \quad (\text{A.31})$$



Conditional CDF and PDF

Definition. Given two jointly distributed random variables X and Y , then the **CONDITIONAL CDF** of X conditioned on $\{Y = y\}$ is given by $F_{X|Y} : \mathbb{R} \rightarrow [0, 1]$, with

$$F_{X|Y}(x|y) = P(\{X \leq x\} | \{Y = y\}), \quad \text{if } Y \text{ only takes discrete values.} \quad (\text{A.32})$$

If Y is a continuous random variable, the probability of $\{Y = y\}$ is zero, and thus the **CONDITIONAL CDF** must be defined alternatively as

$$F_{X|Y}(x|y) = \int_{-\infty}^x \frac{f_{X,Y}(\xi, y)}{f_Y(y)} d\xi, \quad (\text{A.33})$$

assuming the existence of the PDFs and $f_Y(y) > 0$. Obviously, the **CONDITIONAL PDF** is

$$f_{X|Y}(x|y) = \frac{f_{X,Y}(\xi, y)}{f_Y(y)}. \quad (\text{A.34})$$

A.4 Mean and Covariance

Mean Value of Random Variables

Definition. The MEAN VALUE of a continuous real-valued random variable X is

$$\mu_X = E[X] = \int_{-\infty}^{\infty} x f_X(x) dx. \quad (\text{A.35})$$

Properties.

$$E[\alpha X + \beta Y] = \alpha E[X] + \beta E[Y], \quad (\text{A.36})$$

$$X \leq Y \Rightarrow E[X] \leq E[Y], \quad (\text{A.37})$$

for all $\alpha, \beta \in \mathbb{R}$.

Definition. The VARIANCE of random variable X is

$$\sigma_X^2 = \text{Var} [X] = \text{E} \left[(X - \text{E} [X])^2 \right] = \text{E} [X^2] - \text{E} [X]^2, \quad (\text{A.38})$$

Definition. The COVARIANCE and CORRELATION of two random variables X and Y is

$$c_{X,Y} = \text{Cov} [X, Y] = \text{E} [(X - \text{E} [X]) (Y - \text{E} [Y])] = \text{E} [XY] - \text{E} [X] \text{E} [Y], \quad (\text{A.39})$$

$$r_{X,Y} = \text{E} [XY]. \quad (\text{A.40})$$

Properties. Given the random variables X, Y, U and V and $\alpha, \beta, \gamma, \delta \in \mathbb{R}$, it follows

$$\text{Cov} [\alpha X + \beta, \gamma Y + \delta] = \alpha \gamma \text{Cov} [X, Y], \quad (\text{A.41})$$

$$\text{Cov} [X + U, Y + V] = \text{Cov} [X, Y] + \text{Cov} [X, V] + \text{Cov} [U, Y] + \text{Cov} [U, V]. \quad (\text{A.42})$$

Special Case. $\text{Var} [\alpha X + \beta] = \alpha^2 \text{Var} [X]$.

Correlation of Random Variables

Definition. Random variables X and Y are UNCORRELATED if

$$c_{X,Y} = 0 \quad \Leftrightarrow \quad E[XY] = E[X]E[Y], \quad (\text{A.43})$$

otherwise we call X and Y CORRELATED.

Definition. By normalization we obtain the CORRELATION COEFFICIENT $\rho_{X,Y} \in [-1, 1]$

$$\rho_{X,Y} = \frac{c_{X,Y}}{\sigma_X \sigma_Y}. \quad (\text{A.44})$$

Stochastic interrelations of random variables:

$$\text{INDEPENDENT : } F_{XY} = F_X F_Y \quad (\text{A.45})$$

$$\text{UNCORRELATED : } \mu_{XY} = \mu_X \mu_Y \quad (\text{A.46})$$

$$\text{ORTHOGONAL : } \mu_{XY} = 0. \quad (\text{A.47})$$

Example

The joint CDF of n real-valued Gaussian random variables $\mathbf{X} = [X_1, \dots, X_n]^\top$ is given by

$$f_{X_1, \dots, X_n}(x_1, \dots, x_n) = \frac{1}{\sqrt{2\pi}^n \sqrt{\det \mathbf{C}}} \exp \left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \mathbf{C}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right). \quad (\text{A.48})$$

The respective MEAN VECTOR and COVARIANCE MATRIX is given by

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_n \end{bmatrix} = \begin{bmatrix} \mathbb{E}[X_1] \\ \vdots \\ \mathbb{E}[X_n] \end{bmatrix} = \mathbb{E} \begin{bmatrix} X_1 \\ \vdots \\ X_n \end{bmatrix} = \mathbb{E}[\mathbf{X}] \quad (\text{A.49})$$

and

$$\mathbf{C} = \mathbb{E}[(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^\top] = \mathbb{E} \left[\begin{bmatrix} X_1 - \mu_1 \\ \vdots \\ X_n - \mu_n \end{bmatrix} [X_1 - \mu_1, \dots, X_n - \mu_n] \right] = \mathbf{C}^\top,$$

respectively. A common shorthand notation is

$$\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{C}).$$

Covariance matrices of uncorrelated random variables X_1, \dots, X_n are diagonal

$$\mathbf{C} = \text{diag} [\sigma_1^2, \dots, \sigma_n^2], \quad (\text{A.50})$$

with variances $\sigma_i^2 = \text{Var} [X_i]$, $i \in \{1, \dots, n\}$ as main diagonal elements.

Given n uncorrelated Gaussian random variables, the joint PDF is

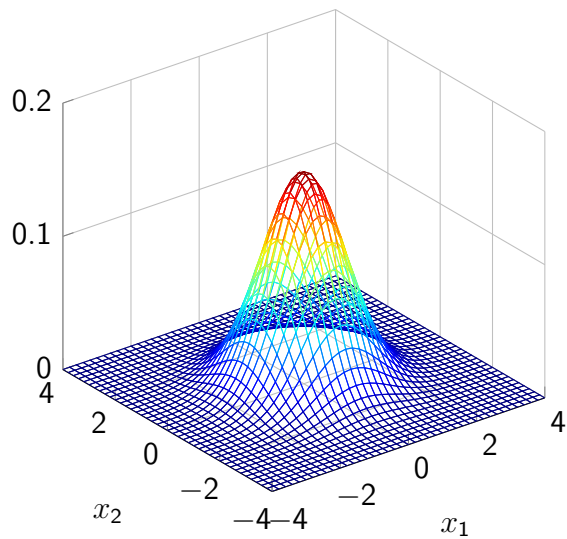
$$\begin{aligned} f_{X_1, \dots, X_n}(x_1, \dots, x_n) &= \frac{1}{\prod_{i=1}^n \sqrt{2\pi} \sqrt{\sigma_i^2}} \exp \left(-\frac{1}{2} \sum_{i=1}^n (x_i - \mu_i) \sigma_i^{-2} (x_i - \mu_i) \right) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi \sigma_i^2}} e^{-\frac{(x_i - \mu_i)^2}{2\sigma_i^2}} = \prod_{i=1}^n f_{X_i}(x_i). \end{aligned}$$

Special Case.

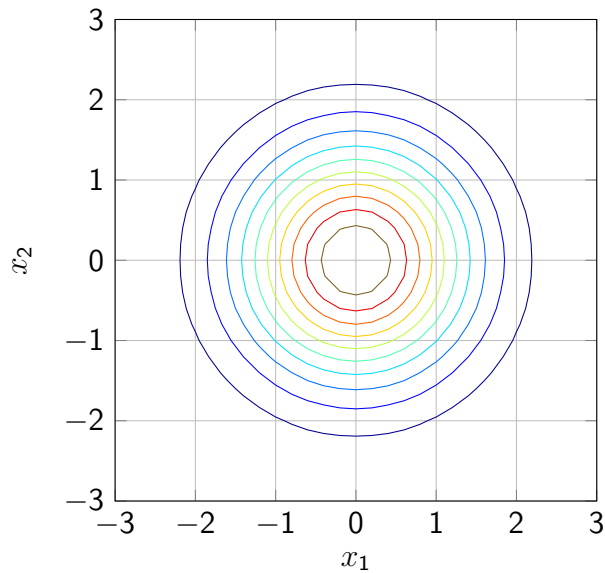
UNCORRELATED Gaussian random variables are always STOCHASTICALLY INDEPENDENT!

