

Vorlesungsskript

35-407

Mathematische Grundlagen  
der  
Nachrichtentechnik

Prof. Dr. J. L. Massey

Institut für Signal- und  
Informationsverarbeitung  
ETH Zürich

Herbst 1992

## Vorwort

Dieses Skript ist eine korrigierte und leicht revidierte Version des Skripts, das während dem Wintersemester 1987/88 laufend verteilt wurde.

Die für diese Vorlesung zuständigen Assistenten (Richard Gut, Andreas Gygi und Markus Hufschmid) haben dieses Skript mit dem Textsystem  $\text{\LaTeX}$  erstellt und viele Sprachfehler von mir verbessert.

Ich bin ihnen nicht nur für ihre ausserordentlich grosse Arbeit am Skript, sondern auch für all ihre Bestrebungen, diese neue Vorlesung zu einer erfolgreichen Lehrveranstaltung zu machen, sehr dankbar.

J. L. Massey, März 1988

In dieser Version wurden einige Fehler korrigiert. Die Seitenumbrüche und -zahlen sind aber nicht kontrolliert!

November 1989

Auch diesen Herbst konnten wieder einige (sicher nicht die letzten!) Fehler eliminiert werden und kleine Verbesserungen angebracht werden — so stimmen zum Beispiel die Querverweise wieder. Neu wurde ein Anhang über verallgemeinerte Funktionen und den Dirac-Stoss angefügt.

Herbst 1990

In dieser Auflage wurde neu das Kapitel 7.5 über die Maximum-Likelihood-Filterung (Viterbi-Algorithmus) eingefügt. Einige Fehler konnten eliminiert werden und gewisse formale Detailänderungen wurden vorgenommen.

Herbst 1991

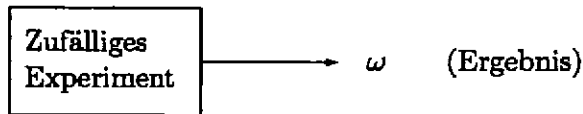
# Inhaltsverzeichnis

<b>1 WAHRSCHEINLICHKEITSTHEORIE</b>	<b>1</b>
1.1 Mathematisches Modell . . . . .	1
1.2 Charakterisierung von Zufallsgrößen . . . . .	5
1.3 Erwartungswerte von Zufallsgrößen . . . . .	8
1.4 Charakterisierung n-dimensionaler Zufallsgrößen . . . . .	10
1.5 Charakterisierung bedingter Ereignisse . . . . .	17
<b>2 ENTSCHEIDUNGSTHEORIE</b>	<b>24</b>
2.1 Problemstellung . . . . .	24
2.2 MAP- und ML-Regel . . . . .	25
2.3 Allgemeines Bayessches Problem . . . . .	26
2.4 Satz von Neyman-Pearson . . . . .	30
2.5 MINIMAX-Regel . . . . .	34
<b>3 SCHÄTZUNGSTHEORIE</b>	<b>36</b>
3.1 Problemstellung . . . . .	36
3.2 Bayessche MMSE-Schätzung . . . . .	36
3.3 Allgemeines Bayessches Schätzproblem . . . . .	38
<b>4 LINEARE ANNÄHERUNG IN EINEM SKALARPRODUKTRAUM</b>	<b>41</b>
4.1 Problemstellung . . . . .	41
4.2 Reeller Skalarproduktraum . . . . .	43
4.3 Orthogonalität . . . . .	45
4.4 Norm für einen reellen Vektorraum . . . . .	48
4.5 Optimale lineare Annäherung . . . . .	49
<b>5 LINEARE MMSE-SCHÄTZUNG</b>	<b>54</b>
5.1 Einführung . . . . .	54
5.2 Menge aller Zufallsgrößen als Vektorraum . . . . .	54
5.3 Zufallsgrößen mit endlichem zweitem Moment . . . . .	57
5.4 Lineare MMSE-Schätzung . . . . .	60
5.5 Lineare MMSE-Schätzung im Gaußschen Fall . . . . .	63
<b>6 ZEITDISKRETE STOCHASTISCHE PROZESSE UND LINEARE ZEITDISKRETE SYSTEME</b>	<b>66</b>
6.1 Stochastische Prozesse . . . . .	66
6.2 Lineare zeitdiskrete Systeme (LDS) . . . . .	67
6.3 Stochastische Prozesse als Eingang eines LDS . . . . .	68
6.4 Die z-Transformation . . . . .	71
<b>7 FILTERUNG STOCHASTISCHER SIGNALE</b>	<b>74</b>
7.1 FIR-Filter und die Wiener-Hopf-Gleichung . . . . .	74
7.2 Das nichtkausale Wiener-Filter . . . . .	78
7.3 Das kausale Wiener-Filter . . . . .	80
7.4 Zeitvariante Filterung: Das Kalman-Filter . . . . .	84
7.5 Maximum-Likelihood-Filterung (Viterbi-Algorithmus) . . . . .	93



# 1 WAHRSCHEINLICHKEITSTHEORIE

## 1.1 Mathematisches Modell



"Wahrscheinlichkeitssystem" =  $(\Omega, \mathcal{A}, P)$

$\Omega =$  Ergebnisraum ("Sample Space").

(oder "Vereinigung der elementaren Ereignisse")

= Menge der möglichen, einander ausschliessenden Ergebnisse.

Bei jedem Versuch des zufälligen Experimentes tritt genau ein Ergebnis  $\omega, \omega \in \Omega$ , ein.

### Beispiele:

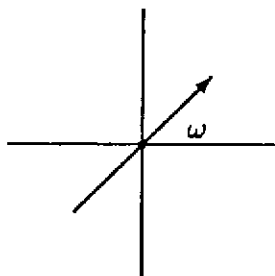
- (1) zufälliges Experiment: Werfen einer Münze.

$$\Omega = \{K, Z\}$$

- (2) zufälliges Experiment: Werfen von zwei nicht identischen Münzen.

$$\Omega = \{Kk, Kz, Zk, Zz\}$$

- (3) zufälliges Experiment: Drehung eines Zeigers.



$$\Omega = \{\omega : 0 \leq \omega < 2\pi\} .$$

$\mathcal{A}$  = Klasse der Ereignisse (oder "zufällige Ereignisse").

Grob gesagt, ist ein Ereignis eine Untermenge von  $\Omega$ .

Genau gesagt, ist ein Ereignis eine Untermenge von  $\Omega$ , die zu  $\mathcal{A}$  gehört.

Die Ereignisse in  $\mathcal{A}$  müssen eine  $\sigma$ -Algebra bilden, d.h. eine endliche oder abzählbar unendliche Reihe von Mengenoperationen von Ereignissen in  $\mathcal{A}$  muss immer ein Ereignis in  $\mathcal{A}$  liefern. Um eine  $\sigma$ -Algebra zu sein, genügt es, dass

(i)  $A \in \mathcal{A} \Rightarrow \bar{A} \in \mathcal{A}$   
wobei  $\bar{A} = \{\omega : \omega \in \Omega, \omega \notin A\}$   
und

(ii)

$$A_1, A_2, A_3, \dots \in \mathcal{A} \Rightarrow \bigcup_{i=1}^{+\infty} A_i \in \mathcal{A}.$$

**Beispiele: (Fortsetzung)**

(1)

$$\mathcal{A} = ( \underbrace{\{K\}, \{Z\}}_{\text{elementare Ereignisse}}, \underbrace{\{\}}_{\emptyset}, \underbrace{\{K, Z\}}_{\Omega} )$$

$A$  = Ereignis, dass entweder  $K$  oder  $Z$  eintritt  
=  $\{K, Z\} = \Omega$ .

$B$  = Ereignis, dass weder  $K$  noch  $Z$  eintritt  
=  $\{\} = \emptyset$ .

$C$  = Ereignis, dass  $K$  eintritt =  $\{K\}$ .

$D$  = Ereignis, dass  $Z$  eintritt =  $\{Z\}$ .

(2)

$$A = (\{\}, \{Kk, Kz, Zk, Zz\}, \{Kk\}, \{Kz\}, \{Zk\}, \{Zz\}, \\ \{Kk, Kz\}, \{Kk, Zk\}, \{Kk, Zz\}, \{Kz, Zk\}, \{Kz, Zz\}, \{Zk, Zz\}, \\ \{Kk, Kz, Zk\}, \{Kk, Kz, Zz\}, \{Kk, Zk, Zz\}, \{Kz, Zk, Zz\}).$$

$A$  = Ereignis, dass genau eine Zahl eintritt

$$= \{Kz, Zk\}.$$

= Ereignis, dass genau ein Kopf eintritt.

$B$  = Ereignis, dass ein oder mehrere Köpfe eintreten

$$= \{Kz, Kk, Zk\}.$$

(3)

$A$  = Klasse von allen Unterintervallen des Intervalls  $[0, 2\pi)$  zusammen mit allen endlichen oder abzählbar unendlichen Summen solcher Intervalle.

$A_1$  = Ereignis, dass der Zeiger im ersten Quadrant landet

$$= \{\omega : 0 \leq \omega < \frac{\pi}{2}\}.$$

Ein Ereignis tritt genau dann ein, wenn das Ergebnis des zufälligen Experimentes zu diesem Ereignis gehört. Also treten bei jedem Versuch des zufälligen Experimentes mehrere Ereignisse ein.

$P$  = Wahrscheinlichkeitsmass ("Probability Measure").

$P$  ist eine reellwertige Funktion mit dem Definitionsbereich  $\mathcal{A}$ , d.h.

$$P : \mathcal{A} \rightarrow \mathbb{R},$$

die die Axiome von Kolmogorov (1903-1987) erfüllen muss, namentlich:

(i)  $A \in \mathcal{A} \Rightarrow 0 \leq P[A] \leq 1;$

(ii)  $P[\Omega] = 1;$

(iii)

$$A_i \cap A_j = \emptyset \quad \Rightarrow \quad \begin{cases} P[\bigcup_{i=1}^n A_i] = \sum_{i=1}^n P[A_i] \\ P[\bigcup_{i=1}^{\infty} A_i] = \sum_{i=1}^{+\infty} P[A_i]. \end{cases}$$

( $A_i \cap A_j$  nennt man auch Verbundereignis)

Es gelten dann zum Beispiel:

$$A \cap \bar{A} = \emptyset \Rightarrow P[A \cup \bar{A}] = P[A] + P[\bar{A}]$$

Deswegen gilt:

$$P[\Omega] = 1 = P[A] + P[\bar{A}].$$

$$P[\bar{A}] = 1 - P[A] \Rightarrow P[\emptyset] = 0.$$

$\Omega$  heisst das sichere Ereignis.

$\emptyset$  heisst das unmögliche Ereignis.

Ein Ereignis und nur ein Ereignis hat eine Wahrscheinlichkeit.

Nichts in der Wahrscheinlichkeitstheorie sagt, was das physikalisch richtige Wahrscheinlichkeitsmass  $P$  sein muss!

Beispiele: (Fortsetzung)

- (1) Wir wählen  $P[\{K\}] = 0.3$   
 $\Rightarrow P[\{Z\}] = 1 - P[\{K\}] = 0.7$   
 $P[\emptyset] = 0$   
 $P[\{Z, K\}] = P[\Omega] = 1.$

- (2) Wir wählen  $P[\{Kk\}] = P[\{Kz\}] = P[\{Zk\}] = \frac{1}{4}$   
 $\Rightarrow P[\{Zz\}] = \frac{1}{4}$   
 $P[A] = P[\{Zk\} \cup \{Kz\}] = P[\{Zk\}] + P[\{Kz\}] = \frac{1}{2}$   
usw.

- (3) Wir wählen  $P[A] = \frac{L(A)}{2\pi}$   $A \in \mathcal{A}$

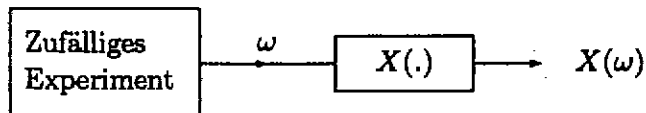
dabei ist  $L(A)$  gleich der Summe der Längen der sich nicht überlappenden Intervalle, die  $A$  bilden.

$$L(A_1) = \frac{\pi}{2} \quad P[A_1] = \frac{\pi/2}{2\pi} = \frac{1}{4}.$$



## 1.2 Charakterisierung von Zufallsgrößen

Eine Zufallsgröße ("random variable")  $X$  ist eine Funktion  $X, X: \Omega \rightarrow \mathbb{R}$ , so, dass für jedes  $x \in \mathbb{R}$  die Menge  $\{\omega: X(\omega) < x\}$  ein Ereignis in  $\mathcal{A}$  ist.



$F_X$  = Verteilungsfunktion von  $X$ .

$F_X: \mathbb{R} \rightarrow \mathbb{R}$

$$F_X(x) \triangleq \underbrace{P[\{\omega: X(\omega) < x\}]}_{\text{gekürzt als } P[X < x]}.$$

Eigenschaften:

(i)  $x_1 < x_2 \implies F_X(x_1) \leq F_X(x_2)$ .

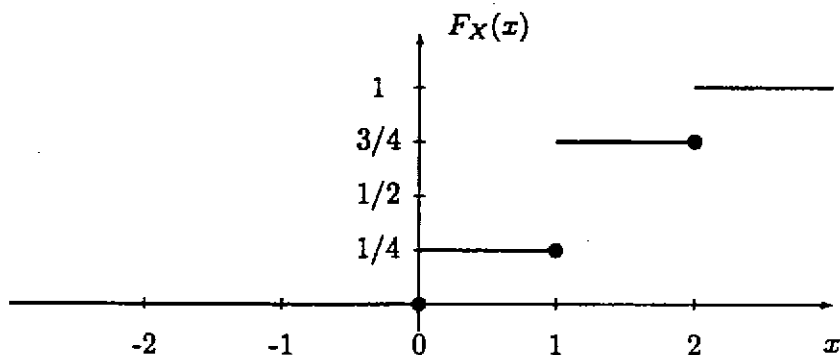
und

(ii)  $\lim_{x \rightarrow -\infty} F_X(x) = 0, \quad \lim_{x \rightarrow \infty} F_X(x) = 1;$

Beispiele: (Fortsetzung)

(2)  $X$  = Anzahl Zahlen, die eintreten

$$P[X = 0] = P[X = 2] = \frac{1}{4} \quad P[X = 1] = \frac{1}{2}$$



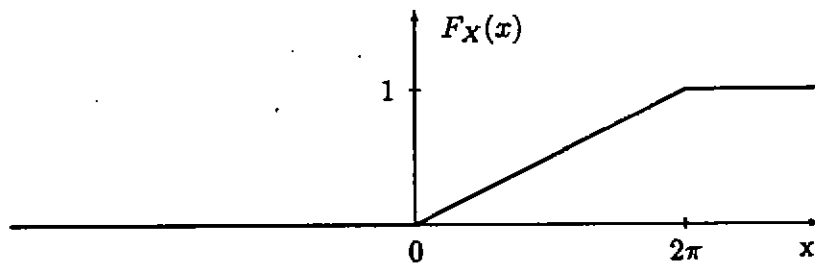
$F_X$  ist immer linksseitig stetig.

- (3)  $X$  = Winkel zwischen dem Zeiger und der horizontalen Achse. Das heisst  $X = \omega$ . Es folgt, dass

$$\{\omega : X(\omega) < x\} = \begin{cases} \emptyset & x \leq 0 \\ \{\omega : 0 \leq \omega < x\} & 0 < x \leq 2\pi \\ \Omega & x > 2\pi \end{cases}$$

Aber  $L(\{\omega : 0 \leq \omega < x\}) = x$ ,  $0 < x \leq 2\pi$ .

$$F_X(x) = P[X < x] = \begin{cases} 0 & x \leq 0 \\ \frac{x}{2\pi} & 0 < x \leq 2\pi \\ 1 & x > 2\pi. \end{cases}$$



$p_X$  : Wahrscheinlichkeitsdichte von  $X$

$$p_X : \mathbb{R} \rightarrow \mathbb{R} \tag{1.1}$$

$$p_X \triangleq \frac{d}{dx} F_X(x). \tag{1.2}$$

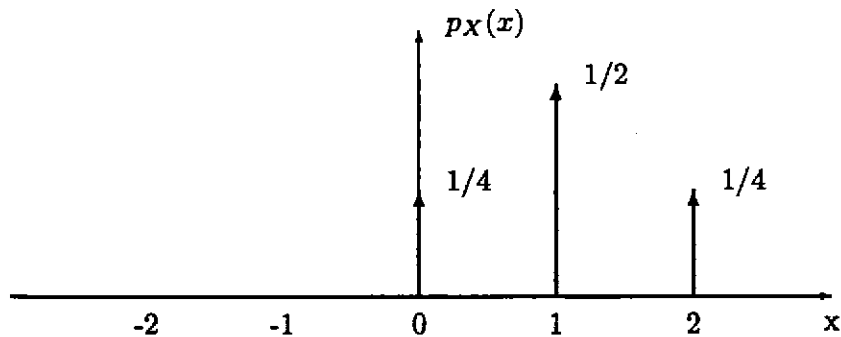
Eigenschaften:

- (i)  $p_X(x) \geq 0$  (oder enthält eine positive Diracfunktion auf dem Punkt  $x$ ), für jedes  $x \in \mathbb{R}$ .
- (ii)  $\int_{-\infty}^{+\infty} p_X(x) dx = 1$ .

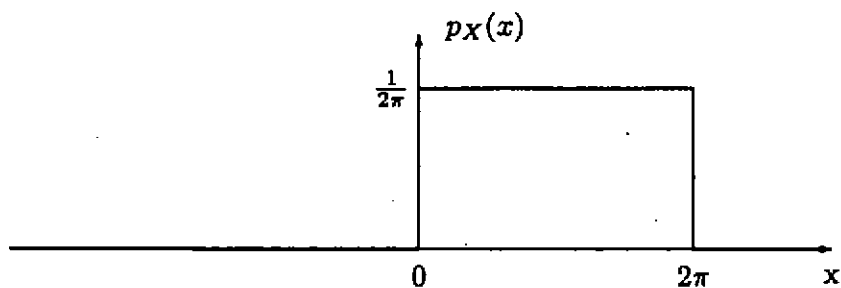
Diese beiden Eigenschaften sind die einzigen, die eine Wahrscheinlichkeitsdichte erfüllen muss.

Beispiele: (Fortsetzung)

$$(2) \quad p_X(x) = \frac{1}{4}\delta(x) + \frac{1}{2}\delta(x-1) + \frac{1}{4}\delta(x-2).$$



(3)



Es folgt aus der Definition von  $p_X(x)$ , dass

$$F_X(x) = \int_{-\infty}^x p_X(z) dz \quad \text{ist.}$$

Wenn  $X(\Omega)$  entweder endlich oder abzählbar unendlich ist, dann wird die Zufallsgröße als **diskret** bezeichnet. Äquivalent ist  $X$  genau dann diskret, wenn  $p_X$  aus einer Reihe von Diracfunktionen besteht. Wenn  $F_X$  stetig ist, dann wird auch die Zufallsgröße  $X$  als **stetig** bezeichnet.

N.B. Es gibt Zufallsgrößen, die weder diskret noch stetig sind. Dazu gehört z.B. eine Zufallsgröße, deren Verteilungsfunktion  $F_X$  folgende Eigenschaften aufweist:

- $F_X$  ist stückweise stetig (d.h. stetig bis auf einige Sprungstellen)
- $F_X$  ist nicht stückweise konstant (da in diesem Fall die entsprechende Zufallsgröße diskret wäre)

**Beispiele:** (Fortsetzung)

- (2)  $X$  ist diskret, weil  $X(\Omega) = \{0, 1, 2\}$  ist.
- (3)  $X$  ist stetig, weil  $F_X$  stetig ist. (N.B. Es spielt keine Rolle, ob  $p_X$  stetig oder nichtstetig ist.)

$P_X$ : Wahrscheinlichkeitsfunktion der diskreten Zufallsgröße  $X$ .  
( $X(\Omega)$  bezeichnet den Wertebereich der Funktion  $X$ .)

$$P_X : X(\Omega) \rightarrow R \quad (1.3)$$

$$P_X(x) \triangleq \underbrace{P[\{\omega : X(\omega) = x\}]}_{\text{gekürzt als } P[X=x]} \quad (1.4)$$

**Beispiel:** (Fortsetzung)

(2)

$x$	$P_X(x)$
0	1/4
1	1/2
2	1/4

Sei  $X$  diskret, dann gilt

$$p_X(x) = \sum_i P_X(x_i) \cdot \delta(x - x_i),$$

wobei  $x_1, x_2, \dots$  die möglichen Werte von  $X$  sind.

### 1.3 Erwartungswerte von Zufallsgrößen

Sei  $f$  eine reellwertige Funktion mit einem Definitionsbereich, der  $X(\Omega)$  einschliesst, dann ist der Erwartungswert von  $f(X)$  wie folgt definiert:

$$E[f(X)] = \int_{-\infty}^{+\infty} f(x) p_X(x) dx.$$

Falls  $X$  diskret ist, kann der Erwartungswert auch folgendermassen berechnet werden:

$$E[f(X)] = \sum_i f(x_i) \cdot P_X(x_i).$$

Sei  $f(x) = c$ , dann gilt offensichtlich:

$$E[c] = \int_{-\infty}^{+\infty} c \cdot p_X(x) dx = c \cdot \int_{-\infty}^{+\infty} p_X(x) dx = c.$$

Für  $f(x) = x^n$  erhält man das sogenannte **n-te Moment von X**, nämlich

$$E[X^n] = \int_{-\infty}^{+\infty} x^n \cdot p_X(x) dx.$$

Im diskreten Fall kann das n-te Moment wie folgt bestimmt werden:

$$E[X^n] = \sum_i (x_i)^n \cdot P_X(x_i).$$

**Beispiele (Fortsetzung)**

$$(2) \quad E[X] = 0 \cdot \frac{1}{4} + 1 \cdot \frac{1}{2} + 2 \cdot \frac{1}{4} = 1$$

$$E[X^2] = 0^2 \cdot \frac{1}{4} + 1^2 \cdot \frac{1}{2} + 2^2 \cdot \frac{1}{4} = \frac{3}{2}.$$

$$(3) \quad E[X] = \int_{-\infty}^{+\infty} x \cdot p_X(x) dx = \int_0^{2\pi} x \cdot \frac{1}{2\pi} dx = \pi$$

$$E[X^2] = \int_{-\infty}^{+\infty} x^2 \cdot p_X(x) dx = \int_0^{2\pi} x^2 \cdot \frac{1}{2\pi} dx = \frac{4}{3} \cdot \pi^2.$$

**Linearität des Erwartungswertoperators:**

$$E[c_1 f_1(X) + c_2 f_2(X)] = c_1 E[f_1(X)] + c_2 E[f_2(X)].$$

Die Varianz von  $X$  ist wie folgt definiert:

$$\text{Var}[X] \triangleq E[(X - m_X)^2],$$

wobei  $m_X = E[X]$ . Eine wichtige Beziehung ist die folgende:

$$\begin{aligned} \text{Var}[X] &= E[X^2 - 2m_X X + (m_X)^2] \\ &= E[X^2] - 2m_X E[X] + (m_X)^2 && \text{(Linearität)} \\ &= E[X^2] - (E[X])^2. \end{aligned}$$

**Beispiele (Fortsetzung)**

$$(2) \quad \text{Var}[X] = E[X^2] - (E[X])^2 = \frac{3}{2} - 1^2 = \frac{1}{2}.$$

$$(3) \quad \text{Var}[X] = \frac{4}{3}\pi^2 - \pi^2 = \frac{\pi^2}{3}.$$

Sei  $a < b$ , dann gilt:

$$\{\omega : X(\omega) < b\} = \underbrace{\{\omega : X(\omega) < a\} \cup \{\omega : a \leq X(\omega) < b\}}_{\text{unvereinbare Ereignisse}}.$$

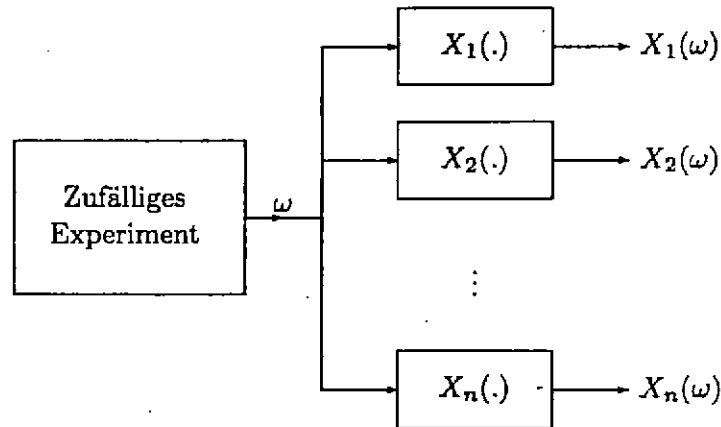
Also

$$\begin{aligned} P[X < b] &= P[X < a] + P[a \leq X < b] \\ P[a \leq X < b] &= P[X < b] - P[X < a] \\ &= F_X(b) - F_X(a) \\ &= \int_a^b p_X(x) dx. \end{aligned}$$

Falls  $\mathcal{R}$  ein Teilgebiet der reellen Achse bezeichnet, dann gilt allgemein:

$$P[X \in \mathcal{R}] = \int_{\mathcal{R}} p_X(x) dx.$$

## 1.4 Charakterisierung n-dimensionaler Zufallsgrößen



$\underline{X} = (X_1, X_2, \dots, X_n)$  heisst ein **n-dimensionaler Zufallsvektor**.

Bemerkung:

In diesem Skript bezeichnet  $(X_1, X_2, \dots, X_n)$  einen **Kolonnenvektor** und  $[X_1, X_2, \dots, X_n]$  einen **Zeilenvektor**.

$F_{\underline{X}}$ : **n-dimensionale Verteilungsfunktion von  $\underline{X}$** .

$$F_{\underline{X}}: R^n \rightarrow R$$

$$F_{\underline{X}}(\underline{x}) \triangleq P\{\{\omega : X_1(\omega) < x_1 \text{ und } \dots \text{ und } X_n(\omega) < x_n\}\}.$$

abgekürzt:  $P[X_1 < x_1, \dots, X_n < x_n]$  oder sogar:  $P[\underline{X} < \underline{x}]$

Eigenschaften:

$$(i) \lim_{x_i \rightarrow -\infty} F_{\underline{X}}(x_1, \dots, x_i, \dots, x_n) = 0$$

$$\lim_{x_i \rightarrow +\infty} F_{\underline{X}}(x_1, \dots, x_i, \dots, x_n) = F_{\underline{X}-}(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n),$$

wobei  $\underline{X}- = (X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n)$ .

(ii)  $F_{\underline{X}}(x_1, \dots, x_n)$  nimmt zu (oder genauer: nimmt nicht ab) mit zunehmenden  $x_i$ .

$p_{\underline{X}}$ : **n-dimensionale Verteilungsdichte (Verbunddichte) von  $\underline{X}$** .

$$p_{\underline{X}}: R^n \rightarrow R$$

$$p_{\underline{X}}(\underline{x}) \triangleq \frac{\partial^n}{\partial x_1 \dots \partial x_n} F_{\underline{X}}(x_1, \dots, x_n).$$

Eigenschaften:

$$(i) p_{\underline{X}}(\underline{x}) \geq 0 \quad \text{für jedes } \underline{x} \in R^n$$

(oder enthält eine positive, mehrdimensionale Diracfunktion auf dem Punkt  $\underline{x}$ ).

$$(ii) \int_{-\infty}^{+\infty} p_{\underline{X}}(x_1, \dots, x_i, \dots, x_n) dx_i = p_{\underline{X}-}(x_1, \dots, x_{i-1}, x_{i+1}, \dots, x_n),$$

wobei  $\underline{X}- = (X_1, \dots, X_{i-1}, X_{i+1}, \dots, X_n)$ .

N.B. Aus (ii) folgt, dass

$$\int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} p_{\underline{X}}(x_1, \dots, x_n) dx_1 \dots dx_n = 1.$$

$P_{\underline{X}}$ : n-dimensionale Wahrscheinlichkeitsfunktion des diskreten Zufallsvektors  $\underline{X}$ .

$$P_{\underline{X}} : \underline{X}(\Omega) \rightarrow R$$

$$P_{\underline{X}}(\underline{x}) \triangleq P[\underline{X} = \underline{x}].$$

Sei  $f$  eine reellwertige Funktion mit einem Definitionsbereich, der  $\underline{X}(\Omega)$  einschliesst, dann ist der Erwartungswert wie folgt definiert:

$$E[f(\underline{X})] \triangleq \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty} f(x_1, \dots, x_n) \cdot p_{\underline{X}}(x_1, \dots, x_n) dx_1 \dots dx_n.$$

Falls  $\underline{X}$  diskret ist, dann gilt:

$$E[f(\underline{X})] = \sum_{\underline{x}_1} \dots \sum_{\underline{x}_n} f(x_1, \dots, x_n) \cdot P_{\underline{X}}(x_1, \dots, x_n).$$

**Linearität des Erwartungswertoperators**

$$E[c_1 f_1(\underline{X}) + c_2 f_2(\underline{X})] = c_1 E[f_1(\underline{X})] + c_2 E[f_2(\underline{X})].$$

Insbesondere gilt immer:

$$E[X_1 + X_2 + \dots + X_n] = E[X_1] + E[X_2] + \dots + E[X_n].$$

Die Zufallsgrössen  $X_1, X_2, \dots, X_n$  heissen **unabhängig**, falls

$$F_{\underline{X}}(x_1, x_2, \dots, x_n) = F_{X_1}(x_1) \cdot F_{X_2}(x_2) \cdot \dots \cdot F_{X_n}(x_n),$$

für alle  $(x_1, \dots, x_n)$  in  $R^n$ .

Oder äquivalent: Die stetigen Zufallsgrössen  $X_1, X_2, \dots, X_n$  sind unabhängig, falls

$$p_{\underline{X}}(x_1, x_2, \dots, x_n) = p_{X_1}(x_1) \cdot p_{X_2}(x_2) \cdot \dots \cdot p_{X_n}(x_n),$$

für alle  $(x_1, \dots, x_n)$  in  $R^n$ .

Ebenfalls äquivalent: Die diskreten Zufallsgrössen  $X_1, X_2, \dots, X_n$  sind unabhängig, falls

$$P_{\underline{X}}(x_1, x_2, \dots, x_n) = P_{X_1}(x_1) \cdot P_{X_2}(x_2) \cdot \dots \cdot P_{X_n}(x_n),$$

für alle  $(x_1, \dots, x_n)$  in  $\underline{X}(\Omega)$ .

Falls  $X$  und  $Y$  unabhängig sind, dann gilt:

$$E[X \cdot Y] = E[X] \cdot E[Y]$$

**Beweis:**

$$\begin{aligned} E[X \cdot Y] &= \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} x \cdot y \cdot p_{XY}(x, y) \, dx \, dy \\ &= \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} x \cdot y \cdot p_X(x) \cdot p_Y(y) \, dx \, dy && \text{(unabhängig)} \\ &= \int_{-\infty}^{+\infty} x \cdot p_X(x) \, dx \int_{-\infty}^{+\infty} y \cdot p_Y(y) \, dy \\ &= E[X] \cdot E[Y]. \end{aligned} \quad \square$$

Es ist nun leicht zu beweisen, dass

$$\text{Var}[X_1 + X_2 + \dots + X_n] = \text{Var}[X_1] + \text{Var}[X_2] + \dots + \text{Var}[X_n].$$

gilt, falls  $X_1, X_2, \dots, X_n$  unabhängig sind. (Tatsächlich gilt diese Gleichung mit der schwächeren Annahme, dass  $X_1, X_2, \dots, X_n$  unkorreliert sind, wie unten definiert wird.)

Sei  $\mathcal{R}$  ein Teilgebiet von  $R^n$ , dann gilt:

$$P[\underline{X} \in \mathcal{R}] = \int \dots \int_{\mathcal{R}} p_{\underline{X}}(\underline{x}) \, dx_1 \dots dx_n.$$

Im diskreten Fall gilt:

$$P[\underline{X} \in \mathcal{R}] = \sum \dots \sum_{\mathcal{R}} P_{\underline{X}}(\underline{x}).$$

Die Zufallsgrößen  $X_1, X_2, \dots, X_n$  heissen i.i.d. (nach dem Englischen "independent and identically distributed"), falls sie unabhängig sind und dieselbe Verteilung besitzen. (vgl. Übungsaufgabe).

Die Kovarianz der Zufallsgrößen  $X$  und  $Y$  ist wie folgt definiert:

$$\text{Cov}(X, Y) = E[(X - m_X)(Y - m_Y)],$$

wobei  $m_X = E[X]$  und  $m_Y = E[Y]$ .

**Bemerkungen:**



(1)  $Cov(X, X) = Var[X].$

(2) Seien  $X$  und  $Y$  unabhängig, dann ist

$$\begin{aligned} Cov(X, Y) &= E[XY - m_X Y - m_Y X + m_X m_Y] \\ &= E[XY] - m_X E[Y] - m_Y E[X] + m_X m_Y \\ (\text{unabhängig}) \Rightarrow &= E[X]E[Y] - m_X m_Y - m_Y m_X + m_X m_Y \\ &= 0. \end{aligned}$$

(3)  $Cov(X, Y) = E[XY] - E[X] \cdot E[Y].$

Wenn  $Cov(X, Y) = 0$ , gilt nicht unbedingt, dass  $X$  und  $Y$  unabhängig sind.

Die Zufallsgrößen  $X_1, X_2, \dots, X_n$  heissen **unkorreliert**, wenn

$$Cov(X_i, X_j) = 0 \text{ für alle } i \neq j.$$

Wenn  $X_1, X_2, \dots, X_n$  unabhängig sind, dann sind sie auch unkorreliert. Die Umkehrung gilt aber im allgemeinen nicht.

Die Kovarianzmatrix des Zufallsvektors  $\underline{X} = (X_1, X_2, \dots, X_n)$  ist die Matrix

$$\Lambda = \begin{bmatrix} \lambda_{11} & \dots & \lambda_{1n} \\ \vdots & & \vdots \\ \lambda_{n1} & \dots & \lambda_{nn} \end{bmatrix}$$

wobei  $\lambda_{ij} = Cov(X_i, X_j)$ .

Seien  $X_1, \dots, X_n$  unkorreliert. Dann ist  $\Lambda$  eine Diagonalmatrix, d.h.  $\lambda_{ij} = 0$  für alle  $i \neq j$ .

Die Zufallsgröße  $X$  heisst **gaussverteilt** (oder **normalverteilt**), falls

$$p_X(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-m)^2}{2\sigma^2}}$$

wobei  $\sigma^2 = Var[X]$  und  $m = E[X]$ .

Der Zufallsvektor  $\underline{X} = (X_1, X_2, \dots, X_n)$  heisst **gaussverteilt** (oder **normalverteilt**), wenn

$$p_{\underline{X}}(\underline{x}) = \frac{1}{(\sqrt{2\pi})^n \sqrt{|\det(\Lambda)|}} e^{-\frac{1}{2}(\underline{x}-\underline{m})^T \Lambda^{-1}(\underline{x}-\underline{m})},$$

wobei  $\Lambda$  die Kovarianzmatrix von  $\underline{X}$  und  $\underline{m} = (E[X_1], \dots, E[X_n])$  ist.

Wenn  $X_1, X_2, \dots, X_n$  unkorreliert sind und  $\underline{X} = (X_1, X_2, \dots, X_n)$  gaussverteilt ist, dann sind  $X_1, X_2, \dots, X_n$  auch unabhängig.

**Beweis:**

$$\begin{aligned} \Lambda &= \text{diag}(\sigma_1^2, \sigma_2^2, \dots, \sigma_n^2) && \text{(unkorreliert)} \\ \Rightarrow \Lambda^{-1} &= \text{diag}(\sigma_1^{-2}, \sigma_2^{-2}, \dots, \sigma_n^{-2}) \\ \Rightarrow (\underline{x} - \underline{m})^T \Lambda^{-1} (\underline{x} - \underline{m}) &= \sum_{i=1}^n (x_i - m_i)^2 / \sigma_i^2. && (m_i = E[X_i]) \end{aligned}$$

Weil  $\underline{X}$  gaussverteilt ist, gilt:

$$\begin{aligned} p_{\underline{X}}(\underline{x}) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma_i} e^{-(x_i - m_i)^2 / 2\sigma_i^2} \\ &= \prod_{i=1}^n p_{X_i}(x_i). \end{aligned}$$

□

Aus obigem Beweis sehen wir, dass die einzelnen Komponenten des Zufallsvektors  $X_1, X_2, \dots, X_n$  gaussverteilt sind. Dies ist ein Sonderfall folgender allgemeinen Tatsache:

Sei der Zufallsvektor  $\underline{X} = (X_1, \dots, X_n)$  gaussverteilt, dann ist auch jeder Untervektor von  $\underline{X}$  gaussverteilt.

z.B. sei  $(X_1, X_2, X_3)$  gaussverteilt  $\Rightarrow$  auch  $(X_1, X_3)$  ist gaussverteilt.

Es ist nicht schwierig, diese Eigenschaft direkt zu beweisen. Wir integrieren dazu die Wahrscheinlichkeitsdichten der unerwünschten Komponenten von  $-\infty$  bis  $\infty$ . Hier wollen wir aber auf die Details eines solchen Beweises verzichten.

Wir stellen jedoch fest, dass diese Eigenschaft einen Spezialfall folgender wichtiger Eigenschaft darstellt:

Sei  $\underline{X} = (X_1, \dots, X_n)$  gaussverteilt und sei  $A$  eine  $m \times n$  reelle Matrix mit Rang  $m$ . Dann ist auch  $\underline{Y} = A\underline{X} = (Y_1, \dots, Y_m)$  gaussverteilt.

(Grob gesagt: eine lineare Operation auf gaussverteilte Zufallsgrößen liefert wieder gaussverteilte Zufallsgrößen!)

Um diese sehr wichtige Eigenschaft leichter beweisen zu können, führen wir folgende Notation ein:

Sei  $\mathbf{G}(\underline{X})$  eine Matrix (oder ein Vektor) wobei jede Komponente  $g_{ij}(\underline{X})$  eine reellwertige Funktion des Zufallsvektors  $\underline{X}$  ist. Dann schreiben wir

$$E[\mathbf{G}(\underline{X})]$$

für die Matrix (oder den Vektor) mit entsprechenden Komponenten  $E[g_{ij}(\underline{X})]$ . Mit dieser Notation können wir für einen beliebig verteilten Zufallsvektor  $\underline{X} = (X_1, \dots, X_n)$  schreiben:

$$E[\underline{X}] = (E[X_1], \dots, E[X_n]) = \underline{m}_X,$$

und

$$E[(\underline{X} - \underline{m}_X)(\underline{X} - \underline{m}_X)^T] = \Lambda_X.$$

Die  $(i, j)$ -te Komponente der Matrix innerhalb des Erwartungswertoperators ist  $(X_i - E[X_i])(X_j - E[X_j])$ .

Eine wichtige Eigenschaft dieser Matrixformulierung ist:

$$E[\mathbf{A}\mathbf{G}(\underline{X})\mathbf{B}] = \mathbf{A}E[\mathbf{G}(\underline{X})]\mathbf{B},$$

was direkt aus der Linearität des Erwartungsoperators folgt, wobei  $\mathbf{A}$  und  $\mathbf{B}$  Matrizen mit den entsprechenden Dimensionen sind.

Sei  $\underline{X} = (X_1, X_2, \dots, X_n)$  ein beliebig verteilter Zufallsvektor mit dem Mittelwert  $\underline{m}_X$  und der Kovarianzmatrix  $\Lambda_X$  und sei  $\mathbf{A}$  eine beliebige reelle  $m \times n$  Matrix, dann ist  $\underline{Y} = (Y_1, Y_2, \dots, Y_m) = \mathbf{A}\underline{X}$  ein Zufallsvektor mit dem Mittelwert

$$\underline{m}_Y = \mathbf{A}\underline{m}_X$$

und der Kovarianzmatrix

$$\Lambda_Y = \mathbf{A}\Lambda_X\mathbf{A}^T.$$

**Beweis:**

$$\underline{Y} = \mathbf{A}\underline{X} \quad \Rightarrow \quad E[\underline{Y}] = \mathbf{A}E[\underline{X}] \quad (\text{Linearität})$$

$$\text{d.h.} \quad \underline{m}_Y = \mathbf{A}\underline{m}_X.$$

$$\begin{aligned} \Lambda_Y &= E[(\underline{Y} - \underline{m}_Y)(\underline{Y} - \underline{m}_Y)^T] \\ &= E[(\mathbf{A}\underline{X} - \mathbf{A}\underline{m}_X)(\mathbf{A}\underline{X} - \mathbf{A}\underline{m}_X)^T] \\ &= E[(\mathbf{A}(\underline{X} - \underline{m}_X))(\mathbf{A}(\underline{X} - \underline{m}_X))^T] \\ &= E[\mathbf{A}(\underline{X} - \underline{m}_X)(\underline{X} - \underline{m}_X)^T\mathbf{A}^T] \\ &= \mathbf{A}E[(\underline{X} - \underline{m}_X)(\underline{X} - \underline{m}_X)^T]\mathbf{A}^T \\ &= \mathbf{A}\Lambda_X\mathbf{A}^T \end{aligned}$$

□

Weiter brauchen wir folgende allgemeine Tatsache:

Sei  $G : R^n \rightarrow R^n$  eine invertierbare Funktion  $\underline{y} = G(\underline{x})$  wobei  $y_i = g_i(x_1, \dots, x_n)$  so, dass  $g_i$  alle partiellen Ableitungen erster Ordnung für  $i = 1, \dots, n$  besitzt. Weiter sei  $\underline{X} = (X_1, \dots, X_n)$  ein stetiger Zufallsvektor. Dann hat  $\underline{Y} = G(\underline{X})$  die Wahrscheinlichkeitsdichte:

$$p_{\underline{Y}}(\underline{y}) = \frac{1}{|\mathcal{J}_G(G^{-1}(\underline{y}))|} p_{\underline{X}}(G^{-1}(\underline{y}))$$

wobei  $\mathcal{J}_G(\underline{x})$  die Jacobische Determinante von  $G(\underline{x})$ , d.h.

$$\mathcal{J}_G(\underline{x}) = \det \begin{bmatrix} \frac{\partial g_1(\underline{x})}{\partial x_1} & \dots & \frac{\partial g_1(\underline{x})}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_n(\underline{x})}{\partial x_1} & \dots & \frac{\partial g_n(\underline{x})}{\partial x_n} \end{bmatrix}$$

ist.

