

# Informationstheorie

## Diskrete Modelle und Verfahren

Professor Dr. Rudolf Mathar

Lehrstuhl für Theoretische Informationstechnik

Rheinisch-Westfälische Technische Hochschule Aachen

<http://www.ti.rwth-aachen.de>

Erschienen im Teubner-Verlag: Informationstheorie, diskrete Modelle und Verfahren. ISBN 3-519-02574-4, Stuttgart, 1996.

## Vorwort

Das vorliegende Buch ist aus Vorlesungen entstanden, die ich an der Universität Augsburg und der RWTH Aachen gehalten habe. Aus der Fülle des vorliegenden Materials im Bereich der Informationstheorie sind in vier Kapiteln wichtige Begriffsbildungen und Ergebnisse zusammengestellt. Ausgangspunkt ist der von Shannon geprägte Begriff der Entropie. Er ermöglicht die theoretischen Untersuchungen zur Kodierung diskreter Quellen und die abstrakte Bewertung der Übertragungsgüte von diskreten Kanälen. Den Abschluß bildet ein Kapitel über fehlerkorrigierende Codes. Diese sehr kurze Einführung in die Kodierungstheorie scheint mir wichtig, um wenigstens einige Konzepte und Verfahren zur Konstruktion fehlerkorrigierender Codes mit kleiner Irrtumswahrscheinlichkeit bereitzustellen, deren Existenz durch den Shannonschen Fundamentalsatz gesichert wird. Wegen ihres hohen Datendurchsatzes und der effizienten Implementierbarkeit sind hierbei die Faltungskodierer mit ihrer Trellis-Darstellung und dem Viterbi-Algorithmus zur Dekodierung besonders behandelenswert.

Nach meiner Erfahrung deckt der Stoff eine vierstündige einführende Vorlesung über Informationstheorie ab, bei der die konstruktive Kodierungstheorie allerdings nur kurz angeschnitten wird. Zum Abschluß jedes Kapitels finden sich Übungsaufgaben, die nach Studium des zugehörigen Stoffs ohne weitere Hilfsmittel gelöst werden können. Insofern eignet sich das vorliegende Buch sowohl als Begleitmaterial zu einer Vorlesung als auch zum Selbststudium. Im Text sind parallel zu den deutschen Fachbegriffen auch die englischen eingeführt, sofern sie sich nicht nur leicht in der Schreibweise unterscheiden. Diese finden sich auch im Index wieder. Ich möchte hiermit Anfängern den Einstieg in die englische Literatur erleichtern.

Wichtige Konzepte zur Modellierung von Nachrichtenquellen und Übertragungskanälen stammen aus der Wahrscheinlichkeitstheorie. Für den vorliegenden Text reichen in den meisten Fällen diskrete Wahrscheinlichkeitsverteilungen aus. Die wichtigsten Definitionen und Sätze werden in einem

einführenden Kapitel bereitgestellt. Dieser Abschnitt ist allerdings nur als kurze Wiederholung gedacht, eine gründliche Erarbeitung des Stoffs sollte in einführenden Vorlesungen oder Büchern über “Stochastik” erfolgen.

Ich hoffe, mit dem vorliegenden Buch eine Stoffauswahl anzubieten, die Informatikern, Mathematikern und auch an Grundlagen interessierten Elektrotechnikern eine gründliche Einführung in diskrete Modelle der Informationstheorie ermöglicht. Auf Exaktheit und Motivation der eingeführten Begriffe wird hierbei besonderer Wert gelegt.

Dr. Gerd Brücks und Dr. Rolf Hager haben durch gewissenhaftes Durchsehen von Vorgängerversionen des Manuskripts einige Unklarheiten und Schreibfehler in der endgültigen Fassung verhindert, wofür ich ihnen herzlich danke. Weiterhin möchte ich Dr. Jürgen Mattfeldt und Dr. Thomas Niessen danken, die einige interessante Übungsaufgaben beigesteuert und durch ihre Kritik zur Verbesserung der endgültigen Fassung beigetragen haben.

Aachen, im SS 1996

Rudolf Mathar

# Inhalt

<b>Vorwort</b>	<b>3</b>
<b>1 Einleitung</b>	<b>7</b>
<b>2 Stochastische Grundlagen</b>	<b>11</b>
2.1 Zufallsvariable und ihre Verteilung . . . . .	11
2.2 Markoff-Ketten . . . . .	16
2.3 Übungsaufgaben . . . . .	19
<b>3 Information und Entropie</b>	<b>22</b>
3.1 Entropie und Transinformation . . . . .	24
3.2 Axiomatische Charakterisierung der Entropie . . . . .	35
3.3 Übungsaufgaben . . . . .	39
<b>4 Kodierung diskreter Quellen</b>	<b>43</b>
4.1 Codes fester Länge . . . . .	47
4.2 Codes variabler Länge . . . . .	52
4.3 Binäre Suchbäume . . . . .	69
4.4 Stationäre Quellen, Markoff-Quellen . . . . .	76
4.5 Übungsaufgaben . . . . .	85
<b>5 Diskrete gedächtnislose Kanäle</b>	<b>88</b>
5.1 Kanalkapazität . . . . .	89
5.2 Kanaldekodierung . . . . .	97
5.3 Der Shannonsche Fundamentalsatz . . . . .	102
5.4 Kaskadenkanäle und Umkehrung des Fundamentalsatzes . . .	115
5.5 Übungsaufgaben . . . . .	122

6 Inhalt

<b>6</b>	<b>Fehlerkorrigierende Codes</b>	<b>126</b>
6.1	Blockcodes und Hamming-Distanz . . . . .	126
6.2	Lineare Codes . . . . .	129
6.3	Faltungskodes und der Viterbi-Algorithmus . . . . .	135
6.4	Übungsaufgaben . . . . .	144
<b>7</b>	<b>Anhang: endliche Körper</b>	<b>146</b>
	<b>Literatur</b>	<b>150</b>
	<b>Index</b>	<b>153</b>

# 1 Einleitung

Mit dem Schlagwort ‘Informationszeitalter’ wird die Wirkung der wichtigsten Ressource der heutigen Zeit charakterisiert. In der Tat scheint es so zu sein, daß der leichte Zugang zu vielfältigen Informationen den wissenschaftlichen Fortschritt zumindest in technischen Bereichen erheblich beschleunigt. Natürlich ist ‘Information’ ein weitspannender Begriff, dessen Bedeutung nur aus dem jeweiligen Kontext hervorgeht. Informationstheorie behandelt nicht den Inhalt oder die Bedeutung von ‘Informationen’, sondern Systeme, Modelle und Methoden zur Beschreibung der Entstehung und Übermittlung von Daten.

Informationsübertragende Systeme lassen sich nach ihrer Funktionsweise in analoge und digitale unterscheiden. Beispiele für analoge, kontinuierliche Informationsübertragung sind Radio, Fernsehen, Schallplatten, Telefon und motorische Reize auf Nervenbahnen bei biologischen Systemen. Digitale, diskrete Übertragung findet zum Beispiel bei Compact Discs, Digital Audio Tapes, Computernetzwerken, digitaler Sprach- und Datenübertragung beim Mobilfunk und bei digitaler Bildübertragung auf Satellitenkanälen statt. Bei einigen dieser Beispiele können die Grenzen jedoch nicht scharf gezogen werden, da moderne Kommunikationssysteme wegen der besseren Qualität, den effizienteren Fehlerkorrekturverfahren und der leichteren Implementierbarkeit zunehmend auf digitaler Übertragungstechnik basieren.

Wir werden uns mit diskreten Systemen bei der Erzeugung und Übertragung von Signalen beschäftigen. Typischerweise hat man bei solchen Systemen eine Quelle, die zu diskreten Zeitpunkten Nachrichten oder Signale aussendet. Diese Signale werden durch einen Quellenkodierer komprimiert. Für Ausgabeblocke von Buchstaben des Quellenkodierers stellt der Kanalkodierer dann eine Menge von Kodewörtern bereit, die fehlerrobust über einen gestörten Kanal übertragen werden können. Nach Übertragung wird dieser Vorgang wieder umgekehrt. Zunächst dekodiert der Kanaldekodierer in die ursprünglichen Blöcke, die der Quellendekodierer beim Ziel in der dort lesbaren Form abliefern.

## 8 1 Einleitung

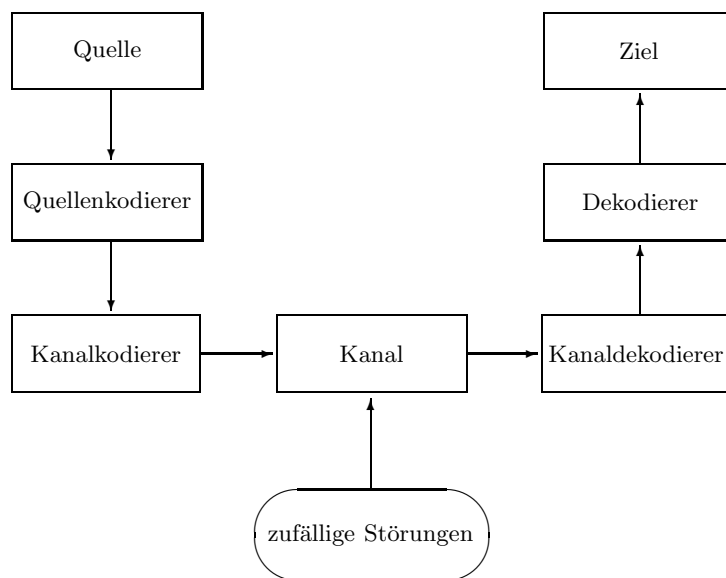


Abb. 1.1 Das Modell der Informationsübertragung

Zwei verschiedene Zielsetzungen der Kodierung kristallisieren sich bei diesem Prozeß heraus. Im Quellenkodierer lautet die Aufgabe, durch Entfernen von Redundanz Information zu verdichten, um mit möglichst wenig Aufwand die angebotenen Quellwörter übertragen zu können. Der Kanalkodierer hat zum Ziel, einen Code bereitzustellen, dessen Kodewörter auch nach Übertragung im gestörten Kanal trotz des Auftretens von Fehlern noch identifiziert werden können. Hier wird gezielt wieder Redundanz hinzugefügt, allerdings so, daß mit den überflüssigen Bits eine möglichst große Zahl von Übertragungsfehlern korrigiert werden kann. Ist die Quelle zum Beispiel ein Terminal, das Buchstaben im ASCII-Kode erzeugt, könnte als Quellenkodierer eines der bekannten Datenkompressionsprogramme verwendet werden. Ein Faltungskodierer würde Blöcke des komprimierten Files mit Redundanzbits zur eventuellen Fehlerkorrektur versehen. Diese vor Fehlern geschützte Information wird übertragen, beim Empfänger dekodiert und wieder dekomprimiert. Abbildung 1.1 zeigt den Ablauf der Übertragung als Blockschaltbild.