Correlation between the random variables X_1, \dots, X_n leads to a non-diagonal covariance matrix \mathbf{C} . Since \mathbf{C} is in general POSITIVE SEMIDEFINITE (PSD), $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{C})$ can be transformed to $\mathbf{Y} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, by the linear transformation

$$\mathbf{Y} = \mathbf{C}^{-\frac{1}{2}} (\mathbf{X} - \boldsymbol{\mu}). \quad (\text{A.51})$$

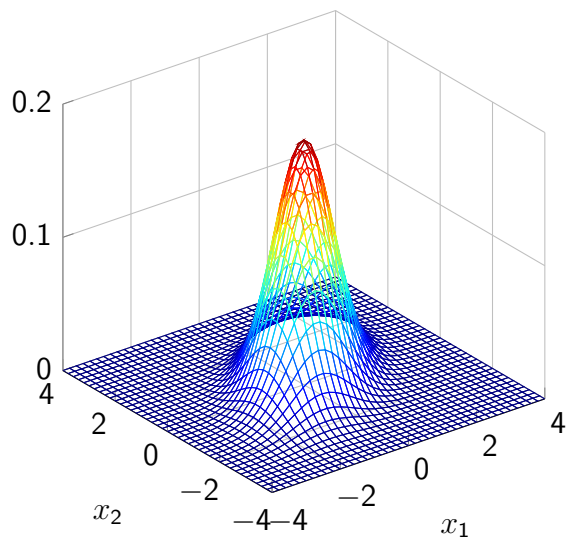


(a) Joint PDF $f_{X_1, X_2}(x_1, x_2)$

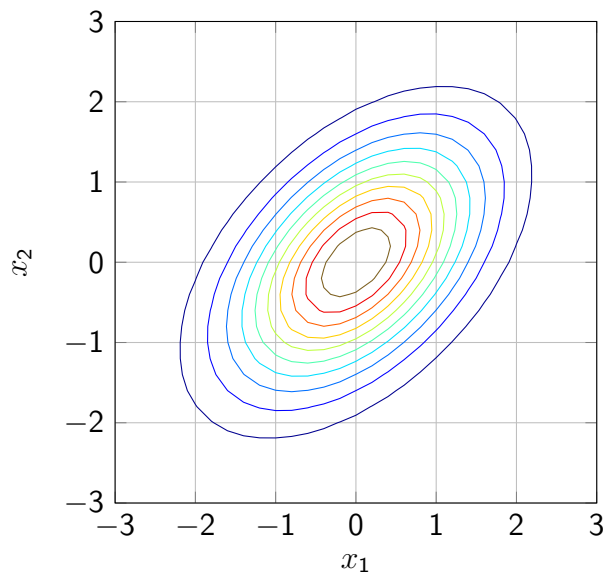


(b) Contour lines of $f_{X_1, X_2}(x_1, x_2)$

Fig. A.6: Joint PDF of the UNCORRELATED bivariate standard normal distribution.



(a) Joint PDF $f_{X_1, X_2}(x_1, x_2)$



(b) Contour lines of $f_{X_1, X_2}(x_1, x_2)$

Fig. A.7: Joint PDF of the CORRELATED bivariate standard normal distribution.

A.5 Conditional Expectation

Definition. The CONDITIONAL EXPECTATION of a random variable X , given an observation y of random variable Y is

$$E[X|Y = y] = \int_{-\infty}^{\infty} x f_{X|Y}(x|y) dx. \quad (\text{A.52})$$

The CONDITIONAL EXPECTATION $E[X|Y]$ is obviously a function of the observation Y , i.e., $E[X|Y]$ is still a random variable due to the random nature of the conditional argument.

The EXPECTATION of $E[X|Y]$ with respect to the random variable Y obtains the expectation of the unconditioned random variable X , i.e.,

$$E[E[X|Y]] = E[X]. \quad (\text{A.53})$$

Proof.

$$\begin{aligned}
 E \left[\int_{-\infty}^{\infty} x f_{X|Y}(x|Y) dx \right] &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x f_{X|Y}(x|y) dx f_Y(y) dy \\
 &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x \underbrace{f_{X|Y}(x|y) f_Y(y)}_{f_{X,Y}(x,y)} dy dx \\
 &= \int_{-\infty}^{\infty} x \underbrace{\int_{-\infty}^{\infty} f_{X,Y}(x,y) dy}_{\text{MARGINALIZATION}} dx \\
 &= \int_{-\infty}^{\infty} x f_X(x) dx \\
 &= E[X].
 \end{aligned}$$

The general version of this results for a given function $g(X, Y)$ is

$$E[E[g(X, Y)|Y]] = E[E[g(X, Y)|X]] = E[g(X, Y)], \quad (\text{A.54})$$

which in the special case above results into $E[E[X|Y]] = E[X]$.

B. Linear Operators – Matrices

B.1 Definition

Definition. A LINEAR OPERATOR is a function $A : \mathbb{S} \rightarrow \mathbb{S}'$ between VECTOR SPACES \mathbb{S} and \mathbb{S}' , such that for all $x, y \in \mathbb{S}$ and $\alpha \in \mathbb{C}$ (or \mathbb{R}) the following holds,

$$\text{ADDITIVITY : } A(x + y) = Ax + Ay \quad (\text{B.1})$$

$$\text{SCALABILITY : } A(\alpha x) = \alpha Ax. \quad (\text{B.2})$$

In this tutorial, we only consider LINEAR OPERATORS between finite-dimensional COMPLEX VECTOR SPACES \mathbb{C}^N and \mathbb{C}^M with $\alpha \in \mathbb{C}$. Consequently, a LINEAR OPERATOR A is represented by a MATRIX $\mathbf{A} : \mathbb{C}^N \rightarrow \mathbb{C}^M$.

Derivation.

Assume $\{u_n\}_{n=1}^N$ and $\{u'_m\}_{m=1}^M$ is an ORTHONORMAL BASIS of an abstract¹ N -dimensional vector space \mathbb{S} and the abstract M -dimensional vector space \mathbb{S}' , respectively, and $A : \mathbb{S} \rightarrow \mathbb{S}'$ is a LINEAR OPERATOR with $y = Ax$ for all $x \in \mathbb{S}$ and $y \in \mathbb{S}'$.

Then, with $\mathbf{x} = [x_1, \dots, x_N]^T$ and $\mathbf{y} = [y_1, \dots, y_M]^T$ representing the coordinates of x and y with respect to the basis $\{u_n\}_{n=1}^N$ and $\{u'_m\}_{m=1}^M$, respectively, we obtain the Matrix \mathbf{A} with $\mathbf{y} = \mathbf{A}\mathbf{x}$ as follows:

$$y_m = \langle y, u'_m \rangle = \langle Ax, u'_m \rangle = \left\langle A \sum_{n=1}^N x_n u_n, u'_m \right\rangle = \sum_{n=1}^N \langle Au_n, u'_m \rangle x_n \quad (\text{B.3})$$

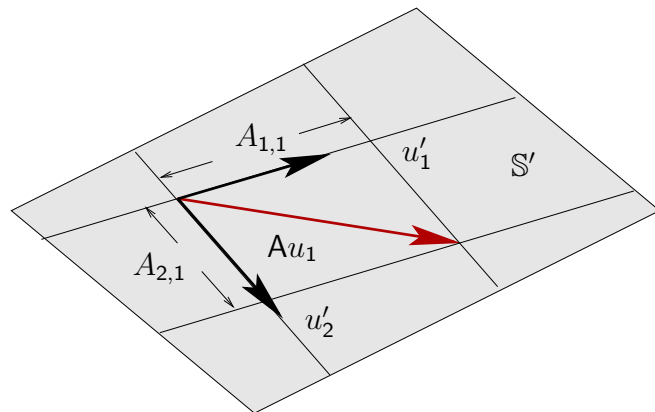
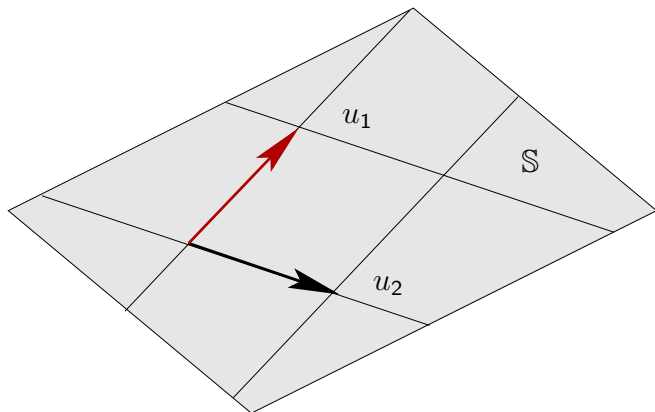
Consequently, the elements of matrix \mathbf{A} are obtained by

$$[\mathbf{A}]_{m,n} \triangleq \langle Au_n, u'_m \rangle \quad \text{and} \quad [\mathbf{A}^T]_{n,m} \triangleq \langle Au_n, u'_m \rangle. \quad (\text{B.4})$$

¹For instance the vector space of bandlimited periodic signals.