**Beweis:**

Sei  $\underline{x}$  ein beliebiger Punkt in  $R^n$  und  $\mathcal{R}_1$  ein beliebig kleines Teilgebiet von  $R^n$  mit dem Volumen  $\Delta_1$ , das  $\underline{x}$  enthält und sei  $\Delta_2$  das Volumen von  $\mathcal{R}_2 = G(\mathcal{R}_1)$ , dann gilt:

$$P[\underline{Y} \in \mathcal{R}_2] = P[\underline{X} \in \mathcal{R}_1]$$

sodass

$$p_{\underline{Y}}(G(\underline{x}))\Delta_2 \approx p_{\underline{X}}(\underline{x})\Delta_1.$$

Aus der Analysis wissen wir, dass

$$\frac{\Delta_2}{\Delta_1} \approx |\mathcal{J}_G(\underline{x})|$$

und dass die beiden Näherungen genau werden, falls  $\Delta_1 \rightarrow 0$ . Somit erhalten wir:

$$p_{\underline{Y}}(G(\underline{x})) \cdot |\mathcal{J}_G(\underline{x})| = p_{\underline{X}}(\underline{x})$$

was äquivalent zu unserer Behauptung ist.

**Bemerkung:**

Manchmal ist  $\mathcal{J}_{G^{-1}}(\underline{y})$  leichter zu berechnen als  $\mathcal{J}_G(G^{-1}(\underline{y}))$ . In diesem Falle ist folgende Tatsache nützlich:

$$\mathcal{J}_{G^{-1}}(\underline{y}) = \frac{1}{\mathcal{J}_G(G^{-1}(\underline{y}))}.$$

Sei  $\underline{X} = (X_1, \dots, X_n)$  ein beliebiger stetiger Zufallsvektor und sei  $A$  eine invertierbare reelle  $n \times n$  Matrix, dann gilt für den Zufallsvektor  $\underline{Y} = A\underline{X}$ :

$$p_{\underline{Y}}(\underline{y}) = \frac{1}{|\det(A)|} p_{\underline{X}}(A^{-1}\underline{y}).$$

**Beweis:**

Definieren wir  $G(\underline{x}) = A\underline{x}$ , dann ist  $G$  invertierbar mit

$$g_i(\underline{x}) = a_{i1}x_1 + \dots + a_{in}x_n,$$

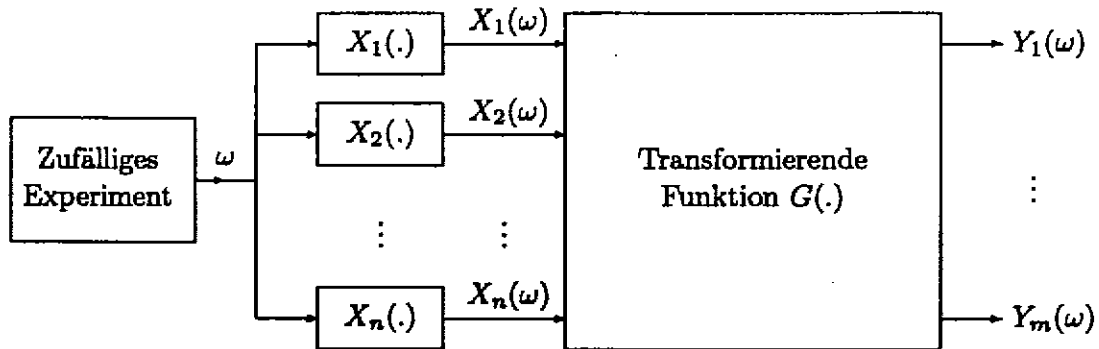
wobei  $a_{ij}$  die  $(i, j)$ -te Komponente von  $A$  ist. Also haben wir

$$\frac{\partial g_i(\underline{x})}{\partial x_j} = a_{ij}$$

so, dass  $J_G(\underline{x}) = \det(A)$  für alle  $\underline{x}$ . □

Der Leser sollte jetzt in der Lage sein, die 'sehr wichtige Eigenschaft' auf Seite 14 selber zu beweisen (vgl. Übungsaufgabe).

**Transformation von Zufallsgrößen:**



## 1.5 Charakterisierung bedingter Ereignisse

Seien  $A$  und  $B$  Ereignisse. Dann ist die **bedingte Wahrscheinlichkeit von  $B$  unter der Bedingung, dass  $A$  eintritt**, wie folgt definiert:

$$P(B|A) \triangleq \frac{P[A \cap B]}{P[A]}, \text{ falls } P[A] \neq 0.$$

Es ist jetzt aber zu bemerken, dass das "P" beim  $P(B|A)$  nicht dasselbe ist, wie das "P" beim  $P[B]$ . Tatsächlich sollte man es folgendermassen formulieren:

Sei  $A$  ein Ereignis mit  $P[A] \neq 0$ . Dann ist  $P(\cdot|A)$ , das **bedingte Wahrscheinlichkeitsmass unter der Bedingung, dass  $A$  eintritt**, die Funktion:

$$P(\cdot|A) : \mathcal{A} \rightarrow \mathcal{R}$$

$$P(B|A) = \frac{P[A \cap B]}{P[A]}.$$

Es folgt nun unmittelbar aus den Eigenschaften für das Wahrscheinlichkeitsmass  $P$ , dass

(i)

$$B \in \mathcal{A} \Rightarrow 0 \leq P(B|A) \leq 1;$$

(ii)

$$P(\Omega|A) = 1;$$

(iii)

$$B_i \cap B_j = \emptyset \Rightarrow \begin{cases} P\left(\bigcup_{i=1}^n B_i|A\right) = \sum_{i=1}^n P(B_i|A) \\ P\left(\bigcup_{i=1}^{\infty} B_i|A\right) = \sum_{i=1}^{\infty} P(B_i|A) \end{cases} \quad i \neq j$$

Wir sehen, dass  $P(\cdot|A)$  die drei Axiome von Kolmogorov erfüllt. Es ist aber zu bemerken, dass  $P(\cdot|A)$  von  $P$  eindeutig bestimmt ist, d.h.  $P(\cdot|A)$  ist eine Art "sekundäres" Wahrscheinlichkeitsmass.

Wenn  $P[A] = 0$  ist, kann man  $P(\cdot|A)$  nicht wie oben definieren. Man sagt dann, "dass  $P(B|A)$  mathematisch nicht definiert ist, wenn  $P[A] = 0$  ist". Aber bei mehreren Anwendungen kommt es vor, dass ein bedingtes Wahrscheinlichkeitsmass  $P(\cdot|A)$  "physikalisch bestimmt" ist, obwohl  $P[A] = 0$ . Solange dieses  $P(\cdot|A)$  die drei Axiome von Kolmogorov erfüllt, kriegt man keine mathematischen Widersprüche, wenn man dieses "physikalisch bestimmte"  $P(\cdot|A)$  anwendet. Vielleicht sollte man am besten sagen, falls  $P[A] = 0$  und  $P(\cdot|A)$  physikalisch nicht bestimmt ist, dass  $P(\cdot|A)$  ein unbekanntes Mass ist, das die Axiome von Kolmogorov erfüllt. Egal ob  $P[A] \neq 0$  oder  $P[A] = 0$  gilt dann immer:

$$P[A \cap B] = P[A]P(B|A).$$

Es gilt dann auch immer die Multiplikationsregel:

$$P[B_1 \cap B_2 \cap \dots \cap B_n] = P[B_1]P(B_2|B_1) \dots P(B_n|B_1 \cap B_2 \cap \dots \cap B_{n-1}).$$

**Beweis:**

$$\begin{aligned} P[B_1 \cap B_2 \cap \dots \cap B_n] &= P[(B_1 \cap \dots \cap B_{n-1}) \cap B_n] \\ &= P[B_1 \cap \dots \cap B_{n-1}]P(B_n|B_1 \cap B_2 \cap \dots \cap B_{n-1}) \\ &\text{usw.} \end{aligned}$$

□

Die Ereignisse  $A_1, A_2, \dots, A_k$  bilden eine komplette (oder vollständige) Klasse von paarweise unvereinbaren Ereignissen, falls gilt:

$$A_i \cap A_j = \emptyset, \text{ für alle } i \neq j$$

und

$$A_1 \cup A_2 \cup \dots \cup A_k = \Omega .$$

Der sogenannte "Satz von der totalen Wahrscheinlichkeit":

Wenn  $A_1, A_2, \dots, A_k$  eine komplette Klasse von paarweise unvereinbaren Ereignissen bilden, dann gilt:

$$P[B] = P[A_1]P(B|A_1) + P[A_2]P(B|A_2) + \dots + P[A_k]P(B|A_k) .$$

Beweis:

$$\begin{aligned} B &= B \cap \Omega \\ &= B \cap (A_1 \cup A_2 \cup \dots \cup A_k) \\ &= \underbrace{(B \cap A_1) \cup (B \cap A_2) \cup \dots \cup (B \cap A_k)}_{\text{eine Summe von paarweise unvereinbaren Ereignissen}} \end{aligned}$$

Aber nach dem dritten Axiom von Kolmogorov gilt:

$$\begin{aligned} P[B] &= \sum_{i=1}^k P[B \cap A_i] \\ &= \sum_{i=1}^k P[A_i]P(B|A_i) . \end{aligned}$$

□

**Korollar 1: (Die Bayessche Formel)**

Wenn  $A_1, A_2, \dots, A_k$  eine komplette Klasse von paarweise unvereinbaren Ereignissen bilden und  $P[B] \neq 0$ , dann gilt:

$$P(A_i|B) = \frac{P[A_i]P(B|A_i)}{\sum_{j=1}^k P[A_j]P(B|A_j)}$$

Beweis: Der Nenner ist gerade  $P[B]$ ! □

Die bedingte Verteilungsfunktion von  $X$  unter der Bedingung, dass  $A$  eintritt, ist die Funktion:

$$F_{X|A} : R \rightarrow R$$

$$F_{X|A}(x) = \underbrace{P(\{\omega : X(\omega) < x\}|A)}_{\text{abgekürzt als } P(X < x|A)}$$

Die bedingte Wahrscheinlichkeitsdichte von  $X$  unter der Bedingung, dass  $A$  eintritt, ist die Funktion:

$$p_{X|A} : R \rightarrow R$$

$$p_{X|A}(x) = \frac{d}{dx} F_{X|A}(x).$$

Eigenschaften:

- (i)  $p_{X|A}(x) \geq 0$  (oder enthält eine positive Diracfunktion im Punkt  $x$ ), für alle  $x \in R$ .
- (ii)

$$\int_{-\infty}^{\infty} p_{X|A}(x) dx = 1.$$

Bedeutung: Sei  $\mathcal{R}$  ein Teilgebiet von  $R$ . Dann ist:

$$\underbrace{P(\{\omega : X(\omega) \in \mathcal{R}\}|A)}_{\text{abgekürzt als } P(X \in \mathcal{R}|A)} = \int_{\mathcal{R}} p_{X|A}(x) dx.$$

Korollar 2:

Wenn  $A_1, A_2, \dots, A_k$  eine komplette Klasse von paarweise unvereinbaren Ereignissen bilden, dann gilt:

$$p_X(x) = P[A_1]p_{X|A_1}(x) + \dots + P[A_k]p_{X|A_k}(x).$$

Beweis:

$$\begin{aligned} F_X(x) &= P[X < x] \\ &= P[A_1]P(X < x|A_1) + \dots + P[A_k]P(X < x|A_k) \\ &= P[A_1]F_{X|A_1}(x) + \dots + P[A_k]F_{X|A_k}(x). \end{aligned}$$



Zum Schluss bleibt nur noch abzuleiten. □

Der bedingte Erwartungswert von  $f(X)$  unter der Bedingung, dass  $A$  eintritt, ist wie folgt definiert:

$$E[f(X)|A] = \int_{-\infty}^{\infty} f(x)p_{X|A}(x)dx .$$

**Korollar 3: ("Satz des totalen Erwartungswertes")**

Wenn  $A_1, A_2, \dots, A_k$  eine komplette Klasse von paarweise unvereinbaren Ereignissen bilden, dann gilt:

$$E[f(X)] = P[A_1]E[f(X)|A_1] + \dots + P[A_k]E[f(X)|A_k] .$$

**Beweis:** Man braucht nur die Formel für  $p_X(x)$  von Korollar 2 in die Definition

$$E[f(X)] = \int_{-\infty}^{\infty} f(x)p_X(x)dx$$

einzusetzen. □

Wir überlassen es nun dem Leser, alle obigen Formeln, die für bedingte Wahrscheinlichkeiten gelten, so zu ergänzen, dass statt einer Zufallsgrösse  $X$  ein **Zufallsvektor**  $\underline{X}$  betrachtet wird.

Falls  $Y$  eine stetige Zufallsgrösse ist, wissen wir, dass  $P[Y = y] = 0$  ist für alle  $y \in R$ . Trotzdem ist es sinnvoll, für ein stetiges  $Y$  und für  $p_Y(y) \neq 0$ , ein bedingtes Wahrscheinlichkeitsmass  $P(\cdot | Y = y)$  wie folgt zu definieren:

$$P(B|Y = y) = \lim_{\epsilon \rightarrow 0^+} P(B|y \leq Y < y + \epsilon) .$$

Es ist dann leicht zu beweisen (vgl. Übungsaufgabe), dass folgendes gilt:

$$P(B|Y = y) = \frac{p_{Y|B}(y)P[B]}{p_Y(y)}$$

falls  $Y$  stetig und  $p_Y(y) \neq 0$  ist.

Eine unmittelbare Folge dieser Tatsache ist die integrale Form des Satzes von der totalen Wahrscheinlichkeit:

$$P[B] = \int_{-\infty}^{\infty} P(B|Y = y)p_Y(y)dy , \text{ für alle } B \in \mathcal{A} .$$

**Beweis:**

$$\begin{aligned} \int_{-\infty}^{\infty} P(B|Y=y)p_Y(y)dy &= \int_{-\infty}^{\infty} \frac{p_{Y|B}(y)P[B]}{p_Y(y)} p_Y(y)dy \\ &= P[B] \int_{-\infty}^{\infty} p_{Y|B}(y)dy = P[B]. \end{aligned}$$

□

Weil  $P(\cdot|Y=y)$  ein echtes Wahrscheinlichkeitsmass ist, ist es nun auch konsequent, die bedingte Verteilungsfunktion  $F_{X|Y=y}$  und die bedingte Wahrscheinlichkeitsdichte  $p_{X|Y=y}$  wie folgt zu definieren:

$$F_{X|Y=y}(x) = P(X < x|Y=y)$$

und

$$p_{X|Y=y}(x) = \frac{d}{dx} F_{X|Y=y}(x).$$

Wir werden jetzt demonstrieren, dass  $p_{X|Y=y}(x)$  sehr einfach durch  $p_{XY}(x,y)$  und  $p_Y(y)$  ausgedrückt werden kann:

$$\begin{aligned} p_{X|Y=y}(x) &= \frac{d}{dx} P(X < x|Y=y) \\ &= \lim_{\delta \rightarrow 0^+} \frac{P(X < x + \delta|Y=y) - P(X < x|Y=y)}{\delta} \\ &= \lim_{\delta \rightarrow 0^+} \frac{P(x \leq X < x + \delta|Y=y)}{\delta} \\ &= \lim_{\delta \rightarrow 0^+} \lim_{\epsilon \rightarrow 0^+} \frac{P(x \leq X < x + \delta | y \leq Y < y + \epsilon)}{\delta} \\ &= \lim_{\delta \rightarrow 0^+} \lim_{\epsilon \rightarrow 0^+} \frac{P[x \leq X < x + \delta, y \leq Y < y + \epsilon]}{\delta \cdot P[y \leq Y < y + \epsilon]} \\ &\quad (\text{Wir erinnern uns hier, dass } P[A, B] \text{ dasselbe wie } P[A \cap B] \text{ bedeutet)} \\ &= \lim_{\delta \rightarrow 0^+} \lim_{\epsilon \rightarrow 0^+} \frac{\int_x^{x+\delta} \int_y^{y+\epsilon} p_{XY}(\alpha, \beta) d\beta d\alpha}{\delta \int_y^{y+\epsilon} p_Y(\gamma) d\gamma} \\ &= \frac{p_{XY}(x, y)}{p_Y(y)} \end{aligned}$$

falls  $p_Y(y) \neq 0$  ist. Anstatt  $p_{X|Y=y}(x)$  schreibt man normalerweise  $p_{X|Y}(x|y)$ , obwohl die erste Bezeichnung präziser ist.

$$p_{X|Y}(x|y) \triangleq p_{X|Y=y}(x).$$

Obige Herleitung zeigt, dass folgendes gilt:

$$p_{X|Y}(x|y) = \frac{p_{XY}(x, y)}{p_Y(y)}$$

falls  $Y$  stetig und  $p_Y(y) \neq 0$  ist.

Falls nun  $p_Y(y) = 0$ , sagen wir, dass  $p_{X|Y}(x|y)$  "mathematisch nicht bestimmt" ist. In jedem Fall gilt aber:

$$p_{XY}(x, y) = p_{X|Y}(x|y)p_Y(y).$$

Weiter ist es sinnvoll, folgendes zu definieren:

$$E[f(X, Y)|Y = y] = \int_{-\infty}^{\infty} f(x, y)p_{X|Y}(x|y)dx$$

falls  $p_Y(y) \neq 0$  ist.

Eine unmittelbare Folge ist die integrale Form des Satzes vom totalen Erwartungswert:

$$E[f(X, Y)] = \int_{-\infty}^{\infty} E[f(X, Y)|Y = y]p_Y(y)dy.$$

**Beweis:**

$$\begin{aligned} \int_{-\infty}^{\infty} E[f(X, Y)|Y = y]p_Y(y)dy &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) \frac{p_{XY}(x, y)}{p_Y(y)} p_Y(y) dx dy \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y) p_{XY}(x, y) dx dy \\ &= E[f(X, Y)]. \end{aligned}$$

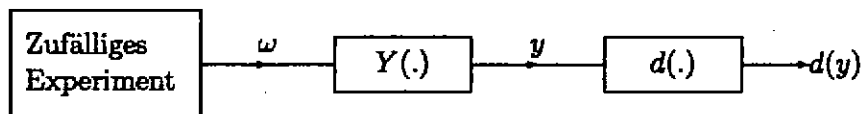
□

Man muss hier wohl kaum erwähnen, dass alle obigen Resultate so ergänzt werden können, dass statt einer Zufallsgrösse  $Y$  ein Zufallsvektor  $\underline{Y}$  betrachtet wird. Wir überlassen es dem Leser, diese Verallgemeinerung durchzuführen.

## 2 ENTSCHEIDUNGSTHEORIE

### 2.1 Problemstellung

Mathematisches Modell:



**Das allgemeine Problem:** Die interessierenden Möglichkeiten  $A_1, A_2, \dots, A_k$  bilden eine komplette Klasse von paarweise unvereinbaren Ereignissen (d.h. bei jedem Versuch des zufälligen Experiments tritt genau eines von diesen  $k$  Ereignissen ein). Gegeben sei die Beobachtung  $Y = y$ . Wir wollen den Entscheid treffen, welche dieser  $k$  Möglichkeiten eingetreten ist. Die Funktion  $d$  ( $d: R \rightarrow \{1, 2, \dots, k\}$ ) ist die **Entscheidungsfunktion**.  $d(y) = i$  bedeute, dass wir entschieden haben, dass  $A_i$  eingetreten ist.

**Bemerkung:** Wir werden mit einer skalaren Beobachtung  $y$  arbeiten. Alle unsere Resultate werden aber sehr leicht zum Fall einer vektoriellen Beobachtung  $\underline{y}$  ergänzbar sein.

**Allgemeine Notation:** Wir schreiben

$$\mathcal{Y}_i = \{y : d(y) = i\}.$$

$\mathcal{Y}_i$  ist der Teil des Beobachtungsraums, worin wir entscheiden, dass  $A_i$  eingetreten ist. Die Spezifizierung von  $\mathcal{Y}_1, \mathcal{Y}_2, \dots, \mathcal{Y}_k$  ist der Spezifizierung von  $d(\cdot)$  äquivalent.

Zunächst müssen wir festhalten,

- (1) was dem Beobachter zum Voraus bekannt ist, und
- (2) was als Kriterium der Güte eines Entscheids gilt.

Die Wahl dieser beiden Dinge kann sehr unterschiedlich sein und liefert uns viele interessante Probleme.

## 2.2 MAP- und ML-Regel

### Erstes Problem:

- (1) Die bedingte Wahrscheinlichkeitsdichte  $p_{Y|A_i}$  für  $i = 1, 2, \dots, k$  als auch  $P[A_i]$  für  $i = 1, 2, \dots, k$  sind dem Beobachter bekannt.
- (2) Der Beobachter will  $P[F]$  minimieren, wobei  $F$  das Ereignis sei, dass der Entscheid falsch ist.

Für dieses Problem ist die optimale Entscheidungsregel leicht zu finden. Zuerst stellen wir fest, dass das Minimieren von  $P[F]$  äquivalent ist zum Maximieren von  $P[\bar{F}]$ . Dabei bedeutet  $\bar{F}$  (das Komplement von  $F$ ) das Ereignis, dass der Entscheid korrekt ist.

Dann folgt, dass

$$P(\bar{F}|Y = y) = P(A_{d(y)}|Y = y),$$

und ferner, aus der integralen Form des Satzes von der totalen Wahrscheinlichkeit, dass

$$P[\bar{F}] = \int_{-\infty}^{\infty} P(A_{d(y)}|Y = y)p_Y(y)dy. \quad (2.1)$$

Von (2.1) können wir schliessen, dass die Regel, die  $P[\bar{F}]$  maximiert, folgende ist: wähle  $d(y)$  als das  $i$ , das  $P[A_i|Y = y]$  maximiert.

Aber im (einzigen interessanten) Falle, wo  $p_Y(y) > 0$  ist, gilt:

$$P(A_{d(y)}|Y = y) = \frac{p_{Y|A_{d(y)}}(y)P[A_{d(y)}]}{p_Y(y)}. \quad (2.2)$$

Weil der Nenner positiv und unabhängig von  $d(y)$  ist, können wir schliessen, dass eine äquivalente Regel gilt: wähle  $d(y)$  als das  $i$ , das  $p_{Y|A_i}(y)P[A_i]$  maximiert.

Wir fassen zusammen:

Die Regel die  $P[F]$  minimiert, ist: Wähle für jedes  $y$ ,  $d(y)$  als das  $i$  (oder als eines der  $i$ , falls es mehr als eines gibt) das

$$p_{Y|A_i}(y)P[A_i]$$

maximiert. (Diese Regel wird als die MAP-Regel bezeichnet, weil diese Regel die maximale "a posteriori"-Wahrscheinlichkeit  $P(A_i|Y = y)$  gibt.)

Für jede Entscheidungsregel (optimale oder nicht optimale) können wir  $P[\bar{F}] = 1 - P[F]$  aus (2.1) mit Hilfe von (2.2) wie folgt berechnen:

$$P[\bar{F}] = \int_{-\infty}^{\infty} p_{Y|A_{d(y)}}(y)P[A_{d(y)}]dy.$$

Aber da  $d(y) = i$  zu  $y \in \mathcal{Y}_i$  äquivalent ist, gilt:

$$P[\bar{F}] = \sum_{i=1}^k \int_{\mathcal{Y}_i} p_{Y|A_i}(y) P[A_i] dy$$

oder

$$P[\bar{F}] = \sum_{i=1}^k \left( P[A_i] \int_{\mathcal{Y}_i} p_{Y|A_i}(y) dy \right).$$

Jetzt stellen wir fest, dass, für den Fall wo  $P[A_i] = 1/k$  ( $i = 1, 2, \dots, k$ ) ist, die MAP Regel sich folgendermassen vereinfacht: wähle  $d(y)$  als das  $i$ , das  $p_{Y|A_i}(y)$  maximiert. Diese Regel (unabhängig angewandt von  $P[A_i]$ ,  $i = 1, 2, \dots, k$ ) heisst die **Maximum-Likelihood- (ML-) Regel**. Diese Regel hat den Vorteil, dass es nicht nötig ist, die "a priori" Wahrscheinlichkeiten  $P[A_i]$ ,  $i = 1, 2, \dots, k$ , zu kennen. Wir haben aber gesehen, dass die ML-Regel  $P[F]$  im Allgemeinen nur dann minimiert, wenn  $P[A_i] = 1/k$ ,  $i = 1, 2, \dots, k$  ist.

**Bemerkung:** Fast immer in der Entscheidungstheorie nimmt man die bedingte Wahrscheinlichkeit  $p_{Y|A_i}$  für  $i = 1, 2, \dots, k$  als gegeben an. Hingegen ist es nur bei gewissen Problemen sinnvoll anzunehmen, dass auch die 'a priori'-Wahrscheinlichkeiten der interessierenden Möglichkeiten (d. h.  $P[A_i]$  für  $i = 1, 2, \dots, k$ ) bekannt sind. Man spricht von einem **Bayesschen Entscheidungsproblem**, falls die 'a priori'-Wahrscheinlichkeiten bekannt sind, da in diesem Fall die 'a posteriori'-Wahrscheinlichkeiten (d. h.  $P(A_i|Y = y)$  für  $i = 1, 2, \dots, k$ ) mit Hilfe der 'Bayesschen Formel' bestimmt werden können (vgl. (2.2) und beachte, dass  $p_Y(y) = P[A_1]p_{Y|A_1}(y) + \dots + P[A_k]p_{Y|A_k}(y)$ ).

### 2.3 Allgemeines Bayessches Problem

Zur Formulierung des allgemeinen Bayesschen Problems, führen wir den Begriff der 'Kosten' eines Entscheids ein. Es ist denkbar, dass es uns viel mehr 'kostet', uns für  $A_1$  zu entscheiden, wenn tatsächlich  $A_2$  eingetreten ist, als uns für  $A_2$  zu entscheiden, wenn tatsächlich  $A_1$  eingetreten ist. Wir führen deshalb eine Zufallsgrösse  $S$  ein, welche den 'Spesen' des Entscheids entspricht. Wir bezeichnen mit  $s_{ij}$  den Wert der Spesen, falls tatsächlich  $A_i$  eingetreten ist und wir uns für  $A_j$  entschieden haben. (Beachte, dass  $s_{ii}$  den Spesen bei einem korrekten Entscheid entspricht). Wir bezeichnen nun  $D_j$  als das Ereignis, dass wir uns für  $A_j$  entschieden haben, oder anders ausgedrückt, dass die Beobachtung  $Y$  dem Entscheidungsgebiet  $\mathcal{Y}_j$  angehört. Dann gilt:

$$S(\omega) = s_{ij}, \quad \text{falls } \omega \in A_i \cap D_j,$$

was die Zufallsgrösse  $S$  vollständig definiert. Ferner gilt offensichtlich:

$$E[S|A_i \cap D_j] = s_{ij}. \quad (2.3)$$

**Zweites Problem:**

- (1) Sowohl  $p_{Y|A_i}$  als auch  $P[A_i]$  für  $i = 1, 2, \dots, k$  sind dem Beobachter bekannt. Zudem sind die  $s_{ij}$  für  $i = 1, 2, \dots, k$ ;  $j = 1, 2, \dots, k$  bekannt.

(2) Der Beobachter will die Entscheidungsfunktion  $d(\cdot)$  so wählen, dass  $E[S]$  minimal wird.

In einem ersten Schritt unserer Herleitung der optimalen Entscheidungsregel, stellen wir fest, dass:

weil sowohl  $A_1, A_2, \dots, A_k$  als auch  $D_1, D_2, \dots, D_k$  eine komplette Klasse von paarweise unvereinbaren Ereignissen bilden, bilden auch die  $k^2$  Ereignisse  $A_i \cap D_j$ ,  $i = 1, 2, \dots, k$ ,  $j = 1, 2, \dots, k$  eine komplette Klasse von paarweise unvereinbaren Ereignissen.

Man darf aus diesem Grund den Satz vom totalen Erwartungswert anwenden:

$$\begin{aligned}
 E[S] &= \sum_{i=1}^k \sum_{j=1}^k E[S|A_i \cap D_j] \cdot P[A_i \cap D_j] \\
 (2.3) \Rightarrow &= \sum_{i=1}^k \sum_{j=1}^k s_{ij} \cdot P[A_i \cap D_j] \\
 &= \sum_{i=1}^k \sum_{j=1}^k s_{ij} \cdot P[A_i] \cdot P(D_j|A_i).
 \end{aligned}$$

Andererseits gilt auch:

$$\begin{aligned}
 P(D_j|A_i) &= P(Y \in \mathcal{Y}_j|A_i) \\
 &= \int_{\mathcal{Y}_j} p_{Y|A_i}(y) dy.
 \end{aligned}$$

Aus diesen beiden Ausdrücken resultiert die Formel:

$$E[S] = \sum_{i=1}^k P[A_i] \sum_{j=1}^k s_{ij} \int_{\mathcal{Y}_j} p_{Y|A_i}(y) dy \quad (2.4)$$

Diese Form ist für die Berechnung von  $E[S]$  für eine beliebige Entscheidungsregel  $d(\cdot)$  gut geeignet, da alle Größen auf der rechten Seite bekannt sind. Hingegen ist (2.4) nicht geeignet, um daraus eine optimale Entscheidungsregel abzuleiten. Wir formen (2.4) deshalb um:

$$E[S] = \sum_{j=1}^k \int_{\mathcal{Y}_j} \sum_{i=1}^k P[A_i] \cdot s_{ij} \cdot p_{Y|A_i}(y) dy.$$

Da  $y \in \mathcal{Y}_j$  äquivalent zu  $d(y) = j$  ist, können wir diese Formel nochmals anders schreiben:

$$E[S] = \int_{-\infty}^{+\infty} \left( \sum_{i=1}^k P[A_i] \cdot s_{id(y)} \cdot p_{Y|A_i}(y) \right) dy \quad (2.5)$$

Für jedes  $y$  können wir  $d(y)$  frei wählen. Um das Integral zu minimieren, wählen wir  $d(y)$  so, dass der Integrand für jedes  $y$  minimal wird. Wir erhalten folgende Regel:

Die Regel, um  $E[S]$  zu minimieren lautet: Wähle für jedes  $y$ ,  $d(y)$  als dasjenige  $j$  (oder als eines derjenigen  $j$ , falls es mehrere gibt), welches

$$\sum_{i=1}^k P[A_i] \cdot s_{ij} \cdot p_{Y|A_i}(y)$$

minimiert.

**Beispiel:** Wir betrachten im Folgenden den Spezialfall, dass es lediglich zwei interessierende Möglichkeiten gibt. Wir bezeichnen diese Möglichkeiten mit  $A_0$ , resp.  $A_1$ . Dies ist für den Fall  $k = 2$  die gebräuchlichere Bezeichnung als  $A_1$  und  $A_2$ . Ferner soll uns das daran erinnern, dass wir die Einschränkung  $k = 2$  gemacht haben.

Wir müssen nun offenbar für jedes  $y$  dasjenige  $j$  wählen, das

$$P[A_0] \cdot s_{0j} \cdot p_{Y|A_0}(y) + P[A_1] \cdot s_{1j} \cdot p_{Y|A_1}(y)$$

minimiert. Der Rand der Entscheidungsgebiete für  $j$  liegt dort, wo

$$P[A_0] \cdot s_{00} \cdot p_{Y|A_0}(y) + P[A_1] \cdot s_{10} \cdot p_{Y|A_1}(y) = P[A_0] \cdot s_{01} \cdot p_{Y|A_0}(y) + P[A_1] \cdot s_{11} \cdot p_{Y|A_1}(y)$$

ist; dort können wir  $j = 0$  oder  $j = 1$  wählen.

Wir können nun unsere Bedingung umformen zu

$$L(y) = \frac{p_{Y|A_0}(y)}{p_{Y|A_1}(y)} = \frac{P[A_1][s_{11} - s_{10}]}{P[A_0][s_{00} - s_{01}]}$$

$L(y)$  entspricht einer "likelihood ratio". Diese ist definiert als:

$$L(y) = \frac{p_{Y|A_0}(y)}{p_{Y|A_1}(y)} \tag{2.6}$$

Das allgemeine Bayessche Problem lässt sich für  $k = 2$  als "Likelihood-ratio"-Test mit der Schwelle

$$T = \frac{P[A_1][s_{11} - s_{10}]}{P[A_0][s_{00} - s_{01}]}$$

auffassen. Die Entscheidungsregel lautet:

$$\begin{aligned} y \in \mathcal{Y}_0, & \quad \text{falls } L(y) > T \\ y \in \mathcal{Y}_1, & \quad \text{falls } L(y) < T \end{aligned}$$

(Wir werden später auf die "likelihood ratio" zurückkommen.) ○

Im folgenden wollen wir einen Spezialfall betrachten:



Alle Spesen für einen falschen Entscheid sind gleich, d.h.

$$s_{ij} = \begin{cases} 1 & \text{falls } i \neq j \\ 0 & \text{falls } i = j \end{cases} \quad (2.7)$$

Mit  $F$  als Ereignis, dass wir einen falschen Entscheid treffen, gilt offensichtlich:

$$S(\omega) = \begin{cases} 1 & \text{falls } \omega \in F \\ 0 & \text{falls } \omega \in \bar{F} \end{cases}$$

Mit dem Satz des totalen Erwartungswerts folgt nun:

$$\begin{aligned} E[S] &= E[S|F] \cdot P[F] + E[S|\bar{F}] \cdot P[\bar{F}] \\ &= 1 \cdot P[F] + 0 \cdot P[\bar{F}] \\ &= P[F]. \end{aligned}$$

Wir erkennen, dass unser zweites Problem für den erwähnten Spezialfall ('alle Fehler kosten dasselbe') identisch mit dem ersten Problem ist. Dies ist nicht ohne weiteres ersichtlich, da die Entscheidungsregeln völlig verschieden aussehen. Dass es sich in Wirklichkeit um die gleichen Entscheidungsregeln handelt, kann man erkennen, wenn man die Spesen in (2.7) um 1 reduziert:

$$s'_{ij} = \begin{cases} 0 & \text{falls } i \neq j \\ -1 & \text{falls } i = j, \end{cases}$$

Dies führt offensichtlich zu einem äquivalenten Problem, weil  $S' = S - 1$  und somit  $E[S'] = E[S] - 1$  gilt. Dann ist

$$\sum_{i=1}^k P[A_i] \cdot s'_{ij} \cdot p_{Y|A_i}(y) = -P[A_j] \cdot p_{Y|A_j}(y),$$

zu minimieren, resp.

$$P[A_j] \cdot p_{Y|A_j}(y)$$

zu maximieren, was mit unserer optimalen Entscheidungsregel aus dem ersten Problem identisch ist.

Wir wollen nun noch untersuchen, welche Folgen es hat, wenn man die Spesen  $S(\Omega)$  folgendermassen modifiziert:

$$s'_{ij} = s_{ij} - \tau_i, \quad (2.8)$$

wobei für jedes  $i$ ,  $\tau_i$  eine beliebige reelle Zahl ist. Für die entsprechende Zufallsgrösse  $S'$  erhält man:

$$E[S'] = \sum_{i=1}^k P[A_i] \cdot E[S'|A_i].$$

Wegen

$$\omega \in A_i \Rightarrow S'(\omega) = S(\omega) - r_i$$

folgt

$$\begin{aligned} E[S'] &= \sum_{i=1}^k P[A_i] \cdot \underbrace{E[S - r_i | A_i]}_{E[S|A_i] - r_i} \\ &= \sum_{i=1}^k P[A_i] \cdot E[S|A_i] - \sum_{i=1}^k r_i \cdot P[A_i] \\ &= E[S] - \underbrace{\sum_{i=1}^k r_i \cdot P[A_i]}_{\text{unabhängig von } d(\cdot)} \end{aligned} \quad (2.9)$$

Die Summe auf der rechten Seite von (2.9) hängt nicht von der Entscheidungsregel ab. Es folgt deshalb: die Minimierung von  $E[S]$  ist äquivalent zur Minimierung von  $E[S']$ . Wählen wir  $r_i = s_{ii}$  für alle  $i$ , dann folgt daraus:  $s_{ii} = 0$  für alle  $i$ . Wir können somit ohne Einschränkung der Allgemeinheit die Spesenfunktion  $s_{ij}$  immer so wählen, dass  $s_{ii} = 0$  für alle  $i$  ist. Da es jedoch manchmal angenehm ist,  $s_{ii} \neq 0$  zu wählen, haben wir darauf verzichtet, diese Beschränkung von Anfang an einzuführen.

## 2.4 Satz von Neyman-Pearson

Wie schon früher erwähnt wurde, spricht man von einem nicht-Bayesschen Entscheidungsproblem, falls die  $p_{Y|A_i}$  für  $i = 1, 2, \dots, k$  bekannt, die 'a priori'-Wahrscheinlichkeiten  $P[A_i]$  jedoch nicht bekannt sind. Wir betrachten im folgenden ein solches nicht-Bayessches Problem. Der Einfachheit halber, aber auch, weil dieser Fall von besonderem Interesse ist, werden wir den Fall behandeln, dass es genau zwei interessierenden Möglichkeiten gibt. Wir bezeichnen diese - wie schon im Beispiel zum Bayesschen Fall - mit  $A_0$ , resp.  $A_1$ .

Zur Erinnerung: Die "likelihood ratio"  $L(y)$  ist definiert als:

$$L(y) = \frac{p_{Y|A_0}(y)}{p_{Y|A_1}(y)}$$

und spielt eine grosse Rolle in der nicht-Bayesschen Entscheidungstheorie (vgl. Übungsaufgabe). Es ist leicht einzusehen, dass die ML-Regel wie folgt ausgedrückt werden kann:

$$\begin{aligned} y &\in \mathcal{Y}_0, & \text{falls } L(y) > 1 \\ y &\in \mathcal{Y}_1, & \text{falls } L(y) < 1 \end{aligned}$$

Falls  $L(y) = 1$ , kann die Zuordnung  $y \in \mathcal{Y}_0$  oder  $y \in \mathcal{Y}_1$  beliebig gewählt werden.

(ML-Entscheidungsregel)

Im nicht-Bayesschen Fall können wir die Fehlerwahrscheinlichkeit **nicht** berechnen (siehe (2.19)). Falls jedoch die Entscheidungsregel bekannt ist, können wir die bedingte Fehlerwahrscheinlichkeit bestimmen:

$$P(F|A_0) = P(Y \in \mathcal{Y}_1|A_0) = \int_{\mathcal{Y}_1} p_{Y|A_0}(y) dy \quad (2.10)$$

$$P(F|A_1) = P(Y \in \mathcal{Y}_0|A_1) = \int_{\mathcal{Y}_0} p_{Y|A_1}(y) dy \quad (2.11)$$

Wir werden nun zunächst zeigen, dass die ML-Regel **robust** ist, in dem Sinne, dass sie immer eine ziemlich kleine Fehlerwahrscheinlichkeit  $P[F]$  liefert. Für die ML-Regel gilt offensichtlich:

$$y \in \mathcal{Y}_0 \quad \Rightarrow \quad \sqrt{L(y)} > 1.$$

Wenn wir den Integranden von (2.11) mit  $\sqrt{L(y)}$  multiplizieren, dann erhalten wir eine obere Grenze für die ML-Regel:

$$P(F|A_1) \leq \int_{\mathcal{Y}_0} \sqrt{p_{Y|A_0}(y) \cdot p_{Y|A_1}(y)} dy$$

Analog können wir den Integranden von (2.10) mit  $1/\sqrt{L(y)}$  multiplizieren und erhalten:

$$P(F|A_0) \leq \int_{\mathcal{Y}_1} \sqrt{p_{Y|A_1}(y) \cdot p_{Y|A_0}(y)} dy$$

Da gilt:  $\mathcal{Y}_0 \cap \mathcal{Y}_1 = \emptyset$  und  $\mathcal{Y}_0 \cup \mathcal{Y}_1$  gleich der gesamten reellen Achse, können wir summieren und erhalten die folgende Grenze für die ML-Regel:

$$P(F|A_0) + P(F|A_1) \leq \int_{-\infty}^{+\infty} \sqrt{p_{Y|A_1}(y) \cdot p_{Y|A_0}(y)} dy \quad (2.12)$$

Ausgehend von (2.12) wollen wir nun eine obere Grenze für die Fehlerwahrscheinlichkeit  $P[F]$  herleiten. Der Satz der totalen Wahrscheinlichkeit sagt aus, dass gilt:

$$P[F] = P(F|A_0) \cdot P[A_0] + P(F|A_1) \cdot P[A_1] \quad (2.13)$$

Die 'a priori'-Wahrscheinlichkeiten sind zwar nicht bekannt, aber wir wissen, dass  $P[A_0] \leq 1$  und  $P[A_1] \leq 1$ . Demnach folgt aus (2.13), dass

$$P[F] \leq P(F|A_0) + P(F|A_1)$$

ist. Mit (2.13) und (2.12) können wir schliessen:

$$P[F] \leq \int_{-\infty}^{+\infty} \sqrt{p_{Y|A_1}(y) \cdot p_{Y|A_0}(y)} dy \quad (2.14)$$

Diese Grenze gilt für die ML-Regel und wird nach ihrem Entdecker **Bhattacharyya-Schranke** genannt. (Eine Verallgemeinerung dieser Grenze für den Fall mit  $k$  Möglichkeiten ist möglich – vgl. Übungsaufgabe).

Die Bhattacharyya-Schranke zeigt, dass die ML-Regel robust ist, sie sagt aber nicht aus, ob die Regel in irgendeiner Weise optimal ist. Tatsächlich haben wir bis jetzt noch kein Gütekriterium für den nicht-Bayesschen Fall eingeführt. Dies wollen wir nun nachholen.