Die Aufgabe lautet nun, geeignete mathematische Modelle zu finden, die die einzelnen Komponenten und ihr Zusammenwirken beschreiben. Häufig wer-

den dazu stochastische Modelle benutzt, insbesondere dann, wenn zufällige Einflüsse wirken. Zufall kommt zum Beispiel ins Spiel, wenn nur eine Wahrscheinlichkeitsverteilung der ausgesendeten Signale bekannt ist (etwa die Verteilung der Zeichen des lateinischen Alphabets in einem deutschen Text) oder zufällige Störungen von Signalen bei der Übertragung im Kanal (Rauschen) auftreten.

Hauptanliegen dieses Textes sind mathematische Modelle von Kommunikationssystemen. Diese orientieren sich unmittelbar an den physikalischen Systemen, die Daten erzeugen oder übertragen. Stochastische Modelle und Methoden sind besonders geeignet, um die Zufallsphänomene bei der Entstehung und Übertragung der Daten zu beschreiben. Dies war der Startpunkt der Theorie, die C. E. Shannon mit seinem 1948 veröffentlichten Aufsatz “A mathematical theory of communication” (siehe [33] oder [34]) begründet hat. Das Hauptaugenmerk lag zunächst auf der funktionellen Beschreibung von Kodierern im Rahmen der Stochastik und Nachweisen ihrer Existenz bzw. Nichtexistenz unter bestimmten Qualitätsanforderungen. Später hinzugekommen sind Anwendungen in der Kryptologie, für die Unsicherheit und Redundanz zur abstrakten Messung von Sicherheit eine zentrale Rolle spielen. Wichtige Auswirkungen hat die Theorie auch auf die Entwicklung von fehlerkorrigierenden Codes, die Untersuchungsgegenstand der Kodierungstheorie sind.

Die Informationstheorie hat sich weiterhin schwunghaft entwickelt. Sie hat sich je nach Setzung der Schwerpunkte

- als Bereich der Wahrscheinlichkeitstheorie, insbesondere bei ergodischen Prozessen, und mit modifizierten Konzepten auch in der Statistik etabliert (siehe [4], [22]),
- auf der algebraischen Seite zu konstruktiven Methoden in der Kodierungstheorie geführt (siehe [16], [30], [37]),
- für die Kryptologie grundlegende Konzepte bereitgestellt (siehe [36]), und
- sie beschäftigt sich auf der angewandten Seite mit Komplexitätsuntersuchungen und Implementierungsfragen (siehe [2], [7]).

Der Text beginnt mit einem einführenden Kapitel in die stochastischen Grundlagen, das eher zur Festlegung der forthin verwendeten Notation als zur Erarbeitung des Stoffs gedacht ist. Anschließend wird das Blockschaltbild aus Abbildung 1.1 stufenweise mit Leben erfüllt. Zur Charakterisierung



des stochastischen Verhaltens von Quellen wird in Kapitel 3 die Entropie mit ihren vielfältigen Eigenschaften eingeführt. Datenkompression und optimale Kodierung sind die Ziele des Quellenkodierers in Kapitel 4. Die anschließende Übertragung in gestörten Kanälen wird in Kapitel 5 behandelt. Algebraische und algorithmische Aspekte treten in den Vordergrund, wenn in Kapitel 6 Codes entwickelt werden, die Fehler korrigieren können.

Mit dem Zeichen ■ wird durchgehend das Ende von Beweisen und Beispielen markiert.

## 2 Stochastische Grundlagen

In den nachfolgenden Abschnitten werden grundlegende Kenntnisse der Wahrscheinlichkeitstheorie als bekannt vorausgesetzt. Insbesondere sollten diskrete Zufallsvariable und Folgen von Zufallsvariablen, der Erwartungswert diskreter Zufallsvariablen und elementare bedingte Verteilungen bekannt sein. Wir benötigen weiterhin das Gesetz großer Zahlen, den Begriff einer stationären Folge von Zufallsvariablen sowie einige elementare Sachverhalte zu Markoff-Ketten.

Die zugehörigen Definitionen werden im folgenden Überblick zusammengestellt. Dieses Kapitel dient eher dazu, durch kurzes Nachlesen die dargestellten Begriffe zu wiederholen als zu einer Erarbeitung der Materie. Eine gründliche Behandlung des Stoffs findet sich zum Beispiel in [24].

### 2.1 Zufallsvariable und ihre Verteilung

Zufallsvariable sind ein bequemes Hilfsmittel, um nicht deterministisch vorhersagbare Ergebnisse zu beschreiben. Die genaue Kenntnis des Wahrscheinlichkeitsraum, auf dem eine Zufallsvariable definiert ist, ist meist nicht von Interesse, von Bedeutung ist oft nur seine Existenz. Ungleich wichtiger in stochastischen Modellen ist die Verteilung von Zufallsvariablen, die in der folgenden Definition eingeführt wird.

**Definition 2.1** (*Zufallsvariable*)

Seien  $(\Omega, \mathcal{A}, P)$  ein Wahrscheinlichkeitsraum und  $(\mathcal{X}, \mathcal{B})$  ein Meßraum. Eine Abbildung  $X : \Omega \rightarrow \mathcal{X}$  heißt Zufallsvariable (ZV), wenn

$$X^{-1}(B) \in \mathcal{A} \text{ für alle } B \in \mathcal{B}. \quad (2.1)$$

Bedingung (2.1) heißt Meßbarkeit der Zufallsvariablen  $X$ . Für meßbare Zufallsvariablen ist für alle  $B \in \mathcal{B}$  wohldefiniert

$$P^X(B) = P(X^{-1}(B)) = P(\{\omega \in \Omega \mid X(\omega) \in B\}) = P(X \in B).$$

$P^X$  ist eine Wahrscheinlichkeitsverteilung auf  $(\mathcal{X}, \mathcal{B})$ , sie heißt Verteilung der Zufallsvariablen  $X$ .

Um hervorzuheben, daß eine Zufallsvariable  $X$  Abbildung von  $\Omega$  nach  $\mathcal{X}$  und zusätzlich meßbar ist, wird die Notation  $X : (\Omega, \mathcal{A}) \rightarrow (\mathcal{X}, \mathcal{B})$  benutzt.

**Definition 2.2** (diskrete Zufallsvariable)

Eine Zufallsvariable  $X$  heißt diskret, wenn eine höchstens abzählbare Menge  $\mathcal{T}$  existiert mit  $P^X(\mathcal{T}) = P(X \in \mathcal{T}) = 1$ .  $\mathcal{T} = \mathcal{T}_X$  heißt Träger von  $X$  bzw. von  $P^X$ .

Besitzt eine Zufallsvariable  $X$  die Verteilung  $P^X$ , so schreiben wir  $X \sim P^X$ . Ist die Verteilung von  $X$  diskret mit  $P(X = x_i) = p_i$ ,  $i = 1, \dots, m$ , und ist  $\mathbf{p} = (p_1, \dots, p_m)$  der zugehörige stochastische Vektor, so bedeutet  $X \sim \mathbf{p}$ , daß  $X$  die durch  $\mathbf{p}$  beschriebene Verteilung besitzt.

**Satz 2.1** (gemeinsame Verteilung, Zufallsvektoren)

Seien  $X_1, \dots, X_n$  Zufallsvariable, definiert auf einem Wahrscheinlichkeitsraum  $(\Omega, \mathcal{A}, P)$ ,  $X_i$  jeweils mit Werten in Meßräumen  $(\mathcal{X}_i, \mathcal{B}_i)$ . Seien  $\mathcal{X} = \times_{i=1}^n \mathcal{X}_i$  und  $\mathcal{B} = \otimes_{i=1}^n \mathcal{B}_i$  die Produkt- $\sigma$ -Algebra. Die gemeinsame Verteilung  $P^{\mathbf{X}} = P^{(X_1, \dots, X_n)}$  des Zufallsvektors  $\mathbf{X} = (X_1, \dots, X_n)$  auf  $(\mathcal{X}, \mathcal{B})$  ist bereits eindeutig bestimmt durch

$$P^{\mathbf{X}}(B_1 \times \dots \times B_n) = P(X_1 \in B_1, \dots, X_n \in B_n)$$

für alle  $B_i \in \mathcal{B}_i$ ,  $i = 1, \dots, n$ .

Analog zu obigem bedeutet  $\mathbf{X} \sim P^{\mathbf{X}}$ , daß der Zufallsvektor  $\mathbf{X}$  die Verteilung  $P^{\mathbf{X}}$  besitzt.

**Satz 2.2** (Folgen von Zufallsvariablen)

Sei  $\{X_n\}_{n \in \mathbb{N}}$  eine Folge von Zufallsvariablen auf einem Wahrscheinlichkeitsraum  $(\Omega, \mathcal{A}, P)$ ,  $X_n$  jeweils mit Werten in Meßräumen  $(\mathcal{X}_n, \mathcal{B}_n)$ . Seien  $\mathcal{X} = \times_{i=1}^{\infty} \mathcal{X}_i$  und  $\mathcal{B} = \otimes_{i=1}^{\infty} \mathcal{B}_i$  die Produkt- $\sigma$ -Algebra. Die gemeinsame Verteilung  $P^{\{X_n\}}$  der Folge auf  $(\mathcal{X}, \mathcal{B})$  ist eindeutig bestimmt durch

$$\begin{aligned} P^{\{X_n\}}(B_1 \times \dots \times B_n \times \mathcal{X}_{n+1} \times \dots) &= P^{(X_1, \dots, X_n)}(B_1 \times \dots \times B_n) \\ &= P(X_1 \in B_1, \dots, X_n \in B_n) \end{aligned}$$

für alle  $n \in \mathbb{N}$ ,  $B_i \in \mathcal{B}_i$ ,  $i = 1, \dots, n$ .  $P^{\{X_n\}}$  wird also bereits durch die endlich-dimensionalen Randverteilungen eindeutig definiert.

**Definition 2.3** (stochastische Unabhängigkeit)

Zufallsvariablen  $X_1, \dots, X_n$  heißen stochastisch unabhängig (s.u.), wenn

$$P^{(X_1, \dots, X_n)}(B_1 \times \dots \times B_n) = P^{X_1}(B_1) \dots P^{X_n}(B_n)$$

für alle  $B_i \in \mathcal{B}_i$  oder äquivalent

$$P(X_1 \in B_1, \dots, X_n \in B_n) = P(X_1 \in B_1) \dots P(X_n \in B_n)$$

für alle  $B_i \in \mathcal{B}_i$ .

Eine Folge von Zufallsvariablen  $\{X_n\}_{n \in \mathbb{N}}$  heißt stochastisch unabhängig, wenn die Zufallsvariablen  $X_1, \dots, X_n$  für alle  $n \in \mathbb{N}$  stochastisch unabhängig sind.

Für diskrete Zufallsvariable vereinfacht sich Definition 2.3 auf die folgende Beziehung.

**Lemma 2.1** Seien  $X_1, \dots, X_n$  diskret mit Trägern  $\mathcal{T}_1, \dots, \mathcal{T}_n$ .  $X_1, \dots, X_n$  sind stochastisch unabhängig genau dann, wenn

$$P(X_1 = x_1, \dots, X_n = x_n) = P(X_1 = x_1) \dots P(X_n = x_n)$$

für alle  $x_i \in \mathcal{T}_i$ ,  $1 \leq i \leq n$ .