$$A : \mathbb{S} \rightarrow \mathbb{S}', \quad u_1 \mapsto Au_1 = A_{1,1}u'_1 + A_{2,1}u'_2$$

$$A_{1,1} \triangleq \langle Au_1, u'_1 \rangle, \quad A_{2,1} \triangleq \langle Au_1, u'_2 \rangle$$



B.2 Elementwise Perspective on Matrices

$$\mathbf{A} \in \mathbb{C}^{M \times N} \text{ with } [\mathbf{A}]_{m,n} = a_{m,n} = \operatorname{Re}\{a_{m,n}\} + \mathrm{j} \operatorname{Im}\{a_{m,n}\} \in \mathbb{C}, \quad (\text{B.5})$$

and

$$\begin{aligned} \mathbf{A}\mathbf{x} &= \begin{bmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,N} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,N} \\ \vdots & \vdots & & \vdots \\ a_{M,1} & a_{M,2} & \cdots & a_{M,N} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} \\ &= \begin{bmatrix} a_{1,1}x_1 + a_{1,2}x_2 + \cdots + a_{1,N}x_N \\ a_{2,1}x_1 + a_{2,2}x_2 + \cdots + a_{2,N}x_N \\ \vdots \\ a_{M,1}x_1 + a_{M,2}x_2 + \cdots + a_{M,N}x_N \end{bmatrix} \\ &= \begin{bmatrix} \sum_{n=1}^N a_{1,n}x_n \\ \sum_{n=1}^N a_{2,n}x_n \\ \vdots \\ \sum_{n=1}^N a_{M,n}x_n \end{bmatrix}. \end{aligned}$$

Definition. MATRIX, TRANSPOSE MATRIX, ADJOINT MATRIX

$$\mathbf{A} = \begin{bmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,N} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,N} \\ \vdots & \vdots & & \vdots \\ a_{M,1} & a_{M,2} & \cdots & a_{M,N} \end{bmatrix},$$

$$\mathbf{A}^{\mathrm{T}} = \begin{bmatrix} a_{1,1} & a_{2,1} & \cdots & a_{M,1} \\ a_{1,2} & a_{2,2} & \cdots & a_{M,2} \\ \vdots & \vdots & & \vdots \\ a_{1,N} & a_{2,N} & \cdots & a_{M,N} \end{bmatrix},$$

$$\mathbf{A}^{\mathrm{H}} = \begin{bmatrix} a_{1,1}^* & a_{2,1}^* & \cdots & a_{M,1}^* \\ a_{1,2}^* & a_{2,2}^* & \cdots & a_{M,2}^* \\ \vdots & \vdots & & \vdots \\ a_{1,N}^* & a_{2,N}^* & \cdots & a_{M,N}^* \end{bmatrix}.$$

The ADJOINT MATRIX $\mathbf{A}^H \in \mathbb{C}^{N \times M}$ of a matrix $\mathbf{A} \in \mathbb{C}^{M \times N}$ has the constituting property

$$\langle \mathbf{A}\mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{A}^H \mathbf{y} \rangle \quad \Leftrightarrow \quad \mathbf{y}^H \mathbf{A}\mathbf{x} = (\mathbf{A}^H \mathbf{y})^H \mathbf{x}, \quad \forall \mathbf{x} \in \mathbb{C}^N, \mathbf{y} \in \mathbb{C}^M. \quad (\text{B.6})$$

MATRIX COLUMNS, MATRIX ROWS

$$\mathbf{A} = [\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_N] \quad \text{with} \quad \mathbf{a}_n = \begin{bmatrix} a_{1,n} \\ a_{2,n} \\ \vdots \\ a_{M,n} \end{bmatrix}$$

$$\mathbf{A}^T = \begin{bmatrix} \mathbf{a}_1^T \\ \mathbf{a}_2^T \\ \vdots \\ \mathbf{a}_N^T \end{bmatrix} \quad \text{with} \quad \mathbf{a}_n^T = [a_{1,n} \ a_{2,n} \ \cdots \ a_{M,n}]$$

$$\mathbf{A}^H = \begin{bmatrix} \mathbf{a}_1^H \\ \mathbf{a}_2^H \\ \vdots \\ \mathbf{a}_N^H \end{bmatrix} \quad \text{with} \quad \mathbf{a}_n^H = [a_{1,n}^* \ a_{2,n}^* \ \cdots \ a_{M,n}^*],$$

with $(\mathbf{A}^T)^T = \mathbf{A}$ and $(\mathbf{A}^H)^H = \mathbf{A}$.

Definition. IDENTITY MATRIX, ZERO MATRIX, INVERSE MATRIX, SELECTION VECTOR

$$\mathbf{I} = \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix} = \text{diag}[1 \ 1 \cdots 1]$$

$$\mathbf{0} = \begin{bmatrix} 0 & 0 & \cdots & 0 \\ 0 & 0 & \cdots & 0 \\ \vdots & \vdots & & \vdots \\ 0 & 0 & \cdots & 0 \end{bmatrix}$$

$$\mathbf{A}^{-1}\mathbf{A} = \mathbf{I} \quad \text{and} \quad \mathbf{A} + (-\mathbf{A}) = \mathbf{0}$$

$$\mathbf{e}_i = [\underbrace{0 \cdots 0}_{i-1 \text{ zeros}} \ 1 \ \underbrace{0 \cdots 0}_{N-i \text{ zeros}}]^\top \in \mathbb{C}^N \quad \text{with}$$

$$\mathbf{A}\mathbf{e}_i = \begin{bmatrix} a_{1,1} & a_{1,2} & \cdots & a_{1,N} \\ a_{2,1} & a_{2,2} & \cdots & a_{2,N} \\ \vdots & \vdots & & \vdots \\ a_{M,1} & a_{M,2} & \cdots & a_{M,N} \end{bmatrix} \begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} = \begin{bmatrix} a_{1,i} \\ a_{2,i} \\ \vdots \\ a_{M,i} \end{bmatrix} = \mathbf{a}_i.$$

B.3 Fundamental Subspace Perspective on Matrices

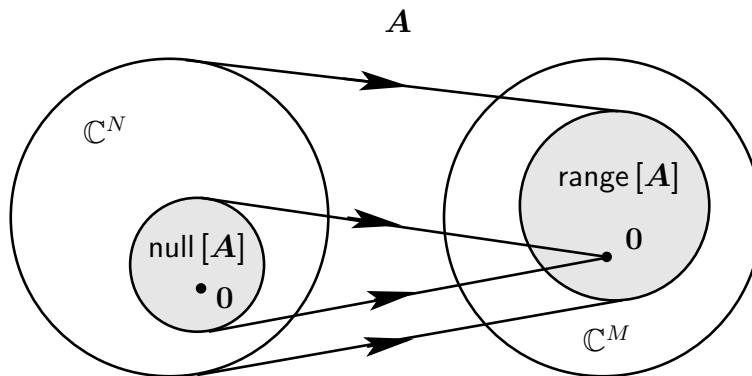
Any LINEAR OPERATOR $\mathbf{A} : \mathbb{C}^N \rightarrow \mathbb{C}^M$ is related to FOUR FUNDAMENTAL SUBSPACES, viz.

$$\text{RANGE SPACE of } \mathbf{A} : \quad \text{range}[\mathbf{A}] = \{ \mathbf{Ax} \mid \mathbf{x} \in \mathbb{C}^N \} \quad (\text{B.7})$$

$$\text{NULL SPACE of } \mathbf{A} : \quad \text{null}[\mathbf{A}] = \{ \mathbf{x} \in \mathbb{C}^N \mid \mathbf{Ax} = \mathbf{0} \} \quad (\text{B.8})$$

$$\text{ORTHOGONAL COMPLEMENT of } \text{range}[\mathbf{A}] : \quad \text{range}[\mathbf{A}]^\perp = \text{null}[\mathbf{A}^H] \quad (\text{B.9})$$

$$\text{ORTHOGONAL COMPLEMENT of } \text{null}[\mathbf{A}] : \quad \text{null}[\mathbf{A}]^\perp = \text{range}[\mathbf{A}^H] . \quad (\text{B.10})$$



B.4 Eigenvectors and Eigenvalues

Definition. An EIGENVECTOR of a matrix $\mathbf{A} \in \mathbb{C}^{M \times M}$ is a nonzero vector $\mathbf{v} \in \mathbb{C}^N$ such that

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v} \quad \text{for some } \lambda \in \mathbb{C}. \quad (\text{B.11})$$

The constant λ is the EIGENVALUE to the EIGENVECTOR \mathbf{v} . The \mathbf{v} and its corresponding λ form a so-called EIGENPAIR of the matrix \mathbf{A} .

For finite dimensional matrices the SET OF EIGENVALUES is discrete, the maximum number of different eigenvalues is M .