**Drittes Problem:**

- (1) Nur  $p_{Y|A_0}$  und  $p_{Y|A_1}$  sind dem Beobachter bekannt.
- (2) Für eine gegebene reelle Zahl  $\alpha$  ( $0 \leq \alpha \leq 1$ ) will der Beobachter  $P(F|A_1)$  minimieren unter der Einschränkung, dass  $P(F|A_0) \leq \alpha$  gilt.

Die Lösung dieses Problems ist in folgendem berühmten Satz enthalten:

**Satz von Neyman-Pearson:** Sei  $T$  eine beliebige positive reelle Zahl, so dass die Entscheidungsregel lautet:

$$\begin{aligned}
 &y \in \mathcal{Y}_0, \quad \text{falls } L(y) > T \\
 &y \in \mathcal{Y}_1, \quad \text{falls } L(y) < T
 \end{aligned}$$

Sei für diese Entscheidungsregel

$$P(F|A_0) = \alpha \text{ und } P(F|A_1) = \beta$$

Jede andere Entscheidungsregel, welche  $P(F|A_0) \leq \alpha$  liefert, hat zur Folge, dass  $P(F|A_1) \geq \beta$  wird. Ebenso muss  $P(F|A_0) \geq \alpha$  gelten, falls  $P(F|A_1) \leq \beta$  ist.

**Bemerkung:**  $T$  ist der Schwellwert (englisch: "threshold") für den "Likelihood ratio test", der im Satz von Neyman-Pearson beschrieben ist.

**Beweis:** Weil wir gleichzeitig über zwei verschiedene Entscheidungskriterien sprechen werden, schreiben wir  $F, \mathcal{Y}_0$  und  $\mathcal{Y}_1$  für den "Likelihood ratio test", welcher im Satz beschrieben ist und  $F', \mathcal{Y}'_0$  und  $\mathcal{Y}'_1$  für die entsprechenden Größen einer beliebigen anderen Entscheidungsregel, welche

$$P(F'|A_0) \leq \alpha$$

liefert.

Mit Formel (2.10) erhalten wir:

$$\begin{aligned}
 0 &\leq \underbrace{P(F|A_0)}_{=\alpha} - \underbrace{P(F'|A_0)}_{\leq \alpha} \\
 &= \int_{\mathcal{Y}_1} p_{Y|A_0}(y) dy - \int_{\mathcal{Y}'_1} p_{Y|A_0}(y) dy
 \end{aligned} \tag{2.15}$$

und

$$P(F'|A_1) - P(F|A_1) = \int_{\mathcal{Y}'_0} p_{Y|A_1}(y) dy - \int_{\mathcal{Y}_0} p_{Y|A_1}(Y) dy \quad (2.16)$$

Wir definieren jetzt die beiden Gebiete:

$$\begin{aligned} B &= \{y : y \in \mathcal{Y}_1 \setminus \mathcal{Y}'_1\} \\ B' &= \{y : y \in \mathcal{Y}'_1 \setminus \mathcal{Y}_1\} . \end{aligned}$$

Wegen  $\mathcal{Y}_1 = \bar{\mathcal{Y}}_0$  und  $\mathcal{Y}'_1 = \bar{\mathcal{Y}}'_0$  gilt dann auch:

$$\begin{aligned} B &= \{y : y \in \mathcal{Y}'_0 \setminus \mathcal{Y}_0\} \\ B' &= \{y : y \in \mathcal{Y}_0 \setminus \mathcal{Y}'_0\} \end{aligned}$$

Mit diesen Definitionen können wir (2.15):

$$0 \leq \int_B p_{Y|A_0}(y) dy - \int_{B'} p_{Y|A_0}(y) dy \quad (2.17)$$

und (2.16):

$$P(F'|A_1) - P(F|A_1) = \int_B p_{Y|A_1}(y) dy - \int_{B'} p_{Y|A_1}(y) dy \quad (2.18)$$

umschreiben. Es folgt nun sofort:

$$y \in B \Rightarrow y \in \mathcal{Y}_1 \Rightarrow L(y) \leq T \Rightarrow p_{Y|A_0} \leq T \cdot p_{Y|A_1}$$

und weiter:

$$y \in B' \Rightarrow y \in \mathcal{Y}_0 \Rightarrow L(y) \geq T \Rightarrow p_{Y|A_0} \geq T \cdot p_{Y|A_1} .$$

Unter Benützung dieser beiden Ungleichungen erhält man mit (2.17):

$$0 \leq T \cdot \int_B p_{Y|A_1}(y) dy - T \cdot \int_{B'} p_{Y|A_1}(y) dy$$

Da  $T > 0$  ist, gilt weiter:

$$\int_B p_{Y|A_1}(y) dy - \int_{B'} p_{Y|A_1}(y) dy \geq 0$$

Man erhält schlussendlich unter Verwendung der Formel (2.18):

$$P(F'|A_1) - P(F|A_1) \geq 0$$

respektive:

$$P(F'|A_1) \geq P(F|A_1) = \beta,$$

was genau gleich der Aussage des Satzes von Neyman-Pearson ist.  $\square$

Leider existiert keine Verallgemeinerung des Satzes von Neyman-Pearson für  $k$  interessierende Möglichkeiten, falls  $k > 2$ .

Für den allgemeinen Fall von  $k$  interessierenden Möglichkeiten gilt:

$$P[F] = P(F|A_1) \cdot P[A_1] + P(F|A_2) \cdot P[A_2] + \dots + P(F|A_k) \cdot P[A_k]. \quad (2.19)$$

Die bedingten Fehlerwahrscheinlichkeiten hängen nur von der Entscheidungsregel und den bedingten Wahrscheinlichkeitsdichten  $p_{Y|A_i}$ ,  $i = 1, 2, \dots, k$  ab, da:

$$P(F|A_i) = 1 - \int_{y_i} p_{Y|A_i}(y) dy.$$

Wenn wir diese bedingten Fehlerwahrscheinlichkeiten als fest und bestimmt durch die Entscheidungsregeln betrachten, erkennen wir, dass es eine schlechteste Wahl für die (unbekannten) 'a priori'-Wahrscheinlichkeiten  $P[A_i]$ ,  $i = 1, 2, \dots, k$  in (2.19) gibt, welche  $P[F]$  maximiert. Es ist dies der Fall, wenn wir  $P[A_i] = 1$  für dasjenige  $i$  wählen, welches  $P(F|A_i)$  maximiert. Für diesen schlechtesten Fall ist  $P[F]$  gleich:

$$P_{Fsf} \triangleq \max_{1 \leq i \leq k} P(F|A_i). \quad (2.20)$$

Wir wissen ferner, dass unter Garantie gilt:

$$P[F] \leq P_{Fsf}$$

für die tatsächlichen 'a priori' Wahrscheinlichkeiten.

## 2.5 MINIMAX-Regel

**Viertes Problem:**

- (1) Nur  $p_{Y|A_i}$  für  $i = 1, 2, \dots, k$  sind dem Beobachter bekannt.
- (2) Der Beobachter will  $P_{Fsf}$  minimieren.

Die Entscheidungsregel, welche dieses Problem löst, heisst **MINIMAX-Regel**, weil dadurch eine Minimierung der maximal möglichen Fehlerwahrscheinlichkeit erreicht wird. Im allgemeinen ist es nicht einfach, die MINIMAX-Regel zu finden, für den Fall  $k = 2$  erhält man die Lösung mit dem Satz von Neyman-Pearson.

**Korollar zum Satz von Neyman-Pearson:** Falls für den im Satz von Neyman-Pearson beschriebenen "Likelihood ratio test" gilt:  $P(F|A_0) = P(F|A_1)$ , dann ist diese Entscheidungsregel auch MINIMAX.

**Beweis:** Falls  $P(F|A_0) = P(F|A_1)$  gilt, dann folgt aus dem Satz von Neyman-Pearson:

$$\max \{P(F'|A_0), P(F'|A_1)\} \geq \max \{P(F|A_0), P(F|A_1)\},$$

was gemäss (2.20) zur Behauptung des Korollars äquivalent ist. □

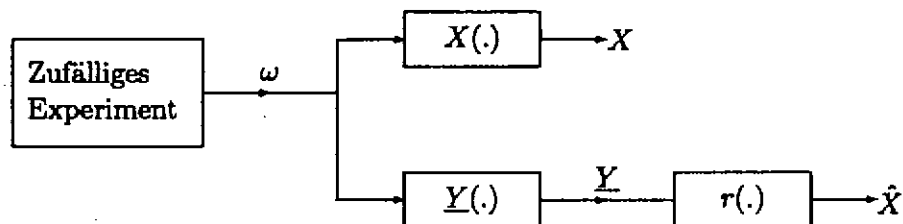
Da die ML-Regel ein "Likelihood ratio test" ist (nämlich mit dem Schwellwert  $T = 1$ ), folgt aus dem Korollar, dass eine ML-Regel auch MINIMAX ist, falls  $P(F|A_0) = P(F|A_1)$  gilt.

**Wichtige Bemerkung:** Alle Resultate des Kapitels über die Entscheidungstheorie lassen sich leicht ergänzen auf den Fall, wo die Beobachtung ein Zufallsvektor  $\underline{Y}$  anstelle einer Zufallsgrösse  $Y$  ist. Man braucht dazu lediglich die Integration über Teilgebiete der reellen Achse  $R$  durch Integration über dem entsprechenden Teilgebiet des  $n$ -dimensionalen Vektorraums  $R^n$  zu ersetzen.

### 3 SCHÄTZUNGSTHEORIE

#### 3.1 Problemstellung

Mathematisches Modell:



dabei ist die Zufallsgrösse  $X$  derjenige Parameter, den wir schätzen wollen, der Zufallsvektor  $\underline{Y} = (Y_1, Y_2, \dots, Y_n)$  die Beobachtung und die Funktion  $r$  ( $r : R^n \rightarrow R$ ) die Schätzungsregel.

[Normalerweise stellen wir uns vor, dass  $X$  die Beobachtung  $\underline{Y}$  in irgendeiner Weise verursacht hat. Diese Interpretation ist aber keine Voraussetzung.] Im Allgemeinen kann  $X$  unendlich viele mögliche Werte annehmen, wobei für jede Schätzungsregel  $P[\hat{X} = X] = 0$  gilt. Es hat dann keinen Sinn, eine Schätzung als "falsch" oder "richtig" zu bezeichnen, sondern man bezeichnet eine Schätzung als "näher" oder "ferner" vom richtigen Wert. Dies ist der wesentliche Unterschied zwischen der Entscheidungstheorie und der Schätzungstheorie. Die beiden Theorien sind jedoch insofern gleich, dass wir auch hier festhalten müssen,

- (1) was dem Beobachter bekannt ist, und
- (2) was als Kriterium der Güte einer Schätzung gilt.

Bei fast allen Schätzungsproblemen ist  $p_{\underline{Y}|X}(\cdot|x)$  bekannt für alle  $x \in X(\Omega)$ . Jedoch ist die a priori Wahrscheinlichkeitsdichte  $p_X$  nicht immer bekannt.

#### 3.2 Bayessche MMSE-Schätzung

Erstes Problem (Bayessche "Minimum mean-squared-error- (MMSE-)" Schätzung):

- (1) Sowohl die bedingte Wahrscheinlichkeitsdichte  $p_{\underline{Y}|X}(\cdot|x)$  für alle  $x \in X(\Omega)$ , als auch  $p_X$  sind bekannt.
- (2) Wir wollen die Schätzungsregel  $r(\cdot)$  finden, die  $E[(X - \hat{X})^2]$  minimiert.



Bemerkung:  $E[(X - \hat{X})^2]$  ist der "mean-squared error (MSE)".

Bevor wir dieses Problem lösen, lohnt es sich, das folgende einfache Lemma zu beweisen.

**Lemma:** Sei  $X$  eine beliebige Zufallsgrösse und  $c$  eine reelle Zahl. Dann gilt:

$$E[(X - c)^2] \geq \text{Var}[X] .$$

Das Gleichheitszeichen gilt dann und nur dann, wenn  $c = E[X]$  ist.

**Beweis:** Mit  $m = E[X]$  können wir schreiben:

$$\begin{aligned} E[(X - c)^2] &= E[(X - m + m - c)^2] \\ &= E[(X - m)^2 + 2(X - m)(m - c) + (m - c)^2] \\ &= E[(X - m)^2] + 2(m - c)(E[X] - m) + (m - c)^2 \\ &= \text{Var}[X] + (m - c)^2 . \end{aligned}$$

□

Die bedingte Varianz von  $X$  gegeben  $\underline{Y} = \underline{y}$  ist definiert als

$$\begin{aligned} \text{Var}[X|\underline{Y} = \underline{y}] &\triangleq E[(X - m')^2|\underline{Y} = \underline{y}] \\ \text{wobei} \quad m' &= E[X|\underline{Y} = \underline{y}] . \end{aligned}$$

**Korollar:** Sei  $X$  eine beliebige Zufallsgrösse,  $\underline{Y}$  ein beliebiger Zufallsvektor und  $c$  eine reelle Zahl, dann gilt:

$$E[(X - c)^2|\underline{Y} = \underline{y}] \geq \text{Var}[X|\underline{Y} = \underline{y}] .$$

Das Gleichheitszeichen gilt dann und nur dann, wenn  $c = E[X|\underline{Y} = \underline{y}]$  ist.

Jetzt sind wir bereit, das erste Problem anzupacken. Weil sowohl  $p_{\underline{Y}|X}(\cdot|x)$  als auch  $p_X$  bekannt sind, stellen wir zuerst fest, dass  $p_{X|\underline{Y}}(\cdot|\underline{y})$  aus der "integralen Form der Bayesschen Formel"

$$p_{X|\underline{Y}}(x|\underline{y}) = \frac{p_X(x)p_{\underline{Y}|X}(\underline{y}|x)}{\int_{-\infty}^{+\infty} p_X(\alpha)p_{\underline{Y}|X}(\underline{y}|\alpha)d\alpha} \quad (3.1)$$

berechnet werden kann.

[Diese Formel folgt aus der Tatsache, dass  $p_{X\underline{Y}}(x, \underline{y}) = p_X(x)p_{\underline{Y}|X}(\underline{y}|x)$  und  $p_{\underline{Y}}(\underline{y}) = \int_{-\infty}^{+\infty} p_{X\underline{Y}}(x, \underline{y})dx$  gilt.]

Wir schliessen daraus für das erste Problem, dass sowohl  $p_{X|\underline{Y}}(\cdot|\underline{y})$  als auch  $E[X|\underline{Y} = \underline{y}] = \int_{-\infty}^{+\infty} xp_{X|\underline{Y}}(x|\underline{y})dx$  vom Beobachter berechnet werden können.

Zunächst nützen wir die integrale Form des Satzes vom totalen Erwartungswert aus und schreiben

$$\begin{aligned}
E[(X - \hat{X})^2] &= E[(X - r(Y))^2] \\
&= \int_{-\infty}^{+\infty} E[(X - r(\underline{Y}))^2 | \underline{Y} = \underline{y}] p_{\underline{Y}}(\underline{y}) d\underline{y}
\end{aligned}$$

wobei  $d\underline{y} \triangleq dy_1 dy_2 \dots dy_n$  und  $\int_{-\infty}^{+\infty} \triangleq \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} \dots \int_{-\infty}^{+\infty}$  ( $n$  mal).

Also haben wir

$$E[(X - \hat{X})^2] = \int_{-\infty}^{+\infty} E[(X - r(\underline{y}))^2 | \underline{Y} = \underline{y}] p_{\underline{Y}}(\underline{y}) d\underline{y} .$$

Aber  $p_{\underline{Y}}(\underline{y})$  ist immer  $\geq 0$ . Daraus können wir schliessen, dass wir um  $E[(X - \hat{X})^2]$  zu minimieren für jedes  $\underline{y}$  den Wert  $r(\underline{y})$  so wählen, dass  $E[(X - r(\underline{y}))^2 | \underline{Y} = \underline{y}]$  minimiert wird. Nun folgt aus dem Korollar, dass wir  $r(\underline{y})$  als  $E[X | \underline{Y} = \underline{y}]$  wählen müssen.

Wir fassen zusammen:

Die Schätzungsregel, die  $E[(X - \hat{X})^2]$  minimiert, ist:

$$r(\underline{y}) = E[X | \underline{Y} = \underline{y}] .$$

Der resultierende MMSE ist:

$$\text{MMSE} = \int_{-\infty}^{+\infty} \text{Var}[X | \underline{Y} = \underline{y}] p_{\underline{Y}}(\underline{y}) d\underline{y} .$$

Nun ist es leicht das verallgemeinerte Bayessche Problem zu lösen.

### 3.3 Allgemeines Bayessches Schätzproblem

**Zweites Problem:**

- (1)  $p_{\underline{Y}|X}(\cdot|x)$  für alle  $x \in X(\Omega)$  als auch  $p_X$  sind bekannt.
- (2) Wir wollen die Schätzungsregel  $r(\cdot)$  finden, die  $E[S(X, \hat{X})]$  minimiert, wobei  $S : R^2 \rightarrow R$  und  $S(\alpha, \beta)$  den Spesen für  $X = \alpha$  und  $\hat{X} = \beta$  entspricht.

Von der integralen Form des Satzes vom totalen Erwartungswert erhalten wir wiederum

$$\begin{aligned}
E[S(X, \hat{X})] &= E[S(X, r(\underline{Y}))] \\
&= \int_{-\infty}^{+\infty} E[S(X, r(\underline{Y})) | \underline{Y} = \underline{y}] p_{\underline{Y}}(\underline{y}) d\underline{y} \\
&= \int_{-\infty}^{+\infty} E[S(X, r(\underline{y})) | \underline{Y} = \underline{y}] p_{\underline{Y}}(\underline{y}) d\underline{y} .
\end{aligned}$$

Weil immer  $p_Y(\underline{y}) \geq 0$  ist, können wir folgendes daraus schliessen:

Die Schätzungsregel die  $E[S(X, \hat{X})]$  minimiert lautet:

Wähle für jedes  $\underline{y}$ ,  $r(\underline{y})$  als das  $\beta$ , das  $E[S(X, \beta)|\underline{Y} = \underline{y}]$  minimiert.

**Bemerkung:** Für  $S(\alpha, \beta) = (\alpha - \beta)^2$  geht das zweite Problem ins erste über.

**Drittes Problem:** Entsprechend dem zweiten Problem für den Sonderfall wo

$$S(\alpha, \beta) = \begin{cases} 0, & |\alpha - \beta| < \Delta \\ 1, & |\alpha - \beta| \geq \Delta \end{cases}$$

und  $\Delta$  eine kleine positive Zahl ist, d.h. kleine Fehler sind kostenlos aber alle andern Fehler kosten gleich viel.

Für dieses Bayessche Problem erhalten wir:

$$\begin{aligned} E[S(X, \beta)|\underline{Y} = \underline{y}] &= \int_{-\infty}^{+\infty} S(x, \beta) p_{X|\underline{Y}}(x|\underline{y}) dx \\ &= 1 + \int_{-\infty}^{+\infty} [S(x, \beta) - 1] p_{X|\underline{Y}}(x|\underline{y}) dx \\ &= 1 - \int_{\beta-\Delta}^{\beta+\Delta} p_{X|\underline{Y}}(x|\underline{y}) dx \\ &\approx 1 - 2\Delta p_{X|\underline{Y}}(\beta|\underline{y}). \end{aligned}$$

Also müssen wir  $\beta$  so wählen, dass  $p_{X|\underline{Y}}(\beta|\underline{y})$  maximiert wird. Aber weil

$$p_{X|\underline{Y}}(\beta|\underline{y}) = \frac{p_X(\beta) p_{\underline{Y}|X}(\underline{y}|\beta)}{p_Y(\underline{y})}$$

und weil der Nenner positiv und unabhängig von  $\beta$  ist, sehen wir, dass es äquivalent ist,  $p_X(\beta) p_{\underline{Y}|X}(\underline{y}|\beta)$  zu maximieren.

Die Schätzungsregel die  $E[S(X, \hat{X})]$  minimiert wenn kleine Fehler kostenlos sind und alle andern Fehler kostengleich, lautet:

Wähle für jedes  $\underline{y}$ ,  $r(\underline{y})$  als das  $\beta$  das  $p_X(\beta) p_{\underline{Y}|X}(\underline{y}|\beta)$  maximiert.

Es ist bei dieser Schätzungsregel zu bemerken, dass  $\hat{X}$  immer zu  $X(\Omega)$  gehört, weil sonst  $p_X(\beta) = 0$  ist.

Als Spezialfall des dritten Problems ergibt sich weiter folgende Schätzungsregel:

Die Schätzungsregel die  $E[S(X, \hat{X})]$  minimiert, wenn kleine Fehler nichts kosten und alle andern Fehler gleichviel kosten und wenn weiter der Parameter  $X$  gleichverteilt über  $X(\Omega)$  ist [ d.h.  $p_X(x) = c$  für alle  $x \in X(\Omega)$  ], lautet:

Wähle für jedes  $\underline{y}$ ,  $r(\underline{y})$  als das  $\beta$  in  $X(\Omega)$  das  $p_{Y|X}(\underline{y}|\beta)$  maximiert.

Diese Schätzungsregel wird als maximum-likelihood- (ML-) Schätzungsregel bezeichnet. Es ist zu bemerken, dass der Beobachter die ML-Schätzungsregel benutzen kann ohne die a priori Wahrscheinlichkeit  $p_X$  zu kennen. Wie in der Entscheidungstheorie ist die ML-Regel auch von grosser Bedeutung in der Schätzungstheorie.

## 4 LINEARE ANNÄHERUNG IN EINEM SKALARPRODUKT- RAUM

### 4.1 Problemstellung

Zuerst wollen wir einige Begriffe aus der linearen Algebra zusammenfassen.

Ein **reeller Vektorraum** ist gegeben durch eine nicht-leere Menge  $V$  (die Elemente in  $V$  heißen **Vektoren**) und zwei Regeln. Eine Regel gilt für die Addition von zwei Vektoren, die zweite Regel gilt für die Multiplikation eines Vektors mit einer reellen Zahl. Diese beiden Regeln müssen folgende Axiome erfüllen.

$$(V0) \quad \underline{u}, \underline{v} \in V \Rightarrow \underline{u} + \underline{v} \in V \quad (\text{Geschlossenheit});$$

$$(V1) \quad \underline{u}, \underline{v}, \underline{w} \in V \Rightarrow \underline{u} + (\underline{v} + \underline{w}) = (\underline{u} + \underline{v}) + \underline{w} \quad (\text{Assoziativität});$$

$$(V2) \quad \text{Es gibt ein eindeutig bestimmtes Element } \underline{0} \text{ in } V \text{ mit } \underline{v} + \underline{0} = \underline{v} \text{ für alle } \underline{v} \in V \text{ (Nullvektor);}$$

$$(V3) \quad \text{Zu jedem } \underline{v} \in V \text{ gibt es ein eindeutig bestimmtes Element } -\underline{v} \text{ in } V \text{ mit } \underline{v} + (-\underline{v}) = \underline{0} \\ (\text{entgegengesetzter Vektor});$$

$$(V4) \quad \underline{u}, \underline{v} \in V \Rightarrow \underline{u} + \underline{v} = \underline{v} + \underline{u} \quad (\text{Kommutativität});$$

$$(R0) \quad a \in R, \underline{v} \in V \Rightarrow a\underline{v} \in V \quad (\text{Geschlossenheit});$$

$$(R1) \quad a \in R, \underline{u}, \underline{v} \in V \Rightarrow a(\underline{u} + \underline{v}) = a\underline{u} + a\underline{v} \quad (\text{erstes Distributivgesetz});$$

$$(R2) \quad a, b \in R, \underline{v} \in V \Rightarrow (a + b)\underline{v} = a\underline{v} + b\underline{v} \quad (\text{zweites Distributivgesetz});$$

$$(R3) \quad a, b \in R, \underline{v} \in V \Rightarrow a(b\underline{v}) = (ab)\underline{v} \quad (\text{Assoziativität});$$

$$(R4) \quad \underline{v} \in V \Rightarrow 1 \cdot \underline{v} = \underline{v}. \quad (\text{Skalierung mit Eins}).$$

**Beispiele:**

$$(1) \quad V = \{(r_1, r_2, \dots, r_n) : r_i \in R, i = 1, 2, \dots, n\} \\ (r_1, r_2, \dots, r_n) + (s_1, s_2, \dots, s_n) \triangleq (r_1 + s_1, r_2 + s_2, \dots, r_n + s_n) \\ a(r_1, r_2, \dots, r_n) \triangleq (ar_1, ar_2, \dots, ar_n).$$

Dieser Vektorraum wird normalerweise mit  $R^n$  bezeichnet.

$$(2) \quad \Omega = \text{eine beliebige nicht-leere Menge.}$$

$V = \{f : f : \Omega \rightarrow R\}$ , d.h.,  $V$  ist die Menge aller reellwertigen Funktionen mit Definitionsbereich  $\Omega$ .

$f_1 + f_2$  ist als die Funktion  $g$  mit

$$g(\omega) = f_1(\omega) + f_2(\omega) \text{ für alle } \omega \in \Omega \text{ definiert.}$$

$af$  ist als die Funktion  $g$  mit

$$g(\omega) = af(\omega) \text{ für alle } \omega \in \Omega \text{ definiert.}$$

$\underline{0}$  ist die Funktion mit

$$\underline{0}(\omega) = 0 \text{ für alle } \omega \in \Omega.$$

(Der Einfachheit halber schreiben wir üblicherweise die Funktion  $\underline{0}$  nur als  $0$ .)

○

Ein **Unterraum**  $U$  des reellen Vektorraums  $V$  ist eine Untermenge  $U$ , die auch ein reeller Vektorraum ist.  $U = V$  ist auch ein Unterraum von  $V$ .  $U = \{\underline{0}\}$  heisst der **triviale Unterraum** von  $V$ .

**Test für einen Unterraum:** Sei  $U$  eine nicht leere Untermenge des reellen Vektorraums  $V$ , so ist  $U$  dann und nur dann ein Unterraum von  $V$ , wenn

(1)  $\underline{u}, \underline{v} \in U \Rightarrow \underline{u} + \underline{v} \in U$  (Geschlossenheit)

und

(2)  $a \in R, \underline{u} \in U \Rightarrow a\underline{u} \in U$  (Geschlossenheit).

In diesem Kapitel werden wir künftig mit  $V$  immer einen reellen Vektorraum bezeichnen. Sei  $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_n$  in  $V$  und  $a_1, a_2, \dots, a_n$  in  $R$ , dann heisst

$$a_1\underline{v}_1 + a_2\underline{v}_2 + \dots + a_n\underline{v}_n, \quad n < \infty$$

eine **Linearkombination** der Vektoren  $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_n$ . Die Menge aller Linearkombinationen der Vektoren  $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_n$  ist ein Unterraum von  $V$ . Dieser Unterraum wird als  $S(\underline{v}_1, \underline{v}_2, \dots, \underline{v}_n)$  bezeichnet. Man sagt, er werde von  $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_n$  **aufgespannt**.

$$S(\underline{v}_1, \underline{v}_2, \dots, \underline{v}_n) = \{a_1\underline{v}_1 + a_2\underline{v}_2 + \dots + a_n\underline{v}_n : a_i \in R, i = 1, 2, \dots, n\}.$$

Die Vektoren  $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_n$  heissen **linear unabhängig**, falls

$$a_1\underline{v}_1 + a_2\underline{v}_2 + \dots + a_n\underline{v}_n = \underline{0} \Rightarrow a_1 = a_2 = \dots = a_n = 0,$$

sonst heissen sie **linear abhängig**. Da  $a \cdot \underline{0} = \underline{0}$  auch für  $a \neq 0$ , ist  $\underline{0}$  immer linear abhängig.

Die Vektoren  $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_n$  heissen eine **Basis** für  $V$ , falls

(1)  $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_n$  linear unabhängig sind

und

(2)  $S(\underline{v}_1, \underline{v}_2, \dots, \underline{v}_n) = V$

gilt. Hat  $V$  eine Basis, dann ist die **Dimension** von  $V$  (bezeichnet mit  $\dim(V)$ ) als die Anzahl Vektoren in einer Basis definiert. Man sagt  $\dim(\{\underline{0}\}) = 0$ . Hat  $V$  keine Basis und sei  $V \neq \{\underline{0}\}$ , dann definiert man  $\dim(V) = \infty$ .

**Beispiele:** (Fortsetzung)

(1)  $\underline{v}_1 = (1, 0, \dots, 0)$ ,  $\underline{v}_2 = (0, 1, \dots, 0)$ , ...,  $\underline{v}_n = (0, 0, \dots, 1)$  sind eine Basis für  $V = \mathbb{R}^n$ .  
Deshalb ist  $\dim(\mathbb{R}^n) = n$ .

(2) Sei  $\Omega$  eine endliche Menge  $\Omega = \{\omega_1, \omega_2, \dots, \omega_n\}$ , dann sind  $f_1, f_2, \dots, f_n$  mit

$$f_i(\omega_j) = \begin{cases} 1, & j = i \\ 0, & j \neq i \end{cases}$$

eine Basis für  $V$ . Deshalb ist  $\dim(V) = n$ . Sei  $\Omega$  eine unendliche Menge, dann ist  $\dim(V) = \infty$ .

○

## 4.2 Reeller Skalarproduktraum

Ein **Skalarprodukt** für einen reellen Vektorraum  $V$  ist eine reellwertige Funktion mit dem Definitionsbereich  $V \times V$  (d.h. die Argumente sind ein Paar von Vektoren), die die folgenden Axiome erfüllt (wobei wir für den Wert der Funktion mit den Argumenten  $\underline{u}$  und  $\underline{v}$  die Abkürzung  $\langle \underline{u}, \underline{v} \rangle$  schreiben):

(S1)  $\langle \underline{u}, \underline{v} \rangle = \langle \underline{v}, \underline{u} \rangle$  für alle  $\underline{u}, \underline{v} \in V$  (Symmetrie);

(S2)  $\langle a\underline{u} + b\underline{v}, \underline{w} \rangle = a\langle \underline{u}, \underline{w} \rangle + b\langle \underline{v}, \underline{w} \rangle$  für alle  $a, b \in \mathbb{R}$ ,  $\underline{u}, \underline{v}, \underline{w} \in V$  (Bilinearität);

(S3)  $\langle \underline{v}, \underline{v} \rangle \geq 0$  mit Gleichheit dann und nur dann, wenn  $\underline{v} = \underline{0}$  ist (Positive Definitheit);

**Beispiele:**

(1)  $V = \mathbb{R}^n$ .

$p_1, p_2, \dots, p_n$  seien positive reelle Zahlen. Dann ist

$$\langle (r_1, r_2, \dots, r_n), (s_1, s_2, \dots, s_n) \rangle \triangleq p_1 r_1 s_1 + p_2 r_2 s_2 + \dots + p_n r_n s_n$$

ein Skalarprodukt für  $\mathbb{R}^n$ .

Der Fall  $p_1 = p_2 = \dots = p_n = 1$  ergibt das sogenannte "dot - Produkt"; in diesem Falle schreibt man  $\underline{u} \cdot \underline{v}$ .

(2)  $V =$  Menge der stetigen reellwertigen Funktionen, deren Definitionsbereich das geschlossene Intervall  $[a, b]$  mit  $a < b$  ist.

$p(x) =$  eine stetige Funktion mit  $p(x) > 0$  für  $a \leq x \leq b$ .

$$\langle f, g \rangle \triangleq \int_a^b p(x) \cdot f(x) \cdot g(x) dx$$

ist ein Skalarprodukt für  $V$ .

Ein reeller Vektorraum  $V$  zusammen mit einem Skalarprodukt für  $V$  heisst reeller **Skalarproduktraum**.

**Cauchy - Schwarz Ungleichung :** In einem reellen Skalarproduktraum ist

$$(\langle \underline{u}, \underline{v} \rangle)^2 \leq \langle \underline{u}, \underline{u} \rangle \langle \underline{v}, \underline{v} \rangle.$$

Das Gleichheitszeichen gilt dann und nur dann, wenn  $\underline{u}$  und  $\underline{v}$  linear abhängig sind.

**Beweis:**

Die Behauptung gilt trivialerweise, wenn entweder  $\underline{u} = \underline{0}$  oder  $\underline{v} = \underline{0}$ . Also können wir annehmen, dass  $\underline{u} \neq \underline{0}$  und  $\underline{v} \neq \underline{0}$  gelte. Wir definieren dann eine reellwertige Funktion  $f(x)$  durch

$$f(x) = \langle \underline{u} - x\underline{v}, \underline{u} - x\underline{v} \rangle \quad x \in \mathbb{R}.$$

Wegen (S3) gilt für alle  $x$ , dass  $f(x) \geq 0$ .

Es folgt aus der Bilinearität des Skalarprodukts, dass

$$\begin{aligned} f(x) &= \langle \underline{u}, \underline{u} - x\underline{v} \rangle - x \langle \underline{v}, \underline{u} - x\underline{v} \rangle \\ &= \langle \underline{u}, \underline{u} \rangle - x \langle \underline{u}, \underline{v} \rangle - x \langle \underline{v}, \underline{u} \rangle + x^2 \langle \underline{v}, \underline{v} \rangle \\ &= \langle \underline{u}, \underline{u} \rangle - 2x \langle \underline{u}, \underline{v} \rangle + x^2 \langle \underline{v}, \underline{v} \rangle. \end{aligned}$$

Durch Ableiten erhalten wir:

$$\begin{aligned} f'(x) &= -2 \langle \underline{u}, \underline{v} \rangle + 2x \langle \underline{v}, \underline{v} \rangle \\ f''(x) &= 2 \langle \underline{v}, \underline{v} \rangle > 0 \quad (\text{weil } \underline{v} \neq \underline{0}). \end{aligned}$$

Es folgt, dass die Lösung von  $f'(x) = 0$  das eindeutige Minimum von  $f(x)$  ergibt, nämlich

$$x = \frac{\langle \underline{u}, \underline{v} \rangle}{\langle \underline{v}, \underline{v} \rangle}.$$

Der minimale Wert von  $f(x)$  ist somit

$$\begin{aligned} f\left(\frac{\langle \underline{u}, \underline{v} \rangle}{\langle \underline{v}, \underline{v} \rangle}\right) &= \langle \underline{u}, \underline{u} \rangle - 2 \frac{(\langle \underline{u}, \underline{v} \rangle)^2}{\langle \underline{v}, \underline{v} \rangle} + \frac{(\langle \underline{u}, \underline{v} \rangle)^2}{\langle \underline{v}, \underline{v} \rangle} \\ &= \frac{\langle \underline{u}, \underline{u} \rangle \langle \underline{v}, \underline{v} \rangle - (\langle \underline{u}, \underline{v} \rangle)^2}{\langle \underline{v}, \underline{v} \rangle} \\ &\geq 0 \end{aligned}$$

wobei die Ungleichheit aus der Tatsache folgt, dass  $f(x) \geq 0$  für alle  $x$  gilt. Weil  $\langle \underline{v}, \underline{v} \rangle > 0$ , folgt nun, dass



$$\langle \underline{u}, \underline{u} \rangle \langle \underline{v}, \underline{v} \rangle - (\langle \underline{u}, \underline{v} \rangle)^2 \geq 0.$$

Das Gleichheitszeichen gilt dann und nur dann, wenn es ein  $x$  gibt, für das  $f(x) = 0$  ist; somit muss aber auch  $\underline{u} - x\underline{v} = \underline{0}$  für ein  $x$  gelten. Das bedeutet aber, dass  $\underline{u}$  und  $\underline{v}$  linear abhängig sind.  $\square$

**Beispiele:**

(1)  $V = \mathbb{R}^n$

$$\langle \underline{u}, \underline{v} \rangle = \underline{u} \cdot \underline{v}$$

[ d.h., wenn  $\underline{u} = (a_1, \dots, a_n)$  und  $\underline{v} = (b_1, \dots, b_n)$ , dann gilt  $\langle \underline{u}, \underline{v} \rangle = \sum_{i=1}^n a_i b_i$  .]

Die Cauchy-Schwarz-Ungleichung gibt

$$\left( \sum_{i=1}^n a_i b_i \right)^2 \leq \sum_{i=1}^n a_i^2 \sum_{j=1}^n b_j^2$$

wobei das Gleichheitszeichen dann und nur dann gilt, wenn  $(a_1, \dots, a_n)$  und  $(b_1, \dots, b_n)$  proportional sind. Diese Form der Ungleichung ist als die **Cauchy-Ungleichung** bekannt.

(2)  $V =$  Menge der stetigen reellwertigen Funktionen mit dem Definitionsbereich  $[a, b]$ , wobei  $a < b$  ist.

$$\langle f, g \rangle = \int_a^b f(x) \cdot g(x) dx.$$

Die Cauchy - Schwarz Ungleichung liefert

$$\left( \int_a^b f(x) \cdot g(x) dx \right)^2 \leq \int_a^b f^2(x) dx \int_a^b g^2(y) dy,$$

wobei das Gleichheitszeichen dann und nur dann gilt, wenn  $f$  und  $g$  proportional sind. Diese Form der Ungleichung ist als die **Schwarz-Ungleichung** bekannt.  $\circ$

### 4.3 Orthogonalität

Die Vektoren  $\underline{u}, \underline{v}$  in einem Skalarproduktraum  $V$  heißen **orthogonal**, falls  $\langle \underline{u}, \underline{v} \rangle = 0$ . Sei  $U$  ein Unterraum von  $V$ , dann schreibt man  $U^\perp$  für die Untermenge von  $V$ , die alle Vektoren  $\underline{v}$  enthält, für welche sich  $\langle \underline{u}, \underline{v} \rangle = 0$  ergibt für alle  $\underline{u} \in U$ . Seien  $\underline{v}_1$  und  $\underline{v}_2$  in  $U^\perp$ , dann gilt  $\langle \underline{u}, \underline{v}_1 + \underline{v}_2 \rangle = \langle \underline{u}, \underline{v}_1 \rangle + \langle \underline{u}, \underline{v}_2 \rangle = 0$  für alle  $\underline{u} \in U$ ; somit ist  $\underline{v}_1 + \underline{v}_2$  in  $U^\perp$ . Sei  $\underline{v}$  in  $U^\perp$  und  $c$  in  $\mathbb{R}$ , dann ist  $\langle \underline{u}, c\underline{v} \rangle = c \langle \underline{u}, \underline{v} \rangle = 0$  für alle  $\underline{u} \in U$ . Es folgt aus dem Test für einen Unterraum, dass  $U^\perp$  ein Unterraum von  $V$  ist.  $U^\perp$  heisst das **Orthogonalkomplement** von  $U$  in  $V$ .

**Lemma 1:**

$$U \cap U^\perp = \{0\}.$$

**Beweis:**

Sei  $\underline{u}$  in  $U$  und auch in  $U^\perp$ , dann muss auch  $\langle \underline{u}, \underline{u} \rangle = 0$  gelten. Aus Axiom (S3) folgt nun, dass  $\underline{u} = \underline{0}$  ist.  $\square$

**Lemma 2:** Kann  $\underline{v} \in V$  als  $\underline{v} = \underline{u}_1 + \underline{u}_2$  mit  $\underline{u}_1 \in U$  und  $\underline{u}_2 \in U^\perp$  geschrieben werden, dann sind  $\underline{u}_1$  und  $\underline{u}_2$  eindeutig bestimmt.

**Beweis:**

Sei  $\underline{v} = \underline{u}_1 + \underline{u}_2 = \underline{u}_3 + \underline{u}_4$  mit  $\underline{u}_1, \underline{u}_3 \in U$  und  $\underline{u}_2, \underline{u}_4 \in U^\perp$ , dann gilt:

$$\underline{u}_1 - \underline{u}_3 = \underline{u}_4 - \underline{u}_2$$

Aber da  $U$  und  $U^\perp$  Vektorräume sind, gilt:  $\underline{u}_1 - \underline{u}_3 \in U$  und  $\underline{u}_4 - \underline{u}_2 \in U^\perp$ . Aus Lemma 1 folgt nun, dass  $\underline{u}_1 - \underline{u}_3 = \underline{u}_4 - \underline{u}_2 = \underline{0}$  ist.  $\square$

**Bemerkung:**

Es ist oft der Fall, dass jedes  $\underline{v}$  in  $V$  als  $\underline{v} = \underline{u}_1 + \underline{u}_2$  mit  $\underline{u}_1 \in U$  und  $\underline{u}_2 \in U^\perp$  geschrieben werden kann. Die pathologische Ausnahme tritt nur dann ein, wenn  $\dim(U) = \infty$  gilt.

Die Vektoren  $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_n$  heissen orthogonal, wenn jedes Paar  $\underline{v}_i, \underline{v}_j$  mit  $i \neq j$  orthogonal ist.

**Lemma 3:** Seien  $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_m$  orthogonale Vektoren in einem reellen Skalarproduktraum  $V$  mit  $\underline{v}_i \neq \underline{0}$  für  $i = 1, 2, \dots, m$  und sei  $\underline{v}$  ein beliebiger Vektor in  $V$ . Dann ist

$$\underline{u} = \underline{v} - \sum_{i=1}^m \frac{\langle \underline{v}, \underline{v}_i \rangle}{\langle \underline{v}_i, \underline{v}_i \rangle} \underline{v}_i$$

orthogonal zu  $\underline{v}_i$  für  $i = 1, 2, \dots, m$ .

**Beweis:**

$$\begin{aligned} \langle \underline{u}, \underline{v}_j \rangle &= \langle \underline{v}, \underline{v}_j \rangle - \sum_{i=1}^m \frac{\langle \underline{v}, \underline{v}_i \rangle}{\langle \underline{v}_i, \underline{v}_i \rangle} \langle \underline{v}_i, \underline{v}_j \rangle \\ &= \langle \underline{v}, \underline{v}_j \rangle - \frac{\langle \underline{v}, \underline{v}_j \rangle}{\langle \underline{v}_j, \underline{v}_j \rangle} \langle \underline{v}_j, \underline{v}_j \rangle \\ &= \langle \underline{v}, \underline{v}_j \rangle - \langle \underline{v}, \underline{v}_j \rangle = 0 \\ &\text{für } j = 1, 2, \dots, m. \end{aligned}$$

$\square$

### Orthogonalbasissatz:

Jeder reelle Skalarproduktraum  $V$  mit  $1 \leq \dim(V) < \infty$  besitzt eine Basis, die orthogonal ist.