Stochastische Unabhängigkeit von Zufallsvariablen ist also äquivalent zur stochastischen Unabhängigkeit der Ereignisse  $\{X_1 \in B_1\}, \dots, \{X_n \in B_n\}$  im Wahrscheinlichkeitsraum  $(\Omega, \mathcal{A}, P)$  für alle  $B_i \in \mathcal{B}_i$  im Sinn der bekannten Durchschnitts-Produkt-Beziehung.

Das folgende Lemma gestattet die Überprüfung der stochastischen Unabhängigkeit induktiv jeweils durch Anhängen einer weiteren unabhängigen Zufallsvariablen an einen Vektor aus stochastisch unabhängigen Komponenten.

**Lemma 2.2** Zufallsvariable  $X_1, \dots, X_n$  sind stochastisch unabhängig genau dann, wenn  $(X_1, \dots, X_j)$  und  $X_{j+1}$  stochastisch unabhängig sind für alle  $j = 1, \dots, n - 1$ .

Stochastische Unabhängigkeit der gesendeten Buchstaben ist für viele Quellen zu restriktiv. In natürlichen Sprachen bestehen offensichtlich Abhängigkeiten zwischen aufeinanderfolgenden Buchstaben. Ein adäquates Modell sind stationäre Folgen von Zufallsvariablen, bei denen die Invarianz aller endlich-dimensionalen Randverteilungen gegen eine Verschiebung des Indexbereichs gefordert wird. Interpretiert man die Indizes als Zeitparameter, ergibt sich die zeitliche Invarianz des stochastischen Verhalten der Quelle, d.h. die Wahrscheinlichkeit für das Auftreten einer bestimmten Buchstabensequenz ab der Zeit  $n = 0$  ist dieselbe wie ab jeder anderen Zeit  $n > 0$ .

**Definition 2.4** (*Stationarität*)

Eine Folge von Zufallsvariablen  $\{X_n\}_{n \in \mathbb{N}}$  heißt stationär, wenn

$$P(X_{i_1}, \dots, X_{i_n}) = P(X_{i_1+s}, \dots, X_{i_n+s})$$

für alle  $1 \leq i_1 < i_2 < \dots < i_n$ ,  $s \in \mathbb{N}$ .

Erwartungswerte und bedingte Verteilungen können für Zufallsvariable mit höchstens abzählbarem Wertebereich elementar erklärt werden. Die folgenden Definitionen beziehen sich nur auf diesen Fall.

**Definition 2.5** (*Erwartungswert diskreter Zufallsvariablen*)

Sei  $X$  eine diskrete Zufallsvariable mit höchstens abzählbarem Träger  $\mathcal{T} = \{t_1, t_2, \dots\} \subset \mathbb{R}$ . Falls  $\sum_{i=1}^{\infty} |t_i| P(X = t_i) < \infty$ , heißt

$$E(X) = \sum_{i=1}^{\infty} t_i P(X = t_i) \quad (2.2)$$

Erwartungswert von  $X$ .

Bei Hintereinanderausführung einer diskreten Zufallsvariablen  $X$  und einer reellwertigen, meßbaren Abbildung  $h$ , also  $(\Omega, \mathcal{A}, P) \xrightarrow{X} (\mathcal{X}, \mathcal{B}) \xrightarrow{h} \mathbb{R}$ , läßt sich der Erwartungswert wie folgt berechnen. Falls  $\sum_{i=1}^{\infty} |h(t_i)| P(X = t_i) < \infty$  ist, gilt

$$E(h(X)) = \sum_{i=1}^{\infty} h(t_i) P(X = t_i).$$

**Definition 2.6** (bedingte Verteilung)

Seien  $X, Y$  diskrete Zufallsvariablen auf einem Wahrscheinlichkeitsraum  $(\Omega, \mathcal{A}, P)$ . Durch

$$P(X = x | Y = y) = \begin{cases} \frac{P(X=x, Y=y)}{P(Y=y)}, & \text{falls } P(Y = y) > 0 \\ P(X = x), & \text{falls } P(Y = y) = 0 \end{cases}, \quad x \in \mathcal{T}_X,$$

wird die bedingte Verteilung von  $X$  unter (der Hypothese)  $\{Y = y\}$  definiert.

Offensichtlich gilt für alle  $x \in \mathcal{T}_X$ , daß

$$P(X = x) = \sum_{y \in \mathcal{T}_Y} P(X = x | Y = y) P(Y = y),$$

unabhängig von der speziellen Festsetzung der bedingten Verteilung im Fall  $P(Y = y) = 0$  in Definition 2.6. Ist  $P(Y = y) = 0$ , kann die bedingte Verteilung beliebig gewählt werden, ohne daß dies Einfluß auf die Berechnung der Verteilung von  $X$  mit obiger Formel hätte.

**Definition 2.7** (bedingter Erwartungswert)

Seien  $X, Y$  diskrete, reellwertige Zufallsvariable auf einem Wahrscheinlichkeitsraum  $(\Omega, \mathcal{A}, P)$ .

$$E(X | Y = y) = \sum_{x \in \mathcal{T}_X} x P(X = x | Y = y)$$

heißt bedingter Erwartungswert von  $X$  unter  $\{Y = y\}$ .

Man beachte, daß  $E(X | Y = y) = h(y)$  eine Funktion von  $y$  ist. Dies definiert eine Zufallsvariable  $h(Y) = E(X | Y)$ .  $h(Y)$  heißt bedingte Erwartung von  $X$  unter  $Y$ . Bildet man den Erwartungswert der Zufallsvariablen  $h(Y)$ , so ergibt sich der Erwartungswert von  $X$ . Dies sieht man an der folgenden Formel.

$$\begin{aligned} E(E(X | Y)) &= \sum_y \left( \sum_x x P(X = x | Y = y) \right) P(Y = y) \\ &= \sum_x x \sum_y P(X = x, Y = y) = \sum_x x P(X = x) = E(X). \end{aligned} \quad (2.3)$$

**Satz 2.3** (Starkes Gesetz großer Zahlen, SGGZ)

$\{X_n\}_{n \in \mathbb{N}}$  sei eine Folge von stochastisch unabhängigen Zufallsvariablen auf einem Wahrscheinlichkeitsraum  $(\Omega, \mathcal{A}, P)$ , identisch verteilt mit  $P^{X_n} = P^{X_1}$  für alle  $n \in \mathbb{N}$ , und es existiere  $E(X_1) = \mu$ . Dann gilt:

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i = \mu \quad P\text{-fast sicher, d.h.}$$

$$P\left(\left\{\omega \mid \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{i=1}^n X_i(\omega) = \mu\right\}\right) = 1.$$

Satz 2.3 besagt, daß das arithmetische Mittel von Zufallsvariablen ‘fast sicher’ oder ‘mit Wahrscheinlichkeit 1’ gegen den Erwartungswert konvergiert. Man sagt dann auch,  $\{X_n\}_{n \in \mathbb{N}}$  genügt dem starken Gesetz großer Zahlen (SGGZ).

Da jede  $P$ -fast sicher konvergente Folge von Zufallsvariablen auch  $P$ -stochastisch konvergiert, folgt unter den Voraussetzungen von Satz 2.3 die stochastische Konvergenz von  $\frac{1}{n} \sum_{i=1}^n X_i$ , also

$$\lim_{n \rightarrow \infty} P\left(\left|\frac{1}{n} \sum_{i=1}^n X_i - \mu\right| > \varepsilon\right) = 0 \quad (2.4)$$

für alle  $\varepsilon > 0$ . Man sagt,  $\{X_n\}_{n \in \mathbb{N}}$  genügt dem schwachen Gesetz großer Zahlen, wenn 2.4 gilt.

## 2.2 Markoff-Ketten

Bei natürlichen Sprachen ist die Wahrscheinlichkeit für das Auftreten eines Buchstabens abhängig davon, welche Buchstaben als Vorgänger verwendet wurden. Im Deutschen ist die Wahrscheinlichkeit sehr klein, daß nach der Kombination ‘al’ am Wortanfang ein ‘d’ oder ‘e’ auftritt. Dies kann dann nur ein Schreibfehler sein, wovon man sich mit Hilfe eines Lexikons überzeugen kann. Markoff-Ketten sind ein geeignetes Modell, um stochastische Abhängigkeit zwischen zeitlich aufeinanderfolgenden Ereignissen zu beschreiben. Zunächst wird lediglich die Abhängigkeit auf den unmittelbaren Vorgängerzeitpunkt berücksichtigt, durch Vergrößern des Zustandsraum

auf kartesische Produkte können jedoch auch weiterreichende Abhängigkeiten modelliert werden. Diese Methode wird in Kapitel 4 bei Markoff-Quellen eingesetzt.

**Definition 2.8** (Markoff-Kette)

Eine Folge von diskreten Zufallsvariablen  $\{X_n\}_{n \in \mathbb{N}_0}$ , alle mit demselben höchstens abzählbaren Träger  $\mathcal{S} = \mathcal{T}_{X_n}$ , heißt Markoff-Kette, wenn

$$\begin{aligned} P(X_n = x_n \mid X_{n-1} = x_{n-1}, \dots, X_0 = x_0) \\ = P(X_n = x_n \mid X_{n-1} = x_{n-1}) \end{aligned} \quad (2.5)$$

für alle  $n \in \mathbb{N}$  und alle  $x_0, \dots, x_n \in \mathcal{S}$  mit  $P(X_{n-1} = x_{n-1}, \dots, X_1 = x_1) > 0$ .  $\mathcal{S}$  heißt Zustandsraum der Markoff-Kette,  $P(X_n = x_n \mid X_{n-1} = x_{n-1})$  Übergangswahrscheinlichkeit von  $x_{n-1}$  nach  $x_n$  im  $n$ -ten Schritt und  $P^{X_0}$  Anfangsverteilung der Markoff-Kette. Eine Markoff-Kette heißt homogen, wenn (2.5) unabhängig von  $n$  ist.

Markoff-Ketten sind die “erste Stufe” einer Verallgemeinerung von stochastisch unabhängigen Folgen von Zufallsvariablen: Abhängigkeiten auf direkte Vorgänger sind zugelassen, weiterreichende jedoch nicht. In Kapitel 4 wird sich zeigen, daß dieses relativ einfache Modell doch recht weittragend ist.