The EIGENVALUES of SELF-ADJOINT MATRICES $\mathbf{A}^H = \mathbf{A}$ are real-valued and the EIGENVECTORS form an ONB.

Proof.

$$(1) \quad \lambda \mathbf{v}^H \mathbf{v} = \mathbf{v}^H (\lambda \mathbf{v}) = \mathbf{v}^H \mathbf{A} \mathbf{v} = \mathbf{v}^H \mathbf{A}^H \mathbf{v} = (\mathbf{A} \mathbf{v})^H \mathbf{v} = (\lambda \mathbf{v})^H \mathbf{v} = \lambda^* \mathbf{v}^H \mathbf{v} \quad \Rightarrow \quad \lambda = \lambda^*$$

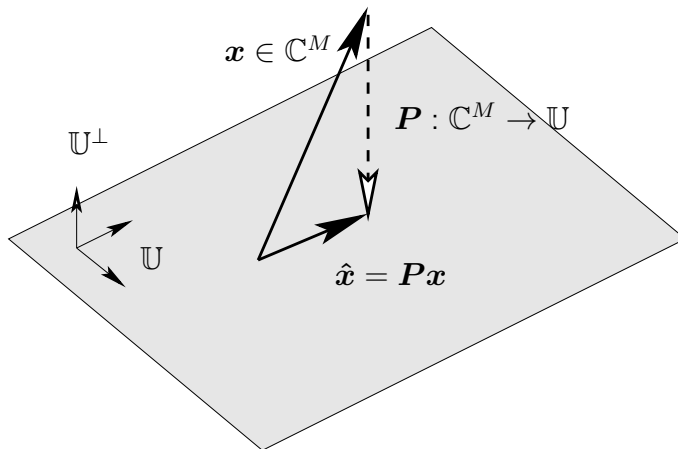
$$(2) \quad \lambda_i \mathbf{v}_j^H \mathbf{v}_i = \mathbf{v}_j^H (\lambda_i \mathbf{v}_i) = \mathbf{v}_j^H \mathbf{A} \mathbf{v}_i = \mathbf{v}_j^H \mathbf{A}^H \mathbf{v}_i = (\mathbf{A} \mathbf{v}_j)^H \mathbf{v}_i = (\lambda_j \mathbf{v}_j)^H \mathbf{v}_i = \lambda_j^* \mathbf{v}_j^H \mathbf{v}_i \quad \Rightarrow \quad \mathbf{v}_j^H \mathbf{v}_i = 0$$

since λ_i and λ_j are real-valued (1st part) and different for $i \neq j$ by assumption.

B.5 Orthogonal Projectors

Best Approximation and Orthogonal Projection

The PROJECTION OPERATOR $\mathbf{P} : \mathbb{C}^M \rightarrow \mathbb{C}^M$ is related to the BEST APPROXIMATION problem, i.e., find $\hat{\mathbf{x}} \in \mathbb{U} \subset \mathbb{C}^M$ that is closest to a given $\mathbf{x} \in \mathbb{C}^M$. As a result, we obtain $\hat{\mathbf{x}} = \mathbf{P}\mathbf{x} \perp \mathbf{x} - \mathbf{P}\mathbf{x}$.



According to the CLOSEST POINT THEOREM² we obtain

$$\text{EXISTENCE : } \exists \hat{\mathbf{x}} \in \mathbb{U} \text{ such that } \|\mathbf{x} - \hat{\mathbf{x}}\| \leq \|\mathbf{x} - \mathbf{s}\| \text{ for every } \mathbf{s} \in \mathbb{U} \quad (\text{B.12})$$

$$\text{ORTHOGONALITY : } \mathbf{x} - \hat{\mathbf{x}} \perp \mathbb{U} \quad (\text{B.13})$$

$$\text{LINEARITY : } \hat{\mathbf{x}} = \mathbf{P}\mathbf{x} \text{ and } \mathbf{P} \text{ only depends on } \mathbb{U} \quad (\text{B.14})$$

$$\text{INDEMPOTENCY : } \mathbf{P}^2 = \mathbf{P} \quad (\text{B.15})$$

$$\text{SELF-ADJOINTNESS : } \mathbf{P}^H = \mathbf{P}. \quad (\text{B.16})$$

Note.

PROJECTION OPERATORS are always INDEMPOTENT.

ORTHOGONAL PROJECTORS must be INDEMPOTENT and SELF-ADJOINT.

OBLIQUE PROJECTORS (non-orthogonal projectors) are only INDEMPOTENT, i.e., $\mathbf{P}^2 = \mathbf{P}$, but SELF-ADJOINTNESS does not hold.

OBLIQUE PROJECTORS are not considered in this tutorial.

²The CLOSEST POINT THEOREM applies for the category of CONVEX SETS which LINEAR SUBSPACES are belonging to.

Orthogonal Projection on Subspaces.

Given a pair of matrices $\mathbf{A} : \mathbb{C}^N \rightarrow \mathbb{C}^M$ and $\mathbf{B} : \mathbb{C}^M \rightarrow \mathbb{C}^N$ and \mathbf{A} is a LEFT-INVERSE of \mathbf{B} , i.e., $\mathbf{AB} = \mathbf{I}$, then

$$\mathbf{BA} \text{ is a PROJECTOR onto range } [\mathbf{B}] \subset \mathbb{C}^N, \quad (\text{B.17})$$

since $(\mathbf{BA})^2 = \mathbf{BABA} = \mathbf{B(AB)A} = \mathbf{BA}$.

If \mathbf{BA} is also SELF-ADJOINT, then

$$\mathbf{BA} \text{ is an ORTHOGONAL PROJECTOR onto range } [\mathbf{B}] \subset \mathbb{C}^N, \quad (\text{B.18})$$

since $(\mathbf{BA})^2 = \mathbf{BABA} = \mathbf{B(AB)A} = \mathbf{BA}$ and $(\mathbf{BA})^H = \mathbf{BA}$.

Example. Assume the column vectors of \mathbf{U} form an ONB of \mathbb{U} , then $\mathbf{P} = \mathbf{UU}^H$ is the ORTHOGONAL PROJECTOR onto \mathbb{U} . For a proof consider the plugins

$$\mathbf{A} = \mathbf{U}^H \quad \text{and} \quad \mathbf{B} = \mathbf{U}. \quad (\text{B.19})$$

Orthogonal Projection via Pseudo-Inverse.

If $(\mathbf{A}^H \mathbf{A})^{-1}$ exists, then $\mathbf{A}_{\text{left}}^+ = (\mathbf{A}^H \mathbf{A})^{-1} \mathbf{A}^H$ is a LEFT-INVERSE of \mathbf{A} , i.e.,

$$(1) \quad \mathbf{A} \mathbf{A}_{\text{left}}^+ = \mathbf{A} (\mathbf{A}^H \mathbf{A})^{-1} \mathbf{A}^H \text{ is the ORTHOGONAL PROJECTOR onto } \text{range} [\mathbf{A}] \quad (\text{B.20})$$

with so-called PSEUDO-INVERSE $\mathbf{A}_{\text{left}}^+$.

Otherwise, if $(\mathbf{A} \mathbf{A}^H)^{-1}$ exists, then $\mathbf{A}_{\text{right}}^+ = \mathbf{A}^H (\mathbf{A} \mathbf{A}^H)^{-1}$ is a RIGHT-INVERSE of \mathbf{A} , i.e.,

$$(2) \quad \mathbf{A}_{\text{right}}^+ \mathbf{A} = \mathbf{A}^H (\mathbf{A} \mathbf{A}^H)^{-1} \mathbf{A} \text{ is the ORTHOGONAL PROJECTOR onto } \text{range} [\mathbf{A}^H] \quad (\text{B.21})$$

with the PSEUDO-INVERSE $\mathbf{A}_{\text{right}}^+$.

Weighted Sum of Disjoint Orthogonal Projectors

Assume rank- R matrices $\mathbf{A} : \mathbb{C}^M \rightarrow \mathbb{C}^M$ that can be represented by means of a WEIGHTED SUM of rank- R_i DISJOINT ORTHOGONAL PROJECTORS $\mathbf{P}_i : \mathbb{C}^M \rightarrow \mathbb{C}^M$,

$$\mathbf{A} = \sum_{i=1}^P w_i \mathbf{P}_i, \quad (\text{B.22})$$

with rank $R = \sum_{i=1}^P R_i$ and $\text{range}[\mathbf{P}_i] \perp \text{range}[\mathbf{P}_j]$ for all $i \neq j$, $i = 1, \dots, P$, i.e., all projectors are mutually orthogonal to each other.

A weighted sum of disjoint orthogonal projectors $\mathbf{A} : \mathbb{C}^M \rightarrow \mathbb{C}^M$ referring to (B.22) is a NORMAL MATRIX with

$$\mathbf{A}\mathbf{A}^H = \mathbf{A}^H\mathbf{A}. \quad (\text{B.23})$$

The respective projectors and weights depend on \mathbf{A} .