**Beweis:**

Seien  $\underline{u}_1, \underline{u}_2, \dots, \underline{u}_n$  eine beliebige Basis für  $V$ , dann setzen wir  $\underline{v}_1 = \underline{u}_1$  und

$$\underline{v}_{m+1} = \underline{u}_{m+1} - \sum_{i=1}^m \frac{\langle \underline{u}_{m+1}, \underline{v}_i \rangle}{\langle \underline{v}_i, \underline{v}_i \rangle} \underline{v}_i \quad (4.1)$$

für  $m = 1, 2, \dots, n-1$ . Da  $\underline{u}_1, \underline{u}_2, \dots, \underline{u}_n$  linear unabhängig sind, ist es nicht möglich, dass  $\underline{v}_i = \underline{0}$  für irgendein  $i$ . Aus Lemma 3 folgt direkt, dass  $\underline{v}_2$  orthogonal zu  $\underline{v}_1$  ist, dass  $\underline{v}_3$  orthogonal zu  $\underline{v}_1$  und  $\underline{v}_2$  ist, usw. Also sind die Vektoren  $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_n$  orthogonal. [ Die Methode (4.1), eine Reihe von Vektoren zu orthogonalisieren, heisst **Gram-Schmidt-Orthogonalisierungsverfahren**. ] Weil wir  $\underline{u}_m$  als eine Linearkombination von  $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_n$  für  $m = 1, 2, \dots, n$  schreiben können, müssen auch  $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_n$  eine Basis für  $V$  sein.  $\square$

**Lemma 4:** Seien  $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_m$  Vektoren in einem reellen Skalarproduktraum  $V$  und sei  $U = \mathcal{S}(\underline{v}_1, \underline{v}_2, \dots, \underline{v}_m)$ , dann ist  $\underline{v} \in V$  ein Vektor in  $U^\perp$  dann und nur dann, wenn  $\langle \underline{v}, \underline{v}_i \rangle = 0$  für  $i = 1, 2, \dots, m$ .

**Beweis:**

Sei  $\underline{v} \in U^\perp$ , dann gilt trivialerweise  $\langle \underline{v}, \underline{v}_i \rangle = 0$  für  $i = 1, 2, \dots, m$ . Sei  $\langle \underline{v}, \underline{v}_i \rangle = 0$  für  $i = 1, 2, \dots, m$  und sei  $\underline{u}$  ein beliebiger Vektor in  $U = \mathcal{S}(\underline{v}_1, \underline{v}_2, \dots, \underline{v}_m)$ , sodass  $\underline{u} = c_1 \underline{v}_1 + c_2 \underline{v}_2 + \dots + c_m \underline{v}_m$  geschrieben werden kann, dann gilt:

$$\begin{aligned} \langle \underline{v}, \underline{u} \rangle &= \langle \underline{u}, \underline{v} \rangle \\ &= \langle c_1 \underline{v}_1 + c_2 \underline{v}_2 + \dots + c_m \underline{v}_m, \underline{v} \rangle \\ &= c_1 \langle \underline{v}, \underline{v}_1 \rangle + c_2 \langle \underline{v}, \underline{v}_2 \rangle + \dots + c_m \langle \underline{v}, \underline{v}_m \rangle \\ &= 0. \end{aligned}$$

Deswegen muss  $\underline{v}$  zu  $U^\perp$  gehören.  $\square$

**Zerlegungssatz für einen Skalarproduktraum :**

Seien  $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_m$  Vektoren in einem reellen Skalarproduktraum  $V$  und sei  $U = \mathcal{S}(\underline{v}_1, \underline{v}_2, \dots, \underline{v}_m)$ . Dann kann jeder Vektor

$$\underline{v} \in V \text{ als } \underline{v} = \underline{u}_1 + \underline{u}_2 \text{ mit } \underline{u}_1 \in U, \underline{u}_2 \in U^\perp$$

geschrieben werden. Die Vektoren  $\underline{u}_1$  und  $\underline{u}_2$  sind eindeutig bestimmt.

**Beweis:**

Wähle eine orthogonale Basis  $\underline{w}_1, \underline{w}_2, \dots, \underline{w}_k$  für  $U$ , wobei  $k = \dim(U) \leq m$ . Sei  $k = 0$ , dann ist  $U = \{0\}$  und  $U^\perp = V$  und der Satz ist trivialerweise erfüllt. Sei  $k > 0$ , dann folgt aus Lemma 3, dass  $\underline{v} \in V$  als die Summe von  $\underline{u}$  und einem Vektor in  $\mathcal{S}(\underline{w}_1, \underline{w}_2, \dots, \underline{w}_k) = U$  geschrieben werden kann, wobei  $\underline{u}$  orthogonal zu  $\underline{w}_i$  für  $i = 1, 2, \dots, k$  ist. Es folgt dann von Lemma 4, dass  $\underline{u} \in U^\perp$  gilt. Dass diese beiden Vektoren ( einer von  $U$  und einer von  $U^\perp$ ) eindeutig bestimmt sind, folgt direkt aus Lemma 2.  $\square$

#### 4.4 Norm für einen reellen Vektorraum

Eine Norm für einen reellen Vektorraum  $V$  ist eine reellwertige Funktion mit Definitionsbereich  $V$ , die folgende Axiome erfüllt (wir schreiben  $\|\underline{v}\|$  für den Wert der Funktion mit dem Argument  $\underline{v}$ ):

(N1)  $\|\underline{v}\| \geq 0$ . Das Gleichheitszeichen gilt dann und nur dann, wenn  $\underline{v} = \underline{0}$  ist (positive Definitheit);

(N2)  $c \in \mathbb{R} \Rightarrow \|c\underline{v}\| = |c|\|\underline{v}\|$  (Skalierung);

(N3)  $\|\underline{v}_1 + \underline{v}_2\| \leq \|\underline{v}_1\| + \|\underline{v}_2\|$  (Dreiecksungleichung).

Eine Norm ist ein Mass für den Betrag eines Vektors. Im Allgemeinen gibt es eine grosse Auswahl von Normen für einen reellen Vektorraum; für einen reellen Skalarproduktraum gibt es eine "natürliche" Norm.

In einem reellen Skalarproduktraum  $V$ , ist  $\|\underline{v}\| = \sqrt{\langle \underline{v}, \underline{v} \rangle}$  eine Norm.

**Bemerkung:**  $\sqrt{\quad}$  bezeichnet die positive Quadratwurzel.

**Beweis:**

Axiome (N1) und (S3) sind äquivalent. Aus (S2) folgt:

$$\langle c\underline{v}, c\underline{v} \rangle = c \langle \underline{v}, c\underline{v} \rangle = c^2 \langle \underline{v}, \underline{v} \rangle$$

Darum gilt (N2). Unter Benützung von (S1) und (S2) merken wir, dass

$$\begin{aligned} \langle \underline{v}_1 + \underline{v}_2, \underline{v}_1 + \underline{v}_2 \rangle &= \langle \underline{v}_1, \underline{v}_1 + \underline{v}_2 \rangle + \langle \underline{v}_2, \underline{v}_1 + \underline{v}_2 \rangle \\ &= \langle \underline{v}_1, \underline{v}_1 \rangle + \langle \underline{v}_1, \underline{v}_2 \rangle + \langle \underline{v}_2, \underline{v}_1 \rangle + \langle \underline{v}_2, \underline{v}_2 \rangle \\ &= \langle \underline{v}_1, \underline{v}_1 \rangle + 2 \cdot \langle \underline{v}_1, \underline{v}_2 \rangle + \langle \underline{v}_2, \underline{v}_2 \rangle \end{aligned}$$

gilt. Es folgt ferner aus der Cauchy-Schwarz Ungleichung, dass

$$\begin{aligned} \langle \underline{v}_1 + \underline{v}_2, \underline{v}_1 + \underline{v}_2 \rangle &\leq \langle \underline{v}_1, \underline{v}_1 \rangle + 2\sqrt{\langle \underline{v}_1, \underline{v}_1 \rangle \langle \underline{v}_2, \underline{v}_2 \rangle} + \langle \underline{v}_2, \underline{v}_2 \rangle \\ &= (\sqrt{\langle \underline{v}_1, \underline{v}_1 \rangle} + \sqrt{\langle \underline{v}_2, \underline{v}_2 \rangle})^2, \end{aligned}$$

und somit gilt auch (N3). □

Jedesmal wenn wir von einer Norm in einem reellen Skalarproduktraum sprechen, werden wir diese "natürliche" Norm meinen.

## 4.5 Optimale lineare Annäherung

Jetzt sind wir endlich in der Lage, das Problem der linearen Annäherung anzupacken.

**Das Problem der linearen Annäherung in einem reellen Skalarproduktraum  $V$ :** Finde für einen beliebigen Vektor  $\underline{v} \in V$  eine lineare Kombination  $\hat{\underline{v}}$  der Vektoren  $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_m$  in  $V$ , sodass  $\|\underline{v} - \hat{\underline{v}}\|^2$  minimal ist.

**Aequivalente Problemstellung:** Finde einen Vektor  $\hat{\underline{v}}$  in  $U = \mathcal{S}(\underline{v}_1, \underline{v}_2, \dots, \underline{v}_m)$ , der  $\|\underline{v} - \hat{\underline{v}}\|^2$  minimiert.

**Bemerkung:**

Das Problem, wie man die Koeffizienten  $c_1, c_2, \dots, c_m$  so finden kann, dass  $\hat{\underline{v}} = c_1\underline{v}_1 + c_2\underline{v}_2 + \dots + c_m\underline{v}_m$  werden wir später betrachten.

**Satz über die lineare Annäherung in einem reellen Skalarproduktraum:**

Seien  $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_m$  Vektoren in einem reellen Skalarproduktraum  $V$  und sei  $\underline{v}$  in  $V$ , dann ist der Vektor  $\hat{\underline{v}}$  in  $\mathcal{S}(\underline{v}_1, \underline{v}_2, \dots, \underline{v}_m) = U$ , welcher  $\|\underline{v} - \hat{\underline{v}}\|^2$  minimiert, eindeutig bestimmt. Falls  $\underline{v}$  als  $\underline{v} = \underline{u}_1 + \underline{u}_2$  mit  $\underline{u}_1 \in U$  und  $\underline{u}_2 \in U^\perp$  geschrieben wird, dann ist

$$\hat{\underline{v}} = \underline{u}_1$$

diese optimale Annäherung.

**Bemerkung:** Den Vektor  $\hat{\underline{v}}$  bezeichnet man als die **Projektion von  $\underline{v}$  auf  $U = \mathcal{S}(\underline{v}_1, \underline{v}_2, \dots, \underline{v}_m)$ .**

**Beweis:**

Aus dem Zerlegungssatz folgt, dass  $\underline{v} = \underline{u}_1 + \underline{u}_2$  mit  $\underline{u}_1 \in U = \mathcal{S}(\underline{v}_1, \underline{v}_2, \dots, \underline{v}_m)$  und  $\underline{u}_2 \in U^\perp$  gilt, wobei  $\underline{u}_1$  und  $\underline{u}_2$  eindeutig bestimmt sind. Für eine beliebige Annäherung  $\hat{\underline{v}} \in \mathcal{S}(\underline{v}_1, \underline{v}_2, \dots, \underline{v}_m)$  erhält man deshalb:

$$\begin{aligned}\|\underline{v} - \hat{\underline{v}}\|^2 &= \|\underline{v} - \underline{u}_1 + \underline{u}_1 - \hat{\underline{v}}\|^2 \\ &= \langle (\underline{v} - \underline{u}_1) + (\underline{u}_1 - \hat{\underline{v}}), (\underline{v} - \underline{u}_1) + (\underline{u}_1 - \hat{\underline{v}}) \rangle \\ &= \langle \underline{v} - \underline{u}_1, (\underline{v} - \underline{u}_1) + (\underline{u}_1 - \hat{\underline{v}}) \rangle + \langle \underline{u}_1 - \hat{\underline{v}}, (\underline{v} - \underline{u}_1) + (\underline{u}_1 - \hat{\underline{v}}) \rangle \\ &= \langle \underline{v} - \underline{u}_1, \underline{v} - \underline{u}_1 \rangle + 2\langle \underline{v} - \underline{u}_1, \underline{u}_1 - \hat{\underline{v}} \rangle + \langle \underline{u}_1 - \hat{\underline{v}}, \underline{u}_1 - \hat{\underline{v}} \rangle \\ &= \|\underline{v} - \underline{u}_1\|^2 + 2\langle \underline{v} - \underline{u}_1, \underline{u}_1 - \hat{\underline{v}} \rangle + \|\underline{u}_1 - \hat{\underline{v}}\|^2.\end{aligned}$$

Aber, da  $\underline{v} - \underline{u}_1 = \underline{u}_2 \in U^\perp$  und  $\underline{u}_1 - \hat{\underline{v}} \in U$  ist, gilt:  $\langle \underline{v} - \underline{u}_1, \underline{u}_1 - \hat{\underline{v}} \rangle = 0$ .

Wir schliessen, dass

$$\|\underline{v} - \hat{\underline{v}}\|^2 = \|\underline{v} - \underline{u}_1\|^2 + \|\underline{u}_1 - \hat{\underline{v}}\|^2.$$

Es folgt, dass

$$\|\underline{v} - \hat{\underline{v}}\|^2 \geq \|\underline{v} - \underline{u}_1\|^2,$$

wobei das Gleichheitszeichen dann und nur dann gilt, wenn  $\hat{\underline{v}} = \underline{u}_1$ .

□

Wir formulieren diesen Satz in einer für die Anwendung besser geeigneten Form.

**Orthogonalitätsprinzip für die lineare Annäherung in einem reellen Skalarproduktraum:**

Seien  $\underline{v}_1, \dots, \underline{v}_m$  und  $\underline{v}$  Vektoren in einem reellen Skalarproduktraum, dann ist  $\hat{\underline{v}} = c_1 \underline{v}_1 + c_2 \underline{v}_2 + \dots + c_m \underline{v}_m$  genau dann der Vektor in  $S(\underline{v}_1, \underline{v}_2, \dots, \underline{v}_m)$  der  $\|\underline{v} - \hat{\underline{v}}\|^2$  minimiert, wenn der Fehler  $\underline{v} - \hat{\underline{v}}$  orthogonal zu  $\underline{v}_i$  ( $i = 1, 2, \dots, m$ ) ist, also wenn

$$\langle \underline{v} - \hat{\underline{v}}, \underline{v}_i \rangle = 0 \quad (i = 1, 2, \dots, m) \quad (4.2)$$

gilt.

Für diese optimale lineare Annäherung gilt:

$$\|\underline{v} - \hat{\underline{v}}\|^2 = \langle \underline{v} - \hat{\underline{v}}, \underline{v} \rangle = \|\underline{v}\|^2 - \|\hat{\underline{v}}\|^2. \quad (4.3)$$

**Beweis:**

Aus dem Satz der linearen Annäherung folgt, dass  $\hat{\underline{v}} \in S(\underline{v}_1, \dots, \underline{v}_m)$  genau dann die optimale lineare Annäherung von  $\underline{v}$  ist, wenn  $\underline{v} - \hat{\underline{v}} \in U^\perp$  ist. Wegen Lemma 4 ist dies zu Formel (4.2) äquivalent.

Sei  $\hat{\underline{v}} \in S(\underline{v}_1, \dots, \underline{v}_m) = U$ , dann gilt:

$$\begin{aligned} \|\underline{v} - \hat{\underline{v}}\|^2 &= \langle \underline{v} - \hat{\underline{v}}, \underline{v} - \hat{\underline{v}} \rangle \\ &= \langle \underline{v} - \hat{\underline{v}}, \underline{v} \rangle - \langle \underline{v} - \hat{\underline{v}}, \hat{\underline{v}} \rangle \end{aligned}$$

Weil  $\langle \underline{v} - \hat{\underline{v}} \rangle \in U^\perp$  und  $\hat{\underline{v}} \in U$  sind, ist  $\langle \underline{v} - \hat{\underline{v}}, \hat{\underline{v}} \rangle = 0$ .

Weiter gilt dann:

$$\begin{aligned} \|\underline{v} - \hat{\underline{v}}\|^2 = \langle \underline{v} - \hat{\underline{v}}, \underline{v} \rangle &= \langle \underline{v}, \underline{v} \rangle - \langle \hat{\underline{v}}, \underline{v} \rangle = \langle \underline{v}, \underline{v} \rangle - \langle \hat{\underline{v}}, \underline{v} - \hat{\underline{v}} + \hat{\underline{v}} \rangle \\ &= \langle \underline{v}, \underline{v} \rangle - \langle \hat{\underline{v}}, \underline{v} - \hat{\underline{v}} \rangle - \langle \hat{\underline{v}}, \hat{\underline{v}} \rangle = \langle \underline{v}, \underline{v} \rangle - \langle \hat{\underline{v}}, \hat{\underline{v}} \rangle. \end{aligned}$$

□

**Bemerkung:**

Die Gleichung (4.2) heisst **Orthogonalitätsgleichung** für die lineare Annäherung in einem reellen Skalarproduktraum.

### Matrixform der Orthogonalitätsgleichung:

Seien  $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_m$  und  $\underline{v}$  Vektoren in einem reellen Skalarproduktraum, dann ist  $\hat{\underline{v}} = c_1 \underline{v}_1 + \dots + c_m \underline{v}_m$  genau dann der Vektor in  $S(\underline{v}_1, \dots, \underline{v}_m)$  der  $\|\underline{v} - \hat{\underline{v}}\|^2$  minimiert, wenn die Koeffizienten  $c_1, c_2, \dots, c_m$  die Gleichung

$$\begin{bmatrix} \langle \underline{v}_1, \underline{v}_1 \rangle & \langle \underline{v}_1, \underline{v}_2 \rangle & \dots & \langle \underline{v}_1, \underline{v}_m \rangle \\ \langle \underline{v}_2, \underline{v}_1 \rangle & \langle \underline{v}_2, \underline{v}_2 \rangle & \dots & \langle \underline{v}_2, \underline{v}_m \rangle \\ \vdots & & & \\ \langle \underline{v}_m, \underline{v}_1 \rangle & \langle \underline{v}_m, \underline{v}_2 \rangle & \dots & \langle \underline{v}_m, \underline{v}_m \rangle \end{bmatrix} \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_m \end{bmatrix} = \begin{bmatrix} \langle \underline{v}, \underline{v}_1 \rangle \\ \langle \underline{v}, \underline{v}_2 \rangle \\ \vdots \\ \langle \underline{v}, \underline{v}_m \rangle \end{bmatrix} \quad (4.4)$$

befriedigen.

### Bemerkung:

Gleichung (4.4) wird als **Matrixform der Orthogonalitätsbeziehung** bezeichnet. Man bemerke, dass  $\langle \underline{v}_i, \underline{v}_j \rangle = \langle \underline{v}_j, \underline{v}_i \rangle$  bedeutet, dass die  $m \times m$  Matrix in (4.4) **symmetrisch** ist.

### Beweis:

Wir stellen fest, dass

$$\begin{aligned} \langle \underline{v} - \hat{\underline{v}}, \underline{v}_j \rangle &= \langle \underline{v}_j, \underline{v} - \hat{\underline{v}} \rangle \\ &= \langle \underline{v}_j, \underline{v} - \sum_{i=1}^m c_i \underline{v}_i \rangle \\ &= \langle \underline{v}_j, \underline{v} \rangle - \sum_{i=1}^m c_i \langle \underline{v}_j, \underline{v}_i \rangle \end{aligned}$$

gilt. Es folgt somit, dass  $\langle \underline{v} - \hat{\underline{v}}, \underline{v}_j \rangle = 0$  äquivalent zu

$$\sum_{i=1}^m \langle \underline{v}_j, \underline{v}_i \rangle c_i = \langle \underline{v}, \underline{v}_j \rangle$$

ist. Diese Gleichung entspricht der  $j$ -ten Zeile in (4.4). □

### Linearitätseigenschaft der optimalen linearen Annäherung:

Seien  $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_m, \underline{w}_1$  und  $\underline{w}_2$  Vektoren in einem reellen Skalarproduktraum, sei  $\hat{\underline{w}}_1$  (bzw.  $\hat{\underline{w}}_2$ ) der Vektor in  $S(\underline{v}_1, \dots, \underline{v}_m)$  der  $\|\underline{w}_1 - \hat{\underline{w}}_1\|^2$  (bzw.  $\|\underline{w}_2 - \hat{\underline{w}}_2\|^2$ ) minimiert und seien  $a$  und  $b$  reelle Zahlen, dann ist

$$\hat{\underline{v}} = a \hat{\underline{w}}_1 + b \hat{\underline{w}}_2$$

der Vektor in  $S(\underline{v}_1, \dots, \underline{v}_m)$ , der  $\|\underline{v} - \hat{\underline{v}}\|^2$  minimiert, wobei

$$\underline{v} = a \underline{w}_1 + b \underline{w}_2$$

ist.

**Beweis:**

Es folgt aus

$$\langle \underline{w}_1 - \hat{w}_1, \underline{v}_i \rangle = 0 \quad i = 1, 2, \dots, m$$

und

$$\langle \underline{w}_2 - \hat{w}_2, \underline{v}_i \rangle = 0 \quad i = 1, 2, \dots, m$$

dass

$$\begin{aligned} \langle (a \underline{w}_1 + b \underline{w}_2) - (a \hat{w}_1 + b \hat{w}_2), \underline{v}_i \rangle &= \langle a (\underline{w}_1 - \hat{w}_1) + b (\underline{w}_2 - \hat{w}_2), \underline{v}_i \rangle \\ &= a \langle \underline{w}_1 - \hat{w}_1, \underline{v}_i \rangle + b \langle \underline{w}_2 - \hat{w}_2, \underline{v}_i \rangle \\ &= 0 \end{aligned}$$

für  $i = 1, 2, \dots, m$  ist. □

**Trennungseigenschaft der optimalen linearen Annäherung:**

Seien  $\underline{v}_1, \dots, \underline{v}_m, \underline{v}$  Vektoren in einem reellen Skalarproduktraum, sei

$$\langle \underline{v}_i, \underline{v}_j \rangle = 0 \quad \text{für } i \in \{1, 2, \dots, k\}, j \in \{k+1, k+2, \dots, m\},$$

und sei  $\hat{\underline{v}} = \hat{\underline{v}}_A + \hat{\underline{v}}_B$  wobei  $\hat{\underline{v}}_A \in \mathcal{S}(\underline{v}_1, \dots, \underline{v}_k)$  und  $\hat{\underline{v}}_B \in \mathcal{S}(\underline{v}_{k+1}, \dots, \underline{v}_m)$ , dann und nur dann ist  $\hat{\underline{v}}$  der Vektor in  $\mathcal{S}(\underline{v}_1, \dots, \underline{v}_m)$ , der  $\|\underline{v} - \hat{\underline{v}}\|^2$  minimiert, wenn  $\hat{\underline{v}}_A$  (bzw.  $\hat{\underline{v}}_B$ ) der Vektor in  $\mathcal{S}(\underline{v}_1, \dots, \underline{v}_k)$  (bzw.  $\mathcal{S}(\underline{v}_{k+1}, \dots, \underline{v}_m)$ ) ist, der  $\|\underline{v} - \hat{\underline{v}}_A\|^2$  (bzw.  $\|\underline{v} - \hat{\underline{v}}_B\|^2$ ) minimiert.

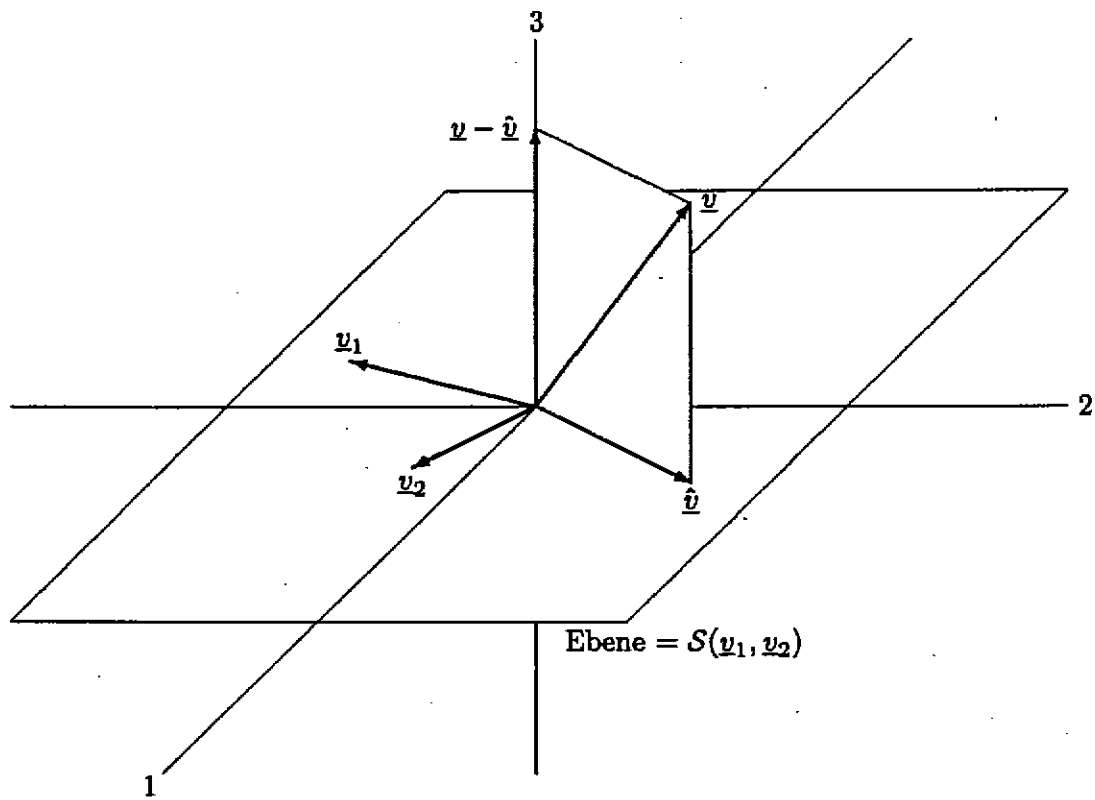
**Beweis:**

$$\begin{aligned} \langle \underline{v} - \hat{\underline{v}}, \underline{v}_j \rangle &= \langle \underline{v} - (\hat{\underline{v}}_A + \hat{\underline{v}}_B), \underline{v}_j \rangle \\ &= \langle \underline{v} - \hat{\underline{v}}_A, \underline{v}_j \rangle - \langle \hat{\underline{v}}_B, \underline{v}_j \rangle \\ &= \langle \underline{v} - \hat{\underline{v}}_A, \underline{v}_j \rangle \quad \text{für } j \in \{1, 2, \dots, k\}, \end{aligned}$$

weil  $\hat{\underline{v}}_B$  und  $\underline{v}_j$  orthogonal sind. Deswegen gilt  $\langle \underline{v} - \hat{\underline{v}}, \underline{v}_j \rangle = 0$  für  $j \in \{1, \dots, k\}$  genau dann, wenn  $\langle \underline{v} - \hat{\underline{v}}_A, \underline{v}_j \rangle = 0$  ist für  $j \in \{1, \dots, k\}$ , d.h. wenn  $\hat{\underline{v}}_A$  der Vektor in  $\mathcal{S}(\underline{v}_1, \dots, \underline{v}_k)$  ist, der  $\|\underline{v} - \hat{\underline{v}}_A\|^2$  minimiert. Analog gilt  $\langle \underline{v} - \hat{\underline{v}}, \underline{v}_j \rangle = 0$  für  $j \in \{k+1, \dots, m\}$  genau dann, wenn  $\hat{\underline{v}}_B$  der Vektor in  $\mathcal{S}(\underline{v}_{k+1}, \dots, \underline{v}_m)$  ist, der  $\|\underline{v} - \hat{\underline{v}}_B\|^2$  minimiert. □

**Bemerkung zum Orthogonalitätsprinzip:** Obwohl  $\hat{\underline{v}}$  immer eindeutig bestimmt ist, sind die Koeffizienten  $c_1, c_2, \dots, c_m$  dann und nur dann eindeutig bestimmt, wenn die Vektoren  $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_m$  linear unabhängig sind. Das bedeutet, dass die Matrix in (4.4) dann und nur dann invertierbar ist, wenn  $\underline{v}_1, \underline{v}_2, \dots, \underline{v}_m$  linear unabhängig sind.





Graphische Darstellung des Orthogonalitätsprinzips in  $V = \mathbb{R}^3$  mit  $\langle \underline{u}, \underline{v} \rangle = \underline{u} \cdot \underline{v}$ .

## 5 LINEARE MMSE-SCHÄTZUNG

### 5.1 Einführung

Wir kehren nun zum Problem der statistischen Schätzung zurück, das heisst, wir beobachten den Zufallsvektor  $\underline{Y} = (Y_1, Y_2, \dots, Y_n)$  und wollen eine optimale Schätzung  $\hat{X}$  des Parameters  $X$  machen. Als Kriterium der Güte der Schätzung wählen wir die Grösse des "mean-squared-error (MSE)"  $E[(X - \hat{X})^2]$ . Wir führen nun jedoch die Beschränkung ein, dass unsere Schätzungsregel linear sein muss. Vorerst nehmen wir an, dass dem Beobachter die gleichen Grössen wie bei der Bayesschen Schätzung bekannt sind.

**Das Problem der linearen MMSE-Schätzung:**

- (1) Bekannt sind  $p_{\underline{Y}|X}(\cdot|x)$  für alle  $x \in X(\Omega)$  und  $p_X$ .
- (2) Wir wollen reelle Zahlen  $c_1, c_2, \dots, c_n$  so finden, dass die Schätzung  $\hat{X} = c_1 Y_1 + c_2 Y_2 + \dots + c_n Y_n$  den MSE  $E[(X - \hat{X})^2]$  minimiert.

Aus  $p_{\underline{Y}|X}$  und  $p_X$  können wir  $p_{X\underline{Y}}$  bestimmen, das heisst, wir wissen alles über die statistischen Verhältnisse der Zufallsgrössen  $X, Y_1, Y_2, \dots, Y_n$ . Wir werden jedoch später erkennen, dass für die optimale lineare Schätzungsregel nicht alle diese Grössen bekannt sein müssen.

Unsere Strategie in diesem Kapitel ist die folgende: Zuerst werden wir zeigen, dass das Problem der linearen MMSE-Schätzung zum Problem der linearen Annäherung in einem Vektorraum äquivalent ist. Demnach können wir die im vorhergehenden Kapitel hergeleiteten Verfahren anwenden. Unser erster Schritt in diese Richtung ist der folgende Satz:

### 5.2 Menge aller Zufallsgrössen als Vektorraum

**Satz 1:** Die Menge aller Zufallsgrössen, die auf einem Ergebnisraum  $\Omega$  definiert werden können, ist ein reeller Vektorraum.

**Beweis:** Wir wissen, dass die Menge aller reellwertigen Funktionen mit Definitionsbereich  $\Omega$  ein Vektorraum ist (vgl. Beispiel (2), Seite 41). Die Menge aller Zufallsgrössen mit Definitionsbereich  $\Omega$  ist eine Untermenge dieser Menge. Es bleibt noch zu beweisen, dass diese Untermenge den Test für einen Unterraum (siehe Seite 42) besteht, das heisst, wir müssen beweisen, dass die Summe zweier Zufallsgrössen wieder eine Zufallsgrösse ist und dass eine Konstante multipliziert mit einer Zufallsgrösse auch wieder eine Zufallsgrösse ist. Im Kapitel I haben wir diese Tatsache als selbstverständlich angenommen. Jetzt lohnt es sich, diese Selbstverständlichkeit durch einen Beweis zu untermauern. Unterbrechen wir hier den Beweis vom Satz 1, um einige nützliche Tatsachen herzuleiten.

**Lemma 1:** Seien  $X$  und  $Y$  Zufallsgrössen mit dem Definitionsbereich  $\Omega$  und sei  $X(\omega) + Y(\omega) < \beta$  für ein gewisses  $\omega \in \Omega$ . Dann gibt es eine rationale Zahl  $q$  mit  $X(\omega) < q$  und  $Y(\omega) < \beta - q$ .

**Bemerkung:** Wir bezeichnen die Menge der rationalen Zahlen mit  $Q$ .

**Beweis:** Da

$$\Delta = \beta - X(\omega) - Y(\omega) > 0$$

ist, gibt es sicher eine rationale Zahl  $q$  mit

$$X(\omega) < q < X(\omega) + \frac{\Delta}{2}. \quad (5.1)$$

Demnach gilt auch:

$$q < X(\omega) + \frac{\Delta}{2} = \frac{\beta + X(\omega) - Y(\omega)}{2} < \frac{\beta + q - Y(\omega)}{2},$$

woraus folgt, dass

$$2q < \beta + q - Y(\omega)$$

gilt, was wiederum zu

$$Y(\omega) < \beta - q \quad (5.2)$$

äquivalent ist. □

**Lemma 2:** Seien  $X$  und  $Y$  Zufallsgrößen mit Definitionsbereich  $\Omega$  und  $\beta$  eine reelle Zahl, dann gilt:

$$\{\omega : X(\omega) + Y(\omega) < \beta\} = \bigcup_{q \in \mathbb{Q}} \{\omega : X(\omega) < q \text{ und } Y(\omega) < \beta - q\}.$$

**Beweis:** Gehört  $\omega$  zur Vereinigung auf der rechten Seite, dann existiert ein  $q \in \mathbb{Q}$  mit  $X(\omega) < q$  und  $Y(\omega) < \beta - q$ , woraus  $X(\omega) + Y(\omega) < \beta$  folgt. Deshalb ist diese Vereinigung eine Untermenge der Menge auf der linken Seite. Gehört  $\omega$  der Menge auf der linken Seite an, dann gilt:  $X(\omega) + Y(\omega) < \beta$  und es folgt mit dem Lemma 1, dass  $\omega$  der Vereinigung auf der rechten Seite angehört. Aus diesem Grund sind die Mengen auf beiden Seiten des Gleichheitszeichens identisch. □

**Korollar:** Seien  $X$  und  $Y$  Zufallsgrößen mit dem Definitionsbereich  $\Omega$ , dann ist auch  $X + Y$  eine Zufallsgröße.

**Beweis:** Wir müssen beweisen, dass die Menge  $\{\omega : X(\omega) + Y(\omega) < \beta\}$  für jedes  $\beta$  zur Klasse der Ereignisse  $\mathcal{A}$  gehört. Durch Umformung von Lemma 2 erhalten wir:

$$\{\omega : X(\omega) + Y(\omega) < \beta\} = \bigcup_{q \in \mathbb{Q}} [\{\omega : X(\omega) < q\} \cap \{\omega : Y(\omega) < \beta - q\}].$$

Ferner erinnern wir uns, dass  $\mathcal{A}$  eine  $\sigma$ -Algebra ist (siehe Seite 2). Da  $X$  und  $Y$  Zufallsgrößen sind, müssen  $\{\omega : X(\omega) < q\}$  und  $\{\omega : Y(\omega) < \beta - q\}$  Ereignisse in  $\mathcal{A}$  sein und es folgt dann, dass auch  $\{\omega : X(\omega) < q\} \cap \{\omega : Y(\omega) < \beta - q\}$  ein Ereignis in  $\mathcal{A}$  sein muss. Es gibt nur abzählbar unendlich viele rationale Zahlen und da die Vereinigung von abzählbar unendlich

vielen Ereignissen in  $\mathcal{A}$  wieder ein Ereignis in  $\mathcal{A}$  ist, folgt, dass  $\{\omega : X(\omega) + Y(\omega) < \beta\}$  ein Ereignis in  $\mathcal{A}$  ist. □

**Bemerkung:**

In diesem Beweis sieht man den Grund für die Hypothese der Wahrscheinlichkeitstheorie, dass  $\mathcal{A}$  eine  $\sigma$ -Algebra sein muss. Wäre eine Vereinigung von abzählbar unendlich vielen Ereignissen nicht auch immer ein Ereignis, dann wäre  $X + Y$  nicht immer eine Zufallsgrösse, was eine sehr unangenehme Tatsache wäre. Man sieht nun auch die Notwendigkeit des dritten Axioms von Kolmogorov, welches die Berechnung der Wahrscheinlichkeit einer solchen Vereinigung von Ereignissen erlaubt.

**Lemma 3:**

Sei  $X$  eine Zufallsgrösse mit Definitionsbereich  $\Omega$  und sei  $\beta$  eine reelle Zahl, dann ist die Menge  $\{\omega : X(\omega) > \beta\}$  ein Ereignis, d.h.  $\{\omega : X(\omega) > \beta\}$  gehört zu  $\mathcal{A}$ .

**Bemerkung:**

Die Definition einer Zufallsgrösse fordert nur, dass die Menge  $\{\omega : X(\omega) < \beta\}$  zu  $\mathcal{A}$  gehört.

**Beweis:**

Weil für jede reelle Zahl  $\alpha$  die Menge  $\{\omega : X(\omega) < \alpha\}$  ein Ereignis ist, muss das Komplement  $\{\omega : X(\omega) \geq \alpha\}$  auch ein Ereignis sein. Zu einem  $\omega$  mit  $X(\omega) > \beta$  gibt es sicher eine rationale Zahl  $q$  mit  $X(\omega) \geq q > \beta$ . Wir sehen, dass

$$\{\omega : X(\omega) > \beta\} = \bigcup_{q \in Q_\beta} \{\omega : X(\omega) \geq q\}$$

gilt, wobei  $Q_\beta$  die Menge von rationalen Zahlen  $q$  mit  $q > \beta$  ist. Es folgt, dass die Menge  $\{\omega : X(\omega) > \beta\}$  eine Vereinigung von abzählbar unendlich vielen Ereignissen in  $\mathcal{A}$  ist. Damit können wir schliessen, dass die Menge  $\{\omega : X(\omega) > \beta\}$  auch ein Ereignis in  $\mathcal{A}$  ist. □

**Korollar:**

Sei  $X$  eine Zufallsgrösse mit Definitionsbereich  $\Omega$  und sei  $c$  eine reelle Zahl. Dann ist  $cX$  auch eine solche Zufallsgrösse.

**Beweis:**

Wir müssen für jede reelle Zahl  $\beta$  zeigen, dass die Menge  $\{\omega : cX(\omega) < \beta\}$  ein Ereignis ist. Dieses gilt trivialerweise, falls  $c \geq 0$  ist. Sei  $c < 0$ , dann gilt

$$\{\omega : cX(\omega) < \beta\} = \{\omega : X(\omega) > \frac{\beta}{c}\}.$$

Aus Lemma 3 können wir schliessen, dass diese Menge ein Ereignis ist. □

**Beweis von Satz 1: (Fortsetzung)**

Die beiden Korollare, die wir jetzt bewiesen haben, zeigen, dass die Menge aller Zufallsgrössen mit Definitionsbereich  $\Omega$  den Test für einen Unterraum bestehen. Deshalb ist diese Menge auch ein Vektorraum.

### 5.3 Zufallsgrößen mit endlichem zweitem Moment

#### Satz 2:

Die Menge  $V$  aller Zufallsgrößen mit Definitionsbereich  $\Omega$ , die ein endliches zweites Moment haben, ist ein Vektorraum.

#### Beweis:

Es genügt zu beweisen, dass diese Menge  $V$  von Zufallsgrößen mit endlichem zweitem Moment den Test für einen Unterraum besteht. Seien  $X$  und  $Y$  Zufallsgrößen mit  $E[X^2] < \infty$  und  $E[Y^2] < \infty$ , dann gelten

$$0 \leq E[(X + Y)^2] = E[X^2] + E[Y^2] + 2E[XY]$$

und

$$0 \leq E[(X - Y)^2] = E[X^2] + E[Y^2] - 2E[XY].$$

Somit gelten auch

$$-2E[XY] \leq E[X^2] + E[Y^2]$$

und

$$2E[XY] \leq E[X^2] + E[Y^2],$$

was zu

$$2|E[XY]| \leq E[X^2] + E[Y^2]$$

äquivalent ist. Es folgt, dass

$$\begin{aligned} E[(X + Y)^2] &\leq E[X^2] + E[Y^2] + 2|E[XY]| \\ &\leq 2E[X^2] + 2E[Y^2] < \infty \end{aligned}$$

gilt. Sei  $X$  eine Zufallsgröße mit  $E[X^2] < \infty$  und sei  $c$  eine reelle Zahl, dann gilt:

$$E[(cX)^2] = E[c^2 X^2] = c^2 E[X^2] < \infty.$$

□

Nun werden wir demonstrieren, dass die Funktion  $\langle \cdot, \cdot \rangle : V \times V \rightarrow R$  mit

$$\langle X, Y \rangle = E[XY] \tag{5.3}$$

„fast“ ein Skalarprodukt für den Vektorraum  $V$  der Zufallsgrößen mit endlichem zweitem Moment ist. Zuerst stellen wir fest, dass

$$\langle X, Y \rangle = E[XY] = E[YX] = \langle Y, X \rangle$$

gilt, das dem Axiom (S1) entspricht. Dann bemerken wir, dass

$$\begin{aligned}
\langle aX + bY, Z \rangle &= E[(aX + bY)Z] \\
&= E[aXZ + bYZ] \\
&= aE[XZ] + bE[YZ] \\
&= a\langle X, Z \rangle + b\langle Y, Z \rangle
\end{aligned}$$

gilt. Das entspricht Axiom (S2). Und zuletzt sehen wir, dass

$$\langle X, X \rangle = E[X^2] \geq 0$$

gilt. Das entspricht Axiom (S3) teilweise.

Die einzige Frage die noch bleibt, lautet: Gilt  $E[X^2] = 0$  nur wenn  $X = 0$  ist? Leider heisst die Antwort "nein"! Um dieses Paradoxon besser zu verstehen, schreiben wir:

$$A_1 = \{\omega : X(\omega) = 0\}$$

und

$$A_2 = \{\omega : X(\omega) \neq 0\}.$$

Weil  $A_1$  und  $A_2$  eine komplette Klasse von paarweise unvereinbaren Ereignissen bilden, folgt mit dem Satz des totalen Erwartungswertes, dass

$$E[X^2] = E[X^2|A_1]P[A_1] + E[X^2|A_2]P[A_2].$$

Aber wenn  $\omega \in A_1$  gilt, dann ist  $X(\omega) = 0$ . Damit sehen wir, dass  $E[X^2|A_1] = 0$  und

$$E[X^2] = E[X^2|A_2]P[A_2] \tag{5.4}$$

ist.