Im folgenden werden nur homogene Markoff-Ketten betrachtet, deren Zustandsraum  $\mathcal{S} = \{s_1, \dots, s_r\}$ ,  $r \in \mathbb{N}$ , darüberhinaus endlich ist. Dieser Fall reicht zur Beschreibung des Verhaltens von Quellen mit endlichem Alphabet aus. Die Übergangswahrscheinlichkeiten lassen sich dann zu einer  $(r \times r)$ -Matrix  $\mathbf{II}$  wie folgt zusammenfassen. Bezeichne

$$p_{ij} = P(X_n = s_j \mid X_{n-1} = s_i),$$

falls ein  $n \in \mathbb{N}$  existiert mit  $P(X_{n-1} = s_i) > 0$ . Anderenfalls setze  $p_{ij} \geq 0$ ,  $j = 1, \dots, r$ , beliebig so, daß  $\sum_{j=1}^r p_{ij} = 1$ . Die Matrix

$$\mathbf{II} = (p_{ij})_{1 \leq i, j \leq r}$$

heißt Übergangsmatrix der Markoff-Kette. Die Übergangsmatrix bzw. Markoff-Kette heißt irreduzibel, wenn von jedem Zustand zu jedem anderen eine Kette aus positiven Übergangswahrscheinlichkeiten besteht. Genauer, für alle  $i, j \in \{1, \dots, r\}$  existieren  $i_1, i_2, \dots, i_k \in \{1, \dots, r\}$  mit

$$p_{i_1} \cdot p_{i_1 i_2} \cdots p_{i_{k-1} i_k} > 0.$$



Im Fall eines endlichen Zustandsraums läßt sich die Randverteilung der Zufallsvariablen  $X_n$  durch einen stochastischen Vektor

$$\mathbf{p}(n) = (p_1(n), \dots, p_r(n)) \quad \text{mit} \quad p_i(n) = P(X_n = s_i), \quad i = 1, \dots, r,$$

beschreiben. Die Menge der stochastischen Vektoren der Länge  $r$  wird im folgenden bezeichnet mit

$$\mathcal{P}_r = \left\{ (p_1, \dots, p_r) \mid p_i \geq 0, \sum_{i=1}^r p_i = 1 \right\}$$

Stochastische Vektoren werden im folgenden stets als Zeilenvektoren geschrieben.

$\mathbf{p}(0) = (p_1(0), \dots, p_r(0))$  repräsentiert entsprechend die Anfangsverteilung. Als Anwendung der Kolmogoroff-Smirnoff-Gleichung (siehe z.B. [24]) folgt

$$\mathbf{p}(n) = \mathbf{p}(n-1)\mathbf{II}, \quad n \in \mathbb{N}.$$

Mit Hilfe der Anfangsverteilung und der Übergangsmatrix  $\mathbf{II}$  kann hiermit die Verteilung von  $X_n$  iterativ berechnet werden zu

$$\mathbf{p}(n) = \mathbf{p}(n-1)\mathbf{II} = \mathbf{p}(n-2)\mathbf{II}^2 = \dots = \mathbf{p}(0)\mathbf{II}^n. \quad (2.6)$$

Besonders wichtig sind Verteilungen, die die Randverteilungen einer Markoff-Kette invariant lassen. Diese werden in der folgenden Definition eingeführt.

**Definition 2.9** (*stationäre Verteilung einer Markoff-Kette*)

$\{X_n\}_{n \in \mathbb{N}_0}$  sei eine homogene Markoff-Kette mit endlichem Zustandsraum  $\mathcal{S} = \{s_1, \dots, s_r\}$ . Ein stochastischer Vektor  $\mathbf{p}^* \in \mathcal{P}_r$  heißt stationäre Verteilung, wenn  $\mathbf{p}^*\mathbf{II} = \mathbf{p}^*$ .

Die Berechnung einer stationären Verteilung bedeutet bei endlichem Zustandsraum, das homogene Gleichungssystem  $(\mathbf{I} - \mathbf{II})'\mathbf{v} = \mathbf{o}$  in der Menge der stochastischen Vektoren  $\mathbf{v}' \in \mathcal{P}_r$  zu lösen.

Ist die Anfangsverteilung  $\mathbf{p}(0)$  einer Markoff-Kette stationär im Sinn von Definition 2.9, so folgt mit (2.6), daß alle Randverteilungen  $\mathbf{p}(n)$  mit  $\mathbf{p}(0)$  übereinstimmen. Mehr noch, in diesem Fall ist sogar die ganze Folge stationär.

**Lemma 2.3** (*stationäre Markoff-Kette*)

$\{X_n\}_{n \in \mathbb{N}_0}$  sei eine homogene Markoff-Kette mit stationärer Anfangsverteilung  $\mathbf{p}(0)$ . Dann ist die Folge  $\{X_n\}_{n \in \mathbb{N}_0}$  stationär im Sinn von Definition 2.4.

Unter gewissen Regularitätsvoraussetzungen stabilisieren sich Markoff-Ketten unabhängig von der Anfangsverteilung in einer stationären Verteilung in dem Sinn, daß die Verteilung von  $X_n$  mit wachsendem  $n$  gegen eine stationäre Verteilung konvergiert. Nach einer genügend langen Startphase einer solchen Markoff-Kette wird man das Auftreten der Zustände entsprechend der stationären Verteilung beobachten. Dies besagt der folgende Satz.

**Satz 2.4** (*stationäre Verteilung und Limesverteilung*)

$\{X_n\}_{n \in \mathbb{N}_0}$  sei eine irreduzible, homogene Markoff-Kette mit endlichem Zustandsraum  $\mathcal{S}$ . Es existiere ein  $m \in \mathbb{N}$  derart, daß  $\mathbf{P}^m$  eine Spalte aus lauter positiven Elementen besitzt. Dann gilt

$$\lim_{n \rightarrow \infty} \mathbf{P}^n = \begin{pmatrix} p_1^* & \cdots & p_r^* \\ \vdots & & \vdots \\ p_1^* & \cdots & p_r^* \end{pmatrix} \quad \text{und} \quad \lim_{n \rightarrow \infty} p_i(n) = p_i^*$$

für alle  $i = 1, \dots, r$ .  $\mathbf{p} = (p_1^*, \dots, p_r^*)$  ist dann die eindeutig bestimmte stationäre Verteilung.

Die Voraussetzung, daß  $\mathbf{P}^m$  eine Spalte mit lauter positiven Elementen besitzt, ist für irreduzible Markoff-Ketten äquivalent zur Aperiodizität. Diese besagt, daß Zustände nicht nur periodisch mit positiver Wahrscheinlichkeit besucht werden, sondern daß nach genügend langer Startphase jeder Zeitpunkt als Besuchszeit in Frage kommt. Für eine genaue Definition sei auf [24] verwiesen.

## 2.3 Übungsaufgaben

**Aufgabe 2.1**  $X$  mit Träger  $T_X = \{x_1, x_2\}$  und  $Y$  mit Träger  $T_Y = \{y_1, y_2\}$  seien diskrete Zufallsvariable auf einem Wahrscheinlichkeitsraum  $(\Omega, \mathcal{A}, P)$ . Es gelte

$$P(\{X = x_1\} \cup \{Y = y_1\}) = 0.58, \\ P(X = x_2) = 0.7 \text{ und } P(X = x_1, Y = y_1) = 0.12.$$

Bestimmen Sie die gemeinsame Verteilung von  $(X, Y)$ . Sind  $X$  und  $Y$  stochastisch unabhängig?

**Aufgabe 2.2** Der diskrete Zufallsvektor  $(X, Y)$  besitze die Zähldichte ( $0 < p < 1$  ein Parameter)

$$f(i, j) = \begin{cases} 2^{-(j+1)}(1-p)^{i-j}p, & \text{falls } i \geq j \\ 0, & \text{sonst} \end{cases}, \quad i, j \in \mathbb{N}_0.$$

Berechnen Sie  $E(X)$ ,  $E(Y)$  und  $E(X \cdot Y)$ .

**Aufgabe 2.3**  $X$  und  $Y$  seien stochastisch unabhängige, je mit Parameter  $\lambda > 0$  poissonverteilte Zufallsvariable, d.h.

$$P(X = k) = P(Y = k) = e^{-\lambda} \frac{\lambda^k}{k!}, \quad k \in \mathbb{N}_0.$$

Bestimmen Sie die bedingte Verteilung von  $X$  unter  $X + Y = k$ , also die Verteilung  $P^X(\cdot | X + Y = k)$ ,  $k \in \mathbb{N}_0$ .

**Aufgabe 2.4**

- a) Für  $n \in \mathbb{N}$  berechne man durch die Wahl eines geeigneten Wahrscheinlichkeitsraums die Summe der Quersummen aller natürlichen Zahlen  $\leq 10^n$ .
- b) Die Eulersche  $\varphi$ -Funktion ist für  $n \in \mathbb{N}$  definiert durch

$$\varphi(n) := \text{Anzahl der zu } n \text{ teilerfremden natürlichen Zahlen } \leq n.$$

Beweisen Sie mittels wahrscheinlichkeitstheoretischer Überlegungen die Formel

$$\varphi(n) = n \prod_{p|n} \left(1 - \frac{1}{p}\right),$$

wobei das Produkt über alle Primzahlen gebildet wird, die  $n$  teilen.

Hinweis: Man betrachte eine Zufallsvariable  $X$  mit  $P(X = i) = \frac{1}{n}$ ,  $1 \leq i \leq n$ , und Ereignisse  $E(p) = \{\omega \mid p \text{ teilt } X(\omega)\}$ .

**Aufgabe 2.5** Eine homogene Markoff-Kette mit zwei Zuständen besitze die Übergangsmatrix  $\mathbf{II} = \begin{pmatrix} \alpha & 1 - \alpha \\ 1 - \beta & \beta \end{pmatrix}$ ,  $0 \leq \alpha, \beta \leq 1$ .

Für welche  $\alpha$  und  $\beta$  existiert eine stationäre Verteilung? Bestimmen Sie diese gegebenenfalls. Berechnen Sie eine Anfangsverteilung  $\mathbf{p}(0)$ , so daß  $\mathbf{p}(n) = \mathbf{p}(0)$  für alle  $n \in \mathbb{N}$  und alle  $0 \leq \alpha, \beta \leq 1$ . Ist diese Anfangsverteilung eindeutig bestimmt?

**Aufgabe 2.6** (Thanks to Bernard van Cutsem for this nice exercise.) In einem Feld mit  $r$  Zellen sind  $r$  verschiedene Objekte abgespeichert. Diese werden zu Zeitpunkten  $n \in \mathbb{N}$  gemäß folgendem Algorithmus in dem Feld gesucht.

Durchsuche das Feld von links nach rechts nach aufsteigenden Indizes. Wird das gesuchte Objekt in Zelle  $k$  gefunden, verschiebe die Zelleninhalte  $1, \dots, k-1$  jeweils um eine Position nach rechts und speichere das gefundene Objekt in Zelle 1. (Der Sinn dieses Algorithmus ist, im Laufe der Zeit häufig gesuchte Elemente schnell in Zellen mit kleinem Index zu finden.)

Beschreiben Sie die sukzessiven Positionen von Objekt 1 unter diesem Algorithmus durch eine Markoff-Kette. Wie lautet die stationäre Verteilung, wenn die Zugriffswahrscheinlichkeit für Objekt 1 den Wert  $a$ ,  $0 < a < 1$ , und für die Objekte  $2, \dots, r$  den Wert  $b = (1 - a)/(r - 1)$  hat?

### 3 Information und Entropie

Ein präziser Begriff für Unsicherheit bzw. Informationsgewinn ist bei der kompakten Kodierung von Quellen und der Beschreibung der Leistungsfähigkeit von gestörten Kanälen von großer Bedeutung. Hierbei ist es gleichgültig, ob eine Zufallsexperiment durch die Unsicherheit über den Ausgang vor Ausführung oder den Informationsgewinn nach Bekanntwerden des Ausgangs beurteilt wird. Beide Größen können durch dieselbe Maßzahl gemessen werden, wie durch das folgende einführende Beispiel motiviert wird.