Note. The weights take COMPLEX VALUES, i.e., $w_i \in \mathbb{C}$.

B.6 Special Matrices

Definition. $\mathbf{A} : \mathbb{C}^M \rightarrow \mathbb{C}^M$ referring to (B.22) is a SELF-ADJOINT MATRIX when

$$\mathbf{A}^H = \mathbf{A}. \quad (\text{B.24})$$

SELF-ADJOINT MATRICES are also NORMAL MATRICES.

Note. The weights strictly take REAL VALUES, i.e., $w_i \in \mathbb{R}$.

Definition. $\mathbf{A} : \mathbb{C}^M \rightarrow \mathbb{C}^M$ referring to (B.22) is a POSITIVE DEFINITE MATRIX when it is SELF-ADJOINT³ and

$$\mathbf{x}^H \mathbf{A} \mathbf{x} \geq 0, \quad \text{for all } \mathbf{x} \in \mathbb{C}^M. \quad (\text{B.25})$$

Note. The weights strictly take NONNEGATIVE VALUES, i.e., $w_i \geq 0$.

³There exist more general definitions of POSITIVE DEFINITE MATRICES.

Eigenpairs of Normal Matrices

NORMAL MATRICES \mathbf{A} are LINEAR COMBINATIONS of ORTHOGONAL PROJECTORS, i.e., $\mathbf{A} = \sum_{i=1}^P w_i \mathbf{P}_i$.

Consequently, any $\mathbf{u}_j \in \text{range}[\mathbf{P}_j]$ and its corresponding weight w_j form an EIGENPAIR of \mathbf{A} , since for every element of $\text{range}[\mathbf{P}_j]$ we obtain $\mathbf{P}_j \mathbf{u}_j = \mathbf{u}_j$ by definition and $\mathbf{P}_i \mathbf{u}_j = \mathbf{0}$ for $i \neq j$ due to the orthogonality of projectors, and thus

$$\mathbf{A} \mathbf{u}_j = \left(\sum_{i=1}^P w_i \mathbf{P}_i \right) \mathbf{u}_j \quad (\text{B.26})$$

$$= \sum_{i=1}^P w_i (\mathbf{P}_i \mathbf{u}_j) = w_j \mathbf{u}_j. \quad (\text{B.27})$$

Obviously, each w_j and \mathbf{u}_j form an EIGENPAIR consisting of an EIGENVECTOR and an EIGENVALUE of the matrix:

$$\mathbf{A} \mathbf{u}_j = w_j \mathbf{u}_j, \quad \text{for all } j = 1, \dots, P. \quad (\text{B.28})$$

Spectral Theorem of Normal Matrices

Since (1) NORMAL MATRICES \mathbf{A} are LINEAR COMBINATIONS of ORTHOGONAL PROJECTORS \mathbf{P}_i and any $\mathbf{u}_j \in \text{range}[\mathbf{P}_j]$ and its corresponding weight w_j form an EIGENPAIR of \mathbf{A} , and

(2) since any \mathbf{P}_i can be represented as a LINEAR COMBINATION of rank-1 projectors $\mathbf{u}_{i,j}\mathbf{u}_{i,j}^H$ with respect to elements of an ONB $\{\mathbf{u}_{i,j}\}_{j=1}^{R_i}$, with

$$\text{range}[\mathbf{P}_i] = \text{range}[\mathbf{U}_i] = \text{span}[\mathbf{u}_{i,1} \cdots \mathbf{u}_{i,R_i}], \quad (\text{B.29})$$

with $\mathbf{U}_i = [\mathbf{u}_{i,1}, \dots, \mathbf{u}_{i,R_i}]$ and R_i is the dimension of $\text{range}[\mathbf{P}_i]$, we obtain the main result of the SPECTRAL THEOREM of NORMAL MATRICES,

$$\mathbf{A} = \sum_{i=1}^P w_i \mathbf{P}_i \quad (\text{B.30})$$

$$= \sum_{i=1}^P w_i \mathbf{U}_i \mathbf{U}_i^H = \sum_{i=1}^P \sum_{j=1}^{R_i} w_i \mathbf{u}_{i,j} \mathbf{u}_{i,j}^H \quad (\text{B.31})$$

$$= [\mathbf{U}_1 \cdots \mathbf{U}_P] \begin{bmatrix} w_1 \mathbf{I}_{R_1 \times R_1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & w_P \mathbf{I}_{R_P \times R_P} \end{bmatrix} \begin{bmatrix} \mathbf{U}_1^H \\ \vdots \\ \mathbf{U}_P^H \end{bmatrix}. \quad (\text{B.32})$$

B.7 Singular Value Decomposition

Derivation and Properties of the SVD

Any COMPACT LINEAR OPERATOR A can be decomposed by means of the SINGULAR VALUE DECOMPOSITION (SVD) and any FINITE DIMENSIONAL MATRIX $A : \mathbb{C}^N \rightarrow \mathbb{C}^M$ is a compact linear operator.

(1) Since the product of matrices $A^H A$ —the so-called GRAMIAN MATRIX—is NORMAL, SELF-ADJOINT and POSITIVE DEFINITE,⁴ the SPECTRAL THEOREM can be applied such that

$$A^H A = \sum_{i=1}^P \lambda_i P_i = \sum_{i=1}^R \lambda_i v_i v_i^H, \quad (\text{B.33})$$

with $\text{rank } R = R_1 + \cdots + R_P \leq \min(M, N)$ and at least P different weights $\lambda_i > 0$.

The vectors $\{v_i\}_{i=1}^R$ form an ONB of range $[A^H]$.

⁴ $(A^H A)^H = A^H A$ and $v^H A^H A v = v^H \lambda v = \lambda v^H v = \lambda \|v\|^2 \geq 0$.

(2) Given the ORTHONORMAL VECTORS $\{\mathbf{v}_i\}_{i=1}^R$, the vectors $\{\mathbf{u}_i\}_{i=1}^R$, with

$$\mathbf{u}_i \triangleq \lambda_i^{-\frac{1}{2}} \mathbf{A} \mathbf{v}_i, \quad (\text{B.34})$$

again constitute an ONB:

$$\begin{aligned} \mathbf{u}_j^H \mathbf{u}_i &= (\lambda_j^{-\frac{1}{2}} \mathbf{A} \mathbf{v}_j)^H \lambda_i^{-\frac{1}{2}} \mathbf{A} \mathbf{v}_i \\ &= \lambda_j^{-\frac{1}{2}} \lambda_i^{-\frac{1}{2}} \mathbf{v}_j^H \mathbf{A}^H \mathbf{A} \mathbf{v}_i \\ &= \lambda_j^{-\frac{1}{2}} \lambda_i^{-\frac{1}{2}} \lambda_i \mathbf{v}_j^H \mathbf{v}_i \\ &= \begin{cases} 1 & ; \quad i = j \\ 0 & ; \quad \text{otherwise} \end{cases} . \end{aligned}$$

Definition.

The elements of an ONB $\{\mathbf{v}_i\}_{i=1}^R$ with EIGENVALUES $\{\lambda_i\}_{i=1}^R$ of a GRAMIAN MATRIX $\mathbf{A}^H \mathbf{A}$ are called RIGHT SINGULAR VECTORS of \mathbf{A} .

The elements of the ONB $\{\mathbf{u}_i = 1/\sqrt{\lambda_i} \mathbf{A} \mathbf{v}_i\}_{i=1}^R$ with EIGENVALUES $\{\lambda_i\}_{i=1}^R$ of a GRAMIAN MATRIX $\mathbf{A}^H \mathbf{A}$ are called corresponding LEFT SINGULAR VECTORS of \mathbf{A} .

The SINGULAR VALUE DECOMPOSITION (SVD) of an arbitrary matrix $\mathbf{A} : \mathbb{C}^N \rightarrow \mathbb{C}^M$ is constituted by

$$\mathbf{A} = \sum_{i=1}^R \sigma_i \mathbf{u}_i \mathbf{v}_i^H, \quad \text{with matrix rank } R \quad (\text{B.35})$$

$$\mathbf{u}_j^H \mathbf{u}_i = \delta_{i,j} \quad \text{and} \quad \mathbf{v}_j^H \mathbf{v}_i = \delta_{i,j} \quad \text{and} \quad \sigma_i > 0 \quad i, j = 1, \dots, Q. \quad (\text{B.36})$$

with $\delta_{i,j} = 1$ when $i = j$ and zero otherwise.