Wenn  $\omega \in A_2$  gilt, dann gilt  $(X(\omega))^2 > 0$ . Es folgt somit, dass

$$E[X^2|A_2] > 0 \tag{5.5}$$

gilt. Ist  $E[X^2] = 0$ , dann folgt aus (5.4) und (5.5), dass  $P[A_2] = 0$  gelten muss. Äquivalent gilt  $P[A_1] = 1 - P[A_2] = 1$ . Da aber

$$P[A_1] = P[X = 0]$$

ist haben wir folgendes Lemma bewiesen:

**Lemma 4:** Sei  $E[X^2] = 0$ , dann gilt  $P[X = 0] = 1$ .

Es ist in der Wahrscheinlichkeitstheorie üblich zu schreiben:

$$X = Y \tag{mW1}$$

(wobei (mW1) "mit Wahrscheinlichkeit 1" bedeutet), wenn

$$P[X = Y] = 1$$

gilt. Also kann Gleichung (5.5) äquivalent als

$$X = 0 \quad (\text{mW1}) \quad (5.6)$$

geschrieben werden. Die Verallgemeinerung von Lemma 4 (dessen leichten Beweis wird dem Leser überlassen) heisst:

**Lemma 5:** Sei  $E[(X - Y)^2] = 0$ , dann gilt

$$X = Y \quad (\text{mW1}).$$

Es ist auch in der Wahrscheinlichkeitstheorie üblich, Zufallsgrössen  $X$  und  $Y$  mit  $X = Y$  (mW1) als dieselben Zufallsgrössen zu betrachten. [ Wir machen etwas ähnliches, wenn wir ein Gleichheitszeichen zwischen eine nichtstetige Funktion und eine unendlich Summe von Fourier Funktionen schreiben; die Gleichung gilt nur für "fast alle" Werte des Arguments]. Um ganz ehrlich zu sein, müssen wir dann sagen, dass die Gleichung  $X = Y$  nur bedeutet, dass  $X$  und  $Y$  zu derselben Äquivalenzklasse von Zufallsgrössen gehören. Wir definieren nun, dass die Zufallsgrössen  $X$  und  $Y$  genau dann äquivalent sind, wenn  $X = Y$  (mW1) gilt. Künftig werden wir bloss  $X=Y$  statt  $X=Y$  (mW1) schreiben, d.h. wir werden stillschweigend annehmen, dass wir tatsächlich mit Äquivalenzklassen von Zufallsgrössen arbeiten, obwohl wir die üblichen Notationen für Zufallsgrössen benützen.

Mit der obigen Abmachung sehen wir, dass

$$\langle X, X \rangle = E[X^2] \geq 0$$

gilt. Die Gleichheit gilt genau dann, wenn  $X = 0$  ist, was dem Axiom (S3) entspricht.

Wir haben somit folgendes bewiesen:

**Der Skalarproduktraum von Zufallsgrössen:** Die Funktion

$$\langle X, Y \rangle = E[XY]$$

ist ein Skalarprodukt für den Vektorraum  $V$  der Zufallsgrössen mit endlichem zweitem Moment.

Die entsprechende Norm heisst:

$$\|X\| = \sqrt{E[X^2]}.$$

Die Zufallsgrössen  $X$  und  $Y$  heissen orthogonal, falls

$$E[XY] = 0$$

gilt.

## 5.4 Lineare MMSE-Schätzung

Wir sehen, dass der zweite Teil des Problems der linearen MMSE-Schätzung wie folgt umschrieben werden kann:

- (2) Wir wollen reelle Zahlen  $c_1, c_2, \dots, c_n$  finden, so dass  $\hat{X} = c_1 Y_1 + c_2 Y_2 + \dots + c_n Y_n$  die Zufallsgrösse in  $\mathcal{S}(Y_1, \dots, Y_n)$  ist, die  $\|X - \hat{X}\|^2 = E[(X - \hat{X})^2]$  minimiert.

Von der Matrixform der Orthogonalitätsgleichung wissen wir, dass es genügt,  $\langle Y_i, Y_j \rangle = E[Y_i Y_j]$  und  $\langle X, Y_i \rangle = E[XY_i]$  für  $i = 1, 2, \dots, n$  und  $j = 1, 2, \dots, n$  zu kennen, um die Koeffizienten  $c_1, \dots, c_n$  zu finden. Deshalb können wir den ersten Teil des Problems der linearen MMSE-Schätzung wie folgt reduzieren, ohne die Lösung zu ändern:

- (1) Sowohl  $\langle Y_i, Y_j \rangle = E[Y_i Y_j]$  als auch  $\langle X, Y_i \rangle = E[XY_i]$  sind für  $i = 1, 2, \dots, n$  und  $j = 1, 2, \dots, n$  bekannt.

**Bemerkung:** Um  $E[(X - \hat{X})^2]$  zu berechnen, muss man auch  $E[X^2]$  kennen.

Wir wissen jedoch schon aus dem vorhergehenden Kapitel, wie man dieses Problem löst.

**Orthogonalitätsprinzip für die lineare MMSE-Schätzung:** Seien  $Y_1, Y_2, \dots, Y_n$  und  $X$  Zufallsgrössen mit endlichem zweiten Moment, dann ist  $\hat{X} = c_1 Y_1 + c_2 Y_2 + \dots + c_n Y_n$  genau dann diejenige Zufallsgrösse in  $\mathcal{S}(Y_1, Y_2, \dots, Y_n)$ , welche  $E[(X - \hat{X})^2]$  minimiert, falls der Fehler  $X - \hat{X}$  orthogonal zu  $Y_i$  ( $i = 1, \dots, n$ ) ist, das heisst, falls

$$E[(X - \hat{X})Y_i] = 0 \quad (i = 1, \dots, n) \quad (5.7)$$

gilt. Für diese optimale lineare Schätzung gilt:

$$\text{MSE} = E[(X - \hat{X})^2] = E[(X - \hat{X})X] = E[X^2] - E[\hat{X}^2]. \quad (5.8)$$

**Matrixform der Orthogonalitätsgleichung für die lineare MMSE-Schätzung:** Seien  $Y_1, \dots, Y_n$  und  $X$  Zufallsgrössen mit endlichem zweiten Moment, dann ist

$$\hat{X} = c_1 Y_1 + c_2 Y_2 + \dots + c_n Y_n$$

genau dann die Zufallsgrösse in  $\mathcal{S}(Y_1, \dots, Y_n)$  die  $E[(X - \hat{X})^2]$  minimiert, wenn die Koeffizienten  $c_1, \dots, c_n$  die Gleichung

$$\begin{bmatrix} E[Y_1^2] & E[Y_1 Y_2] & \dots & E[Y_1 Y_n] \\ E[Y_2 Y_1] & E[Y_2^2] & \dots & E[Y_2 Y_n] \\ \vdots & & & \\ E[Y_n Y_1] & E[Y_n Y_2] & \dots & E[Y_n^2] \end{bmatrix} \cdot \begin{bmatrix} c_1 \\ c_2 \\ \vdots \\ c_n \end{bmatrix} = \begin{bmatrix} E[XY_1] \\ E[XY_2] \\ \vdots \\ E[XY_n] \end{bmatrix} \quad (5.9)$$

befriedigen.

Die Gleichung (5.7) heisst Orthogonalitätsgleichung für die lineare MMSE-Schätzung.



**Linearitätseigenschaft der linearen MMSE-Schätzung:** Seien  $Y_1, Y_2, \dots, Y_n, W_1$  und  $W_2$  Zufallsgrößen mit endlichem zweiten Moment, sei  $\hat{W}_1$  (bzw.  $\hat{W}_2$ ) diejenige Zufallsgröße in  $\mathcal{S}(Y_1, Y_2, \dots, Y_n)$ , welche  $E[(W_1 - \hat{W}_1)^2]$ , (bzw.  $E[(W_2 - \hat{W}_2)^2]$ ) minimiert und seien  $a_1$  und  $a_2$  reelle Zahlen, dann ist

$$\hat{X} = a_1 \hat{W}_1 + a_2 \hat{W}_2$$

diejenige Zufallsgröße in  $\mathcal{S}(Y_1, Y_2, \dots, Y_n)$ , welche  $E[(X - \hat{X})^2]$  minimiert, wobei

$$X = a_1 W_1 + a_2 W_2$$

ist.

**Trennungseigenschaft der linearen MMSE-Schätzung:** Seien  $Y_1, \dots, Y_k, Y_{k+1}, \dots, Y_n$  und  $X$  Zufallsgrößen mit endlichem zweiten Moment, sei

$$E[Y_i Y_j] = 0 \quad \text{für } i \in \{1, \dots, k\}, j \in \{k+1, \dots, n\},$$

und sei  $\hat{X} = \hat{X}_A + \hat{X}_B$  mit  $\hat{X}_A \in \mathcal{S}(Y_1, \dots, Y_k)$  und  $\hat{X}_B \in \mathcal{S}(Y_{k+1}, \dots, Y_n)$ , dann ist  $\hat{X}$  genau dann diejenige Zufallsgröße in  $\mathcal{S}(Y_1, Y_2, \dots, Y_n)$ , welche  $E[(X - \hat{X})^2]$  minimiert, falls  $\hat{X}_A$  (bzw.  $\hat{X}_B$ ) die Zufallsgröße in  $\mathcal{S}(Y_1, \dots, Y_k)$  (bzw.  $\mathcal{S}(Y_{k+1}, \dots, Y_n)$ ) ist, welche  $E[(X - \hat{X}_A)^2]$  (bzw.  $E[(X - \hat{X}_B)^2]$ ) minimiert.

Wir wollen nun zeigen, dass es für die lineare Schätzung von Vorteil ist, wenn man annimmt, dass dem Beobachter die triviale Zufallsgröße  $Y_{n+1} = 1$  bekannt ist.

**Lemma 6:** Seien  $X$  und  $\hat{X}$  Zufallsgrößen mit endlichem zweiten Moment, dann ist der Fehler  $X - \hat{X}$  genau dann orthogonal zu der Zufallsgröße 1, wenn  $E[\hat{X}] = E[X]$  ist.

**Beweis:**  $\langle X - \hat{X}, 1 \rangle = E[(X - \hat{X})1] = E[X] - E[\hat{X}].$  □

Das Lemma 6 zeigt, dass die lineare MMSE-Schätzung von  $X$  in  $\mathcal{S}(Y_1, Y_2, \dots, Y_n, 1)$  genau dann identisch mit der linearen MMSE-Schätzung  $\hat{X}$  von  $X$  in  $\mathcal{S}(Y_1, Y_2, \dots, Y_n)$  ist, falls  $E[X] = E[\hat{X}]$  gilt. Ist  $E[X] \neq E[\hat{X}]$  dann liefert die lineare MMSE-Schätzung von  $X$  in  $\mathcal{S}(Y_1, Y_2, \dots, Y_n, 1)$  einen kleineren MSE als die Schätzung  $\hat{X}$ .

Der nachfolgende Satz zeigt, dass die lineare MMSE-Schätzung von  $X$  in  $\mathcal{S}(Y_1, Y_2, \dots, Y_n, 1)$  durch die Lösung von nur  $n$  (statt  $n + 1$ ) linearen Gleichungen gefunden werden kann.

**Satz über die lineare MMSE-Schätzung mit dem zusätzlichen Datum 1:** Seien  $Y_1, Y_2, \dots, Y_n$  und  $X$  Zufallsgrößen mit endlichem zweiten Moment und sei  $\hat{X}'$  die lineare MMSE-Schätzung von  $X' = X - E[X]$  in  $\mathcal{S}(Y'_1, \dots, Y'_n)$  mit  $Y'_i = Y_i - E[Y_i]$  ( $i = 1, \dots, n$ ), dann ist

$$\hat{X} = \hat{X}' + E[X]$$

die lineare MMSE-Schätzung von  $X$  in  $\mathcal{S}(Y_1, Y_2, \dots, Y_n, 1)$ .

**Beweis:** Da  $\hat{X}'$  eine Linearkombination von  $Y'_1, \dots, Y'_n$  ist und da  $E[Y'_i] = 0$  für  $i = 1, \dots, n$  ist, folgt daraus, dass  $E[\hat{X}'] = 0$  gilt. Deshalb gilt:

$$E[X - \hat{X}] = E[X - \hat{X}' - E[X]] = 0,$$

das heisst,  $X - \hat{X}$  ist orthogonal zu  $Y_{n+1} = 1$ . Es bleibt noch zu beweisen, dass  $X - \hat{X} = X - \hat{X}' - E[X]$  auch orthogonal zu  $Y_i$  ist für  $i = 1, \dots, n$ :

$$\begin{aligned} E[(X - \hat{X}' - E[X])Y_i] &= E[(X - E[X] - \hat{X}')(Y_i - E[Y_i] + E[Y_i])] \\ &= E[(X' - \hat{X}')(Y_i' + E[Y_i])] \\ &= E[(X' - \hat{X}')Y_i'] + E[(X' - \hat{X}')]E[Y_i] \\ &= 0 + 0 \\ &= 0 \end{aligned} \quad \text{für } i = 1, \dots, n.$$

□

Folgende Beziehungen sind für die Anwendung des obigen Satzes wichtig:

$$\begin{aligned} \langle Y_i', Y_j' \rangle &= E[Y_i' Y_j'] \\ &= E[(Y_i - E[Y_i])(Y_j - E[Y_j])] \\ &= \text{Cov}(Y_i, Y_j). \end{aligned}$$

Analog gilt:

$$\langle X', Y_i' \rangle = \text{Cov}(X, Y_i).$$

Es folgt, dass die Gleichung (5.7) als

$$\Lambda_{\underline{Y}} \cdot \underline{c} = \text{Cov}(X, \underline{Y})$$

geschrieben werden kann, wobei  $\underline{c} = (c_1, c_2, \dots, c_n)$  und

$$\text{Cov}(X, \underline{Y}) = \begin{bmatrix} \text{Cov}(X, Y_1) \\ \text{Cov}(X, Y_2) \\ \vdots \\ \text{Cov}(X, Y_n) \end{bmatrix} \quad (5.10)$$

sind. Ferner wollen wir uns merken, dass  $\hat{X}' = c_1 Y_1' + c_2 Y_2' + \dots + c_n Y_n'$  als

$$\hat{X}' = \underline{c}^T \underline{Y}' = \underline{c}^T (\underline{Y} - E[\underline{Y}])$$

geschrieben werden kann.

**Korollar zum Satz über die lineare MMSE-Schätzung mit dem zusätzlichen Datum 1:** Seien  $Y_1, \dots, Y_n$  und  $X$  Zufallsgrößen mit endlichem zweitem Moment. Dann ist

$$\hat{X} = \underline{c}^T (\underline{Y} - E[\underline{Y}]) + E[X] \quad (5.11)$$

genau die Zufallsgröße in  $\mathcal{S}(Y_1, \dots, Y_n, 1)$ , die  $E[(X - \hat{X})^2]$  minimiert, wenn  $\underline{c}$  die Gleichung

$$\Lambda_{\underline{Y}} \cdot \underline{c} = \text{Cov}(X, \underline{Y}) \quad (5.12)$$

befriedigt, wobei  $\underline{Y} = (Y_1, \dots, Y_n)$  ist.

## 5.5 Lineare MMSE-Schätzung im Gaußschen Fall

Bei der Anwendung kommt es oft vor (aber nicht immer!), dass  $Y_1, \dots, Y_n$  und  $X$  gemeinsam gaußverteilt sind. Der folgende wichtige Satz behandelt diesen Fall und zeigt den Grund dafür, dass die lineare MMSE-Schätzung von besonderer Bedeutung ist.

**Satz über die lineare MMSE-Schätzung im Gaußschen Fall:** Seien  $Y_1, \dots, Y_n$  und  $X$  gemeinsam gaußverteilt. Dann ist die lineare MMSE-Schätzung  $\hat{X}$  von  $X$  in  $\mathcal{S}(Y_1, \dots, Y_n, 1)$  sowohl identisch mit der Bayesschen MMSE-Schätzung von  $X$  für die Beobachtung  $\underline{Y} = (Y_1, \dots, Y_n)$  als auch mit der (sogenannten MAP-) Schätzung, die  $E[S(X, \hat{X})]$  minimiert, wenn kleine Fehler kostenlos und alle anderen Fehler kostengleich sind. Ferner gilt in diesem Gaußschen Fall:

$$E[\hat{X} | \underline{Y} = \underline{y}] = E[X | \underline{Y} = \underline{y}] \quad \text{für alle } \underline{y}. \quad (5.13)$$

Das mittlere Fehlerquadrat ist unabhängig von der Beobachtung, d.h. es gilt:

$$E[(X - \hat{X})^2 | \underline{Y} = \underline{y}] = E[(X - \hat{X})^2] \quad \text{für alle } \underline{y}. \quad (5.14)$$

**Bemerkung:** Die triviale Zufallsgröße 1 nützt bei der Bayesschen Schätzung wie auch bei der ML-Schätzung nichts, d.h. die Bayessche (bzw. ML-) Schätzung von  $X$  für die Beobachtung  $Y_1, \dots, Y_n$  und die Beobachtung  $Y_1, \dots, Y_n, 1$  sind immer identisch.

**Beweis:** Es sei  $Y'_i = Y_i - E[Y_i]$  ( $i = 1, \dots, n$ ) und  $X' = X - E[X]$ . Dann gilt, dass  $Y'_1, \dots, Y'_n$  und  $X'$  auch gemeinsam gaußverteilt sind. Wir können  $\hat{X}$  als  $\hat{X} = \hat{X}' + E[X]$  schreiben, wobei  $\hat{X}' = \underline{c}^T \underline{Y}'$  die lineare MMSE-Schätzung von  $X'$  in  $\mathcal{S}(Y'_1, \dots, Y'_n)$  ist. Es folgt, dass  $Y'_1, \dots, Y'_n$  und  $X' - \hat{X}'$  gemeinsam gaußverteilt sind, weil sie durch lineare Transformation aus  $Y'_1, \dots, Y'_n, X'$  hervorgehen. Bedingt durch das Orthogonalitätsprinzip und mit  $\underline{Y}' = (Y'_1, \dots, Y'_n)$  gilt:

$$\text{Cov}(X' - \hat{X}', \underline{Y}') = \underline{0}, \quad (5.15)$$

weil  $E[(X' - \hat{X}')Y'_i] = 0$  für  $i = 1, \dots, n$  ist. Aber Gaußsche Zufallsgrößen und Zufallsvektoren sind genau dann unabhängig, wenn ihre Kovarianz verschwindet. Es folgt von (5.15), dass  $X' - \hat{X}'$  und  $\underline{Y}'$  unabhängig sind, d.h., dass

$$p_{X' - \hat{X}', \underline{Y}'}(x, \underline{y}') = p_{X' - \hat{X}'}(x) \cdot p_{\underline{Y}'}(\underline{y}') \quad \text{für alle } x, \underline{y}' \quad (5.16)$$

ist.

Ferner gilt dann

$$p_{X' - \hat{X}' | \underline{Y}'}(x | \underline{y}') = p_{X' - \hat{X}'}(x) \quad \text{für alle } x, \underline{y}'. \quad (5.17)$$

Aus (5.17) folgt

$$E[X' - \hat{X}' | \underline{Y}' = \underline{y}'] = E[X' - \hat{X}'] = 0,$$

weil  $E[X'] = E[\hat{X}'] = 0$  ist. Es gilt auch, dass

$$\begin{aligned} E[X' - \hat{X}' | \underline{Y}' = \underline{y}'] &= E[X' | \underline{Y}' = \underline{y}'] - E[\hat{X}' | \underline{Y}' = \underline{y}'] \\ &= E[X' | \underline{Y}' = \underline{y}'] - \underline{c}^T \underline{y}' \end{aligned}$$

ist. Daraus können wir schliessen, dass

$$E[X' | \underline{Y}' = \underline{y}'] = \underline{c}^T \underline{y}' \quad (5.18)$$

ist. Dann sehen wir ferner, dass

$$\begin{aligned} E[X | \underline{Y} = \underline{y}] &= E[X' + E[X] | \underline{Y}' + E[\underline{Y}] = \underline{y}] \\ &= E[X' | \underline{Y}' = \underline{y} - E[\underline{Y}]] + E[X] \\ &= \underline{c}^T (\underline{y} - E[\underline{Y}]) + E[X] \end{aligned} \quad (5.19)$$

gilt. Jetzt erkennen wir, dass  $E[X | \underline{Y} = \underline{y}]$  die Bayessche MMSE-Schätzung von  $X$  für die Beobachtung  $\underline{Y} = \underline{y}$  ist und  $\underline{c}^T (\underline{y} - E[\underline{Y}]) + E[X]$  die lineare MMSE-Schätzung von  $X$  für die Beobachtung  $(Y_1, \dots, Y_n, 1) = (\underline{Y}, 1)$  ist, d.h. dass diese zwei Schätzungen gleich sind.

Zunächst wollen wir die Schätzung von  $X$  finden, die  $E[S(X, \hat{X})]$  minimiert, wenn kleine Fehler kostenlos und alle anderen Fehler kostengleich sind. Wir wollen also das  $x$  finden, das  $p_{X|\underline{Y}}(x|\underline{y})$  maximiert (siehe Seite 39). [Wir nennen hier diese Schätzung die MAP-Schätzung von  $X$ .] Weil

$$p_{X|\underline{Y}}(x|\underline{y}) = p_{X' - \hat{X}' | \underline{Y}'}(x - E[X] - \underline{c}^T (\underline{y} - E[\underline{Y}]) | \underline{y} - E[\underline{Y}])$$

ist, folgt aus (5.17), dass

$$p_{X|\underline{Y}}(x|\underline{y}) = p_{X' - \hat{X}'}(x - E[X] - \underline{c}^T (\underline{y} - E[\underline{Y}])) \quad (5.20)$$

gilt. Weil  $X' - \hat{X}'$  gaussverteilt mit Mittelwert 0 ist, so, dass  $x' = 0$  die WSK-Dichte  $p_{X' - \hat{X}'}$  maximiert, folgt aus (5.20), dass

$$x = E[X] + \underline{c}^T (\underline{y} - E[\underline{Y}])$$

der Wert ist, der  $p_{X|\underline{Y}}(x|\underline{y})$  maximiert. Das heisst, die MAP-Schätzung und die lineare MMSE-Schätzung sind identisch.

Aus (5.20) folgt unmittelbar, dass

$$E[X|Y = \underline{y}] = E[X] + \underline{c}^T(\underline{y} - E[Y]) = E[\hat{X}|Y = \underline{y}]$$

und

$$\text{Var}[X - \hat{X}|Y = \underline{y}] = \text{Var}[X - \hat{X}]$$

gelten. Die Gleichungen (5.13) und (5.14) folgen daraus.  $\square$

**Bemerkung:** Die maximum-likelihood- (ML-) Schätzung von  $X$  für die Beobachtung  $Y$  ist mit der linearen MMSE-Schätzung von  $X$  nicht identisch. Als Gegenbeispiel sei  $Y = X + Z$ . Dabei seien  $X$  und  $Z$  unabhängige Gaußsche Zufallsgrößen mit  $E[X] = E[Z] = 0$  und  $\text{Var}[X] = \text{Var}[Z] = 1$ . Dann ist die lineare MMSE-Schätzung  $\hat{X}$  von  $X$

$$\hat{X} = \frac{E[XY]}{E[Y^2]}Y = \frac{1}{2}Y.$$

Aber es gilt:

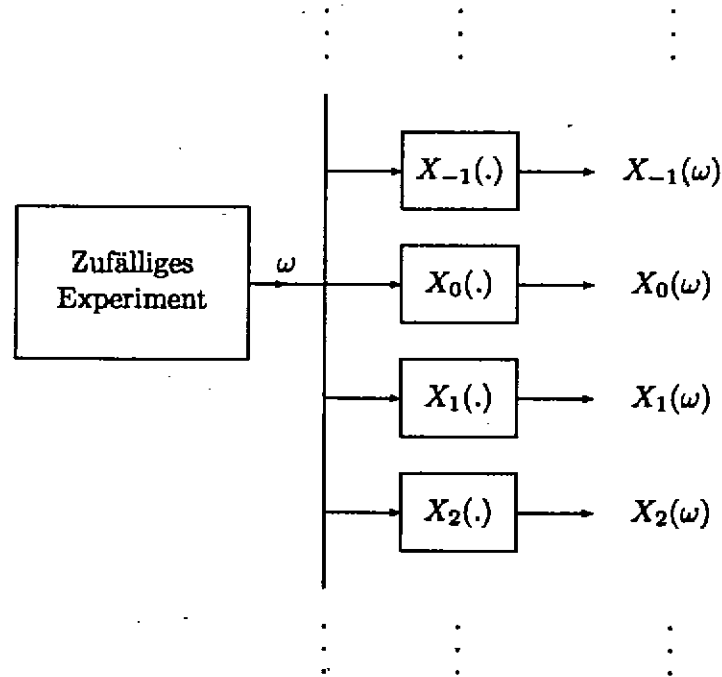
$$p_{Y|X}(y|x) = p_Z(y - x).$$

Somit ist  $x = y$  der Wert, der  $p_{Y|X}(y|x)$  maximiert, d.h.  $\hat{X}_{\text{ML}} = Y$  ist die ML-Schätzung von  $X$ .

## 6 ZEITDISKRETE STOCHASTISCHE PROZESSE UND LINEARE ZEITDISKRETE SYSTEME

### 6.1 Stochastische Prozesse

Mathematisches Modell:



Für jeden Ausgang des zufälligen Experiments erhalten wir eine zweiseitige Sequenz  $\dots, x_{-1}, x_0, x_1, x_2, \dots$  als Wert der "Zufallssequenz"  $\dots, X_{-1}, X_0, X_1, X_2, \dots$ . Interpretieren wir  $X_i$  als eine Zufallsgrösse, die zum Zeitpunkt  $i$  auftritt, dann nennt man eine solche zweiseitige Zufallssequenz einen **zeitdiskreten stochastischen Prozess**. Um dies hervorzuheben, schreibt man gewöhnlich  $X[i]$  statt  $X_i$ , obwohl man dann für den zum Ereignis  $\omega$  gehörenden Wert  $X[i]$  den unschönen Ausdruck  $X[i](\omega)$  verwenden muss. Für den stochastischen Prozess schreibt man  $X[.]$ , was deshalb sinnvoll ist, weil ein stochastischer Prozess als eine Funktion, die jeder ganzen Zahl  $i$  eine Zufallsgrösse  $X[i]$  zuordnet, betrachtet werden kann.

Der stochastische Prozess  $X[.]$  heisst **stationär**, falls für jede positive ganze Zahl  $n$  und für jede beliebige ganze Zahl  $i$ , die Zufallsvektoren  $(X[0], X[1], \dots, X[n-1])$  und  $(X[i], X[i+1], \dots, X[i+n-1])$  die gleichen WSK-Dichten besitzen.  $X[.]$  ist genau dann stationär, falls alle seine statistischen Eigenschaften unabhängig von der absoluten Zeit sind. Die Eigenschaft, dass ein stochastischer Prozess stationär ist, ist eine strenge Einschränkung, die oft schwierig zu testen ist. Für viele Anwendungen genügt eine schwächere Bedingung für die "Unabhängigkeit von der absoluten Zeit": Ein stochastischer Prozess  $X[.]$  heisst **schwach-stationär (s.s.)**, falls

- (1)  $E[X[k]]$  unabhängig von  $k$  und

(2)  $E[X[k] \cdot X[k+i]]$  für alle  $i$  unabhängig von  $k$

ist.

Ist  $X[\cdot]$  schwach-stationär, dann schreibt man

$$m_X = E[X[k]]$$

und

$$R_X[i] = E[X[k] \cdot X[k+i]] \quad \text{für alle } i.$$

Die zweiseitige Sequenz  $R_X[\cdot]$  heisst **Autokorrelationssequenz** (oder "Autokorrelationsfunktion") des stochastischen Prozesses  $X[\cdot]$ . Die Konstante  $m_X$  ist der Mittelwert des stochastischen Prozesses  $X[\cdot]$ . Es ist zu bemerken, dass

$$\text{Cov}(X[k], X[k+i]) = R_X[i] - m_X^2 \quad (6.1)$$

gilt, falls  $X[\cdot]$  s.s. ist.

Eine sehr nützliche zweiseitige Sequenz ist die **Kronecker-delta-Sequenz**  $\delta[\cdot]$ :

$$\delta[k] = \begin{cases} 0 & k \neq 0 \\ 1 & k = 0. \end{cases}$$

$\delta[\cdot]$  kann als "Einheitssequenz" interpretiert werden. Ein stochastischer Prozess  $X[\cdot]$  wird als **weisses Rauschen** bezeichnet, wenn er schwach-stationär mit  $m_X = 0$  und  $R_X[\cdot] = L \cdot \delta[\cdot]$  ist. Aus (6.1) folgt, dass  $X[\cdot]$  genau dann weisses Rauschen ist, wenn er s.s. mit  $m_X = 0$  ist und seine Komponenten unkorrelierte Zufallsgrössen sind. Das bedeutet, dass sich weisses Rauschen schnell mit der Zeit ändert.

Ein stochastischer Prozess heisst **Gaussisch**, falls für alle ganzen Zahlen  $i$  und  $n$  der Zufallsvektor  $(X[i], X[i+1], \dots, X[i+n-1])$  gaussverteilt ist. Da die Gaussverteilung völlig durch den Mittelwert und die Kovarianzen der betreffenden Zufallsgrössen bestimmt ist, folgt:

Ein Gausscher stochastischer Prozess ist genau dann stationär, wenn er schwach-stationär ist.

Selbstverständlich ist ein stationärer Prozess immer auch schwach-stationär, die Umkehrung gilt im Allgemeinen nicht.

## 6.2 Lineare zeitdiskrete Systeme (LDS)

Ein **lineares zeitdiskretes System** (LDS) ist ein zeitinvariantes, lineares System, dessen Ein- und Ausgänge zweiseitige Sequenzen sind. Die **Gewichtssequenz**  $h[\cdot]$  des LDS ist die Antwort (Ausgangssignal) des Systems auf eine Kronecker-delta-Sequenz  $\delta[\cdot]$ . Jede Sequenz  $x[\cdot]$  lässt sich als

$$x[k] = \sum_{i=-\infty}^{+\infty} x[i] \cdot \delta[k-i] \quad \text{für alle } k \quad (6.2)$$

schreiben. Definiert man  $\delta[. - i]$  als die um  $i$  Zeiteinheiten verzögerte Kronecker-delta-Sequenz, kann (6.2) wie folgt umgeschrieben werden:

$$x[.] = \sum_{i=-\infty}^{+\infty} x[i] \cdot \delta[. - i]. \quad (6.3)$$

Aus der Zeitinvarianz des LDS folgt, dass das Eingangssignal  $\delta[. - i]$  das Ausgangssignal  $h[. - i]$  verursacht. Ferner folgt aus der Linearität des LDS, dass das Eingangssignal  $x[.]$  in Formel (6.3) das Ausgangssignal

$$y[.] = \sum_{i=-\infty}^{+\infty} x[i] \cdot h[. - i] \quad (6.4)$$

verursacht. Durch Umformulieren erhalten wir

$$y[k] = \sum_{i=-\infty}^{\infty} x[i] \cdot h[k - i] \quad \text{für alle } k. \quad (6.5)$$

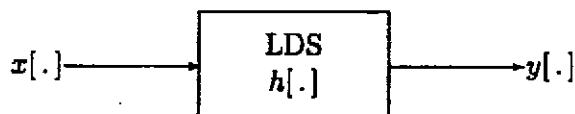
Die Summe in (6.5) heisst **Faltungssumme**. Wir werden (6.4) symbolisch als

$$y[.] = x[.] * h[.] \quad (6.6)$$

schreiben und sagen, dass  $y[.]$  die **Faltung** von  $x[.]$  und  $h[.]$  ist. Es ist leicht zu beweisen (setze dazu  $j = k - i$  in (6.5)), dass

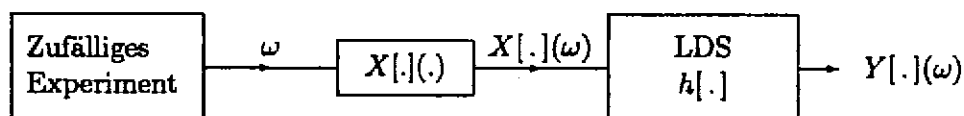
$$h[.] * x[.] = x[.] * h[.] \quad (6.7)$$

gilt, d.h. ein Beobachter des Ausgangssignals  $y[.]$  im folgenden Bild wird es nicht merken, wenn  $x[.]$  und  $h[.]$  vertauscht werden.



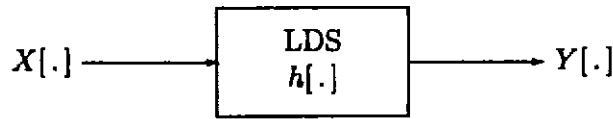
### 6.3 Stochastische Prozesse als Eingang eines LDS

Wir interessieren uns jetzt für folgende Anordnung:





Vereinfacht zeichnet man diese Anordnung üblicherweise wie folgt:



Man darf aber nicht vergessen, dass jeder Versuch des zufälligen Experiments (ein Versuch dauert ewig!) eine und nur eine Eingangssequenz  $x[.]$  für das LDS liefern wird. Man spricht manchmal von einem "stochastischen Eingangssignal"  $X[.]$ . Durch den gewählten Sprachgebrauch sollte man sich aber nicht verwirren lassen.

Aus (6.5) und (6.7) sehen wir, falls  $X[.]$  den Wert  $x[.]$  annimmt, dass  $Y[.]$  den Wert  $y[.]$  mit

$$y[k] = \sum_{i=-\infty}^{\infty} x[k-i] \cdot h[i] \quad \text{für alle } k \quad (6.8)$$

annimmt. Deswegen können wir schreiben

$$Y[k] = \sum_{i=-\infty}^{\infty} X[k-i] \cdot h[i] \quad \text{für alle } k. \quad (6.9)$$

Sei  $X[.]$  schwach stationär, dann ist

$$\begin{aligned} E[Y[k]] &= \sum_{i=-\infty}^{\infty} E[X[k-i]] \cdot h[i] \\ &= m_X \sum_{i=-\infty}^{\infty} h[i] \end{aligned}$$

unabhängig von  $k$ . Ferner ist

$$\begin{aligned} E[Y[k] \cdot Y[k+i]] &= E \left[ \sum_{m=-\infty}^{\infty} X[k-m] \cdot h[m] \sum_{n=-\infty}^{\infty} X[k+i-n] \cdot h[n] \right] \\ &= E \left[ \sum_{m=-\infty}^{\infty} \sum_{n=-\infty}^{\infty} X[k+i-n] \cdot h[n] \cdot X[k-m] \cdot h[m] \right] \\ &= \sum_{m=-\infty}^{\infty} \sum_{n=-\infty}^{\infty} h[m] \cdot h[n] \cdot E[X[k-m] \cdot X[k+i-n]] \\ &= \sum_{m=-\infty}^{\infty} \sum_{n=-\infty}^{\infty} h[m] \cdot h[n] \cdot R_X[i-n+m] \end{aligned}$$

ebenfalls unabhängig von  $k$ .

Wir fassen zusammen:

Sei das Eingangssignal eines LDS mit Gewichtssequenz  $h[\cdot]$  ein schwach stationärer stochastischer Prozess  $X[\cdot]$ , dann ist das Ausgangssignal  $Y[\cdot]$  auch schwach stationär, und zwar mit dem Mittelwert

$$m_Y = m_X \sum_{i=-\infty}^{\infty} h[i] \quad (6.10)$$

und mit der Autokorrelationssequenz

$$R_Y[k] = \sum_{m=-\infty}^{\infty} \sum_{n=-\infty}^{\infty} h[m] \cdot h[n] \cdot R_X[k - n + m] \quad \text{für alle } k. \quad (6.11)$$

Zwei stochastische Prozesse  $X[\cdot]$  und  $Y[\cdot]$  heißen **gemeinsam schwach stationär**, wenn

- (1) sowohl  $X[\cdot]$  als auch  $Y[\cdot]$  schwach stationär sind und
- (2)  $E[X[k] \cdot Y[k + i]]$  für alle  $i$  unabhängig von  $k$  ist.

Man schreibt dann

$$R_{XY}[i] = E[X[k] \cdot Y[k + i]] \quad \text{für alle } i,$$

und nennt  $R_{XY}[\cdot]$  die **Kreuzkorrelationssequenz** von  $X[\cdot]$  und  $Y[\cdot]$ .

Betrachten wir wieder den Fall, wo das Eingangssignal eines LDS ein s.s. Prozess  $X[\cdot]$  ist. Wir wissen schon, dass das Ausgangssignal  $Y[\cdot]$  auch ein s.s. Prozess ist. Ferner sehen wir, dass

$$\begin{aligned} E[X[k] \cdot Y[k + i]] &= E \left[ X[k] \sum_{m=-\infty}^{\infty} X[k + i - m] \cdot h[m] \right] \\ &= \sum_{m=-\infty}^{\infty} h[m] \cdot E[X[k] \cdot X[k + i - m]] \\ &= \sum_{m=-\infty}^{\infty} h[m] \cdot R_X[i - m] \end{aligned}$$

unabhängig von  $k$  ist, wobei wir auch bemerken, dass die letzte Summe genau der Faltungssumme entspricht.

Wir fassen zusammen:

Sei das Eingangssignal eines LDS mit Gewichtssequenz  $h[\cdot]$  ein schwach stationärer Prozess  $X[\cdot]$  und sei  $Y[\cdot]$  das Ausgangssignal, dann sind  $X[\cdot]$  und  $Y[\cdot]$  gemeinsam schwach stationär und zwar mit

$$R_{XY}[k] = \sum_{i=-\infty}^{\infty} h[i] \cdot R_X[k-i] \quad \text{für alle } k, \quad (6.12)$$

d.h.

$$R_{XY}[\cdot] = h[\cdot] * R_X[\cdot].$$

Der Leser soll sich selber überzeugen, dass (6.11) äquivalent als

$$R_Y[\cdot] = h[-\cdot] * h[\cdot] * R_X[\cdot]$$

geschrieben werden kann.

## 6.4 Die z-Transformation

Sei  $f[\cdot]$  eine zweiseitige, reelle oder komplexe Sequenz mit

$$\sum_{k=-\infty}^{\infty} |f[k]| < \infty, \quad (6.13)$$

dann heisst die analytische Funktion  $F$  einer komplexen Variablen  $z$ , die für  $z$  auf dem Einheitskreis in der komplexen Ebene durch

$$F(z) = \sum_{k=-\infty}^{\infty} f[k] \cdot z^{-k} \quad (6.14)$$

definiert ist, die z-Transformierte von  $f[\cdot]$ . Es ist leicht zu beweisen (vgl. Übungsaufgabe), dass

$$f[k] = \frac{1}{2\pi} \int_0^{2\pi} F(e^{j\Omega}) \cdot e^{jk\Omega} d\Omega \quad \text{für alle } k \quad (6.15)$$

die inverse z-Transformation ist. Die komplexe Zahl  $z = e^{j\Omega}$  liegt auf dem Einheitskreis in der komplexen Ebene. Man sieht aus (6.15), dass die Werte von  $F(z)$  für  $z$  auf dem Einheitskreis genügen, um sowohl  $f[\cdot]$  als auch  $F(z)$  vollständig zu bestimmen.

Ist  $h[\cdot]$  die Gewichtssequenz eines LDS, dann nennen wir  $H(z)$  die Übertragungsfunktion des LDS. Wir nennen den Einheitskreis in der komplexen Ebene den Frequenzkreis und wir nennen  $\Omega$  für einen Punkt  $z = e^{j\Omega}$  auf diesem Kreis, die normierte Frequenz. Wir nennen  $H(e^{j\Omega})$ , betrachtet als Funktion von  $\Omega$  für  $0 \leq \Omega < 2\pi$ , den Frequenzgang des LDS.

Sei das Eingangssignal  $x[\cdot]$  eines LDS mit der Übertragungsfunktion  $H(z)$  eine Sequenz, die (6.13) befriedigt (was fast nie der Fall ist, wenn  $x[\cdot]$  der Wert eines stochastischen Prozesses  $X[\cdot]$  ist), dann folgt aus (6.5) für das Ausgangssignal  $y[\cdot]$ , dass

$$\begin{aligned}
\sum_{k=-\infty}^{\infty} y[k]z^{-k} &= \sum_{k=-\infty}^{\infty} \sum_{i=-\infty}^{\infty} x[i] \cdot h[k-i] \cdot z^{-k} \\
&= \sum_{i=-\infty}^{\infty} \sum_{k=-\infty}^{\infty} x[i] \cdot z^{-i} \cdot h[k-i] \cdot z^{-k+i} \\
&= \sum_{i=-\infty}^{\infty} x[i] \cdot z^{-i} \sum_{k=-\infty}^{\infty} h[k-i] \cdot z^{-k+i} \\
&= X(z) \cdot H(z)
\end{aligned}$$

gilt, d.h.

$$y[\cdot] = x[\cdot] * h[\cdot] \Rightarrow Y(z) = X(z) \cdot H(z) \quad \text{für alle komplexen Zahlen } z. \quad (6.16)$$

Obwohl der Wert  $x[\cdot]$  eines s.s. stochastischen Prozesses  $X[\cdot]$  im Allgemeinen keine  $z$ -Transformierte hat, wird die Autokorrelationsfunktion  $R_X[\cdot]$  (betrachtet als eine Zeitsequenz) normalerweise (6.13) befriedigen und darum eine  $z$ -Transformierte haben. Wir schreiben  $S_X(\cdot)$  für die  $z$ -Transformierte von  $R_X[\cdot]$ , d.h.

$$S_X(z) = \sum_{k=-\infty}^{\infty} R_X[k] \cdot z^{-k} \quad (6.17)$$

für alle  $z \in \mathbb{C}$ , für welche obige Summe absolut konvergiert.