Zwei Zufallsexperimente mit jeweils drei möglichen Ausgängen werden durchgeführt. Die zugehörigen Wahrscheinlichkeiten der Ausgänge werden jeweils durch die stochastischen Vektoren

$$\mathbf{p} = (0.3, 0.4, 0.3) \quad \text{und} \quad \mathbf{q} = (0.9, 0.05, 0.05)$$

beschrieben. Der erste Vektor hat eine größere Unbestimmtheit des Ausgangs, da alle Ergebnisse ungefähr dieselbe Wahrscheinlichkeit haben. Nach Beobachten des Ausgangs ist also der Informationsgewinn beim ersten Experiment größer. Demgegenüber kann man beim zweiten Experiment den Ausgang mit hoher Wahrscheinlichkeit richtig vorhersagen. Die Unsicherheit hierüber ist gering, auch der Informationsgewinn nach Beobachten des Experiments ist klein.

Unser Ziel ist es, eine Maßzahl für die Unbestimmtheit oder den Informationsgewinn zu gewinnen. Hierzu wird zunächst ein direkter Vergleich von Wahrscheinlichkeitsvektoren durchgeführt. Es bezeichne

$$\mathcal{P}_m = \left\{ (p_1, \dots, p_m) \mid p_i \geq 0, i = 1, \dots, m, \sum_{i=1}^m p_i = 1 \right\}$$

die Menge der Wahrscheinlichkeitsvektoren der Länge (Dimension)  $m \in \mathbb{N}$ .  $\mathcal{P}_m$  repräsentiert bei festem Träger  $\mathcal{X} = \{x_1, \dots, x_m\}$  durch  $P(\{x_i\}) = p_i$ ,  $i = 1, \dots, m$ , die Menge aller diskreten Wahrscheinlichkeitsverteilungen  $P$  auf  $(\mathcal{X}, \mathfrak{P}(\mathcal{X}))$ .

Durch die folgende Ordnungsrelation auf  $\mathcal{P}_m$  können gewisse  $\mathbf{p}, \mathbf{q} \in \mathcal{P}_m$  bezüglich ihrer Unbestimmtheit verglichen werden. Offensichtlich ist es sinnvoll, nicht auf die Reihenfolge der einzelnen  $p_i$  in  $(p_1, \dots, p_m)$  zu achten, so daß wir uns auf Ordnungen beschränken können, die invariant sind unter Permutationen der Komponenten. Als Repräsentanten werden die Vektoren mit aufsteigenden Komponenten gewählt.

Für  $\mathbf{p} = (p_1, \dots, p_m) \in \mathcal{P}_m$  bezeichne hierzu  $\mathbf{p}_\uparrow = (p_{(1)}, \dots, p_{(m)})$  den zugehörigen Vektor mit aufsteigend geordneten Komponenten, also  $p_{(1)} \leq p_{(2)} \leq \dots \leq p_{(m)}$  und es existiert eine Permutation  $\sigma \in S_m$  (der Menge der Permutationen der Zahlen  $1, \dots, m$ ) mit  $p_{(i)} = p_{\sigma(i)}$  für alle  $i = 1, \dots, m$ .

**Definition 3.1** (Majorisierung)

Seien  $\mathbf{p} = (p_1, \dots, p_m), \mathbf{q} = (q_1, \dots, q_m) \in \mathcal{P}_m$ . Man sagt,  $\mathbf{q}$  majorisiert  $\mathbf{p}$ , bezeichnet mit  $\mathbf{p} \prec \mathbf{q}$ , wenn

$$\sum_{i=1}^k p_{(i)} \geq \sum_{i=1}^k q_{(i)} \text{ für alle } k = 1, \dots, m-1 \text{ und } \sum_{i=1}^m p_{(i)} = \sum_{i=1}^m q_{(i)}.$$

Man prüft leicht nach, daß durch “ $\prec$ ” eine Präordnung auf  $\mathcal{P}_m$  definiert wird, d.h. “ $\prec$ ” ist eine reflexive ( $\mathbf{p} \prec \mathbf{p}$ ) und transitive ( $\mathbf{p} \prec \mathbf{q}$  und  $\mathbf{q} \prec \mathbf{r} \Rightarrow \mathbf{p} \prec \mathbf{r}$ ) Relation auf  $\mathcal{P}_m$ . Antisymmetrie ( $\mathbf{p} \prec \mathbf{q}$  und  $\mathbf{q} \prec \mathbf{p} \Rightarrow \mathbf{p} = \mathbf{q}$ ) gilt lediglich auf der Menge der stochastischen Vektoren mit aufsteigend geordneten Komponenten. Auf dieser Menge liegt also eine partielle Ordnung vor.

Für die eingangs angegebenen stochastischen Vektoren der Länge 3 ergibt sich beispielsweise  $\mathbf{p}_\uparrow = (0.3, 0.3, 0.4)$  und  $\mathbf{q}_\uparrow = (0.05, 0.05, 0.9)$ . Für die Partialsummen gilt:

$k$	1	2		3	
$\mathbf{p}$	0.3	0.3+0.3	= 0.6	0.3+0.3+0.4	= 1.0
$\mathbf{q}$	0.05	0.05+0.05	= 0.1	0.05+0.05+0.9	= 1.0

Folglich wird  $\mathbf{p}$  von  $\mathbf{q}$  majorisiert, also  $\mathbf{p} \prec \mathbf{q}$ .

Stellt man die Partialsummen in Definition 3.1 graphisch dar, indem man die Punkte  $(\frac{k}{m}, \sum_{i=1}^k p_{(i)})$ ,  $k = 0, \dots, m$ , als Koordinaten auffaßt und benachbarte Punkte durch Geraden verbindet, so erhält man die sogenannte *Lorenzkurve*. Hierbei wird die leere Summe als Null definiert, d.h.  $\sum_{i=1}^0 p_{(i)} = 0$ .

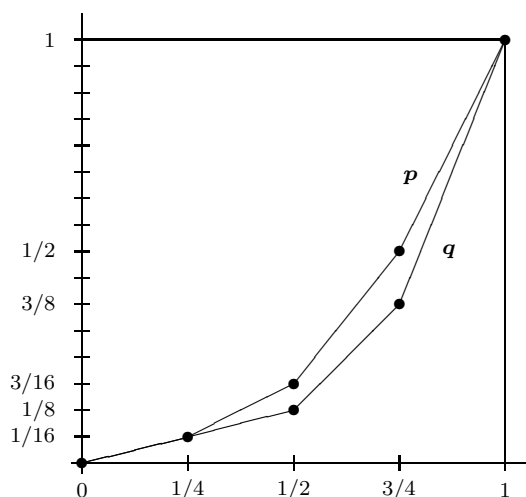


Abb. 3.1 Lorenzkurven,  $\mathbf{p} = (\frac{1}{16}, \frac{1}{8}, \frac{5}{16}, \frac{1}{2}) \prec \mathbf{q} = (\frac{1}{16}, \frac{1}{16}, \frac{1}{4}, \frac{5}{8})$ .

Für das Beispiel  $\mathbf{p} = (\frac{1}{16}, \frac{1}{8}, \frac{5}{16}, \frac{1}{2})$  und  $\mathbf{q} = (\frac{1}{16}, \frac{1}{16}, \frac{1}{4}, \frac{5}{8})$  ergibt sich die Darstellung aus Abbildung 3.1.

Ein Experiment ist umso bestimmter, die zugehörige Wahrscheinlichkeitsverteilung umso konzentrierter auf wenige Ereignisse, je weiter die Kurve nach rechts unten durchgebogen ist. Für  $\mathbf{p}, \mathbf{q} \in \mathcal{P}_m$  liegt die zu  $\mathbf{q}$  gehörige Kurve offensichtlich unterhalb der zu  $\mathbf{p}$  gehörigen genau dann, wenn  $\mathbf{p} \prec \mathbf{q}$ . Dies läßt nun die Interpretation zu, daß das zu  $\mathbf{p}$  gehörige Experiment unbestimmter als das zu  $\mathbf{q}$  gehörige ist, wenn  $\mathbf{p} \prec \mathbf{q}$ .

Nicht alle stochastischen Vektoren sind jedoch bezüglich der Majorisierung vergleichbar. Betrachte etwa  $\mathbf{p} = (0.1, 0.4, 0.5)$  und  $\mathbf{q} = (0.2, 0.2, 0.6)$ . Dann gilt  $\mathbf{p} \not\prec \mathbf{q}$  und  $\mathbf{q} \not\prec \mathbf{p}$ . Daraus folgt, daß  $\prec$  auf  $\mathcal{P}_m$  keine Totalordnung definiert.

### 3.1 Entropie und Transinformation

Um alle stochastischen Vektoren bezüglich ihrer Unbestimmtheit vergleichen zu können, wird ein reelles Funktional  $H$ , eine Maßzahl für Unbestimmtheit, eingeführt. Ein solches  $H$  wird auf der Menge  $\mathcal{P} = \bigcup_{m=1}^{\infty} \mathcal{P}_m$

der stochastischen Vektoren beliebiger endlicher Länge definiert. Nach den vorangegangenen Überlegungen ist eine Minimalforderung an  $H$ , daß für  $\mathbf{p}, \mathbf{q} \in \mathcal{P}_m$  mit  $\mathbf{p} \prec \mathbf{q}$  die Eigenschaft  $H(\mathbf{p}) \geq H(\mathbf{q})$  folgt, d.h. die Antitonie von  $H$  bezüglich “ $\prec$ ”. Eine Funktion, die diese Eigenschaft besitzt, ist die im folgenden definierte Entropie.

**Definition 3.2** (*Entropie*)

Die Abbildung

$$H : \mathcal{P} = \bigcup_{m=1}^{\infty} \mathcal{P}_m \rightarrow \mathbb{R} : (p_1, \dots, p_m) \mapsto - \sum_{i=1}^m p_i \log p_i$$

heißt Entropie. Hierbei wird  $0 \cdot \log 0 = 0$  gesetzt.

Synonym wird  $H(X) = H(p_1, \dots, p_m)$  bezeichnet, falls  $X$  eine diskrete Zufallsvariable ist, deren Verteilung durch  $(p_1, \dots, p_m)$  mit  $P(X = x_i) = p_i$ ,  $x_i \in \mathcal{T}_X$ ,  $i = 1, \dots, m$ , beschrieben wird. Es gilt also

$$H(X) = - \sum_{i=1}^m P(X = x_i) \cdot \log P(X = x_i). \quad (3.1)$$

Man beachte, daß die Definition der Entropie nur von den Wahrscheinlichkeiten, nicht aber von den Trägerpunkten der entsprechenden Verteilung abhängt. Dies ist von einem Maß für Unbestimmtheit zu erwarten. Die Basis des Logarithmus “log” geht in der Definition als multiplikative Konstante ein, da  $\log_b x = \log_a x / \log_a b$  für alle  $a, b > 1$ ,  $x > 0$  gilt. Die Entropie  $H$  ist unter Zusatzbedingungen sogar das einzige vernünftige Maß für Unbestimmtheit oder Informationsgewinn, wie später gezeigt wird.