The set of LEFT SINGULAR VECTORS $\{\mathbf{u}_i\}_{i=1}^R$ of nonzero SINGULAR VALUES $\sigma_i > 0$ provides an ONB for the IMAGE SPACE ($\text{range}[\mathbf{A}] \in \mathbb{C}^M$) of the matrix \mathbf{A} , i.e.,

$$\text{range}[\mathbf{A}] = \text{span}[\mathbf{u}_1 \ \mathbf{u}_2 \ \cdots \ \mathbf{u}_R]. \quad (\text{B.37})$$

Correspondingly, the RIGHT SINGULAR VECTORS $\{\mathbf{v}_i\}_{i=1}^R$ of nonzero SINGULAR VALUES $\sigma_i > 0$ provides an ONB for the IMAGE SPACE ($\text{range}[\mathbf{A}^H] \in \mathbb{C}^N$) of the matrix \mathbf{A}^H , which is equal to the COMPLEMENT of the NULL SPACE ($\text{null}[\mathbf{A}]^\perp \equiv \text{range}[\mathbf{A}^H] \in \mathbb{C}^N$) of the matrix \mathbf{A} , i.e.,

$$\text{range}[\mathbf{A}^H] = \text{null}[\mathbf{A}]^\perp = \text{span}[\mathbf{v}_1 \ \mathbf{v}_2 \ \cdots \ \mathbf{v}_R]. \quad (\text{B.38})$$

Derivation.

Applying the matrix \mathbf{A} to a vector \mathbf{x} can be equivalently expressed by

$$\mathbf{A}\mathbf{x} = \mathbf{A}(\mathbf{x}_- + \mathbf{x}_\perp) = \mathbf{A}\mathbf{x}_-,$$

with $\mathbf{x}_- \in \text{range}[\mathbf{A}^H] \equiv \text{null}[\mathbf{A}]^\perp$ and $\mathbf{x}_\perp \in \text{null}[\mathbf{A}]$.

Since $\{\mathbf{v}_i\}_{i=1}^R$ forms an ONB for $\text{range}[\mathbf{A}^H] \equiv \text{null}[\mathbf{A}]^\perp$, we obtain

$$\begin{aligned}\mathbf{A}\mathbf{x} &= \mathbf{A}\mathbf{x}_- \\ &= \mathbf{A} \sum_{i=1}^R \mathbf{v}_i \mathbf{v}_i^H \mathbf{x}_- \\ &= \sum_{i=1}^R \mathbf{A} \mathbf{v}_i \mathbf{v}_i^H \mathbf{x}_- \\ &= \sum_{i=1}^R \sqrt{\lambda_i} \mathbf{u}_i \mathbf{v}_i^H \mathbf{x}_- \\ &= \sum_{i=1}^R \sqrt{\lambda_i} \mathbf{u}_i \mathbf{v}_i^H \mathbf{x},\end{aligned}$$

with $\sigma_i = \sqrt{\lambda_i}$.

Sorted Matrix Representation

By sorting the SINGULAR VALUES according to

$$\sigma_1 \geq \sigma_2 \geq \cdots \geq \sigma_R \quad \text{with} \quad R = \text{rank}[\mathbf{A}],$$

we obtain the SORTED MATRIX REPRESENTATION of the SVD,

$$\mathbf{A} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^H, \quad \text{with} \quad (\text{B.39})$$

$$\mathbf{U} = [\mathbf{u}_1 \ \mathbf{u}_2 \cdots \mathbf{u}_R] \in \mathbb{C}^{M \times R}, \quad (\text{B.40})$$

$$\mathbf{\Sigma} = \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_R \end{bmatrix} = \text{diag} [\sigma_1 \ \sigma_2 \cdots \sigma_R] \in \mathbb{C}^{R \times R} \quad (\text{B.41})$$

$$\mathbf{V}^H = \begin{bmatrix} \mathbf{v}_1^H \\ \vdots \\ \mathbf{v}_R^H \end{bmatrix} \in \mathbb{C}^{R \times N}. \quad (\text{B.42})$$

Extended Sorted Matrix Representation. By accomplishing the ONB of the LEFT SINGULAR VECTORS (column vectors of \mathbf{U}) and RIGHT SINGULAR VECTORS (column vectors of \mathbf{V}) by their ORTHOGONAL COMPLEMENTS, and taking into account that $Q = \min(M, N)$,

$$\mathbf{U}_{R < Q}^\perp = [\mathbf{u}_{R+1} \ \mathbf{u}_{R+2} \cdots \mathbf{u}_Q] \quad \text{and} \quad \mathbf{U}_{Q < M}^\perp = [\mathbf{u}_{Q+1} \ \mathbf{u}_{Q+2} \cdots \mathbf{u}_M] \quad (\text{B.43})$$

$$\mathbf{V}_{R < Q}^\perp = [\mathbf{v}_{R+1} \ \mathbf{v}_{R+2} \cdots \mathbf{v}_Q] \quad \text{and} \quad \mathbf{V}_{Q < N}^\perp = [\mathbf{v}_{Q+1} \ \mathbf{v}_{Q+2} \cdots \mathbf{v}_N] \quad (\text{B.44})$$

$$\text{range}[\mathbf{U}] \oplus \text{range}[\mathbf{U}_{R < Q}^\perp] \oplus \text{range}[\mathbf{U}_{Q < M}^\perp] = \mathbb{C}^M \quad \text{and} \quad (\text{B.45})$$

$$\text{range}[\mathbf{V}] \oplus \text{range}[\mathbf{V}_{R < Q}^\perp] \oplus \text{range}[\mathbf{V}_{Q < N}^\perp] = \mathbb{C}^N, \quad (\text{B.46})$$

we obtain the EXTENDED SORTED MATRIX REPRESENTATION of the **SVD of matrices with maximum rank** ($R = Q$)

$$\mathbf{A} = \begin{cases} \begin{bmatrix} [\mathbf{U}, \mathbf{U}_{Q < M}^\perp] & \begin{bmatrix} \boldsymbol{\Sigma} \\ \mathbf{0} \end{bmatrix} & \mathbf{V}^H \end{bmatrix} & ; \quad M \geq N \quad (\text{TALL}) \\ \begin{bmatrix} \mathbf{U} & \boldsymbol{\Sigma} & \mathbf{V}^H \end{bmatrix} & ; \quad M = N \quad (\text{SQUARE}) \\ \begin{bmatrix} \mathbf{U} & [\boldsymbol{\Sigma} \ \mathbf{0}] & \begin{bmatrix} \mathbf{V}^H \\ \mathbf{V}_{Q < N}^{\perp, H} \end{bmatrix} \end{bmatrix} & ; \quad M \leq N \quad (\text{WIDE}). \end{cases} \quad (\text{B.47})$$

SVD of Matrices without maximum Rank.

For TALL MATRICES $\mathbf{A}_{M \geq N}$ with linear dependent columns or WIDE MATRICES $\mathbf{A}_{M \leq N}$ with linear dependent rows, i.e., with

$$\text{rank}[\mathbf{A}] < Q = \min\{M, N\}, \quad (\text{B.48})$$

we obtain the cases

$$\mathbf{A}_{M > N} = [\mathbf{U}, \mathbf{U}_{R < Q}^\perp, \mathbf{U}_{Q < M}^\perp] \begin{bmatrix} \begin{bmatrix} \boldsymbol{\Sigma}_{R \times R} & \mathbf{0}_{R \times N-R} \\ \mathbf{0}_{N-R \times R} & \mathbf{0}_{N-R \times N-R} \end{bmatrix} \\ \begin{bmatrix} \mathbf{0}_{M-N \times R} & \mathbf{0}_{M-N \times N-R} \end{bmatrix} \end{bmatrix} \begin{bmatrix} \mathbf{V}^H \\ \mathbf{V}_{R < Q}^{\perp, H} \end{bmatrix} \quad (\text{B.49})$$

$$\mathbf{A}_{M=N} = [\mathbf{U}, \mathbf{U}_{R < Q}^\perp] \begin{bmatrix} \boldsymbol{\Sigma}_{R \times R} & \mathbf{0}_{R \times N-R} \\ \mathbf{0}_{N-R \times R} & \mathbf{0}_{N-R \times N-R} \end{bmatrix} \begin{bmatrix} \mathbf{V}^H \\ \mathbf{V}_{R < Q}^{\perp, H} \end{bmatrix} \quad (\text{B.50})$$

$$\mathbf{A}_{M < N} = [\mathbf{U}, \mathbf{U}_{R < Q}^\perp] \begin{bmatrix} \begin{bmatrix} \boldsymbol{\Sigma}_{R \times R} & \mathbf{0}_{R \times M-R} \\ \mathbf{0}_{M-R \times R} & \mathbf{0}_{M-R \times M-R} \end{bmatrix} \\ \begin{bmatrix} \mathbf{0}_{R \times N-M} \\ \mathbf{0}_{M-R \times N-M} \end{bmatrix} \end{bmatrix} \begin{bmatrix} \mathbf{V}^H \\ \mathbf{V}_{R < Q}^{\perp, H} \\ \mathbf{V}_{Q < N}^{\perp, H} \end{bmatrix}. \quad (\text{B.51})$$