Nehmen wir jetzt an, dass der s.s. Prozess  $X[\cdot]$  das Eingangssignal eines LDS mit der Übertragungsfunktion  $H(z)$  ist und dass  $S_X(z)$  existiert. Dann finden wir aus (6.11), dass

$$\begin{aligned}
\sum_{k=-\infty}^{\infty} R_Y[k] \cdot z^{-k} &= \sum_{k=-\infty}^{\infty} \sum_{m=-\infty}^{\infty} \sum_{n=-\infty}^{\infty} h[m] \cdot h[n] \cdot R_X[k-n+m] \cdot z^{-k} \\
&= \sum_{m=-\infty}^{\infty} \sum_{n=-\infty}^{\infty} \sum_{k=-\infty}^{\infty} h[m] \cdot z^{+m} \cdot h[n] \cdot z^{-n} \cdot R_X[k-n+m] \cdot z^{-k+n-m} \\
&= \sum_{m=-\infty}^{\infty} h[m] \cdot z^{+m} \sum_{n=-\infty}^{\infty} h[n] \cdot z^{-n} \sum_{k=-\infty}^{\infty} R_X[k-n+m] \cdot z^{-k+n-m} \\
&= H(z^{-1}) \cdot H(z) \cdot S_X(z)
\end{aligned}$$

gilt.

Wir fassen zusammen:

Sei das Eingangssignal eines LDS mit Übertragungsfunktion  $H(z)$  ein schwach stationärer stochastischer Prozess  $X[\cdot]$  mit existierendem  $S_X(\cdot)$ . Dann gilt:

$$S_Y(z) = H(z) \cdot H(z^{-1}) \cdot S_X(z) \quad \text{für alle } z \in C. \quad (6.18)$$

Insbesondere gilt auch:

$$S_Y(e^{j\Omega}) = |H(e^{j\Omega})|^2 \cdot S_X(e^{j\Omega}), \quad \text{für } 0 \leq \Omega < 2\pi. \quad (6.19)$$

Die Gleichung (6.18) folgt aus der Tatsache, dass mit  $z = e^{j\Omega}$  auch  $z^{-1} = e^{-j\Omega} = \bar{z}$  gilt, was  $H(e^{-j\Omega}) = \overline{H(e^{j\Omega})}$  entspricht.

Aus der Definition von  $R_Y[\cdot]$  bemerken wir jetzt, dass

$$R_Y[0] = E[Y[k] \cdot Y[k]] = E[Y^2[k]]$$

gilt, d.h., dass  $R_Y[0]$  die mittlere Leistung von  $Y[\cdot]$  ist. Wir sehen auch aus (6.15), dass

$$R_Y[0] = \frac{1}{2\pi} \int_0^{2\pi} S_Y(e^{j\Omega}) d\Omega$$

ist. Es folgt dann, dass

$$E[Y^2[k]] = \frac{1}{2\pi} \int_0^{2\pi} |H(e^{j\Omega})|^2 \cdot S_X(e^{j\Omega}) d\Omega \quad (6.20)$$

gilt.

Wollen wir die mittlere Leistung von  $X[\cdot]$  in einem gewissen normierten Frequenzband messen, dann werden wir  $X[\cdot]$  als Eingangssignal eines LDS mit  $H(e^{j\Omega}) = 1$  (resp. 0) für  $\Omega$  innerhalb (resp. nicht innerhalb) dieses Frequenzbandes einspeisen und die mittlere Leistung des Ausgangssignals  $Y[\cdot]$  messen. Aus (6.19) sehen wir, dass diese mittlere Leistung dem Integral von  $\frac{1}{2\pi} S_X(e^{j\Omega})$  über diesem Frequenzband entspricht. Wir haben damit folgenden Satz bewiesen.

**Wiener-Khintchine-Satz:**

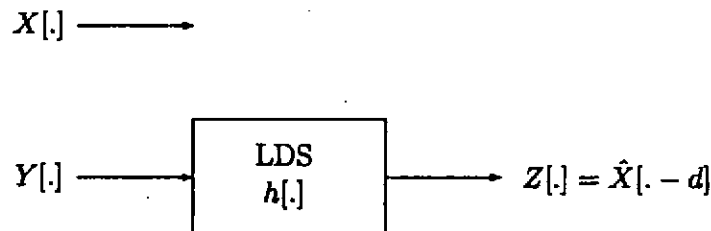
Die Leistungsdichte eines schwach stationären Prozesses  $X[\cdot]$  wird durch die Funktion  $\frac{1}{2\pi} S_X(e^{j\Omega})$  im Definitionsbereich  $0 \leq \Omega < 2\pi$  beschrieben.  $S_X(\cdot)$  ist dabei die  $z$ -Transformierte der Autokorrelationssequenz  $R_X[\cdot]$ .

Dieser Satz liefert uns die Interpretation von  $S_X(\cdot)$ .

## 7 FILTERUNG STOCHASTISCHER SIGNALE

### 7.1 FIR-Filter und die Wiener-Hopf-Gleichung

Mathematisches Modell:



wobei

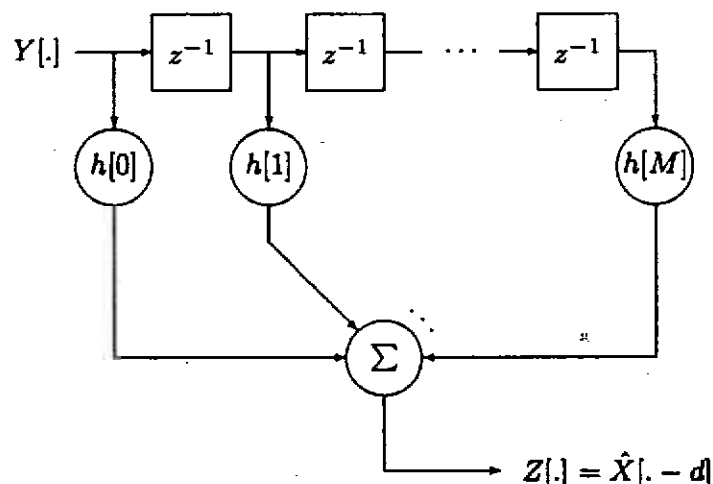
- (1)  $X[.]$  und  $Y[.]$  gemeinsam schwach-stationäre Prozesse sind,
- (2)  $X[.]$  das Signal ist, welches wir schätzen wollen,
- (3)  $Y[.]$  der Prozess ist, den wir beobachten,
- (4) das LDS mit Gewichtssequenz  $h[.]$  unser wählbares Schätzungsfilter ist und
- (5)  $d$  die Verzögerung ist, die für unsere Schätzung erlaubt ist.

Sei  $d = 0$ , dann spricht man von einem **Filterungsproblem**. Für  $d > 0$  nennt man das Problem ein **Glättungsproblem** ("smoothing problem") und für  $d < 0$  hat man es mit einem **Vorhersageproblem** ("prediction problem") zu tun. Für alle  $k$  gilt:  $Z[k] = \hat{X}[k - d]$ , wobei  $Z[.]$  der Prozess am Ausgang unseres Filters ist. Unser Ziel ist es, dass  $\hat{X}[k] = Z[k + d]$  die lineare MMSE-Schätzung von  $X[k]$  für alle  $k$  ist. Im Allgemeinen lassen wir jedoch nicht beliebige Gewichtssequenzen  $h[.]$  zu, das heisst, wir schränken die Menge der erlaubten Sequenzen ein.

Ein LDS mit Gewichtssequenz  $h[.]$  heisst **kausal**, falls  $h[k] = 0$  für alle  $k < 0$  gilt. Wollen wir unser Filter in Echtzeit benützen, dann müssen wir eines mit einer kausalen Gewichtssequenz wählen. Ein kausales Filter mit  $h[M] \neq 0$  und  $h[k] = 0$  für alle  $k > M$ , heisst **FIR-Filter  $M$ -ter Ordnung** (FIR = "finite impulse response". Diese Terminologie ist zwar allgemein gebräuchlich, gibt aber immer wieder zu Missverständnissen Anlass.)  $M$  ist die Länge des Gedächtnisses des FIR-Filters. Ist unser Schätzungsfilter ein FIR-Filter  $M$ -ter Ordnung, dann gilt:

$$Z[k] = \hat{X}[k - d] = \sum_{i=0}^M h[i] \cdot Y[k - i] \quad \text{für alle } k. \quad (7.1)$$

Aus (7.1) folgt, dass ein FIR-Schätzungsfilter wie folgt realisiert werden kann:



wobei die Kästchen mit  $z^{-1}$  eine Verzögerung des Signals um eine Zeiteinheit modellieren. (Eine Verzögerungseinheit wird mit  $z^{-1}$  bezeichnet, weil  $x[.] = y[. - 1]$  durch z-Transformation in die Gleichung  $X(z) = z^{-1} \cdot Y(z)$  übergeht.)

**Bemerkung:** Sei  $m_X \neq 0$  oder/und  $m_Y \neq 0$ , dann wissen wir, dass es von Vorteil ist, die Beobachtung  $Y[.]$  durch  $Y'[.] = Y[.] - m_Y$  zu ersetzen, dann  $X'[.] = X[.] - m_X$  zu schätzen, um schliesslich  $\hat{X}[.] = \hat{X}'[.] + m_X$  als Resultat zu erhalten. Der Einfachheit halber werden wir diese Verbesserung nicht explizit betrachten, sie sollte aber bei der Anwendung nicht vergessen werden!

FIR-Filter sind wahrscheinlich die für die Anwendung wichtigsten zeitinvarianten Schätzungsfilter. Zudem ist es sehr einfach, das optimale FIR-Schätzungsfilter zu finden. Aus (7.1) folgt sofort, dass die Schätzung  $\hat{X}[.]$  von  $X[.]$  eine Linearkombination von  $Y[k], Y[k-1], \dots, Y[k-M]$  ist. Das Orthogonalitätsprinzip sagt aus, dass die lineare Schätzung in diesem Fall genau dann optimal ist, falls der Fehler  $X[k-d] - \hat{X}[k-d]$  orthogonal zu jedem Datum  $Y[k-i]$  ist, das heisst, genau dann, falls:

$$E \left[ Y[k-i] \cdot \left( X[k-d] - \hat{X}[k-d] \right) \right] = 0 \quad \text{für } i = 0, 1, \dots, M$$

gilt. Diese Gleichung lässt sich folgendermassen umschreiben:

$$E \left[ Y[k-i] \cdot \hat{X}[k-d] \right] = R_{YX}[i-d] \quad \text{für } i = 0, 1, \dots, M \quad (7.2)$$

Mit Hilfe von (7.1) sehen wir ferner, dass gilt:

$$\begin{aligned} E \left[ Y[k-i] \cdot \hat{X}[k-d] \right] &= E \left[ Y[k-i] \sum_{j=0}^M h[j] \cdot Y[k-j] \right] \\ &= \sum_{j=0}^M h[j] \cdot R_Y[i-j]. \end{aligned}$$

Wir haben somit folgendes bewiesen:

Das FIR-Filter  $h[\cdot]$  mit Gedächtnislänge  $M$  (oder weniger) liefert zum Zeitpunkt  $k$  genau dann die lineare MMSE-Schätzung  $\hat{X}[k-d]$  von  $X[k-d]$  für die Beobachtung  $Y[k], Y[k-1], \dots, Y[k-M]$ , falls  $h[\cdot]$  die Gleichung

$$\sum_{j=0}^M R_Y[i-j] \cdot h[j] = R_{YX}[i-d] \quad \text{für } i = 0, 1, \dots, M \quad (7.3)$$

befriedigt.

Die Gleichung (7.3) ist die Wiener-Hopf-Gleichung für das FIR-Schätzungsproblem. Sie lässt sich in Matrixform folgendermassen ausdrücken:

$$\begin{bmatrix} R_Y[0] & R_Y[-1] & \dots & R_Y[-M] \\ R_Y[1] & R_Y[0] & \dots & R_Y[1-M] \\ \vdots & \vdots & & \vdots \\ R_Y[M] & R_Y[M-1] & \dots & R_Y[0] \end{bmatrix} \cdot \begin{bmatrix} h[0] \\ h[1] \\ \vdots \\ h[M] \end{bmatrix} = \begin{bmatrix} R_{YX}[-d] \\ R_{YX}[1-d] \\ \vdots \\ R_{YX}[M-d] \end{bmatrix}. \quad (7.4)$$

Da  $R_Y[i] = R_Y[-i]$  gilt, ist die Matrix in (7.4) symmetrisch.

Aus dem Orthogonalitätsprinzip folgt, dass der minimale MSE wie folgt berechnet werden kann:

$$\begin{aligned} \text{MSE} &= E \left[ (X[k-d] - \hat{X}[k-d]) \cdot X[k-d] \right] \\ &= R_X[0] - E \left[ \hat{X}[k-d] \cdot X[k-d] \right] \\ &= R_X[0] - E \left[ Z[k] \cdot X[k-d] \right] \\ &= R_X[0] - E \left[ \sum_{i=0}^M h[i] \cdot Y[k-i] \cdot X[k-d] \right] \\ &= R_X[0] - \sum_{i=0}^M h[i] \cdot R_{YX}[i-d], \end{aligned}$$

wobei wir auch (7.1) benützt haben. Ebenfalls vom Orthogonalitätsprinzip wissen wir, dass

$$\begin{aligned} \text{MSE} &= E \left[ X^2[k-d] \right] - E \left[ \hat{X}^2[k-d] \right] \\ &= R_X[0] - E \left[ Z^2[k] \right] \\ &= R_X[0] - \frac{1}{2\pi} \int_0^{2\pi} |H(e^{j\Omega})|^2 \cdot S_Y(e^{j\Omega}) \, d\Omega \end{aligned}$$

gilt, wobei wir (6.20) benützt haben.



Sei das FIR-Filter  $h[\cdot]$  mit der Gedächtnislänge  $M$  (oder weniger) dasjenige Filter, welches zum Zeitpunkt  $k$  die lineare MMSE-Schätzung von  $X[k-d]$  für die Beobachtung  $Y[k], Y[k-1], \dots, Y[k-M]$  liefert. Dann gilt sowohl

$$\text{MSE} = R_X[0] - \sum_{i=0}^M h[i] \cdot R_{YX}[i-d] \quad (7.5)$$

als auch

$$\text{MSE} = R_X[0] - \frac{1}{2\pi} \int_0^{2\pi} |H(e^{j\Omega})|^2 \cdot S_Y(e^{j\Omega}) d\Omega. \quad (7.6)$$

**Beispiel 1:** Sei

$$Y[k] = X[k] + V[k],$$

wobei  $X[\cdot]$  und  $V[\cdot]$  unabhängige Prozesse sind.  $V[\cdot]$  sei ein weisses Rauschen mit  $R_V[\cdot] = L \cdot \delta[\cdot]$  und  $X[\cdot]$  sei s.s. mit  $m_X = 0$  und  $R_X[k] = \left(\frac{1}{2}\right)^{|k|}$  für alle  $k$ . Wir möchten das FIR-Filter  $h[\cdot]$  mit  $M = 2$  (oder weniger) finden, das zum Zeitpunkt  $k$  die lineare MMSE-Schätzung von  $X[k-1]$  liefert (das heisst  $d = 1$ , also eine um eine Zeiteinheit verzögerte Schätzung). Man kann zeigen (vgl. Übungsaufgabe), dass für alle  $k$  gilt:

$$\begin{aligned} R_Y[k] &= R_X[k] + R_V[k] \\ &= \left(\frac{1}{2}\right)^{|k|} + L \cdot \delta[k]. \end{aligned}$$

Ferner sehen wir, dass

$$R_{YX}[k] = R_X[k] = \left(\frac{1}{2}\right)^{|k|}$$

für alle  $k$  gilt. Die Wiener-Hopf-Gleichung (7.4) liefert:

$$\begin{bmatrix} L+1 & \frac{1}{2} & \frac{1}{4} \\ \frac{1}{2} & L+1 & \frac{1}{2} \\ \frac{1}{4} & \frac{1}{2} & L+1 \end{bmatrix} \cdot \begin{bmatrix} h[0] \\ h[1] \\ h[2] \end{bmatrix} = \begin{bmatrix} \frac{1}{2} \\ 1 \\ \frac{1}{2} \end{bmatrix}.$$

Wir nehmen nun an, dass  $L = 1$  ist (das heisst  $E[V^2[k]] = L = 1$ ). Die Lösung der Wiener-Hopf-Gleichung liefert dann:

$$h[0] = \frac{1}{8}, \quad h[1] = \frac{7}{16}, \quad h[2] = \frac{1}{8}.$$

Der resultierende MMSE lässt sich aus (7.5) als

$$\begin{aligned} \text{MSE} &= R_X[0] - \sum_{i=0}^2 h[i] \cdot R_{YX}[i-1] \\ &= 1 - \frac{1}{8} R_{YX}[-1] - \frac{7}{16} R_{YX}[0] - \frac{1}{8} R_{YX}[1] \\ &= 1 - \frac{1}{8} \cdot \frac{1}{2} - \frac{7}{16} \cdot 1 - \frac{1}{8} \cdot \frac{1}{2} = \frac{7}{16} \end{aligned}$$

bestimmen.

○

## 7.2 Das nichtkausale Wiener-Filter

Jetzt betrachten wir den Fall, wo es keine Beschränkung der  $h[\cdot]$  gibt, d.h.  $h[\cdot]$  darf nicht-kausal sein. Dieser Fall ist im Prinzip besonders einfach. Vom Orthogonalitätsprinzip wissen wir, dass die lineare Schätzung  $\hat{X}[k-d]$  genau dann optimal ist, wenn der Fehler orthogonal zu  $Y[k-i]$  ist, für jedes  $i$  mit  $-\infty < i < \infty$ , d.h.

$$E[Y[k-i](X[k-d] - \hat{X}[k-d])] = 0 \quad \text{für } -\infty < i < \infty.$$

Wir schreiben diese Gleichung als

$$E[Y[k-i]\hat{X}[k-d]] = R_{YX}[i-d] \quad \text{für } -\infty < i < \infty \quad (7.7)$$

um und erkennen dadurch die Ähnlichkeit zu (7.2).

Es gilt ferner

$$\begin{aligned} E[Y[k-i]\hat{X}[k-d]] &= E\left[Y[k-i] \sum_{j=-\infty}^{+\infty} h[j] \cdot Y[k-j]\right] \\ &= \sum_{j=-\infty}^{+\infty} h[j] \cdot R_Y[i-j]. \end{aligned}$$

Damit folgt, dass (7.7) als

$$\sum_{j=-\infty}^{+\infty} h[j] \cdot R_Y[i-j] = R_{YX}[i-d] \quad \text{für } -\infty < i < \infty \quad (7.8)$$

oder noch durchsichtiger als

$$h[\cdot] * R_Y[\cdot] = R_{YX}[\cdot - d] \quad (7.8')$$

geschrieben werden kann. Die Gleichung (7.8) [oder äquivalent (7.8')] heisst die nicht-kausale Wiener-Hopf-Gleichung.

Wenden wir nun die  $z$ -Transformation auf (7.8') an, so bekommen wir:

$$H(z) \cdot S_Y(z) = z^{-d} \cdot S_{YX}(z).$$

Wir fassen zusammen:

Das LDS mit der Übertragungsfunktion

$$H(z) = z^{-d} \frac{S_{YX}(z)}{S_Y(z)} \quad (7.9)$$

ist das nicht-kausale Wiener-Filter, das zum Zeitpunkt  $k$  die lineare MMSE-Schätzung von  $X[k-d]$  für die Beobachtung  $Y[k-i]$ ,  $-\infty < i < \infty$ , liefert.

Das folgende Resultat wird bei den Anwendungen nützlich sein.

Sei  $\alpha$  eine komplexe Zahl mit  $0 < |\alpha| < 1$  und sei  $f[k] = \alpha^{|k|}$  für  $-\infty < k < \infty$ , dann gilt für die z-Transformierte von  $f[\cdot]$ :

$$F(z) = -\frac{(\alpha^{-1} - \alpha)z}{z^2 - (\alpha^{-1} + \alpha)z + 1} \quad (7.10)$$

**Beweis:**

$$\begin{aligned} F(z) &= \sum_{k=-\infty}^{+\infty} f[k] \cdot z^{-k} \\ &= \sum_{k=-\infty}^{+\infty} \alpha^{|k|} \cdot z^{-k} \\ &= 1 + \sum_{k=1}^{+\infty} \alpha^k \cdot z^{-k} + \sum_{k=-1}^{-\infty} \alpha^{-k} \cdot z^{-k} \\ &= 1 + \sum_{k=1}^{+\infty} (\alpha z^{-1})^k + \sum_{i=1}^{+\infty} (\alpha z)^i \\ &= 1 + \frac{\alpha z^{-1}}{1 - \alpha z^{-1}} + \frac{\alpha z}{1 - \alpha z} \\ &= 1 + \frac{\alpha}{z - \alpha} - \frac{z}{z - \alpha^{-1}} \end{aligned}$$

Nach ein bisschen mehr Algebra kommt man schliesslich auf (7.10). □

**Beispiel 2:** Für dieselbe Anordnung wie in Beispiel 1 erhalten wir

$$\begin{aligned} S_{YX}(z) &= S_X(z) \\ &= -\frac{(2 - \frac{1}{2})z}{z^2 - (2 + \frac{1}{2})z + 1} \\ &= -\frac{\frac{3}{2}z}{z^2 - \frac{5}{2}z + 1} \end{aligned}$$

wobei wir (7.10) mit  $\alpha = \frac{1}{2}$  benützt haben. Ferner haben wir

$$\begin{aligned} S_Y(z) &= S_X(z) + S_V(z) \\ &= S_X(z) + 1 \\ &= \frac{z^2 - 4z + 1}{z^2 - \frac{5}{2}z + 1} \end{aligned}$$

Es folgt nun aus (7.9) mit  $d = 1$ , dass

$$H(z) = z^{-1} \frac{(-\frac{3}{2}z)}{z^2 - 4z + 1}$$

die Übertragungsfunktion des nicht-kausalen Wiener-Filters ist. Zunächst stellen wir fest, dass  $\alpha + \alpha^{-1} = 4$  mit  $|\alpha| < 1$  den Werten  $\alpha = 2 - \sqrt{3}$ ,  $\alpha^{-1} = 2 + \sqrt{3}$  und  $\alpha^{-1} - \alpha = 2\sqrt{3}$  entspricht.

Weil

$$H(z) = z^{-1} \left( \frac{\sqrt{3}}{4} \right) \frac{(-2\sqrt{3}z)}{z^2 - 4z + 1}$$

ist, sehen wir nun aus (7.10), dass

$$h[k] = \frac{\sqrt{3}}{4} (2 - \sqrt{3})^{|k-1|} \quad -\infty < k < \infty$$

gilt. Es ist interessant, die Werte  $h[0] = .116$ ,  $h[1] = .432$ ,  $h[2] = .116$  des nicht-kausalen Wiener-Filters mit den entsprechenden Werten  $\frac{1}{8} = .125$ ,  $\frac{7}{16} = .438$ ,  $\frac{1}{8} = .125$  des FIR-Filters aus Beispiel 1 zu vergleichen. Die anderen Werte im nicht-kausalen Wiener-Filter fallen sehr schnell ab. In diesem Fall ist das einfache FIR-Filter fast so gut wie das hochkomplizierte (nicht machbare) nicht-kausale Wiener-Filter. ○

### 7.3 Das kausale Wiener-Filter

Als unser letztes Problem der Wiener-Filterung betrachten wir den Fall, wo wir  $h[\cdot]$  nur insoweit beschränken, als dass das Filter kausal sein soll, d.h. wir werden das kausale Wiener-Filter finden. Dasselbe Vorgehen, welches uns die Gleichungen (7.2) und (7.7) lieferte, gibt uns jetzt die kausale Wiener-Hopf-Gleichung

$$\sum_{j=0}^{+\infty} h[j] R_Y[i-j] = R_{YX}[i-d] \quad \text{für } 0 \leq i < \infty. \quad (7.11)$$

Diese Gleichung ist viel leichter aufzuschreiben als zu lösen! Der Mathematiker Norbert Wiener ein Verfahren entwickelte, (7.11) zu lösen, jedoch war dieses sehr schwierig zu verstehen und zu benützen. Den beiden Ingenieuren Heinrich Bode und Claude Shannon gelang es jedoch, einen Trick zu finden, der dieses Verfahren im Verständnis und in der Anwendung unheimlich viel einfacher macht.

Zu einer beliebigen Gewichtssequenz  $h[\cdot]$  definieren wir den kausalen Teil  $h_{KT}[\cdot]$  wie folgt:

$$h_{KT}[k] = \begin{cases} h[k] & \text{für } k \geq 0 \\ 0 & \text{für } k < 0. \end{cases}$$

**Bode-Shannon-Trick (erste Hälfte):** Sei der beobachtete Prozess  $Y[\cdot]$  weisses Rauschen, dann ist das kausale Wiener-Filter genau der kausale Teil  $h_{KT}[\cdot]$  des nicht-kausalen Wiener-Filters  $h[\cdot]$ .

**Beweis:** Sei  $Y[\cdot]$  weisses Rauschen, dann gilt:

$$\begin{aligned} E\{[Y[i]Y[j]]\} &= R_Y[j-i] \\ &= L\delta[j-i] = 0 \quad \text{für } i \neq j, \end{aligned}$$

d.h.  $Y[i]$  und  $Y[j]$  sind orthogonal, wenn  $i \neq j$  ist. Insbesondere ist jede Zufallsgrösse unter  $Y[k], Y[k-1], Y[k-2], \dots$  orthogonal zu jeder Zufallsgrösse unter  $Y[k+1], Y[k+2], Y[k+3], \dots$ . Aber für das nicht-kausale Wiener-Filter gilt:

$$\begin{aligned} \hat{X}[k-d] = Z[k] &= \sum_{i=-\infty}^{+\infty} h[i] \cdot Y[k-i] \\ &= \sum_{i=0}^{+\infty} h_{KT}[i] \cdot Y[k-i] + \sum_{j=1}^{+\infty} h[-j] \cdot Y[k+j]. \end{aligned}$$

Es folgt dann direkt aus der Trennungseigenschaft der linearen MMSE-Schätzung, dass

$$\sum_{i=0}^{+\infty} h_{KT}[i] \cdot Y[k-i] = \sum_{i=-\infty}^{+\infty} h_{KT}[i] \cdot Y[k-i]$$

die lineare MMSE-Schätzung von  $X[k-d]$  für die Beobachtung  $Y[k], Y[k-1], Y[k-2], \dots$  ist.  $\square$

Leider handelt es sich beim beobachteten Prozess  $Y[\cdot]$  normalerweise nicht um weisses Rauschen. Die kreative Idee von Bode und Shannon war,  $Y[\cdot]$  zu weissem Rauschen umzuformen. Das Filter  $g[\cdot]$  heisst ein "Whitening-Filter" für den s.s. Prozess  $Y[\cdot]$ , wenn

- (1)  $g[\cdot]$  kausal ist,
- (2) das inverse Filter mit der Übertragungsfunktion  $1/G(z)$  auch kausal ist und
- (3)  $S_Y(z) \cdot G(z) \cdot G(z^{-1}) = 1$  (7.12)

gilt. Aus (6.18) und (7.12) sehen wir, dass das Ausgangssignal des Whitening-Filters weisses Rauschen ist, falls  $Y[\cdot]$  das Eingangssignal ist.

Die folgende Tatsache (für den Beweis vgl. Übungsaufgabe) wird uns helfen, ein Whitening-Filter zu finden:

Eine rationale  $z$ -Transformierte entspricht genau dann einem kausalen Filter  $g[\cdot]$ , wenn alle Pole von  $G(z)$  im Inneren des Einheitskreises der komplexen Ebene liegen und auch  $\lim_{z \rightarrow \infty} G(z)$  existiert. Deswegen sind sowohl  $G(z)$  als auch  $1/G(z)$  genau dann Übertragungsfunktionen von kausalen Filtern, wenn alle Pole und alle Nullstellen von  $G(z)$  im Inneren des Einheitskreises liegen und wenn ferner  $\lim_{z \rightarrow \infty} G(z)$  existiert und nicht 0 ist.

**Beispiel 3:**  $S_Y(z)$  aus Beispiel 2 lässt sich wie folgt zerlegen:

$$S_Y(z) = \frac{z^2 - 4z + 1}{z^2 - \frac{5}{2}z + 1} = \frac{(z - (2 - \sqrt{3}))(z - (2 + \sqrt{3}))}{(z - \frac{1}{2})(z - 2)}$$

Weil  $2 - \sqrt{3}$  und  $\frac{1}{2}$  im Inneren des Einheitskreises liegen, folgt, dass die Übertragungsfunktion des Whitening-Filter für  $Y[\cdot]$  gleich

$$G(z) = c \frac{z - \frac{1}{2}}{z - (2 - \sqrt{3})}$$

ist. Da

$$\lim_{z \rightarrow \infty} G(z) \cdot G(z^{-1}) = c^2 \frac{1 \cdot (-\frac{1}{2})}{1 - (2 - \sqrt{3})} = c^2 \frac{2 + \sqrt{3}}{2}$$

und

$$\lim_{z \rightarrow \infty} S_Y(z) = 1$$

gelten, folgt aus (7.12), dass  $c = \sqrt{\frac{2}{2 + \sqrt{3}}}$  ist.

Deshalb ist

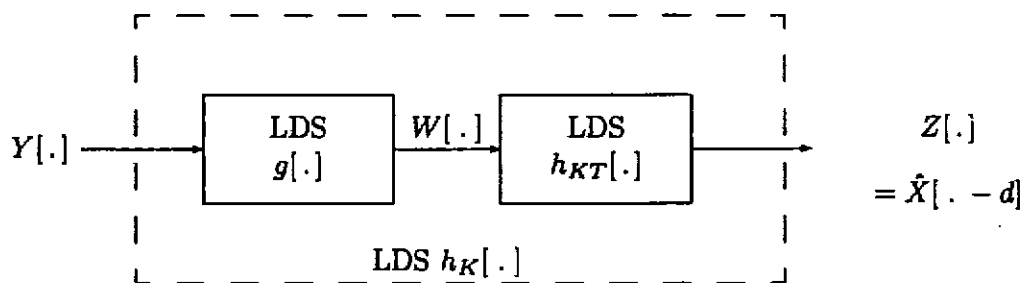
$$G(z) = \sqrt{\frac{2}{2 + \sqrt{3}}} \cdot \frac{z - \frac{1}{2}}{z - (2 - \sqrt{3})}$$

ein Whitening-Filter für  $Y[\cdot]$ .

**Bode-Shannon-Trick (zweite Hälfte):** Sei  $g[\cdot]$  ein Whitening-Filter für den beobachteten Prozess  $Y[\cdot]$ , sei  $h[\cdot]$  das nicht-kausale Wiener-Filter für das resultierende weisse Rauschen  $W[\cdot] = Y[\cdot] * g[\cdot]$  und sei  $h_{KT}[\cdot]$  der kausale Teil von  $h[\cdot]$ , dann ist  $h_K[\cdot] = g[\cdot] * h_{KT}[\cdot]$  das kausale Wiener-Filter für den beobachteten Prozess  $Y[\cdot]$ . Insbesondere gilt für dieses nicht-kausale Filter:

$$H(z) = z^{-d} \cdot G(z^{-1}) \cdot S_{YX}(z) \quad (7.13)$$

In Bildern:



**Beweis:** Zuerst beweisen wir, dass das kausale lineare MMSE Filter, das auf  $W[\cdot]$  operiert, genau so gut ist wie das kausale lineare MMSE-Filter, das auf  $Y[\cdot]$  operiert. Dies folgt direkt aus der Tatsache, dass das Whitening-Filter  $g[\cdot]$  ein kausales inverses Filter hat. Operieren wir zuerst auf  $W[\cdot]$  mit diesem inversen Filter, dann bekommen wir wieder  $Y[\cdot]$ , sodass wir dann

mit dem kausalen Filter für  $Y[\cdot]$  operieren können um dasselbe Ausgangssignal  $Z[\cdot] = \hat{X}[\cdot - d]$  zu erhalten.

Es bleibt nur noch zu zeigen, dass (7.13) gilt. Wir sehen, dass

$$\begin{aligned} R_{WX}[k] &= E[W[i] \cdot X[i+k]] \\ &= E \left[ \sum_{j=-\infty}^{+\infty} g[j] \cdot Y[i-j] \cdot X[i+k] \right] \\ &= \sum_{j=-\infty}^{+\infty} g[j] \cdot R_{YX}[k+j] \\ &= \sum_{i=-\infty}^{+\infty} g[-i] \cdot R_{YX}[k-i] \end{aligned}$$

für alle  $k$  gilt. Somit ist

$$R_{WX}[\cdot] = g[-\cdot] * R_{YX}[\cdot]$$

und auch

$$S_{WX}(z) = G(z^{-1}) \cdot S_{YX}(z).$$

Weil wegen dem Whitening-Filter  $S_W(z) = 1$  ist, folgt nun aus (7.9), dass das nicht-kausale Wiener-Filter für den beobachteten Prozess  $W[\cdot]$  genau

$$H(z) = z^{-d} \cdot G(z^{-1}) \cdot S_{YX}(z)$$

ist. □

**Beispiel 4:** Wir wollen das kausale Wiener-Filter  $h_k[\cdot]$  für den beobachteten Prozess  $Y[\cdot]$  aus Beispiel 1 finden. Aus Beispiel 2 und 3 wissen wir, dass

$$S_{YX}(z) = -\frac{3}{2} \frac{z}{z^2 - \frac{5}{2}z + 1} = -\frac{3}{2} \frac{z}{(z - \frac{1}{2})(z - 2)}$$

und

$$G(z) = \sqrt{\frac{2}{2 + \sqrt{3}}} \cdot \frac{z - \frac{1}{2}}{z - (2 - \sqrt{3})}$$

sind. Wir finden, dass

$$G(z^{-1}) = \sqrt{\frac{2 + \sqrt{3}}{2}} \cdot \frac{z - 2}{z - (2 + \sqrt{3})}$$

ist. Aus (7.13) folgt nun, dass das nicht-kausale Wiener-Filter für das beobachtete weiße Rauschen  $W[\cdot]$

$$H(z) = z^{-1} \left( -\frac{3}{2} \sqrt{\frac{2 + \sqrt{3}}{2}} \right) \frac{z}{(z - \frac{1}{2})(z - (2 + \sqrt{3}))}$$

ist. Es ist im Prinzip einfach (in der Ausführung jedoch etwas mühsam), den kausalen Teil  $h_{KT}[\cdot]$  der entsprechenden Sequenz  $h[\cdot]$  zu finden. Das gewünschte Filter  $h_K[\cdot]$  kann danach als

$$h_K[\cdot] = g[\cdot] * h_{KT}[\cdot]$$

oder als

$$H_K(z) = G(z) \cdot H_{KT}(z)$$

berechnet werden. Wir ersparen dem Leser die entsprechenden, langweiligen Details. ○

#### 7.4 Zeitvariante Filterung: Das Kalman-Filter

Wir modifizieren nun das Filterproblem wesentlich. Bisher standen uns zum Zeitpunkt  $k$  prinzipiell alle Beobachtungen  $Y[k-i]$  für  $i = 0, 1, 2, \dots$  zur Verfügung, und zwar auch dann, wenn wir auf einem kausalen Filter bestanden. Wir haben also bisher angenommen, dass die ganze Vergangenheit zugreifbar ist. Jetzt fordern wir, zusätzlich zur kausalen Beschränkung, dass die Beobachtungen zu einem gewissen Zeitpunkt ( $k = 0$ ) begonnen haben, d.h.

zum Zeitpunkt  $k$  stehen nur die Beobachtungen  
 $Y[0], Y[1], \dots, Y[k]$  zur Verfügung.

Diese neue Annahme bedeutet, dass die Menge der Daten, die wir ausnützen können, um  $X[k]$  zu schätzen, mit  $k$  wächst. Das entsprechende optimale Filter wird dann nicht mehr zeitinvariant sein. Um ein solches Filter praktisch anwenden zu können, muss es rekursiv sein, d.h.

wir wollen die optimale Schätzung von  $X[k-1]$  mit  
Hilfe der Beobachtung  $Y[k]$  zur optimalen Schätzung  
von  $X[k]$  erweitern.

Weil unser Filterproblem nun im Grunde genommen zeitvariant ist, bringt unsere frühere Annahme, dass  $X[\cdot]$  und  $Y[\cdot]$  gemeinsam s.s. sind, keine wesentliche Vereinfachung. Deswegen verzichten wir auf diese Annahme, d.h.

weder  $X[\cdot]$  noch  $Y[\cdot]$  müssen schwach-stationäre  
Prozesse sein.

Wir führen jetzt einige Notationen ein, die uns helfen werden, unsere Herleitung des optimalen Filters anschaulicher zu machen. Sei  $Z[k]$  der Wert eines beliebigen stochastischen Prozesses  $Z[\cdot]$  zum Zeitpunkt  $k$ , dann schreiben wir  $\hat{Z}[k|k-1]$  (resp.  $\hat{Z}[k|k]$ ) für die lineare MMSE-Schätzung



von  $Z[k]$  für die Beobachtungen  $Y[0], Y[1], \dots, Y[k-1]$  (resp.  $Y[0], Y[1], \dots, Y[k]$ ). Der Trick beim Kalman-Filter liegt nun darin, statt direkt mit  $Y[k]$ , mit

$$\tilde{Y}[k] = Y[k] - \hat{Y}[k|k-1], \quad \text{für } k \geq 1 \quad (7.14)$$

zu arbeiten. Wir definieren

$$\tilde{Y}[0] = Y[0], \quad (7.15)$$

sodass  $\tilde{Y}$  für alle  $k \geq 0$  definiert ist. Die Zufallsgrösse  $\tilde{Y}[k]$  heisst die **Innovation** von  $Y[\cdot]$  zum Zeitpunkt  $k$ , d.h. es ist der Teil von  $Y[\cdot]$ , der nicht linear schätzbar ist.

Zunächst stellen wir fest, dass  $Y[0], \dots, Y[k-1], \tilde{Y}[k]$  und  $Y[0], \dots, Y[k-1], Y[k]$  äquivalente Daten für die lineare Schätzung sind, d.h., dass

$$S(Y[0], \dots, Y[k-1], \tilde{Y}[k]) = S(Y[0], \dots, Y[k-1], Y[k]) \quad (7.16)$$

gilt. Dies ist eine Folge der Tatsache, dass  $\tilde{Y}[k]$  eine Linearkombination von  $Y[0], \dots, Y[k-1], Y[k]$  ist, und dass auch  $Y[k] = \tilde{Y}[k] + \hat{Y}[k|k-1]$  eine Linearkombination von  $Y[0], \dots, Y[k-1], \tilde{Y}[k]$  ist. Wir sehen aus (7.14), dass  $\tilde{Y}[k]$  der Fehler bei der linearen MMSE-Schätzung von  $Y[k]$  für die Beobachtung  $Y[0], \dots, Y[k-1]$  ist. Es folgt dann direkt aus dem Orthogonalitätsprinzip, dass  $\tilde{Y}[k]$  orthogonal zu  $Y[i]$  für  $i = 0, 1, \dots, k-1$  ist, d.h.

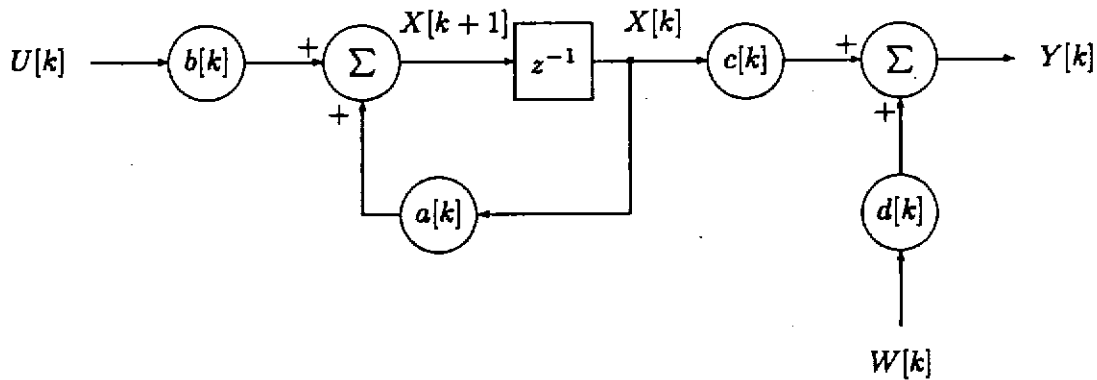
$$E[\tilde{Y}[k]Y[i]] = 0, \quad \text{für } i = 0, 1, \dots, k-1. \quad (7.17)$$

Aus der Trennungseigenschaft der linearen MMSE-Schätzung und mit Hilfe von (7.16) können wir jetzt schliessen, dass

$$\hat{Z}[k|k] = \hat{Z}[k|k-1] + \mathcal{E}(Z[k]|\tilde{Y}[k]) \quad (7.18)$$

gilt, wobei  $\mathcal{E}(Z[k]|\tilde{Y}[k])$  die lineare MMSE-Schätzung von  $Z[k]$  für die Beobachtung  $\tilde{Y}[k]$  bezeichnet. Gleichung (7.18) ist der Schlüssel für die rekursive Lösung des zeitvarianten Filterproblems.

Das spezifische Problem, das wir anpacken wollen, lässt sich in einem Blockdiagramm für den Zeitpunkt  $k$  wie folgt darstellen:



wobei

- (i)  $U[\cdot]$  und  $W[\cdot]$  unabhängige weisse Rauschprozesse mit

$$R_U[\cdot] = R_W[\cdot] = \delta[\cdot]$$

sind.  $X[0]$  ist unabhängig von  $U[k]$  und auch von  $W[k]$  für alle  $k \geq 0$ , und  $E[X^2[0]]$  ist bekannt.

- (ii)  $a[k], b[k], c[k]$  und  $d[k]$  sind bekannt für alle  $k \geq 0$ .