**Lemma 3.1**  $H$  ist antiton bezüglich  $\prec$ , d.h.  $\mathbf{p} \prec \mathbf{q} \Rightarrow H(\mathbf{p}) \geq H(\mathbf{q})$  für alle  $\mathbf{p}, \mathbf{q} \in \mathcal{P}_m$  und  $m \in \mathbb{N}$ .

**Beweis.** Wir benutzen (ohne ihn zu beweisen) den folgenden Satz aus der Theorie der Majorisierung (siehe [23], Seite 22):

Seien  $\mathbf{p}, \mathbf{q} \in \mathcal{P}_m$ . Dann gilt  $\mathbf{p} \prec \mathbf{q}$  genau dann, wenn eine doppelt stochastische Matrix  $\mathbf{S}$  existiert mit  $\mathbf{p} = \mathbf{q}\mathbf{S}$ .

Eine Matrix  $\mathbf{S} = (s_{ij})_{1 \leq i, j \leq m}$  heißt doppelt stochastisch, wenn  $s_{ij} \geq 0$ , ferner alle Zeilensummen  $\sum_{\ell=1}^m s_{i\ell} = 1$  und alle Spaltensummen  $\sum_{\ell=1}^m s_{\ell j} =$



1 sind für alle  $i, j = 1, \dots, m$ .

Weiterhin wird benutzt, daß  $f(x) = -x \ln x$  für  $x \in [0, 1]$  konkav ist, d.h.

$$f\left(\sum_{i=1}^m \alpha_i x_i\right) \geq \sum_{i=1}^m \alpha_i f(x_i)$$

für alle  $x_1, \dots, x_m \in [0, 1]$ ,  $\alpha_i \geq 0$  mit  $\sum_{i=1}^m \alpha_i = 1$ , wobei "ln" den natürlichen Logarithmus (zur Basis  $e$ ) bezeichnet. Dies folgt aus der Stetigkeit von  $f$  in  $[0, 1]$  und der Tatsache, daß  $f''(x) = (-\ln x - 1)' = -\frac{1}{x} < 0$  für alle  $x \in (0, 1]$ .

Seien nun  $\mathbf{p}, \mathbf{q} \in \mathcal{P}_m$  mit  $\mathbf{p} \prec \mathbf{q}$ . Dann existiert eine doppelt stochastische Matrix  $\mathbf{S} = (s_{ij})_{1 \leq i, j \leq m}$  mit  $\mathbf{p} = \mathbf{q}\mathbf{S}$ , d.h.  $p_j = \sum_{i=1}^m q_i s_{ij}$  für alle  $j = 1, \dots, m$ . Es folgt

$$\begin{aligned} H(\mathbf{p}) &= -\sum_{j=1}^m p_j \ln p_j = \sum_{j=1}^m \left( -\left(\sum_{i=1}^m q_i s_{ij}\right) \ln \left(\sum_{i=1}^m q_i s_{ij}\right) \right) \\ &\geq \sum_{j=1}^m \sum_{i=1}^m -s_{ij} q_i \ln q_i = -\sum_{i=1}^m q_i \ln q_i \underbrace{\left(\sum_{j=1}^m s_{ij}\right)}_{=1} = H(\mathbf{q}). \end{aligned}$$

Die Verwendung von Logarithmen zu einer anderen Basis wirkt sich lediglich als positive multiplikative Konstante aus und erhält die  $\geq$ -Beziehung. ■

Für eine diskrete Zufallsvariable  $X$  mit Träger  $\mathcal{X} = \{x_1, \dots, x_m\}$  ist nach (3.1)  $H(X) = -\sum_{i=1}^m P(X = x_i) \log P(X = x_i)$ . Dieser Ausdruck läßt sich als Erwartungswert schreiben. Hierzu definieren wir

$$I_X : \mathcal{X} \rightarrow \mathbb{R} : x_j \mapsto -\log P(X = x_j). \quad (3.2)$$

$I_X$  ist eine Zufallsvariable auf dem Wahrscheinlichkeitsraum  $(\mathcal{X}, \mathfrak{P}(\mathcal{X}), P^X)$ . Mit der Konvention  $0 \cdot \log 0 = 0$  ist  $I_X$  integrierbar, und es gilt mit (2.2)

$$E(I_X) = -\sum_{j=1}^m P(X = x_j) \cdot \log P(X = x_j) = H(X).$$

$I_X(x_j) = -\log P(X = x_j) = \log(1/P(X = x_j))$  heißt hierbei Informationsgehalt des Ereignisses  $\{X = x_j\}$ ,  $j = 1, \dots, m$ .

Die Definition der Entropie eines diskreten Zufallsvektors ist ein Spezialfall von (3.1). Liegen zwei diskrete Zufallsvariable  $X$  und  $Y$  auf einem Wahrscheinlichkeitsraum  $(\Omega, \mathcal{A}, P)$  mit Träger

$$\mathcal{X} = \{x_1, \dots, x_m\} \text{ bzw. } \mathcal{Y} = \{y_1, \dots, y_n\} \quad (3.3)$$

vor, so besitzt der Zufallsvektor  $(X, Y)$  den Träger  $\mathcal{X} \times \mathcal{Y}$ .  $(X, Y)$  ist selbst wieder eine diskrete Zufallsvariable, so daß gemäß Definition 3.2 für die Entropie von  $(X, Y)$

$$H(X, Y) = - \sum_{i=1}^m \sum_{j=1}^n P(X = x_i, Y = y_j) \cdot \log P(X = x_i, Y = y_j)$$

gilt.

Die Entropie der bedingten Verteilung von  $X$  unter  $\{Y = y_j\}$  hat nach Definition 3.2 die Gestalt

$$H(X | Y = y_j) = - \sum_{i=1}^m P(X = x_i | Y = y_j) \cdot \log P(X = x_i | Y = y_j)$$

$H(X | Y = y_j)$  heißt bedingte Entropie von  $X$ , gegeben  $Y = y_j$ . Die bedingte Entropie von  $X$  unter  $Y$  erhält man durch Erwartungswertbildung der bedingten Entropie von  $X$  unter  $\{Y = y_j\}$ , aufgefaßt als Zufallsvariable mit Argument  $y_j \in \mathcal{Y}$  (vgl. (2.3)).

**Definition 3.3** (*bedingte Entropie*)

$X, Y$  seien diskrete Zufallsvariable mit Träger  $\mathcal{X}$  bzw.  $\mathcal{Y}$  aus (3.3).

$$H(X | Y) = \sum_{j=1}^n P(Y = y_j) H(X | Y = y_j) \quad (3.4)$$

heißt *bedingte Entropie von  $X$  unter  $Y$  oder bedingte Entropie von  $X$  gegeben  $Y$* .

Durch einfaches Umrechnen erhält man die Darstellung

$$\begin{aligned} H(X | Y) &= - \sum_{j=1}^n \sum_{i=1}^m P(Y = y_j) \frac{P(X = x_i, Y = y_j)}{P(Y = y_j)} \log P(X = x_i | Y = y_j) \\ &= - \sum_{i=1}^m \sum_{j=1}^n P(X = x_i, Y = y_j) \log P(X = x_i | Y = y_j). \end{aligned}$$

Mit Hilfe der bedingten Entropie wird nun die Transinformation von Zufallsvariablen definiert.

**Definition 3.4** (*Transinformation*)

$X, Y$  seine diskrete Zufallsvariable mit Träger  $\mathcal{X}$  bzw.  $\mathcal{Y}$  aus (3.3).

$$I(X, Y) = H(X) - H(X | Y)$$

heißt *Transinformation oder Synentropie von  $X$  und  $Y$* .

Die Transinformation kann folgendermaßen interpretiert werden. Beschreiben  $X$  und  $Y$  zwei Experimente, deren Ausgänge sich gegenseitig beeinflussen, so mißt  $I(X, Y)$ , um wieviel die Unbestimmtheit im Mittel kleiner wird, wenn man das Ergebnis von  $Y$  kennt.

Für drei Zufallsvariable  $X, Y, Z$  wird die bedingte Transinformation von  $X$  und  $Z$  unter  $Y$ , mit Hilfe der bedingten Entropie definiert als

$$I(X, Z | Y) = H(X | Y) - H(X | Y, Z). \quad (3.5)$$

Wenn keine Unklarheit über die bedingenden Zufallsvariablen besteht, werden die Klammern bei den entsprechenden Zufallsvektoren weggelassen.  $H(X | Y, Z)$  ist also zu lesen als  $H(X | (Y, Z))$ .

Mit der Beziehung  $P(X = x_i) = \sum_j P(X = x_i, Y = y_j)$  erhält man die folgende Darstellung der Transinformation.

$$\begin{aligned} I(X, Y) &= H(X) - H(X | Y) \\ &= - \sum_i P(X = x_i) \log P(X = x_i) \end{aligned}$$

$$\begin{aligned}
& + \sum_{i,j} P(X = x_i, Y = y_j) \log P(X = x_i | Y = y_j) \\
& = \sum_{i,j} P(X = x_i, Y = y_j) \log \frac{P(X = x_i | Y = y_j)}{P(X = x_i)} \quad (3.6)
\end{aligned}$$

Auch die Transinformation besitzt eine Darstellung als Erwartungswert. Definiert man analog zu (3.2) die Zufallsvariable

$$I_{X,Y} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R} : (x_i, y_j) \mapsto \log \frac{P(X = x_i | Y = y_j)}{P(X = x_i)}$$

auf dem Wahrscheinlichkeitsraum  $(\mathcal{X} \times \mathcal{Y}, \mathfrak{P}(\mathcal{X} \times \mathcal{Y}), P^{(X,Y)})$ , so gilt mit (3.6)

$$E(I_{X,Y}) = I(X, Y).$$

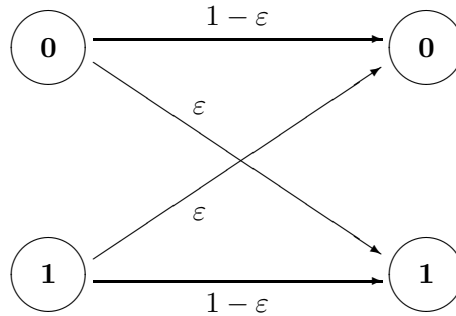
Der Ausdruck  $I_{X,Y}(x_i, y_j) = \log \frac{P(X=x_i|Y=y_j)}{P(X=x_i)}$  heißt hierbei wechselseitige Information der Ereignisse  $\{X = x_i\}$  und  $\{Y = y_j\}$ .

Im folgenden Beispiel werden die oben definierten Begriffe für ein einfaches Kanalmodell untersucht.

**Beispiel 3.1** (binärer symmetrischer Kanal, BSC, engl.: binary symmetric channel)

Durch einen Kanal werden Bits übertragen und ausgegeben. Als Eingabe- und Ausgabealphabet dient  $\mathcal{X} = \{0, 1\} = \mathcal{Y}$ . Die Zufallsvariable  $X$  repräsentiert die (zufällige) Eingabe von Zeichen in den Kanal,  $Y$  entsprechend die Ausgabe.