Fundamental Subspaces (continued)

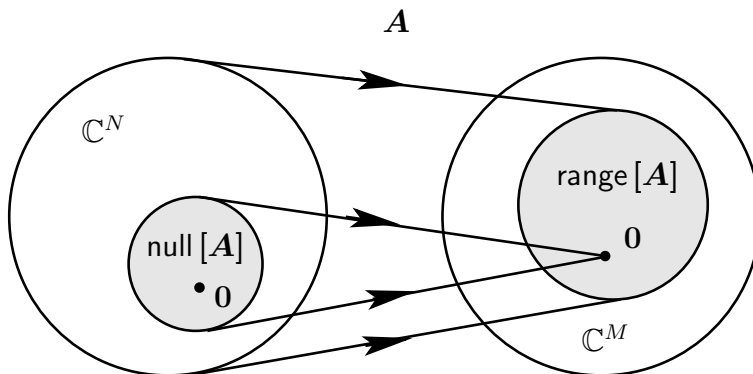
Any LINEAR OPERATOR $\mathbf{A} : \mathbb{C}^N \rightarrow \mathbb{C}^M$ is related to FOUR FUNDAMENTAL SUBSPACES, viz.

$$\text{IMAGE SPACE of } \mathbf{A} : \quad \text{range}[\mathbf{A}] = \{ \mathbf{Ax} \mid \mathbf{x} \in \mathbb{C}^N \} \quad (\text{B.52})$$

$$\text{NULL SPACE of } \mathbf{A} : \quad \text{null}[\mathbf{A}] = \{ \mathbf{x} \in \mathbb{C}^N \mid \mathbf{Ax} = \mathbf{0} \} \quad (\text{B.53})$$

$$\text{ORTHOGONAL COMPLEMENT of } \text{range}[\mathbf{A}] : \quad \text{range}[\mathbf{A}]^\perp = \text{null}[\mathbf{A}^H] \quad (\text{B.54})$$

$$\text{ORTHOGONAL COMPLEMENT of } \text{null}[\mathbf{A}] : \quad \text{null}[\mathbf{A}]^\perp = \text{range}[\mathbf{A}^H] . \quad (\text{B.55})$$



ONB.

The column vectors of $\mathbf{U}_{\text{ext}} = [\mathbf{U}, \mathbf{U}^\perp]$ with $\mathbf{U}^\perp = [\mathbf{U}_{R < Q}^\perp, \mathbf{U}_{Q < M}^\perp]$ and $\mathbf{V}_{\text{ext}} = [\mathbf{V}, \mathbf{V}^\perp]$ with $\mathbf{V}^\perp = [\mathbf{V}_{R < Q}^\perp, \mathbf{V}_{Q < N}^\perp]$ from the EXTENDED SVD of a rank- R matrix $\mathbf{A} \in \mathbb{C}^{M \times N}$ constitute the ORTHONORMAL BASES for the FOUR FUNDAMENTAL SUBSPACES $\text{range}[\mathbf{A}]$, $\text{range}[\mathbf{A}]^\perp$, $\text{null}[\mathbf{A}]$, and $\text{null}[\mathbf{A}]^\perp$, viz.

$$\text{range}[\mathbf{A}] = \text{span}[\mathbf{U}_{\text{ext}} \mathbf{e}_i \mid i = 1, \dots, R] \quad (\text{B.56})$$

$$\text{range}[\mathbf{A}]^\perp = \text{span}[\mathbf{U}_{\text{ext}} \mathbf{e}_i \mid i = R + 1, \dots, Q] \oplus \text{span}[\mathbf{U}_{\text{ext}} \mathbf{e}_i \mid i = Q + 1, \dots, M] \quad (\text{B.57})$$

$$\text{null}[\mathbf{A}] = \text{span}[\mathbf{V}_{\text{ext}} \mathbf{e}_i \mid i = R + 1, \dots, Q] \oplus \text{span}[\mathbf{V}_{\text{ext}} \mathbf{e}_i \mid i = Q + 1, \dots, N] \quad (\text{B.58})$$

$$\text{null}[\mathbf{A}]^\perp = \text{span}[\mathbf{V}_{\text{ext}} \mathbf{e}_i \mid i = 1, \dots, R]. \quad (\text{B.59})$$

The FOUR FUNDAMENTAL SUBSPACES of the matrices $\mathbf{A} : \mathbb{C}^N \rightarrow \mathbb{C}^M$ and $\mathbf{A}^H : \mathbb{C}^M \rightarrow \mathbb{C}^N$ form ORTHOGONAL COMPLEMENTS of the N - and M -dimensional VECTOR SPACES:

$$\mathbb{C}^N = \text{null}[\mathbf{A}] \oplus \text{range}[\mathbf{A}^H], \quad (\text{B.60})$$

$$\mathbb{C}^M = \text{range}[\mathbf{A}] \oplus \text{null}[\mathbf{A}^H]. \quad (\text{B.61})$$

C. Kalman Filter

C.1 Recursive Computation of Conditional Means and Covariances

- a) Exploiting the CHAPMAN-KOLMOGOROV EQUATION (9.11) and the LINEAR STATE SPACE MODEL in Eq. (10.1), we obtain

$$\begin{aligned}
 \hat{\mathbf{x}}_{n|n-1} &= \int_{\mathbb{X}} \mathbf{x}_n f_{\mathbf{x}_n|\mathbf{y}_{(n-1)}}(\mathbf{x}_n|\mathbf{y}_{(n-1)}) d\mathbf{x}_n \\
 &= \int_{\mathbb{X}} \int_{\mathbb{X}} \mathbf{x}_n f_{\mathbf{x}_n|\mathbf{x}_{n-1}}(\mathbf{x}_n|\mathbf{x}_{n-1}) f_{\mathbf{x}_{n-1}|\mathbf{y}_{(n-1)}}(\mathbf{x}_{n-1}|\mathbf{y}_{(n-1)}) d\mathbf{x}_n d\mathbf{x}_{n-1} \\
 &= \int_{\mathbb{X}} \underbrace{\int_{\mathbb{X}} \mathbf{x}_n f_{\mathbf{x}_n|\mathbf{x}_{n-1}}(\mathbf{x}_n|\mathbf{x}_{n-1}) d\mathbf{x}_n}_{\text{cond. expectation of } \mathbf{x}_n \text{ (refer to 10.1)}} f_{\mathbf{x}_{n-1}|\mathbf{y}_{(n-1)}}(\mathbf{x}_{n-1}|\mathbf{y}_{(n-1)}) d\mathbf{x}_{n-1} \\
 &= G_n \underbrace{\int_{\mathbb{X}} \mathbf{x}_{n-1} f_{\mathbf{x}_{n-1}|\mathbf{y}_{(n-1)}}(\mathbf{x}_{n-1}|\mathbf{y}_{(n-1)}) d\mathbf{x}_{n-1}}_{\text{cond. expectation of } \mathbf{x}_{n-1}} = G_n \mu_{\mathbf{x}_{n-1}|\mathbf{y}_{(n-1)}=\mathbf{y}_{(n-1)}}. \tag{C.1}
 \end{aligned}$$

The respective CONDITIONAL STATE COVARIANCE MATRIX of \mathbf{X}_n is obtained as

$$\begin{aligned}
\mathbf{C}_{\mathbf{X}_{n|n-1}} &= \mathbb{E} \left[(\mathbf{X}_n - \hat{\mathbf{x}}_{n|n-1}) (\mathbf{X}_n - \hat{\mathbf{x}}_{n|n-1})^\top \middle| \mathbf{Y}_{(n-1)} = \mathbf{y}_{(n-1)} \right] \\
&= \mathbb{E} \left[(\mathbf{G}_n(\mathbf{X}_{n-1} - \hat{\mathbf{x}}_{n-1|n-1}) + \mathbf{V}_n) (\mathbf{G}_n(\mathbf{X}_{n-1} - \hat{\mathbf{x}}_{n-1|n-1}) + \mathbf{V}_n)^\top \middle| \mathbf{Y}_{(n-1)} = \mathbf{y}_{(n-1)} \right] \\
&= \mathbb{E} \left[(\mathbf{G}_n(\mathbf{X}_{n-1} - \hat{\mathbf{x}}_{n-1|n-1})) (\mathbf{G}_n(\mathbf{X}_{n-1} - \hat{\mathbf{x}}_{n-1|n-1}))^\top \middle| \mathbf{Y}_{(n-1)} = \mathbf{y}_{(n-1)} \right] + \\
&\quad \mathbb{E} \left[(\mathbf{V}_n) (\mathbf{V}_n)^\top \middle| \mathbf{Y}_{(n-1)} = \mathbf{y}_{(n-1)} \right] \\
&= \mathbf{G}_n \mathbf{C}_{\mathbf{X}_{n-1|n-1}} \mathbf{G}_n^\top + \mathbf{C}_{\mathbf{V}_n}, \tag{C.2}
\end{aligned}$$

since $\mathbf{X}_{n-1} - \hat{\mathbf{x}}_{n-1|n-1}$ and \mathbf{V}_n are obviously stochastically independent.