Wir wollen  $\hat{X}[k|k]$  finden, d.h. wir wollen die lineare MMSE-Schätzung von  $X[k]$  für die Beobachtungen  $Y[0], Y[1], \dots, Y[k]$  finden. Wir sehen, dass  $X[k]$  der Zustand eines linearen Systems mit dem Eingangssignal  $U[\cdot]$  ist. Das Signal  $d[\cdot]W[\cdot]$  gilt als Beobachtungsrauschen. Diesem Blockdiagramm entsprechen folgende Gleichungen:

$$Y[k] = c[k]X[k] + d[k]W[k], \quad \text{für } k \geq 0 \quad (7.19)$$

und

$$X[k] = a[k-1]X[k-1] + b[k-1]U[k-1], \quad \text{für } k \geq 1. \quad (7.20)$$

Wir sehen aus (7.19) und (7.20), dass  $Y[i]$  eine Linearkombination von  $X[0], W[0], W[1], \dots, W[i], U[0], U[1], \dots, U[i-1]$  ist. Alle diese Linearkombinationen von Zufallsgrößen sind wegen der Annahme (i) orthogonal zu  $W[k]$  für  $k > i$  (resp. zu  $U[k]$  für  $k \geq i$ ). Deswegen sind sowohl  $W[k]$  als auch  $U[k-1]$  orthogonal zu  $Y[0], Y[1], \dots, Y[k-1]$ . Dank dem Orthogonalitätsprinzip sehen wir nun, dass sowohl

$$\hat{W}[k|k-1] = 0, \quad \text{für } k \geq 0 \quad (7.21)$$

als auch

$$\hat{U}[k-1|k-1] = 0, \quad \text{für } k \geq 1 \quad (7.22)$$

gelten.

Um weiter zu kommen, eignen wir uns die folgende Strategie an. Wir nehmen an, dass wir  $\hat{X}[k-1|k-1]$  schon wissen und dass wir versuchen (mit Hilfe von  $Y[k]$ ),  $\hat{X}[k|k]$  zu finden.

Aus der Linearitätseigenschaft der linearen MMSE-Schätzung und aus (7.20) folgt, dass

$$\hat{X}[k|k-1] = a[k-1]\hat{X}[k-1|k-1] + b[k-1]\hat{U}[k-1|k-1]$$

für  $k \geq 1$  gilt. Ferner haben wir dann mit (7.22)

$$\hat{X}[k|k-1] = a[k-1]\hat{X}[k-1|k-1], \quad \text{für } k \geq 1 \quad (\text{KF-1})$$

und wir stellen fest, dass  $\hat{X}[k|k-1]$  die **Einschrittvorhersage** von  $X[k]$  ist. In ähnlicher Weise liefert (7.19) die Beziehung

$$\hat{Y}[k|k-1] = c[k]\hat{X}[k|k-1] + d[k]\hat{W}[k|k-1],$$

was mit Hilfe von (7.21) zu

$$\hat{Y}[k|k-1] = c[k]\hat{X}[k|k-1], \quad \text{für } k \geq 1 \quad (\text{KF-2})$$

reduziert werden kann. Die Definition (7.14) gibt uns

$$\tilde{Y}[k] = Y[k] - \hat{Y}[k|k-1], \quad \text{für } k \geq 1. \quad (\text{KF-3})$$

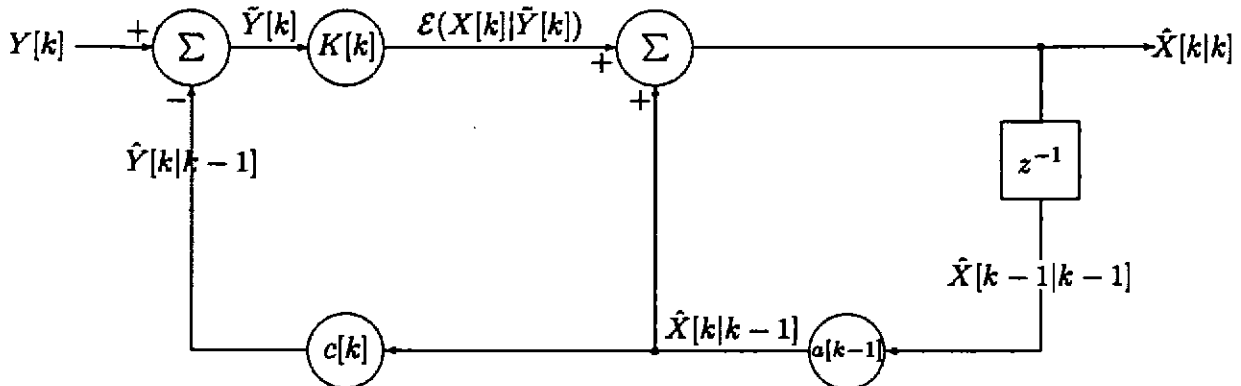
Jetzt müssen wir  $\mathcal{E}(X[k]|\tilde{Y}[k])$  (d.h. die lineare MMSE-Schätzung von  $X[k]$ ) betrachten. Wir wissen, dass diese Schätzung die Form

$$\mathcal{E}(X[k]|\tilde{Y}[k]) = K[k]\tilde{Y}[k], \quad \text{für } k \geq 0 \quad (\text{KF-4})$$

haben muss, wobei  $K[k]$  irgendeine reelle Zahl ist, die wir aus der Orthogonalitätsgleichung finden können.  $K[k]$  heisst der **“Kalman-Gain”** des Kalman-Filters. Wir schieben die Berechnung von  $K[k]$  auf einen späteren Zeitpunkt hinaus, und nehmen an, dass wir  $K[k]$  irgendwie gefunden haben. Dann können wir die Schlüsselbeziehung (7.18) (die der Trennungseigenschaft entspricht) ausnützen, um

$$\hat{X}[k|k] = \hat{X}[k|k-1] + \mathcal{E}(X[k]|\tilde{Y}[k]), \quad \text{für } k \geq 1 \quad (\text{KF-5})$$

zu schreiben. Die Gleichungen (KF-1) bis (KF-5) bilden den sogenannten **Kalman-Filter-Algorithmus**, der uns erlaubt,  $\hat{X}[k|k]$  aus  $\hat{X}[k-1|k-1]$  zu bestimmen. Dieser Algorithmus kann durch das folgende Blockdiagramm dargestellt werden:



Es muss nun noch der Anfangswert  $\hat{X}[-1|-1]$  der Rekursion spezifiziert werden. Da  $\tilde{Y}[0] = Y[0]$  gilt, folgt sofort, dass  $\hat{X}[0|0] = K[0]\tilde{Y}[0]$ . Offensichtlich ist

$$\hat{X}[-1|-1] = 0 \quad (7.23)$$

der gesuchte Anfangswert, was bedeutet, dass der Anfangszustand der Verzögerungselemente in unserem Blockdiagramm gleich Null ist.

Mit diesem Anfangswert gilt (KF-1) bis (KF-5) sowohl für  $k=0$ , als auch für  $k \geq 1$ .

Es bleibt nun noch das Problem der rekursiven Bestimmung des Kalman-Gains zu untersuchen. Zuerst definieren wir den mittleren quadratischen Einschrittvorhersagefehler als

$$v[k] = E \left[ \left( X[k] - \hat{X}[k|k-1] \right)^2 \right], \quad \text{für } k \geq 0 \quad (7.24)$$

und den mittleren quadratischen Schätzungsfehler als

$$s[k] = E \left[ \left( X[k] - \hat{X}[k|k] \right)^2 \right], \quad \text{für } k \geq 0. \quad (7.25)$$

Dabei ist  $s[k]$  diejenige Grösse, welche wir minimieren wollen. Wir werden im folgenden zeigen, dass die Kenntnis von  $v[k]$  ausreicht, um  $K[k]$ ,  $s[k]$  und  $v[k+1]$  berechnen zu können. Diese Tatsache ist nicht ohne weiteres ersichtlich.

Vom Orthogonalitätsprinzip wissen wir, dass  $K[k]$  die Lösung der Gleichung

$$E [\tilde{Y}^2[k]] K[k] = E [X[k]\tilde{Y}[k]] \quad (7.26)$$

ist. Betrachten wir vorerst die rechte Seite von (7.26). Wegen (KF-2) und (7.18) gilt:

$$\tilde{Y}[k] = c[k] (X[k] - \hat{X}[k|k-1]) + d[k]W[k]. \quad (7.27)$$

Durch Einsetzen in (7.26) erhalten wir

$$\begin{aligned} E [X[k]\tilde{Y}[k]] &= c[k]E [X[k] (X[k] - \hat{X}[k|k-1])] + d[k]E [X[k]W[k]] \\ &= c[k]E [X[k] (X[k] - \hat{X}[k|k-1])], \end{aligned} \quad (7.28)$$

da  $X[k]$  und  $W[k]$  orthogonal sind. Aus dem Orthogonalitätsprinzip folgt nun, dass

$$E [X[k] (X[k] - \hat{X}[k|k-1])] = v[k], \quad \text{für } k \geq 0 \quad (7.29)$$

gilt, da dieser Erwartungswert gleich dem MSE für die Einschrittvorhersage von  $X[k]$  ist. Deshalb reduziert sich (7.28) wie folgt:

$$E [X[k]\tilde{Y}[k]] = c[k]v[k], \quad \text{für } k \geq 0. \quad (7.30)$$

Analog findet man mit (7.26), dass

$$E [\tilde{Y}^2[k]] = c^2[k]v[k] + d^2[k], \quad \text{für } k \geq 0 \quad (7.31)$$

gilt. Durch Einsetzen von (7.30) und (7.31) in (7.26) erhalten wir

$$K[k] = \frac{c[k] \cdot v[k]}{c^2[k] \cdot v[k] + d^2[k]} \quad \text{für } k \geq 0. \quad (\text{KG-1})$$

Direkt aus dem Orthogonalitätsprinzip folgt:

$$s[k] = E [X[k] (X[k] - \hat{X}[k|k])],$$

woraus sich mit (KF-4) und (KF-5) ergibt:

$$\begin{aligned} s[k] &= E [X[k] (X[k] - \hat{X}[k|k-1] - K[k]\tilde{Y}[k])] \\ &= E [X[k] (X[k] - \hat{X}[k|k-1])] - K[k]E [X[k]\tilde{Y}[k]]. \end{aligned}$$

Diese Gleichung lässt sich mit Hilfe von (7.29) und (7.30) umformen:

$$s[k] = (1 - K[k]c[k])v[k], \quad \text{für } k \geq 0. \quad (\text{KG-2})$$

Mit (7.20) und (KF-1) folgt nun weiter, dass

$$X[k] - \hat{X}[k|k-1] = a[k-1] (X[k-1] - \hat{X}[k-1|k-1]) + b[k-1]U[k-1]$$

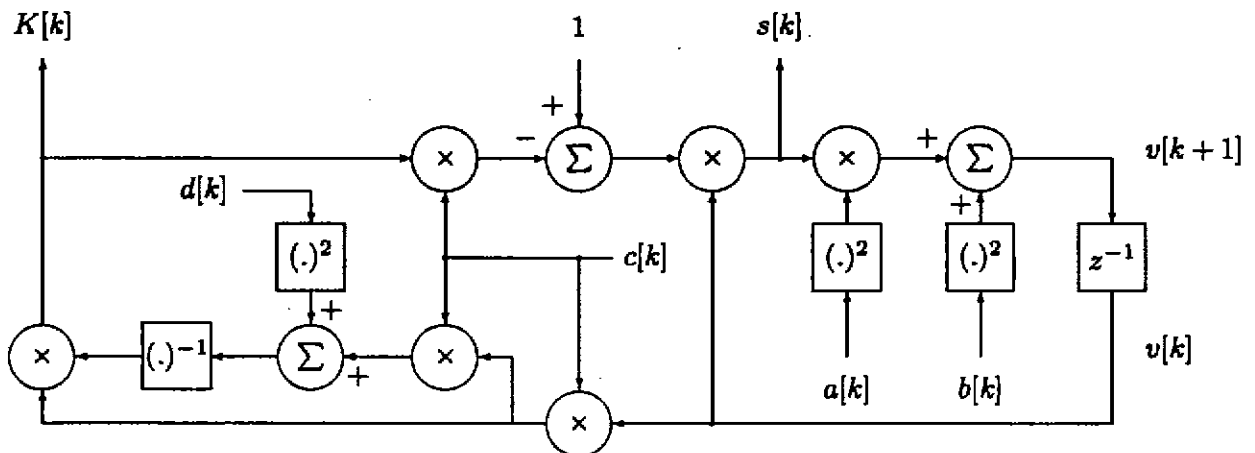
gilt. Da  $U[k-1]$  orthogonal zu  $X[k-1]$  und  $\hat{X}[k-1|k-1]$  ist, erhält man für  $v[k]$ :

$$\begin{aligned} v[k] &= E \left[ (X[k] - \hat{X}[k|k-1])^2 \right] \\ &= a^2[k-1]s[k-1] + b^2[k-1], \quad \text{für } k \geq 1. \end{aligned}$$

Äquivalent erhalten wir:

$$v[k+1] = a^2[k]s[k] + b^2[k], \quad \text{für } k \geq 0. \quad (\text{KG-3})$$

Die Gleichungen (KG-1), (KG-2) und (KG-3) bilden den sogenannten **Kalman-Gain-Algorithmus**, der es uns erlaubt,  $v[k+1]$  aus  $v[k]$  zu bestimmen und daneben sowohl  $K[k]$  als auch  $s[k]$  zu erhalten. Dieser Algorithmus kann durch folgendes Blockdiagramm dargestellt werden:

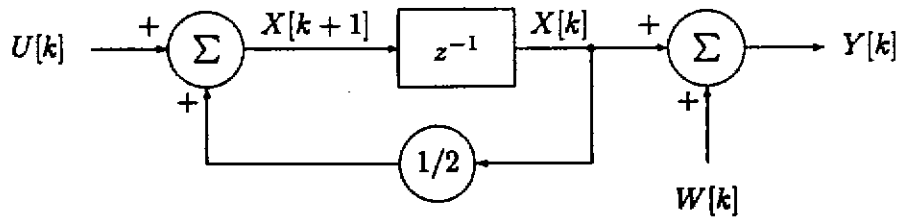


Es bleibt nun noch die Spezifikation des Anfangswerts dieser Rekursion. Da  $\hat{X}(0|-1) = 0$  gilt, erhalten wir für den Anfangswert:

$$v[0] = E [X^2[0]] \quad (7.32)$$

Da sowohl die Eingangssequenzen  $a[\cdot]$ ,  $b[\cdot]$ ,  $c[\cdot]$  und  $d[\cdot]$  als auch der Anfangszustand  $v[0] = E [X^2[0]]$  im Voraus bekannt sind, kann der Kalman-Gain-Algorithmus "off-line" abgearbeitet werden, um so die Kalman-Gain-Sequenz  $K[0], K[1], K[2], \dots$  zu finden. Selbstverständlich ist es auch möglich, diesen Algorithmus in Echtzeit parallel zum Kalman-Filter-Algorithmus abzuarbeiten.

Beispiel:



Der Einfachheit halber haben wir ein zeitinvariantes, lineares System mit

$$a[k] = \frac{1}{2}, \quad b[k] = c[k] = d[k] = 1, \quad \text{für alle } k$$

gewählt.

Wir nehmen an, dass  $E[X^2[0]] = 4/3$  spezifiziert wird. Dadurch erhalten wir mit (7.32)

$$v[0] = \frac{4}{3}$$

als Anfangswert des Kalman-Gain-Algorithmus. Die Gleichungen (KG-1), (KG-2) und (KG-3) lassen sich wie folgt reduzieren:

$$\begin{aligned} K[k] &= \frac{v[k]}{v[k] + 1} \\ s[k] &= (1 - K[k])v[k] = K[k] \\ v[k+1] &= \frac{1}{4}s[k] + 1. \end{aligned}$$

Unter Zuhilfenahme dieser Gleichungen liefert der Algorithmus die in der untenstehenden Tabelle angegebenen Werte:

$k$	$v[k]$	$K[k]$	$s[k]$
0	4/3	4/7	4/7
1	8/7	8/15	8/15
2	17/15	17/32	17/32

Der Anfangswert des Kalman-Filters ist (wie immer)

$$\hat{X}[-1|-1] = 0.$$

Für unser Beispiel reduzieren sich die Gleichungen (KF-1) und (KF-5) wie folgt:

$$\begin{aligned} \hat{X}[k|k-1] &= \frac{1}{2}\hat{X}[k-1|k-1] \\ \hat{Y}[k|k-1] &= \hat{X}[k|k-1] \\ \tilde{Y}[k] &= Y[k] - \hat{Y}[k|k-1] = Y[k] - \hat{X}[k|k-1] \\ \varepsilon(X[k]|\tilde{Y}[k]) &= K[k]\tilde{Y}[k] \\ \hat{X}[k|k] &= \hat{X}[k|k-1] + \varepsilon(X[k]|\tilde{Y}[k]). \end{aligned}$$

Für die Eingangssequenz  $Y[0], Y[1], Y[2], \dots = 0.8, -1.2, 0.1, \dots$  rechnet unser Kalman-Filter wie folgt:

$k$	$\hat{X}[k-1 k-1]$	$\hat{X}[k k-1]$	$Y[k]$	$\tilde{Y}[k]$	$K[k]$	$\varepsilon(X[k]Y[k])$
0	0.00	0.00	0.80	0.80	4/7	0.46
1	0.46	0.23	-1.20	-1.43	8/15	-0.76
2	-0.53	-0.27	0.10	0.37	17/32	0.20
3	-0.07					

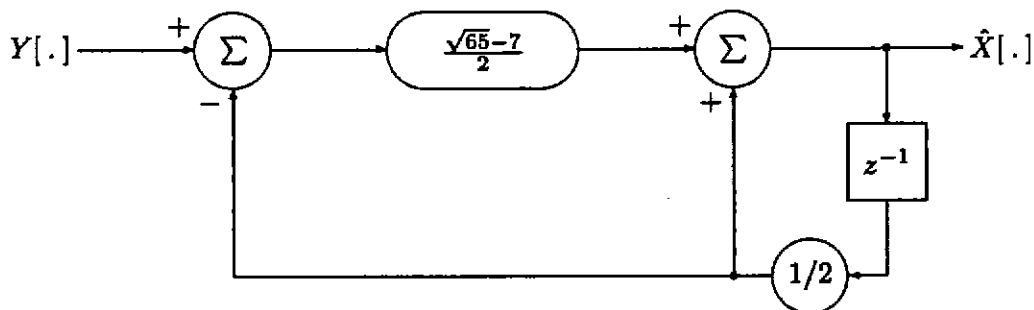
Eingangssignale

vom Kalman-Gain-Algorithmus

Für dieses zeitinvariante Beispiel ist es leicht, zu zeigen, dass

$$\lim_{k \rightarrow \infty} K[k] = \frac{\sqrt{65} - 7}{2} \left( \approx \frac{1}{2} \right),$$

das heisst, das Kalman-Filter konvergiert gegen das folgende zeitinvariante Filter:



mit der Gewichtssequenz

$$h[k] = \begin{cases} 0, & k < 0 \\ \frac{\sqrt{65}-7}{2} \left( \frac{9-\sqrt{65}}{4} \right)^k, & k \geq 0. \end{cases}$$

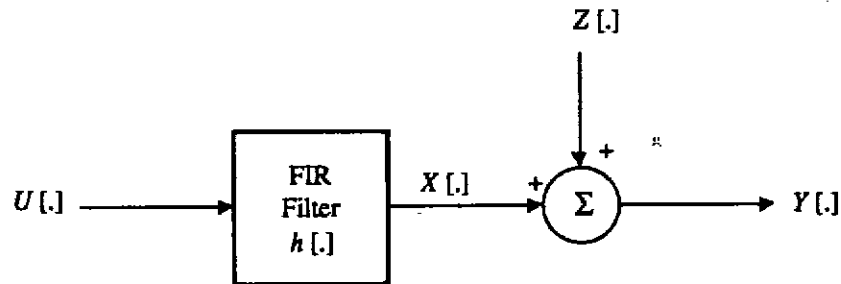
Dieses "asymptotische" Kalman-Filter muss gleich dem kausalen Wiener-Filter des entsprechenden zeitinvarianten Filterungsproblems sein. (vgl. Übungsaufgabe).

Was hier über das Kalman-Filter geschrieben wurde, lässt sich leicht für den Fall verallgemeinern, wo die Zufallsgrößen  $X[k]$ ,  $U[k]$ ,  $W[k]$  und  $Y[k]$  durch Zufallsvektoren  $\underline{X}[k]$ ,  $\underline{U}[k]$ ,  $\underline{W}[k]$  und  $\underline{Y}[k]$  ersetzt werden. Viele praktische Probleme weisen diese Form auf, weshalb das Kalman-Filter in der Praxis weit verbreitet ist. Der Leser, welcher die dem Kalman-Filter zugrunde liegenden Prinzipien gut verstanden hat, sollte in der Lage sein, die Erweiterung auf den vektoriellen Fall selbständig durchzuführen. Ansonsten wird dieses Problem in nahezu jedem Buch über stochastische Filter behandelt. (Leider sind in den meisten Lehrbüchern die Grundlagen entweder nicht erwähnt oder aber sehr gut versteckt!).



## 7.5 Maximum-Likelihood-Filterung (Viterbi-Algorithmus)

Mathematisches Modell:



wobei

- (1)  $U[\cdot]$  ein binärer Prozess mit  $U[k] \in \{-1, +1\}$  für alle  $k$  ist,
- (2) das FIR-Filter ein Gedächtnis  $M$  hat, d.h.  $H(z) = h[0] + h[1]z^{-1} + \dots + h[M]z^{-M}$ ,
- (3)  $Z[\cdot]$  weisses Gaussches Rauschen mit Varianz  $\sigma^2$  ist, das unabhängig von  $U[\cdot]$  (und deshalb auch unabhängig von  $X[\cdot]$ ) ist, und
- (4)  $Y[\cdot]$  der Prozess ist, den wir beobachten.

Wir wollen den Maximum-Likelihood- (ML-) Entscheid für die Sequenz

$$\underline{U} = (U[0], U[1], \dots, U[L-1]) \quad (7.33)$$

treffen unter der Annahme, dass

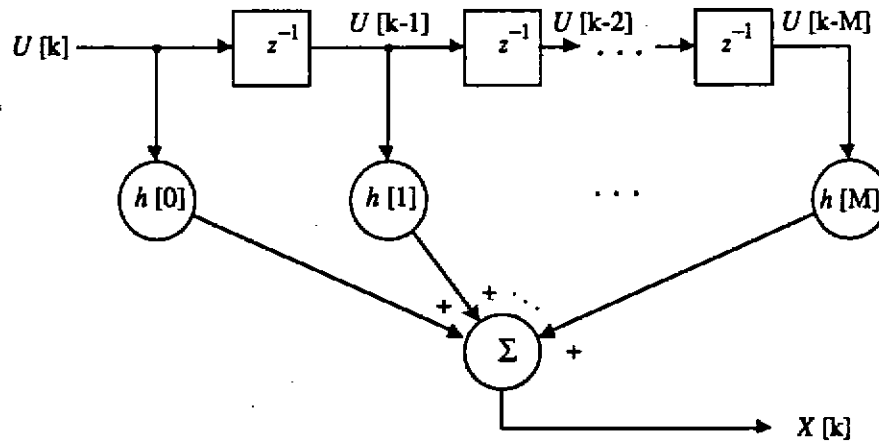
$$U[-1] = U[-2] = \dots = U[-M] = +1 \quad (7.34)$$

und

$$U[L+M-1] = U[L+M-2] = \dots = U[L] = +1 \quad (7.35)$$

gilt.

Aus folgendem Diagramm des FIR-Filters ersieht man,



dass der Zustand des Filters zum Zeitpunkt  $k$  durch

$$\underline{S}[k] = (U[k-1], U[k-2], \dots, U[k-M]) \quad (7.36)$$

gegeben ist. Unsere Annahme (7.34) und (7.35) für die Sequenz  $U[\cdot]$  bedeutet, dass  $\underline{S}[0] = \underline{S}[L+M] = (+1, +1, \dots, +1)$  gilt.

Bevor wir unser Problem anpacken, erwähnen wir ein Prinzip, das sowohl hier als auch bei anderen Anwendungen von Nutzen ist:

**Invertierbarkeitsprinzip für die ML-Entscheidung (bzw. die ML-Schätzung):**  
 Sei  $\underline{X} = f(\underline{U})$ , wobei  $f : \underline{U}(\Omega) \rightarrow \underline{X}(\Omega)$  eine invertierbare Funktion ist, und sei  $\underline{Y}$  der beobachtete Zufallsvektor. Dann ist  $\tau(\cdot)$  mit  $\hat{\underline{X}} = \tau(\underline{Y})$  die ML-Entscheidungsregel (bzw. die ML-Schätzungsregel) genau dann, wenn  $f^{-1}(\tau(\cdot))$  mit  $\hat{\underline{U}} = f^{-1}(\tau(\underline{Y})) = f^{-1}(\hat{\underline{X}})$  die ML-Entscheidung (bzw. die ML-Schätzung) für  $\underline{U}$  ist.

Um dieses Prinzip zu beweisen, bemerken wir, dass  $\underline{X} = \underline{x}$  genau dann gilt, wenn  $\underline{U} = f^{-1}(\underline{x})$  ist. Es folgt

$$p_{\underline{Y}|\underline{X}}(\underline{y}|\underline{x}) = p_{\underline{Y}|\underline{U}}(\underline{y}|f^{-1}(\underline{x})) ,$$

was direkt zum Invertierbarkeitsprinzip führt.

Wir machen den Leser darauf aufmerksam, dass das Invertierbarkeitsprinzip **nicht für andere Entscheidungsregeln (bzw. Schätzungsregeln)**, wie zum Beispiel die MMSE-Schätzungsregel, gilt. Weiter bemerken wir, dass die Frage, ob ein Problem ein ML-Entscheidungsproblem oder ein ML-Schätzungsproblem ist, nur davon abhängt, ob  $\underline{U}$  ein diskreter oder ein kontinuierlicher Zufallsvektor ist. Für unser Problem ist  $\underline{U}$  ein diskreter Zufallsvektor, der  $2^L$  mögliche Werte hat.

Definieren wir jetzt

$$\underline{X} = (X[0], X[1], \dots, X[L+M-1]) . \quad (7.37)$$

Wir wissen vom FIR-Filter, dass

$$X[k] = h[0]U[k] + h[1]U[k-1] + \dots + h[M]U[k-M] \quad (7.38)$$

für alle  $k$  gilt. Mit der Annahme, dass  $h[0], h[1], \dots, h[M]$  nicht alle Null sind, folgt dann aus (7.38) mit Hilfe von (7.34) und (7.35), dass  $\underline{X} = f(\underline{U})$  gilt, wobei  $f$  eine invertierbare Funktion ist. Das heisst also: Aus  $X[0], X[1], \dots, X[L+M-1]$  können wir mit (7.38) und mit Hilfe von (7.34) und (7.35) wieder  $U[0], U[1], \dots, U[L-1]$  finden. Deshalb reduziert sich unser Problem auf das Finden der ML-Entscheidung  $\hat{\underline{X}}$  für  $\underline{X}$  mit anschliessender Inversion  $\hat{\underline{U}} = f^{-1}(\hat{\underline{X}})$ . Weil es genau  $2^L$  mögliche Werte für  $\underline{U}$  gibt, existieren auch genau  $2^L$  mögliche Werte für  $\underline{X}$ .

Analog zu (7.37) definieren wir jetzt

$$\underline{Y} = (Y[0], Y[1], \dots, Y[L+M-1]) \quad (7.39)$$

und

$$\underline{Z} = (Z[0], Z[1], \dots, Z[L+M-1]) \quad (7.40)$$

wobei wir bemerken, dass die Komponenten von  $\underline{Z}$  i.i.d. Gaussische Zufallsgrössen mit Mittelwert 0 und Varianz  $\sigma^2$  sind. Weiter gilt

$$\underline{Y} = \underline{X} + \underline{Z} \quad (7.41)$$

wobei  $\underline{X}$  und  $\underline{Z}$  unabhängige Zufallsvektoren sind. Damit haben wir unser Problem auf das bekannte Problem der ML-Entscheidung eines Signals  $\underline{X}$  in additivem, weissem Gaussischem Rauschen (AWGN: Additive White Gaussian Noise)  $\underline{Z}$  zurückgeführt. Die Lösung (siehe Übungsaufgabe) lautet:

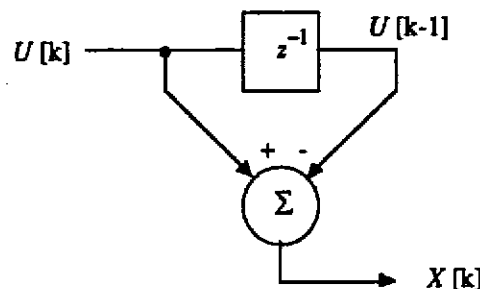
Wähle für die gegebene Beobachtung  $\underline{Y} = \underline{y}$  als  $\hat{\underline{X}}$  dasjenige  $\underline{x}$  in  $\underline{X}(\Omega)$ , welches die Euklidische Distanz zwischen  $\underline{x}$  und  $\underline{y}$  minimiert.

Unsere nächste Aufgabe ist es nun, alle  $2^L$  Sequenzen in  $\underline{X}(\Omega)$  in benutzerfreundlicher Art und Weise zu beschreiben.

Zum besseren Verständnis des Verfahrens arbeiten wir jetzt mit einem besonders einfachen FIR-Filter. Es ist dann aber leicht, in den allgemeinen Fall überzugehen. Das FIR-Filter, das wir betrachten werden, ist durch

$$H(z) = 1 - z^{-1} = \frac{(z-1)}{z} \quad (7.42)$$

gegeben und in folgender Zeichnung abgebildet:



Der Zustand  $\underline{S}[k]$  von (7.36) reduziert sich zu der einzigen Zufallsgrösse

$$S[k] = (U[k-1]) \quad (7.43)$$

Obwohl dieses Beispiel sehr einfach ist, hat es doch praktische Bedeutung. Die Übertragungsfunktion  $H(z) = \frac{(z-1)}{z}$  entspricht dem rauschfreien Teil eines Kanals, der in der Praxis sehr oft vorkommt. Der Grund dafür ist, dass  $H(z) = \frac{(z-1)}{z}$  eine Nullstelle für  $z = 1$  aufweist, was der Gleichspannungskomponente entspricht. (Für  $\Omega = 0$  (Gleichspannung) gilt  $z = e^{j\Omega} = 1$ .) Bei der magnetischen Speicherung, wie auch bei einigen Anwendungen in der Telefonie, strebt man danach, die Übertragungsfunktion  $H(z) = \frac{(z-1)}{z}$  im Kanal zu realisieren, weil es schwierig oder gar unmöglich ist, Gleichspannung zu übertragen. Aus (7.38) oder aus obigem Diagramm sieht man, dass

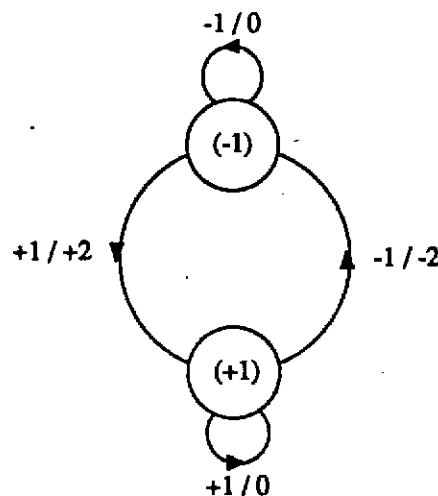
$$X[k] = U[k] - U[k-1] \quad (7.44)$$

ist. Man spricht auch von Intersymbolinterferenz (ISI) im Kanal, weil jedes empfangene Symbol

$$Y[k] = X[k] + Z[k] = U[k] - U[k-1] + Z[k]$$

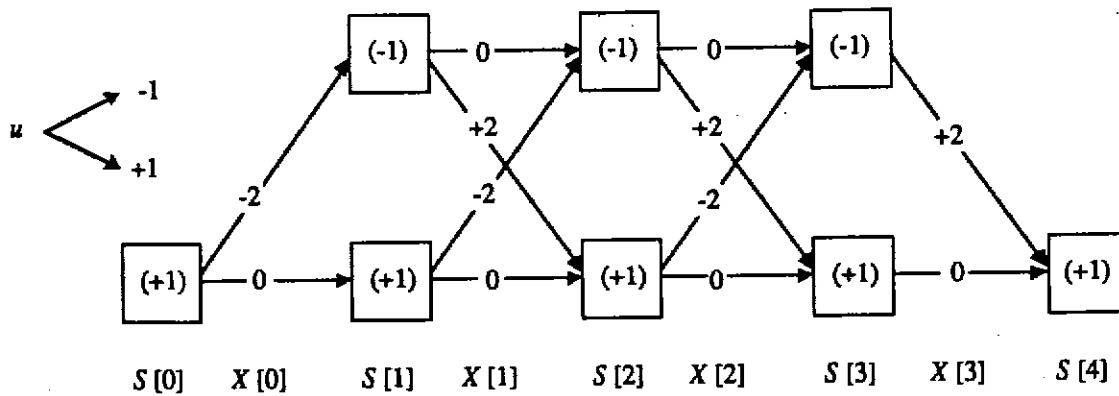
von zwei Datensymbolen abhängt.

Dank der Tatsache, dass  $U[k] \in \{-1, +1\}$  gilt, kann der Zustand des FIR-Filters in unserem Beispiel nur zwei mögliche Werte annehmen, nämlich  $(-1)$  oder  $(+1)$ . Solch eine endliche Zustandsmaschine lässt sich leicht durch ein Zustandsübergangsdiagramm darstellen, das für unser Beispiel wie folgt aussieht:



Im Zustandsübergangsdiagramm entsprechen die Knoten den möglichen Zuständen. Die Übergänge sind mit  $u/x$  beschriftet, wobei  $u$  das Eingangssymbol ist, das den entsprechenden Übergang verursacht, und  $x$  das entsprechende Ausgangssymbol ist. Für den Übergang vom Zustand  $s$  zum Zustand  $s'$  bedeutet zum Beispiel  $S[k] = s$  und  $U[k] = u$ , dass  $X[k] = x$  und  $S[k+1] = s'$  ist.

Nehmen wir nun an, die Länge  $L$  unserer Datensequenz  $\underline{U}$  von (7.33) sei  $L = 3$ . Dann sind  $U[0]$ ,  $U[1]$  und  $U[2]$  frei wählbar, aber  $U[-1] = +1$  und  $U[3] = +1$  sind durch (7.34) und (7.35) bestimmt, was  $S[0] = (+1)$  und  $S[4] = (+1)$  entspricht. Wir können das obige Zustandsdiagramm benutzen, um die möglichen Zeitverläufe unserer endlichen Zustandsmaschine wie folgt darzustellen:

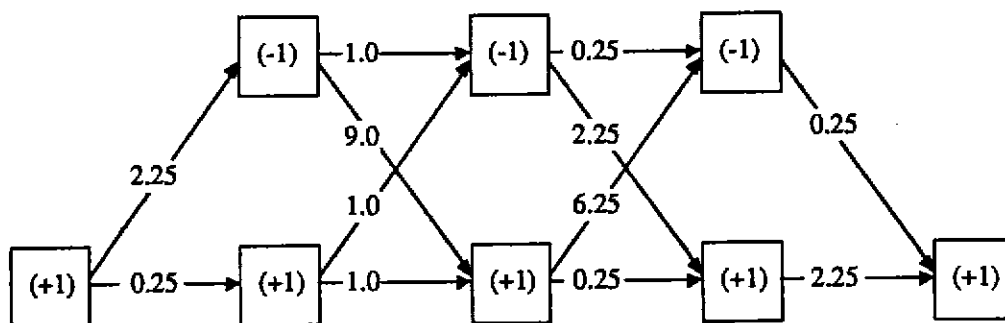


Die Zweige im Diagramm sind mit den Werten des entsprechenden Ausgangssymbols beschriftet. Die Symbole auf einem Pfad von Anfang bis zum Ende des Diagramms stellen eine mögliche Ausgangssequenz  $\underline{X} = (X[0], X[1], X[2], X[3])$  dar. Weiter stellen die Knoten entlang des Pfades die entsprechende Zustandssequenz  $(S[0] = +1, S[1], S[2], S[3], S[4] = +1)$  dar. Wir haben das Diagramm so gestaltet, dass der obere (bzw. der untere) der beiden Zweige, die einen Zustand verlassen, einem Datensymbol  $U[k]$  mit dem Wert  $-1$  (bzw.  $+1$ ) entspricht. Gibt es nur einen Zweig, der einen Zustand verlässt, dann entspricht dieser nach (7.35) immer einem Datensymbol mit dem Wert  $+1$ . Deshalb brauchen wir die Datensymbole in unserem Diagramm nicht explizit einzutragen. Ein solches Diagramm heisst "Trellis" (deutsch: "Spalier"), weil es einem Spalier in einem Rosengarten ähnlich sieht. Der Trellis liefert uns die gewünschte benutzerfreundliche Darstellung aller  $2^L$  möglichen Datensequenzen in  $\underline{X}(\Omega)$ .

Unsere nächste Aufgabe besteht darin, für eine gegebene Empfangssequenz  $\underline{y}$  denjenigen Pfad  $\underline{x}$  im Trellis zu finden, der die minimale Euklidische Distanz zu  $\underline{y}$  aufweist. Dieses Problem löst uns der Viterbi-Algorithmus, den wir nun mit Hilfe eines Beispiels beschreiben. Nehmen wir an, die empfangene Sequenz sei

$$\underline{y} = [-0.5, -1.0, +0.5, +1.5] . \quad (7.45)$$

Wir lassen diese Sequenz eine "Metrik" auf jedem Zweig des Trellis bestimmen, wobei diese Metrik dem Quadrat der Euklidischen Distanz zwischen dem Symbol  $x$  auf dem Zweig im Trellis und dem entsprechenden Symbol in  $\underline{y}$  entspricht. Damit erhalten wir folgendes Trellis-Diagramm:

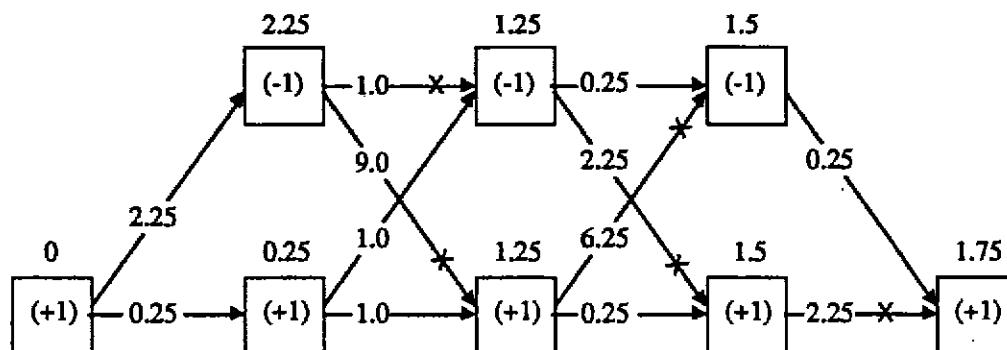


Weil das Quadrat der Euklidischen Distanz zwischen zwei Sequenzen die Summe der Quadrate der Euklidischen Distanz zwischen den einzelnen Symbolen der Sequenzen ist, reduziert sich

unser ML-Entscheidungsproblem wie folgt: **Finde den Pfad  $\underline{x}$  durch obigen Trellis, für den die Summe der Zweigmetriken minimal ist.** Der Viterbi-Algorithmus für die Lösung dieses Problems lässt sich folgendermassen beschreiben:

- (1) Weise dem Anfangsknoten des Trellis die Metrik 0 zu.
- (2) Gehe einen Zweig tiefer im Trellis.
- (3) Führe für jeden Knoten in dieser Tiefe folgende Schritte durch:
  - (i) Addiere bei jedem zu diesem Knoten führenden Zweig die Zweigmetrik und die Metrik des Knotens am Anfang des Zweiges.
  - (ii) Weise die kleinste dieser neu berechneten Metriken diesem Knoten zu und entferne alle Zweige ausser demjenigen, der diese Metrik aufweist. Wenn mehrere Zweige diese kleinste Metrik aufweisen, wähle einen beliebigen Zweig mit kleinster Metrik als "Überlebenden".
- (4) Höre auf, wenn das Ende des Trellis erreicht ist, sonst gehe zu Schritt (2) zurück.

Die einzelnen Schritte des Viterbi-Algorithmus für unser Beispiel sind im folgenden Diagramm dargestellt. Die Zweige, die in Schritt (3)(ii) entfernt wurden, sind dabei mit einem Kreuz markiert.



Wenn wir versuchen, rückwärts durch den Trellis zu gehen, sehen wir, dass **nach Anwendung des Viterbi-Algorithmus nur ein Pfad im Trellis "überlebt"**. Dieser "überlebende" Pfad ist der optimale Pfad, wie folgende Tatsache zeigt:

Wenn ein Zweig des Trellis im Schritt (3)(ii) des Viterbi Algorithmus entfernt wurde, liegt dieser Zweig nicht auf dem Pfad (bzw. nicht auf dem einzigen Pfad) durch den Trellis, der die minimale Summe der Zweigmetriken aufweist.

Diese Aussage gilt, weil die Entfernung eines solchen Zweiges äquivalent zur Vernachlässigung eines Teilpfads ist. Dieser Teilpfad führt vom Anfang des Trellis zu jenem Knoten, zu dem wir einen besseren (bzw. nicht schlechteren) Teilpfad kennen. Weil man beide Teilpfade identisch zu einem Pfad durch den ganzen Trellis erweitern kann, kann der schlechtere (bzw. nicht bessere)

Teilpfad nicht Teil des optimalen Pfads (bzw. nicht Teil eines einzigen optimalen Pfads) durch den Trellis sein.