Die Inputverteilung sei eine Gleichverteilung, d.h.  $P(X = 0) = P(X = 1) = \frac{1}{2}$ . Jedes Bit wird mit der Wahrscheinlichkeit  $(1 - \varepsilon)$  richtig übertragen, mit der Wahrscheinlichkeit  $\varepsilon$  falsch,  $0 \leq \varepsilon \leq 1$ . Dies ergibt die bedingten Wahrscheinlichkeiten  $P(Y = 0 | X = 0) = P(Y = 1 | X = 1) = 1 - \varepsilon$  und  $P(Y = 1 | X = 0) = P(Y = 0 | X = 1) = \varepsilon$ . Der Kanal wird durch die folgende Graphik veranschaulicht.



Als nächstes wird die gemeinsame Verteilung von  $(X, Y)$  berechnet. Es gilt

$$P(X = 0, Y = 0) = P(X = 1, Y = 1) = \frac{1 - \varepsilon}{2},$$

$$P(X = 0, Y = 1) = P(X = 1, Y = 0) = \frac{\varepsilon}{2}.$$

Hieraus erhalten wir die Verteilung von  $Y$  als Gleichverteilung mit  $P(Y = 0) = P(Y = 1) = \frac{1}{2}$  und die bedingten Wahrscheinlichkeiten  $P(X = x_i | Y = y_j)$  zu

$$P(X = 0 | Y = 0) = P(X = 1 | Y = 1) = 1 - \varepsilon,$$

$$P(X = 1 | Y = 0) = P(X = 0 | Y = 1) = \varepsilon.$$

Es folgt bei Verwendung von Logarithmen zur Basis 2:

$$H(X) = H(Y) = -\frac{1}{2} \log_2 \frac{1}{2} - \frac{1}{2} \log_2 \frac{1}{2} = 1$$

$$H(X, Y) = -2 \frac{1 - \varepsilon}{2} \log_2 \frac{1 - \varepsilon}{2} - 2 \frac{\varepsilon}{2} \log_2 \frac{\varepsilon}{2}$$

$$= 1 - (1 - \varepsilon) \log_2(1 - \varepsilon) - \varepsilon \log_2 \varepsilon$$

$$H(X | Y) = -2 \frac{1 - \varepsilon}{2} \log_2(1 - \varepsilon) - 2 \frac{\varepsilon}{2} \log_2 \varepsilon$$

$$= -(1 - \varepsilon) \log_2(1 - \varepsilon) - \varepsilon \log_2 \varepsilon$$

$$I(X, Y) = H(X) - H(X|Y) = 1 + (1 - \varepsilon) \log_2(1 - \varepsilon) + \varepsilon \log_2 \varepsilon$$

Graphisch ergeben sich die in Abbildung 3.2 dargestellten Funktionen von  $\varepsilon$ . Die Interpretation dieses Kanalmodells ist klar. Für  $\varepsilon = 0$  liegt ein ungestörter Kanal vor, für  $\varepsilon = \frac{1}{2}$  ist der Kanal vollständig gestört. In diesem Fall wird jedes gesendete Symbol mit gleicher Wahrscheinlichkeit richtig

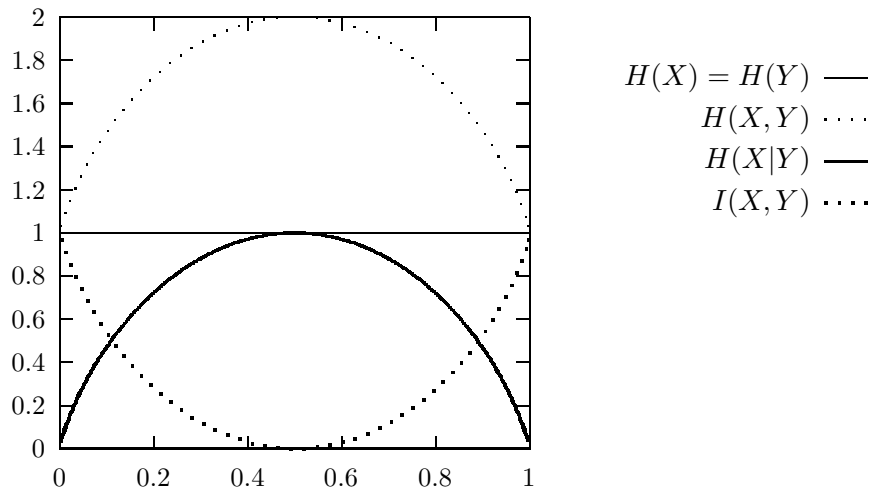


Abb. 3.2 Entropien und Transinformation beim BSC

bzw. falsch übertragen. Obige Kurven spiegeln dieses Verhalten genau wider. Die Entropie (Unbestimmtheit) von  $(X, Y)$  und  $X$  unter  $Y$  sind für  $\varepsilon = \frac{1}{2}$  maximal, die Transinformation von  $X$  und  $Y$  ist hier minimal. Für  $\varepsilon = 1$  überträgt der Kanal wieder ungestört, nur werden 0 und 1 systematisch ausgetauscht. ■

Durch das folgende Lemma werden Zusammenhänge zwischen Entropie, bedingter Entropie und der Transinformation hergestellt.

**Lemma 3.2** Für diskrete Zufallsvariable  $X$  und  $Y$  gilt

- a)  $H(X, Y) = H(X) + H(Y | X) = H(Y) + H(X | Y)$
- b)  $I(X, Y) = H(X) - H(X | Y) = H(Y) - H(Y | X)$   
 $= H(X) + H(Y) - H(X, Y)$

Aus b) folgt insbesondere die Symmetrie der Transinformation in  $X$  und  $Y$ .

**Beweis.** a) Nach Addition einer  $0 = -\log P(X = x_i) + \log P(X = x_i)$  folgt

$$H(X, Y) = - \sum_{i,j} P(X = x_i, Y = y_j) (\log P(X = x_i, Y = y_j))$$

$$\begin{aligned}
& -\log P(X = x_i) + \log P(X = x_i)) \\
= & -\sum_{i,j} P(X = x_i, Y = y_j) \log P(Y = y_j | X = x_i) \\
& -\underbrace{\sum_i \sum_j P(X = x_i, Y = y_j) \log P(X = x_i)}_{=P(X=x_i)} \\
= & H(X) + H(Y | X)
\end{aligned}$$

Die zweite Gleichung folgt analog.

b) läßt sich aus der Definition von  $I(X, Y)$  und a) wie folgt schließen.

$$\begin{aligned}
I(X, Y) &= H(X) - H(X | Y) = H(X) + H(Y) - H(X, Y) \\
&= H(Y) - H(Y | X).
\end{aligned}$$

■

Gleichung a) aus Lemma 3.2 formalisiert die folgende einleuchtende Eigenschaft eines Maßes für Unbestimmtheit. Die Unbestimmtheit eines zusammengesetzten Experiments entsteht additiv aus der Maßzahl für die Unbestimmtheit des ersten und der für die des zweiten, gegeben das erste.

Wir wenden uns nun einigen wichtigen Ungleichungen der Entropie und Transinformation zu. Interessant sind die notwendigen und hinreichenden Bedingungen für Gleichheit, die sich in natürlicher Weise für ein Maß für Unbestimmtheit ergeben.

**Satz 3.1** (Ungleichungen)

$X, Y, Z$  seien diskrete Zufallsvariable, Träger von  $X$  bzw.  $Y$  seien  $\mathcal{X} = \{x_1, \dots, x_m\}$  bzw.  $\mathcal{Y} = \{y_1, \dots, y_n\}$ . Dann gilt:

a)  $0 \leq H(X) \leq \log m,$

wobei Gleichheit in der linken Ungleichung genau dann gilt, wenn ein  $i$  existiert mit  $P(X = x_i) = 1$  (Einpunktverteilung), und Gleichheit rechts genau dann, wenn  $P(X = x_i) = \frac{1}{m}$  für alle  $i = 1, \dots, m$  (Gleichverteilung).

$$\text{b)} \quad 0 \leq I(X, Y) \leq H(X),$$

wobei Gleichheit links genau dann gilt, wenn  $X$  und  $Y$  stochastisch unabhängig sind, und Gleichheit rechts genau dann, wenn für alle  $i, j$  mit  $P(X = x_i, Y = y_j) > 0$  gilt, daß  $P(X = x_i | Y = y_j) = 1$  ( $X$  ist total abhängig von  $Y$ ).

$$\text{c)} \quad 0 \leq H(X | Y) \leq H(X),$$

wobei Gleichheit links genau dann gilt, wenn  $X$  total abhängig von  $Y$  ist, und Gleichheit rechts genau dann, wenn  $X$  und  $Y$  stochastisch unabhängig sind. Ungleichung c) heißt Shannonsche Ungleichung.

$$\text{d)} \quad H(X) \leq H(X, Y) \leq H(X) + H(Y),$$

mit Gleichheit links genau dann, wenn  $Y$  total abhängig von  $X$  ist, und Gleichheit rechts genau dann, wenn  $X$  und  $Y$  stochastisch unabhängig sind.

$$\text{e)} \quad H(X | (Y, Z)) \leq \min \{H(X | Y), H(X | Z)\}.$$

**Beweis.** Wir verwenden die Abkürzungen  $p_i = P(X = x_i)$ ,  $p_{ij} = P(X = x_i, Y = y_j)$  und  $p_{i|j} = P(X = x_i | Y = y_j)$ . Unter Teil (i) werden jeweils die linken Ungleichungen, unter (ii) jeweils die rechten bewiesen.

a)(i) Es gilt  $H(X) = -\sum_i p_i \log p_i \geq 0$ , mit Gleichheit genau dann, wenn  $p_i \log p_i = 0$ , d.h.  $p_i \in \{0, 1\}$  für alle  $i = 1, \dots, m$ . Aus  $\sum_{i=1}^m p_i = 1$  folgt die Behauptung.

(ii) Hier wird die bekannte Ungleichung  $\ln z \leq z - 1$ , falls  $z > 0$ , verwendet. Es gilt

$$\begin{aligned} H(X) - \log m &= \sum_i p_i \log \frac{1}{p_i} - \sum_i p_i \log m = \sum_i p_i \log \frac{1}{mp_i} \\ &= (\log e) \sum_i p_i \ln \frac{1}{mp_i} = (\log e) \sum_{i:p_i>0} p_i \ln \frac{1}{mp_i} \\ &\leq (\log e) \sum_{i:p_i>0} p_i \left( \frac{1}{mp_i} - 1 \right) = (\log e) \left( \sum_{i:p_i>0} \frac{1}{m} - 1 \right) \leq 0, \end{aligned}$$

wobei Gleichheit genau dann gilt, wenn  $p_i = \frac{1}{m}$  für alle  $i = 1, \dots, m$ .

b)(i) Die Darstellung (3.6) der Transinformation und obige Ungleichung



$\ln z \leq z - 1$ ,  $z > 0$ , liefern

$$\begin{aligned} -I(X, Y) &= \sum_{i,j} p_{ij} \log \frac{p_i}{p_{i|j}} = (\log e) \sum_{i,j} p_{ij} \ln \frac{p_i}{p_{i|j}} \\ &\leq (\log e) \sum_{i,j:p_{ij}>0} p_{ij} \left( \frac{p_i}{p_{i|j}} - 1 \right) \\ &= (\log e) \left( \sum_{i,j:p_{ij}>0} p_i p_j - 1 \right) \leq 0 \end{aligned}$$

Gleichheit gilt genau dann, wenn die Bedingungen  $p_{ij} = 0 \Rightarrow p_i p_j = 0$  und  $p_{ij} > 0 \Rightarrow p_i = p_{i|j}$  beide erfüllt sind. Dies ist äquivalent zur stochastischen Unabhängigkeit von  $X$  und  $Y$ .