b) Since both sides of the BAYES RULE in (9.12) are Gaussian PDFs, a comparison of the exponents,¹

$$\log \left(f_{\mathbf{X}_n | \mathbf{Y}_{(n)}}(\mathbf{x}_n | \mathbf{y}_{(n)}) \right) \propto \log \left(f_{\mathbf{Y}_n | \mathbf{X}_n}(\mathbf{y}_n | \mathbf{x}_n) \right) + \log \left(f_{\mathbf{X}_n | \mathbf{Y}_{(n-1)}}(\mathbf{x}_n | \mathbf{y}_{(n-1)}) \right) + \dots,$$

and taking into account the LINEAR OBSERVATION MODEL in Eq. (10.2) results in

$$\hat{\mathbf{x}}_{n|n} = \hat{\mathbf{x}}_{n|n-1} + \mathbf{K}_n (\mathbf{y}_n - \mathbf{H}_n \hat{\mathbf{x}}_{n|n-1}), \tag{C.3}$$

$$\mathbf{C}_{\mathbf{X}_{n|n}} = \mathbf{C}_{\mathbf{X}_{n|n-1}} - \mathbf{K}_n \mathbf{H}_n \mathbf{C}_{\mathbf{X}_{n|n-1}}, \tag{C.4}$$

with the KALMAN GAIN MATRIX

$$\mathbf{K}_n = \mathbf{C}_{\mathbf{X}_{n|n-1}} \mathbf{H}_n^\top \left(\mathbf{H}_n \mathbf{C}_{\mathbf{X}_{n|n-1}} \mathbf{H}_n^\top + \mathbf{C}_{\mathbf{W}_n} \right)^{-1}. \tag{C.5}$$

¹Note, that we are using (9.12) instead of (9.13).

C.2 Proof of the Multivariate Version

From a comparison of the exponents of both sides of (9.12), we obtain

$$\begin{aligned}
 & (x_n - \mu_{x_{n|n}})^\top C_{x_{n|n}}^{-1} (x_n - \mu_{x_{n|n}}) \propto \\
 & (\underbrace{y_n - \mu_{y_n|x_n=x_n}}_{\text{refer to 10.2}})^\top C_{y_n|x_n=x_n}^{-1} (y_n - \mu_{y_n|x_n=x_n}) + (x_n - \mu_{x_{n|n-1}})^\top C_{x_{n|n-1}}^{-1} (x_n - \mu_{x_{n|n-1}}) = \\
 & (y_n - H_n x_n)^\top C_{y_n|x_n=x_n}^{-1} (y_n - H_n x_n) + (x_n - \mu_{x_{n|n-1}})^\top C_{x_{n|n-1}}^{-1} (x_n - \mu_{x_{n|n-1}}),
 \end{aligned}$$

and also

$$\begin{aligned}
 & x_n^\top C_{x_{n|n}}^{-1} x_n - 2x_n^\top C_{x_{n|n}}^{-1} \mu_{x_{n|n}} + \dots \propto \\
 & x_n^\top \left(C_{x_{n|n-1}}^{-1} + H_n^\top C_{y_n|x_n=x_n}^{-1} H_n \right) x_n - 2x_n^\top \left(H_n^\top C_{y_n|x_n=x_n}^{-1} y_n + C_{x_{n|n-1}}^{-1} \mu_{x_{n|n-1}} \right) + \dots
 \end{aligned}$$

By comparing the respective terms, we eventually obtain

$$\begin{aligned}
 C_{x_{n|n}} &= \left(C_{x_{n|n-1}}^{-1} + H_n^\top C_{y_n|x_n=x_n}^{-1} H_n \right)^{-1} \\
 \mu_{x_{n|n}} &= C_{x_{n|n}} \left(H_n^\top C_{y_n|x_n=x_n}^{-1} y_n + C_{x_{n|n-1}}^{-1} \mu_{x_{n|n-1}} \right).
 \end{aligned}$$

Finally, the MATRIX INVERSION LEMMA² yields

$$\begin{aligned}
 C_{\mathbf{x}_{n|n}} &= \left(C_{\mathbf{x}_{n|n-1}}^{-1} + H_n^\top C_{\mathbf{y}_n|\mathbf{x}_n=\mathbf{x}_n}^{-1} H_n \right)^{-1} \\
 &= C_{\mathbf{x}_{n|n-1}} - \underbrace{C_{\mathbf{x}_{n|n-1}} H_n^\top \left(\underbrace{C_{\mathbf{y}_n|\mathbf{x}_n=\mathbf{x}_n} + H_n C_{\mathbf{x}_{n|n-1}} H_n^\top}_{C_{W_n}} \right)^{-1} H_n C_{\mathbf{x}_{n|n-1}}}_{\text{Kalman Gain: } K_n}
 \end{aligned}$$

(*)

and

$$\begin{aligned}
 \mu_{\mathbf{x}_{n|n}} &= \left(C_{\mathbf{x}_{n|n-1}} - K_n H_n C_{\mathbf{x}_{n|n-1}} \right) \left(H_n^\top C_{\mathbf{y}_n|\mathbf{x}_n=\mathbf{x}_n}^{-1} \mathbf{y}_n + C_{\mathbf{x}_{n|n-1}}^{-1} \mu_{\mathbf{x}_{n|n-1}} \right) \\
 &= C_{\mathbf{x}_{n|n-1}} H_n^\top C_{\mathbf{y}_n|\mathbf{x}_n=\mathbf{x}_n}^{-1} \mathbf{y}_n + \mu_{\mathbf{x}_{n|n-1}} - K_n H_n C_{\mathbf{x}_{n|n-1}} H_n^\top C_{\mathbf{y}_n|\mathbf{x}_n=\mathbf{x}_n}^{-1} \mathbf{y}_n - K_n H_n \mu_{\mathbf{x}_{n|n-1}} \\
 &= K_n C_{\mathbf{y}_{n|n-1}} C_{\mathbf{y}_n|\mathbf{x}_n=\mathbf{x}_n}^{-1} \mathbf{y}_n + \mu_{\mathbf{x}_{n|n-1}} - K_n H_n C_{\mathbf{x}_{n|n-1}} H_n^\top C_{\mathbf{y}_n|\mathbf{x}_n=\mathbf{x}_n}^{-1} \mathbf{y}_n - K_n H_n \mu_{\mathbf{x}_{n|n-1}} \\
 &= \mu_{\mathbf{x}_{n|n-1}} + K_n \underbrace{\left(C_{\mathbf{y}_{n|n-1}} - H_n C_{\mathbf{x}_{n|n-1}} H_n^\top \right)}_{C_{\mathbf{y}_n|\mathbf{x}_n=\mathbf{x}_n} \text{ (refer to *)}} C_{\mathbf{y}_n|\mathbf{x}_n=\mathbf{x}_n}^{-1} \mathbf{y}_n - K_n H_n \mu_{\mathbf{x}_{n|n-1}} \\
 &= \mu_{\mathbf{x}_{n|n-1}} + K_n \left(\mathbf{y}_n - H_n \mu_{\mathbf{x}_{n|n-1}} \right).
 \end{aligned}$$

² $(A + BCD)^{-1} = A^{-1} - A^{-1}B(C^{-1} + DA^{-1}B)^{-1}DA^{-1}.$

References

- A. Papoulis and S. U. Pillai. *Probability, Random Variables and Stochastic Processes*, 4th edition, Mc Graw Hill, 2002.
- H. Stark and J. W. Woods. *Probability, Random Processes, and Estimation Theory for Engineers*, 2nd edition, Prentice Hall, 1994.
- W. Utschick. *Stochastische Signale*, Aktuelles Manuskript zur gleichnamigen Vorlesung an der Technischen Universität München, 2012.
- G. Strang. *Linear Algebra and Its Applications*, Harcourt Brace Jovanovich, Publishers, 1988.
- L. N. Trefethen and D. Bau III. *Numerical Linear Algebra*, Society for Industrial and Applied Mathematic, 1997.
- W. Utschick. *Mathematische Methoden der Signalverarbeitung*, Aktuelles Manuskript zur gleichnamigen Vorlesung an der Technischen Universität München, 2012.