Aus dem ersten Trellis erkennen wir, dass für unser Beispiel

$$\underline{x} = [0, -2, 0, 2]$$

die ML-Entscheidung für  $\underline{X}$  ist. Die entsprechende ML-Entscheidung  $\underline{u} = f^{-1}(\underline{x})$  für die Datensequenz  $\underline{U}$  ergibt sich aus dem selben Trellis zu

$$\underline{u} = [+1, -1, -1] .$$

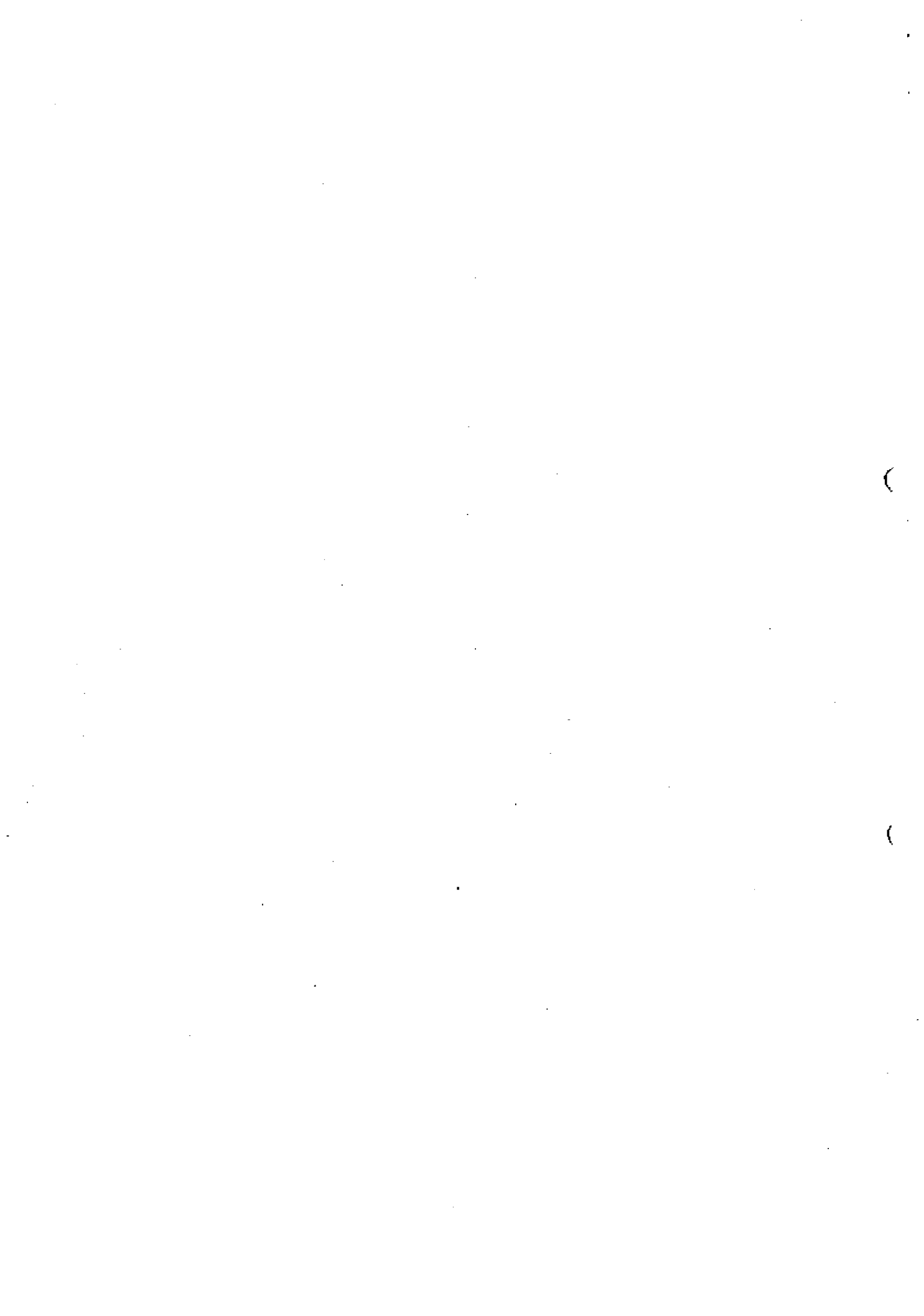
Folgende Punkte sollten dem Leser nun klar sein:

Für die Beschreibung aller möglichen Ausgangssequenzen einer endlichen Zustandsmaschine mit bekanntem Anfangszustand und bekanntem Endzustand, die mit einer Sequenz mit einem endlichen Alphabet gespiesen wird, existiert ein Trellis.

und:

Der Viterbi-Algorithmus kann benutzt werden, um die ML-Entscheidung unter denjenigen Sequenzen, die durch den Trellis bestimmt werden, genau dann zu treffen, wenn man eine Metrik zwischen den beobachteten Symbolen und den Symbolen auf den Trelliszweigen definieren kann mit folgender Eigenschaft: Der Pfad mit minimaler (oder maximaler) Summe der Zweigmetriken entspricht der ML-Entscheidung.

Viele Probleme der modernen Nachrichtentechnik lassen sich in obiger Form darstellen. Deshalb ist in den letzten zwei Jahrzehnten der Viterbi-Algorithmus zu einem sehr geschätzten Werkzeug des Kommunikationsingenieurs geworden. Es gibt viele Tricks, die man bei spezifischen Anwendungen des Viterbi-Algorithmus benutzen kann. Man sieht jedoch allgemein, dass die Komplexität proportional zur Anzahl Zustände der entsprechenden endlichen Zustandsmaschine wächst. Ist diese Zahl kleiner als ca. 1024 ( $2^{10} = 1024$  Zustände), kann der Viterbi-Algorithmus als gutes Lösungsmittel für eine ML-Entscheidung in Betracht gezogen werden.





## A VERALLGEMEINERTE FUNKTIONEN UND DER DIRAC-STOSS

Dem Leser ist sicher schon bekannt, dass der Dirac-Stoss (oder "Delta-Funktion")  $\delta(\cdot)$  eigentlich keine Funktion ist. Andererseits weiss er, dass man mit  $\delta(\cdot)$  fast so arbeiten kann, wie wenn  $\delta(\cdot)$  tatsächlich eine Funktion wäre. Um den Dirac-Stoss mathematisch sauber zu behandeln, braucht man eine Verallgemeinerung des Begriffs "Funktion". Diese Verallgemeinerung kann in verschiedener Art und Weise durchgeführt werden. Wir geben hier eine Verallgemeinerung an, die besonders einfach ist, die aber auch die wichtigsten Aspekte jeder Verallgemeinerung enthält. In diesem Anhang meinen wir mit "Funktion" eine reellwertige Funktion einer reellen Variablen, d.h. eine Funktion  $f : R \rightarrow R$ . Der Leser weiss von den Vorlesungen in Analysis, dass man auch "pathologische" Funktionen definieren kann, z.B. die Funktion

$$f(x) = \begin{cases} -1 & \text{falls } x \text{ rational} \\ 1 & \text{falls } x \text{ irrational.} \end{cases}$$

[Für diese Funktion gilt  $|f(x)| = 1$  für alle  $x$  und damit ergibt sich  $\int_a^b |f(x)| dx = b - a$  für alle  $a, b \in R$ . Aber  $\int_a^b f(x) dx$  existiert nicht als Riemannsches Integral, im Widerspruch zu Behauptung 7 in §3.1.7.2 des Taschenbuchs der Mathematik [1]! Der Leser sollte lernen, jeder mathematischen Behauptung missträuisch gegenüber zu stehen, für die er keinen Beweis gesehen hat. Damit sich  $\int_a^b f(x) dx = b - a$  ergibt, muss man eine Verallgemeinerung des Riemanschen Integrals benützen, z.B. das Lebesguesche Integral.]

Solche pathologischen Funktionen bereiten grosse Mühe in der Analysis, obwohl sie in praktischen Anwendungen kaum vorkommen. Ein Hauptziel einer Theorie der verallgemeinerten Funktionen ist es, diese pathologischen Funktionen für immer zu eliminieren, sodass man mit einfachen Mitteln sauber mathematisch arbeiten kann. Als Vorbereitung für eine solche Theorie führen wir nun zwei Klassen nicht pathologischer Funktionen ein. Aber zuerst eine Bemerkung zur Notation: Man schreibt

$$f(t_0+), f(t_0-), f(+\infty), f(-\infty)$$

für

$$\lim_{t \rightarrow t_0+} f(t), \lim_{t \rightarrow t_0-} f(t), \lim_{t \rightarrow \infty} f(t), \lim_{t \rightarrow -\infty} f(t).$$

$f^{(i)}$  bezeichnet die  $i$ -te Ableitung von  $f$ , wobei  $f^{(0)} = f$ . Man schreibt auch  $f' = f^{(1)}$ ,  $f'' = f^{(2)}$ , usw.

**Definition:** Eine Funktion  $f$  heisst **glatt**, wenn  $f$  beliebig oft differenzierbar ist, und wenn

$$f^{(i)}(+\infty) = f^{(i)}(-\infty) = 0 \quad \text{für } i = 0, 1, 2, \dots \quad (\text{A.1})$$

gilt.

**Beispiel 1:** Die Funktionen  $e^{-t^2}$  und  $\frac{\sin(t)}{t}$  sind glatt. Die Funktion  $e^{-|t|}$  ist nicht glatt, weil sie in  $t = 0$  nicht differenzierbar ist. Die Funktion  $f(t) = t$  ist nicht glatt, weil  $f(+\infty)$  und  $f(-\infty)$  nicht existieren.

Die folgenden Tatsachen liegen auf der Hand:

Sei  $f$  eine glatte Funktion, dann sind sowohl  $cf$  für jedes  $c \in \mathbb{R}$  als auch  $f^{(i)}$  für jede positive, ganze Zahl  $i$  ebenfalls glatte Funktionen.  
Seien  $f_1$  und  $f_2$  glatte Funktionen, dann ist auch  $f_1 + f_2$  eine glatte Funktion.

**Definition:** Eine Funktion  $f$  heisst **brav**, wenn

- (i)  $f$  stetig ist;
- (ii)  $f(t) = 0$  für alle  $t$  ausserhalb eines endlichen Intervalls;
- (iii) es eine ganze Zahl  $m$  gibt, so dass  $f$  beliebig oft differenzierbar ist ausser in höchstens  $m$  Punkten innerhalb jedes Intervalls der Länge 1; und
- (iv) in jedem Punkt  $t_0$ , wo  $f$  nicht beliebig oft differenzierbar ist,  $f^{(i)}(t_0+)$  und  $f^{(i)}(t_0-)$  für jede positive, ganze Zahl  $i$  existieren.

Wieder liegen folgende Tatsachen auf der Hand:

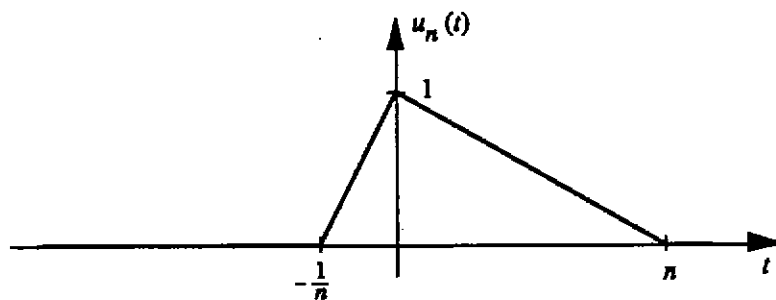
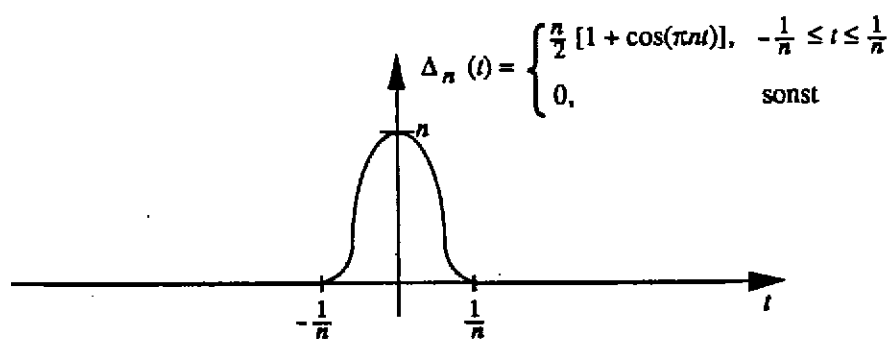
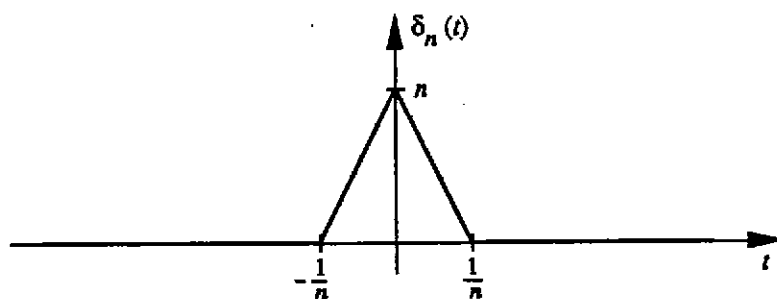
Sei  $f$  eine brave Funktion, dann ist auch  $cf$  für jedes  $c \in \mathbb{R}$  eine brave Funktion.  
Seien  $f_1$  und  $f_2$  brave Funktionen, dann ist auch  $f_1 + f_2$  eine brave Funktion.  
Sei  $f$  eine brave Funktion, die zudem differenzierbar ist, dann ist auch  $f'$  eine brave Funktion.

Für eine brave Funktion  $f$  gilt trivialerweise:

$$f^{(i)}(+\infty) = f^{(i)}(-\infty) = 0 \quad (\text{A.2})$$

für jede nichtnegative ganze Zahl  $i$ .

**Beispiel 2:** Für jede positive, ganze Zahl  $n$  ist jede der folgenden drei Funktionen brav:



Wenn wir eine Funktion  $f$  in der realen Welt beobachten, können wir  $f(t_1)$  und  $f(t_2)$  nicht unterscheiden, falls  $t_1$  und  $t_2$  genügend nahe beieinander liegen. Tatsächlich ist jede Beobachtung irgendwie ein Durchschnitt vieler Werte der Funktion  $f$ . Wir können dann jede Beobachtung von  $f$  als ein Integral  $\int_{-\infty}^{+\infty} f(t)w(t)dt$  für irgendeine Fensterfunktion  $w$  betrachten. Aber die Fensterfunktionen, die physikalisch sinnvoll sind, müssen "glatt" sein. In einem gewissen Sinn können wir dann sagen, dass die Funktionen  $f_1$  und  $f_2$  "physikalisch identisch" sind, falls gilt:

$$\int_{-\infty}^{+\infty} f_1(t)w(t)dt = \int_{-\infty}^{+\infty} f_2(t)w(t)dt$$

für jede glatte Funktion  $w$ . Diese Betrachtungen sind die Basis für die folgende Definition:

**Definition:** Eine Sequenz  $g_1, g_2, g_3, \dots$  von braven Funktionen mit der Eigenschaft, dass

$$\lim_{n \rightarrow \infty} \int_{-\infty}^{+\infty} g_n(t)w(t)dt$$

für jede glatte Funktion  $w$  existiert, spezifiziert eine verallgemeinerte Funktion (**v-Funktion**)  $g$  durch die Regel

$$\int_{-\infty}^{+\infty} g(t)w(t)dt = \lim_{n \rightarrow \infty} \int_{-\infty}^{+\infty} g_n(t)w(t)dt, \quad (\text{A.3})$$

wobei  $w$  eine glatte Funktion ist.

**Bemerkung:** Wir erinnern den Leser daran, dass ein Limes genau dann "existiert" (z.B.  $\lim_{n \rightarrow \infty} a_n = a$ ), wenn dieser Limes einen bestimmten reellen Wert hat (d.h.  $-\infty < a < \infty$ ). Gilt  $\lim_{n \rightarrow \infty} a_n = +\infty$  bzw.  $-\infty$ , dann "existiert" der Limes nicht.

Spezifizieren die Sequenzen von braven Funktionen  $g_n(t)$  und  $\gamma_n(t)$  ( $n = 1, 2, 3, \dots$ ) die v-Funktionen  $g(t)$  und  $\gamma(t)$ , dann spezifizieren die folgenden Sequenzen ebenfalls v-Funktionen:

- (i)  $g_n(t - t_0)$  für beliebige  $t_0 \in R$ ,
- (ii)  $cg_n(t)$  für beliebige  $c \in R$ , und
- (iii)  $g_n(t) + \gamma_n(t)$  ( $n = 1, 2, 3, \dots$ ).

Man schreibt sie als:

- (i)  $g(t - t_0)$ ,
- (ii)  $cg(t)$ , und
- (iii)  $g(t) + \gamma(t)$ .

**Beispiel 3:** Die Funktionen  $\delta_n$  ( $n = 1, 2, 3, \dots$ ) in Beispiel 2 spezifizieren die v-Funktion  $\delta(\cdot)$  in der Weise, dass gilt:

$$\begin{aligned} \int_{-\infty}^{+\infty} \delta(t)w(t)dt &= \lim_{n \rightarrow \infty} \int_{-\infty}^{+\infty} \delta_n(t)w(t)dt = \\ &= \lim_{n \rightarrow \infty} \int_{-\frac{1}{n}}^{\frac{1}{n}} \delta_n(t)w(t)dt = \\ &= w(0) \lim_{n \rightarrow \infty} \int_{-\frac{1}{n}}^{\frac{1}{n}} \delta_n(t)dt = \\ &= w(0) \end{aligned}$$

für jede glatte Funktion  $w$ . Diese einfache v-Funktion  $\delta(\cdot)$  ist nichts anderes als der berühmte Dirac-Stoss.

Wenn wir uns nun an unsere Betrachtungen über "physikalisch identische" Funktionen erinnern, so müssen wir zu folgender Definition gelangen:

Zwei Sequenzen  $g_1, g_2, g_3, \dots$  und  $\gamma_1, \gamma_2, \gamma_3, \dots$  von braven Funktionen spezifizieren dieselbe  $v$ -Funktion, falls

$$\lim_{n \rightarrow \infty} \int_{-\infty}^{+\infty} g_n(t)w(t)dt = \lim_{n \rightarrow \infty} \int_{-\infty}^{+\infty} \gamma_n(t)w(t)dt$$

für jede glatte Funktion  $w$  gilt.

**Beispiel 4:** Für die Funktionen  $\Delta_n$  ( $n = 1, 2, 3, \dots$ ) in Beispiel 2 gilt

$$\begin{aligned} \lim_{n \rightarrow \infty} \int_{-\infty}^{\infty} \Delta_n(t)w(t)dt &= \lim_{n \rightarrow \infty} \int_{-\frac{1}{n}}^{\frac{1}{n}} \Delta_n(t)w(t)dt = \\ &= w(0) \lim_{n \rightarrow \infty} \int_{-\frac{1}{n}}^{\frac{1}{n}} \Delta_n(t)dt = \\ &= w(0) \end{aligned}$$

für jede glatte Funktion  $w$ . Es folgt dann aus Beispiel 3, dass die Sequenz  $\Delta_1, \Delta_2, \Delta_3, \dots$  ebenfalls den Dirac-Stoss spezifiziert.

Man sieht, dass die ganze Geschichte des Dirac-Stosses in der Beziehung

$$\int_{-\infty}^{+\infty} \delta(t)w(t)dt = w(0) \quad (\text{A.4})$$

liegt, die für jede glatte Funktion  $w$  gilt. Man kann aber (A.4) nicht als Definition des Dirac-Stosses benutzen (obwohl genau das in vielen nachlässigen Herleitungen des Dirac-Stosses getan wird), es sei denn, es gebe eine echte Funktion  $\delta(\cdot)$ , so dass (A.4) für jede glatte Funktion  $w$  gilt. Leider gibt es keine solche Funktion, und genau deshalb führen wir überhaupt den Begriff der verallgemeinerten Funktion ein.

Obwohl für eine  $v$ -Funktion  $g$  der Grenzwert in (A.3) nur dann immer existiert, wenn  $w$  eine glatte Funktion ist, können wir unsere Regel (A.3) erweitern zu

$$\int_a^b g(t)f(t)dt \triangleq \lim_{n \rightarrow \infty} \int_a^b g_n(t)f(t)dt \quad (\text{A.5})$$

für alle  $a, b \in \mathbb{R}$  und für jede Funktion  $f$ , deren Grenzwert nach (A.5) existiert.

**Beispiel 5:** Aus Beispiel 3 oder Beispiel 4 ersehen wir, dass für den Dirac-Stoss

$$\int_{-\infty}^{+\infty} \delta(t)f(t)dt = f(0) \quad (\text{A.6})$$

gilt für jede Funktion  $f$ , die stetig ist im Punkt  $t = 0$ . Diese Eigenschaft (A.6) heisst die Siebungs-Eigenschaft des Dirac-Stosses.

**Beispiel 6:** Die Sequenz  $u_1, u_2, u_3, \dots$  von braven Funktionen in Beispiel 2 spezifiziert die  $v$ -Funktion  $u$ . Dabei gilt:

$$\begin{aligned} \int_{-\infty}^{+\infty} u(t)w(t)dt &= \lim_{n \rightarrow \infty} \int_{-\infty}^{+\infty} u_n(t)w(t)dt \\ &= \lim_{n \rightarrow \infty} \int_{-\frac{1}{n}}^n u_n(t)w(t)dt \\ &= \int_0^{\infty} w(t)dt \end{aligned} \quad (\text{A.7})$$

für jede glatte Funktion  $w$ .

Für die Einheitsschrittfunktion  $u(\cdot)$ , die wie folgt definiert ist:

$$u(t) = \begin{cases} 1, & t > 0 \\ 0, & t < 0 \\ \frac{1}{2}, & t = 0 \end{cases}$$

(wobei der Wert  $u(0)$  willkürlich gewählt wurde; oft setzt man auch  $u(0) = 1$ ), gilt auch

$$\int_{-\infty}^{+\infty} u(t)w(t)dt = \int_0^{\infty} w(t)dt$$

für jede glatte Funktion  $w$ . Also können wir sagen, dass die  $v$ -Funktion  $u$  in Beispiel 6 identisch ist mit der Einheitsschrittfunktion  $u(\cdot)$ . Eine  $v$ -Funktion kann also auch eine gewöhnliche Funktion sein. Tatsächlich gibt es für jede Funktion  $\gamma(\cdot)$ , die "physikalisch sinnvoll" ist, eine  $v$ -Funktion  $g$ , so dass  $g$  identisch ist mit der Funktion  $\gamma$ .

Andererseits gibt es  $v$ -Funktionen (z.B. den Dirac-Stoß), die mit keiner gewöhnlichen Funktion identisch sind. In diesem Sinne ist eine  $v$ -Funktion eine echte Verallgemeinerung einer Funktion.

Man kann weiter in diese Richtung gehen. Sei  $g$  eine  $v$ -Funktion und sei  $f$  eine Funktion, so dass

$$\int_a^b g(t)w(t)dt = \int_a^b f(t)w(t)dt \quad (\text{A.8})$$

für jede glatte Funktion  $w$  gilt, wobei  $a, b \in \mathbb{R}$  mit  $a < b$ ; dann sagt man, dass  $g$  und  $f$  im geschlossenen Intervall  $[a, b]$  identisch sind. Falls  $f$  stetig ist in  $[a, b]$ , sagt man sogar, dass

$$g(t) = f(t), \quad t \in [a, b]. \quad (\text{A.9})$$

In diesem Sinne (und nur in diesem Sinne) kann man manchmal (aber nicht immer!) vom Wert der  $v$ -Funktion  $g$  auf einem Punkt  $t$  sprechen.

**Beispiel 7:** Für  $a < b < 0$  gilt

$$\int_a^b \delta(t)w(t)dt = 0 = \int_a^b 0 \cdot w(t)dt.$$

Also gilt

$$\delta(t) = 0, \quad t < 0. \quad (\text{A.10})$$

Analog sieht man, dass

$$\delta(t) = 0, \quad t > 0 \quad (\text{A.11})$$

gilt. Aber es ist nicht möglich, den Wert " $\delta(0)$ " zu definieren. Sogar " $\delta(0) = +\infty$ " ist nicht sinnvoll, wie wir bald sehen werden.

[Gleichungen (A.10) und (A.11) zusammen mit (A.6) bilden die "Definition" des Dirac-Stosses in ungenauen Behandlungen dieser  $v$ -Funktion.]

Sei  $g_1, g_2, g_3, \dots$  eine Sequenz braver Funktionen, die eine v-Funktion  $g$  spezifizieren, und seien weiterhin diese braven Funktionen auch differenzierbar. (Somit ist auch  $g'_1, g'_2, g'_3, \dots$  eine Sequenz braver Funktionen.) Dann spezifiziert die Sequenz  $g'_1, g'_2, g'_3, \dots$  eine v-Funktion  $g'$ , welche die **Ableitung der v-Funktion  $g$**  heisst. Um zu sehen, dass  $g'_1, g'_2, g'_3, \dots$  tatsächlich eine v-Funktion spezifizieren, gehen wir wie folgt vor: Es gilt

$$\begin{aligned} \int_{-\infty}^{+\infty} g'_n(t)w(t)dt &= g_n(t)w(t) \Big|_{-\infty}^{+\infty} - \int_{-\infty}^{+\infty} g_n(t)w'(t)dt \\ &= - \int_{-\infty}^{+\infty} g_n(t)w'(t)dt \end{aligned} \quad (\text{A.12})$$

für jede glatte Funktion  $w$ , wobei wir (A.1) und (A.2) benutzt haben. Die Existenz des Grenzwertes (A.12) für  $n \rightarrow \infty$  folgt aus der Tatsache, dass  $g$  eine v-Funktion und  $w'(t)$  eine glatte Funktion ist.

Wir haben damit bewiesen:

Sei  $g'$  die Ableitung der v-Funktion  $g$ , dann ist  $g'$  die v-Funktion, für die gilt:

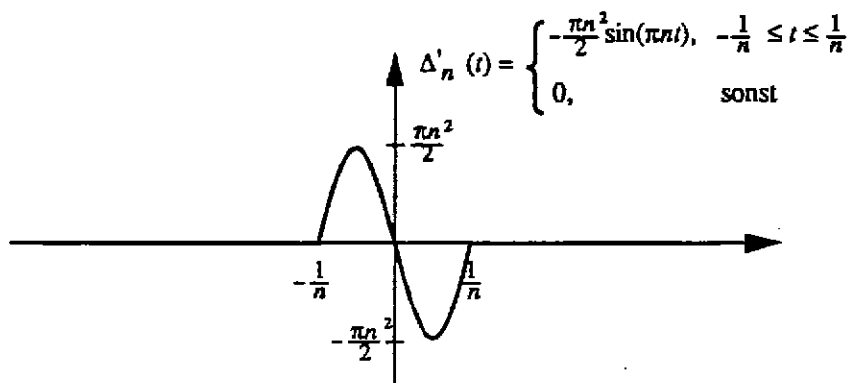
$$\int_{-\infty}^{+\infty} g'(t)w(t)dt = - \int_{-\infty}^{+\infty} g(t)w'(t)dt \quad (\text{A.13})$$

**Beispiel 8:** Für die Ableitung  $\delta'$  des Dirac-Stosses  $\delta$  gilt

$$\begin{aligned} \int_{-\infty}^{+\infty} \delta'(t)w(t)dt &= - \int_{-\infty}^{+\infty} \delta(t)w'(t)dt \\ &= -w'(0) \end{aligned} \quad (\text{A.14})$$

für jede glatte Funktion  $w$ , wobei wir (A.4) benutzt haben.

Der aufmerksame Leser bemerkt, dass wir (A.14) hergeleitet haben, ohne eine Sequenz von braven Funktionen zu finden, die den Dirac-Stoss  $\delta$  spezifiziert. Tatsächlich ist die brave Funktion  $\delta_n$  in Beispiel 2 nicht differenzierbar in den Punkten  $t = -\frac{1}{n}$ ,  $t = 0$  und  $t = \frac{1}{n}$ . Doch die brave Funktion  $\Delta_n$  in Beispiel 2 ist differenzierbar, und zwar wie folgt:



Also können wir  $\Delta'_1, \Delta'_2, \Delta'_3, \dots$  als eine Sequenz von braven Funktionen wählen, welche die v-Funktion  $\delta'$  spezifiziert.

Wie (A.10) und (A.11) aus (A.4) folgen, so folgt aus (A.14), dass

$$\delta'(t) = 0, \text{ alle } t \neq 0. \quad (\text{A.15})$$

Aber man kann  $\delta'(0)$  nicht sinnvoll festlegen, ebensowenig wie  $\delta(0)$ .

Wir überlassen dem Leser den (nichttrivialen) Beweis, dass man für jede v-Funktion  $g$  eine Sequenz  $g_1, g_2, g_3, \dots$  von differenzierbaren braven Funktionen finden kann, die  $g$  spezifiziert.

[Hinweis: Sei  $\gamma_1, \gamma_2, \gamma_3, \dots$  eine Sequenz von braven Funktionen, die nicht alle differenzierbar sind. Zeige, dass man eine Sequenz  $f_1, f_2, f_3, \dots$  von braven Funktionen finden kann mit folgenden Eigenschaften:

- (i)  $f_n$  ist nur in den (höchstens endlich vielen Punkten) nicht beliebig oft differenzierbar, wo auch  $\gamma_n$  nicht beliebig oft differenzierbar ist.
- (ii) In diesen Punkten  $t$  muss gelten:  $f'_n(t-) = -\gamma'_n(t-)$  und  $f'_n(t+) = -\gamma'_n(t+)$ .
- (iii) Ausserdem soll  $\lim_{n \rightarrow \infty} f_n(t) = 0$  für alle  $t \in R$  sein.

Es folgt dann, dass  $g_n = f_n + \gamma_n$  eine differenzierbare brave Funktion ist. Weiter folgt für alle  $a, b \in R$  mit  $a < b$ , dass  $\lim_{n \rightarrow \infty} \int_a^b f_n(t)w(t)dt = 0$  für jede glatte Funktion  $w$  ist, weil  $w$  stetig ist im geschlossenen Intervall  $[a, b]$ . Also gilt  $\lim_{n \rightarrow \infty} \int_{-\infty}^{+\infty} f_n(t)w(t)dt = 0$  für jede glatte Funktion  $w$ , und wir können schliessen, dass die v-Funktion  $f$  identisch ist mit der Nullfunktion. Daraus ergibt sich, dass  $g_1, g_2, g_3, \dots$  die gewünschte Sequenz ist.]

Wir schliessen daraus:

Die Ableitung  $g'$  einer v-Funktion  $g$  existiert immer.

Aus (A.13) und (A.14) folgt allgemein für die  $i$ -te Ableitung  $\delta^{(i)}$  des Dirac-Stosses  $\delta$ , dass

$$\int_{-\infty}^{+\infty} \delta^{(i)}(t)w(t)dt = (-1)^i w^{(i)}(0) \quad (\text{A.16})$$

für jede glatte Funktion  $w$  gilt. Analog zur Herleitung von (A.15) folgt auch, dass

$$\delta^{(i)}(t) = 0, \text{ } t \neq 0 \quad (\text{A.17})$$

für die  $i$ -te Ableitung von  $\delta$  gilt. Nun müsste auch ganz klar sein, dass  $\delta^{(i)}(0)$  nicht sinnvoll definiert werden kann!

**Beispiel 9:** Für die v-Funktion  $u$  (die identisch mit der Einheitsschrittfunktion ist) folgt aus (A.13), (A.7) und (A.1), dass

$$\begin{aligned} \int_{-\infty}^{+\infty} u'(t)w(t)dt &= - \int_{-\infty}^{+\infty} u(t)w'(t)dt = \\ &= - \int_0^{+\infty} w'(t)dt = \\ &= -w(t)|_{t=0}^{\infty} = \\ &= w(0) \end{aligned}$$

für jede glatte Funktion  $w$  gilt. Es folgt nun aus (A.4), dass  $\delta = u'$ , d.h.:



Der Dirac-Stoss ist die Ableitung der Einheitsschrittfunktion.

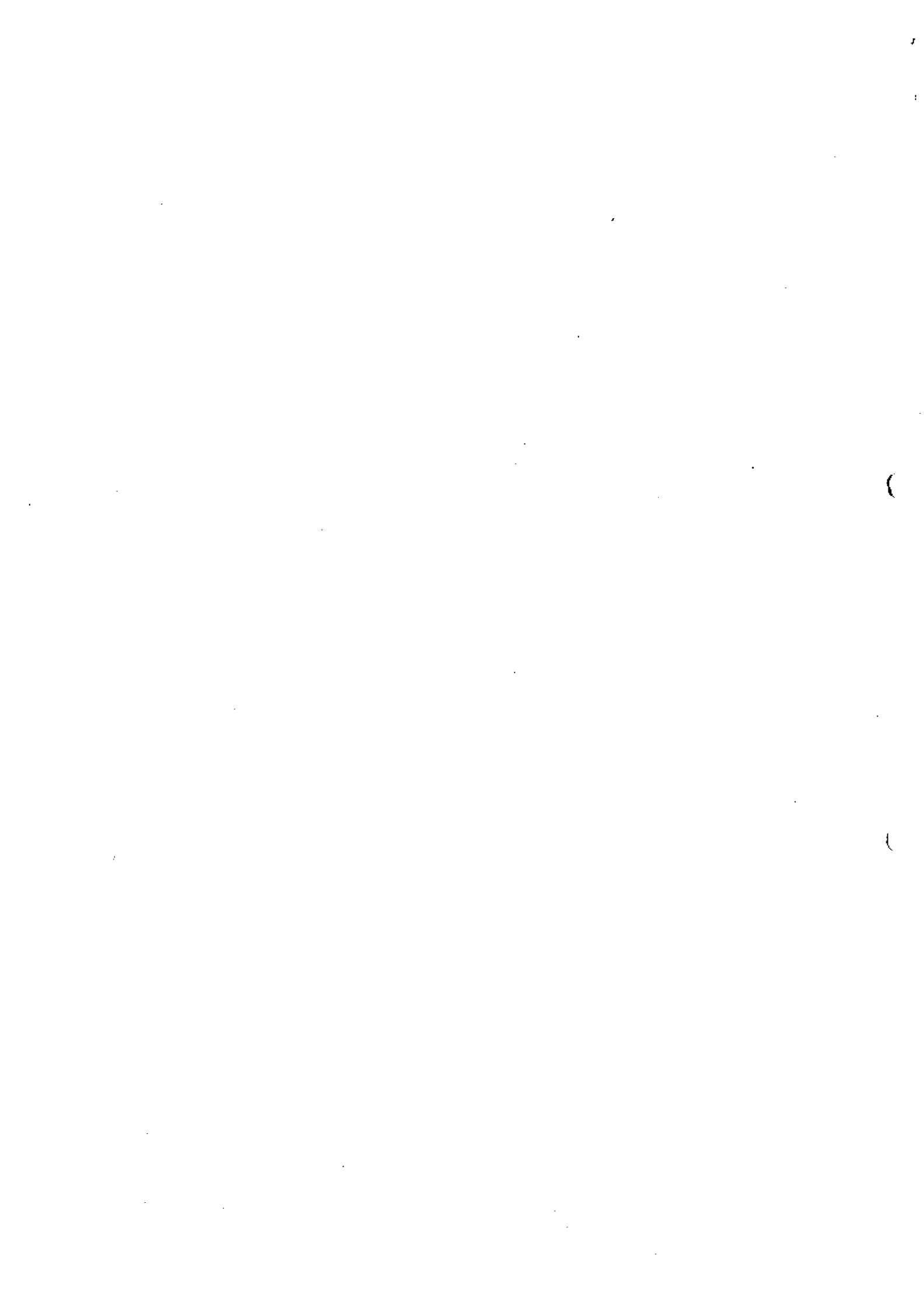
Der Leser hat sicher diese Tatsache bereits intuitiv beim Arbeiten mit dem Dirac-Stoss verwendet.

**Bemerkung:** Der Franzose *Laurent Schwartz* war der erste, der in seinem 1950-1951 veröffentlichten Buch [2] eine mathematisch saubere Behandlung des Dirac-Stosses und ähnlicher "verallgemeinerter Funktionen" vorgeführt hat. Seine Theorie wurde von zwei Engländern (*G. Temple* und *M.J. Lighthill*) sukzessive vereinfacht. Die obige Behandlung folgt im Geist dem Text von *Lighthill* (siehe [3]), enthält aber zusätzliche Vereinfachungen.

Der Leser sollte sich im Klaren sein, dass "glatte" und "brave" Funktionen keine Standardbegriffe in der Mathematik darstellen. Wir hoffen aber, dass diese zwei Begriffe ihm von Nutzen sind.

## Literatur

- [1] I.N. Bronstein, K.A. Semendjajew; "Taschenbuch der Mathematik"; Verlag Harri Deutsch, Thun, 1984.
- [2] L. Schwartz; "Théorie des distributions"; Hermann, Paris, 1950-1951; 2 Bände.
- [3] M.J. Lighthill; "Fourier Analysis and Generalized Functions"; Cambridge University Press, 1962.



## Sachverzeichnis

- $\sigma$ -Algebra 2, 55
- Äquivalenzklasse 59
- Autokorrelationsfunktion 72
- Autokorrelationssequenz 67
- AWGN 95
- Bayessche Formel 19, 37
- Bayessche Schätzung 63
- Beobachtung 24
- Bhattacharyya-Schranke 32
- Bode-Shannon-Trick 80, 82
- Cauchy-Schwarz-Ungleichung 44
- Cauchy-Ungleichung 45
- Delta-Funktion A-1
- Dimension 42
- Dirac-Stoss A-1
- Dot-Produkt 43
- Entscheidungsfunktion 24, 27
- Entscheidungsgebiet 26
- Entscheidungsproblem
  - Bayessches 26
  - nicht-Bayessches 30
- Entscheidungsregel 27, 28
  - optimale 25, 29
- Ereignis 2
  - sicheres 4
  - unmögliches 4
  - unvereinbare -se 18
- Ergebnisraum 1
- Erwartungswert 8, 11
  - bedingter 21
  - Satz vom totalen 27, 29, 38
  - totaler 21
- Faltungssumme 68, 70
- Faltung 68
- Fehlerwahrscheinlichkeit
  - bedingte 31
- Filter
  - FIR- 74, 93
  - inverses 82
  - kausales 81
  - rekursives 84
- Filterung
  - zeitinvariante 75
  - zeitvariante 84
- Filterungsproblem 74
- FIR-Filter 74, 93
- Frequenz, normierte 71
- Frequenzgang 71
- Frequenzkreis 71
- Funktion
  - brave A-2
  - glatte A-1
  - verallgemeinerte A-4
- Gaussverteilung 13, 14
- Gewichtssequenz 67
- Glättungsproblem 74
- Gram-Schmidt-Orthogonalisierungsverfahren 47
- i.i.d. 12
- Innovation 85
- Intersymbolinterferenz 96
- inverses Filter 82
- ISI 96
- Jacobische Determinante 16
- Kalman-Filter 84
  - Algorithmus 88
- Kalman-Gain 87
  - Algorithmus 90
- kausales Filter 81
- kausal 74
- Kolmogorov, Axiome von 3
- Kosten 26
- Kovarianzmatrix 13, 14
- Kovarianz 12
- Kreuzkorrelationssequenz 70
- Kronecker-delta-Sequenz 67
- LDS 67
- Leistungsdichte 73
- Leistung 73
- likelihood ratio 28, 30
  - test 28, 32, 34
- lineare Annäherung 49, 50
  - Linearitätseigenschaft 51
  - Trennungseigenschaft 52
- linear unabhängig 42
- lineares zeitdiskretes System 67
- MAP-Regel 25
- MAP-Schätzung 64
- Maximum-Likelihood-Regel 26
- mean-squared error 37

MINIMAX-Regel 34  
 minimum mean-squared error 35  
 Mittelwert 67  
 ML-Entscheidung 93, 95, 98, 99  
     Invertierbarkeitsprinzip 94  
 ML-Regel 26, 30  
     Grenze 31  
 ML-Schätzungsregel 40  
 ML-Schätzung 63  
 MMSE-Schätzung 36  
     lineare 54, 76  
     lineare im Gaussischen Fall 63  
     lineare mit zus. Datum 1 61  
     Linearitätseigenschaft der linearen 61  
     Orthog.gleichung für lineare 60  
     Orthog.prinzip für die lineare 60  
     Trennungseigenschaft der linearen 61  
 MMSE 38  
 MSE 37, 54, 77  
 Multiplikationsregel 18  
 Neyman-Pearson, Satz von 32, 34  
 Normalverteilung 13  
 Norm 48  
 Orthogonalbasissatz 46  
 Orthogonalitätsgleichung 50, 51  
 Orthogonalkomplement 45  
 orthogonal 45, 46, 50, 59  
 positive Definitheit 43, 48  
 Projektion 49  
 rekursives Filter 84  
 s.s. 66  
 schwach stationär 66  
     gemeinsam 70  
 Schätzung, lineare MMSE- 76  
 Schätzungsregel 36, 38  
     maximum likelihood 40  
 Siebungs-Eigenschaft A-5  
 Skalarproduktraum 44, 59  
     Zerlegungssatz für einen 47  
 Skalarprodukt 43  
 Spesen 26, 29  
 stationär 66  
 stochastischer Prozess  
     Gaussischer 67  
     zeitdiskreter 66  
 Trellis 97  
 Übertragungsfunktion 71  
 unabhängig 11, 12, 13  
 unkorreliert 12, 13, 14  
 Unterraum 42, 45  
     Test für einen 42  
 v-Funktion A-4, A-7  
     Ableitung einer A-7  
 Varianz 9  
     bedingte 37  
 Vektorraum 41, 56  
     reeller 54  
 Verteilungsdichte 10  
 Verteilungsfunktion 5, 10  
     bedingte 19  
 Viterbi-Algorithmus 98, 99  
     Komplexität 99  
 Vorhersageproblem 74  
 Wahrscheinlichkeit  
     a priori 26, 30  
     a posteriori 25  
     bedingte 17  
     Satz von der totalen 19, 25  
     totale 21  
 Wahrscheinlichkeitsdichte 6  
     bedingte 20, 22  
 Wahrscheinlichkeitsfunktion 7, 11  
     bedingte 22  
 Wahrscheinlichkeitsmass 3  
     bedingtes 17  
 weisses Rauschen 67  
 Whitening-Filter 81, 82  
 Wiener-Filter  
     kausales 80, 83  
     nicht-kausales 78  
 Wiener-Hopf-Gleichung 76  
     kausale 80  
     nicht-kausale 78  
 Wiener-Khintchine-Satz 73  
 z-Transformation 71, 78  
 z-Transformierte 81  
 Zufallsgrösse 5  
     Transformation von  $\tilde{z}$ -n 17  
 Zufallsvektor 10  
 Zustandsmaschine, endliche 96  
 Zustandsübergangsdiagramm 96