(ii) Es gilt  $I(X, Y) = H(X) - H(X | Y) \leq H(X)$ , da  $H(X | Y) \geq 0$ . Gleichheit liegt genau dann vor, wenn  $H(X | Y) = 0$ . Die notwendige und hinreichende Bedingung ergibt sich aus der entsprechenden Bedingung in c)(i).

c)(i) Die Nichtnegativität folgt aus der Darstellung

$$H(X|Y) = - \sum_{i,j} p_{ij} \log p_{i|j} = - \sum_{i,j:p_{ij}>0} p_{ij} \log p_{i|j} \geq 0,$$

mit Gleichheit genau dann, wenn  $p_{i|j} = 1$  für alle  $i, j$  mit  $p_{ij} > 0$ .

(ii) Da wegen b)(i)  $0 \leq I(X, Y) = H(X) - H(X | Y)$ , folgt  $H(X) \geq H(X | Y)$ , mit der Bedingung für Gleichheit aus b)(ii).

d)(i) Da mit Lemma 3.2 a)  $H(X, Y) = H(X) + H(Y | X) \geq 0$  gilt, folgt  $H(X) \leq H(X, Y)$ . Analog zu c)(i) schließt man, daß Gleichheit genau dann gilt, wenn  $H(Y | X) = 0$ . Dies ist äquivalent zu  $p_{j|i} = 1$  für alle  $i, j$  mit  $p_{ij} > 0$ .

(ii) Mit Lemma 3.2 b) folgt  $H(X) + H(Y) - H(X, Y) = I(X, Y) \geq 0$ . Gleichheit gilt genau dann, wenn  $I(X, Y) = 0$ , also  $X$  und  $Y$  stochastisch unabhängig sind.

e) Den Beweis der letzten Ungleichung erhält man nach Einsetzen der Definitionen wieder aus der Ungleichung  $\ln z \leq z - 1$ ,  $z > 0$ .

$$\begin{aligned} -H(X|Y) + H(X|Y, Z) &= \sum_{i,j,k} p_{ijk} \log p_{i|j} - \sum_{i,j,k} p_{ijk} \log p_{i|jk} \\ &= - \sum_{i,j,k} p_{ijk} \log \frac{p_{i|jk}}{p_{i|j}} = \sum_{i,j,k:p_{ijk}>0} p_{ijk} \log \frac{p_{i|j}}{p_{i|jk}} \end{aligned}$$

$$\begin{aligned}
 &\leq (\log e) \left( \sum_{i,j,k:p_{ijk}>0} p_{ijk} \frac{p_{i|j}}{p_{i|jk}} - 1 \right) \\
 &= (\log e) \left( \sum_{i,j,k:p_{ijk}>0} \frac{p_{ijk} p_{ij} p_{jk}}{p_j p_{ijk}} - 1 \right) \\
 &= (\log e) \left( \sum_{i,j} p_{ij} - 1 \right) = 0
 \end{aligned}$$

Die Ungleichung  $H(X | Y, Z) \leq H(X | Z)$  folgt analog. Damit sind alle Behauptungen gezeigt. ■

## 3.2 Axiomatische Charakterisierung der Entropie

Das nächste Ziel ist nun eine axiomatische Charakterisierung der Entropie durch möglichst wenige der bisher gezeigten Eigenschaften. Drei Eigenschaften werden ausgewählt, und es wird gezeigt, daß eine Funktion  $H : \mathcal{P} \rightarrow \mathbb{R}$ , die diesen genügt, notwendigerweise bis auf einen konstanten Faktor mit der Entropie aus Definition 3.2 übereinstimmt.

Satz 3.1 a) besagt, daß  $H(X) \leq \log m$  mit Gleichheit genau dann, wenn  $X$  gleichverteilt ist. Für stochastische Vektoren ergibt sich die folgende Formulierung.

(E1) Für alle  $m \in \mathbb{N}$ ,  $(p_1, \dots, p_m) \in \mathcal{P}_m$  gilt

$$H(p_1, \dots, p_m) \leq H\left(\frac{1}{m}, \dots, \frac{1}{m}\right).$$

Es bezeichne  $p_{ij} = P(X = x_i, Y = y_j)$ ,  $p_{i\cdot} = \sum_j p_{ij} = P(X = x_i)$  und  $p_{j|i} = P(Y = y_j | X = x_i) = \frac{p_{ij}}{p_{i\cdot}}$ . Die Beziehung  $H(X, Y) = H(X) + H(Y | X)$  aus Lemma 3.2 a) besagt dann für stochastische Vektoren:

(E2) Für alle  $m, n \in \mathbb{N}$ ,  $(p_{11}, \dots, p_{1n}, p_{21}, \dots, p_{2n}, \dots, p_{m1}, \dots, p_{mn}) \in \mathcal{P}_{m \cdot n}$  gilt

$$H(p_{11}, \dots, p_{mn}) = H(p_{1\cdot}, \dots, p_{m\cdot}) + \sum_{i=1}^m p_{i\cdot} H(p_{1|i}, \dots, p_{n|i}).$$

Mit der Konvention  $0 \cdot \log 0 = 0$  gilt weiterhin:

(E3) Für alle  $m \in \mathbb{N}$ ,  $\ell \in \{0, \dots, m\}$ ,  $(p_1, \dots, p_m) \in \mathcal{P}_m$  gilt

$$H(p_1, \dots, p_\ell, 0, p_{\ell+1}, \dots, p_m) = H(p_1, \dots, p_m),$$

d.h. eine eingefügte 0 verändert die Entropie nicht.

Obige Eigenschaften lassen sich folgendermaßen interpretieren. (E1) besagt, daß die Unbestimmtheit bzw. der Informationsgewinn bei einer Gleichverteilung am größten ist. In (E2) wird die Unbestimmtheit eines Experiments mit zwei Komponenten additiv zusammengesetzt aus der Unbestimmtheit des ersten und der bedingten Unbestimmtheit des zweiten, gegeben das erste. (E3) verlangt, daß Ausgänge eines Experiments, die nur mit der Wahrscheinlichkeit 0 auftreten, die Unbestimmtheit nicht ändern.

**Satz 3.2** (*Axiomatische Charakterisierung der Entropie, Chintchin 1953*)  
 Sei  $H : \mathcal{P} = \bigcup_{m=1}^{\infty} \mathcal{P}_m \rightarrow \mathbb{R}$ ,  $H \not\equiv 0$ , eine reellwertige Abbildung auf der Menge aller stochastischen Vektoren, die den Eigenschaften (E1), (E2) und (E3) genügt. Ferner sei  $H|_{\mathcal{P}_m}$  (die Einschränkung von  $H$  auf  $\mathcal{P}_m$ ) stetig für alle  $m \in \mathbb{N}$ . Dann existiert eine Konstante  $c > 0$  mit

$$H(p_1, \dots, p_m) = -c \sum_{i=1}^m p_i \log p_i.$$

**Beweis.** Setze  $H(\frac{1}{m}, \dots, \frac{1}{m}) = f(m)$ ,  $m \in \mathbb{N}$ . Zunächst wird gezeigt, daß ein  $c > 0$  existiert mit  $f(m) = c \log m$  für alle  $m \in \mathbb{N}$ . Da wegen (E2) und (E1)  $f(m) = H(\frac{1}{m}, \dots, \frac{1}{m}, 0) \leq H(\frac{1}{m+1}, \dots, \frac{1}{m+1}) = f(m+1)$  gilt, ist  $f$  monoton steigend.

Seien nun  $r, s \in \mathbb{N}$ . Setzt man in (E2)  $m = r$ ,  $n = r^{s-1}$  und  $p_{ij} = p_i \cdot p_j$ , wobei  $p_i = \frac{1}{r}$ ,  $p_j = \frac{1}{r^{s-1}}$  für alle  $i, j$ , so folgt  $p_{ij} = \frac{1}{r^s}$ , und die gemeinsame Verteilung bildet eine Gleichverteilung auf  $\{1, \dots, r\} \times \{1, \dots, r^{s-1}\}$  mit stochastisch unabhängigen Komponenten, wie in der folgenden Tabelle skizziert.

$$\underbrace{\begin{array}{ccc|c} 1/r^s & \dots & 1/r^s & 1/r \\ \vdots & & \vdots & \vdots \\ 1/r^s & \dots & 1/r^s & 1/r \end{array}}_{r^{s-1}}$$

Dies eingesetzt in (E2) liefert

$$\begin{aligned} H\left(\frac{1}{r^s}, \dots, \frac{1}{r^s}\right) &= H\left(\frac{1}{r}, \dots, \frac{1}{r}\right) + \frac{1}{r} \sum_{i=1}^r H\left(\frac{1}{r^{s-1}}, \dots, \frac{1}{r^{s-1}}\right) \\ &= H\left(\frac{1}{r}, \dots, \frac{1}{r}\right) + H\left(\frac{1}{r^{s-1}}, \dots, \frac{1}{r^{s-1}}\right). \end{aligned}$$

Mit vollständiger Induktion folgt, daß für alle  $r, s \in \mathbb{N}$

$$f(r^s) = H\left(\frac{1}{r^s}, \dots, \frac{1}{r^s}\right) = sH\left(\frac{1}{r}, \dots, \frac{1}{r}\right) = sf(r). \quad (3.7)$$

Insbesondere ist  $f(1) = sf(1)$  für alle  $s \in \mathbb{N}$ , also gilt  $f(1) = 0$ .

Angenommen  $f(2) = 0$ . Dann ist  $f(2^s) = sf(2) = 0$  für alle  $s \in \mathbb{N}$ , wegen der Monotonie von  $f$  also  $f \equiv 0$ , was nach Voraussetzung ausgeschlossen ist. Also gilt  $f(s) > 0$  für alle  $s \geq 2$ .

Für alle  $r, s, n \in \mathbb{N}$ ,  $r, s \geq 2$  existiert nun  $m \in \mathbb{N}_0$  mit

$$r^m \leq s^n < r^{m+1}. \quad (3.8)$$

Dies äquivalent umgeformt liefert

$$m \log r \leq n \log s < (m+1) \log r \quad \text{bzw.} \quad \frac{m}{n} \leq \frac{\log s}{\log r} < \frac{m+1}{n}. \quad (3.9)$$

Aus der Monotonie von  $f$  folgt mit (3.8), daß  $f(r^m) \leq f(s^n) \leq f(r^{m+1})$ . Mit (3.7) erhält man

$$mf(r) \leq nf(s) \leq (m+1)f(r) \quad \text{bzw.} \quad \frac{m}{n} \leq \frac{f(s)}{f(r)} \leq \frac{m+1}{n}, \quad (3.10)$$

so daß wegen (3.9) und (3.10) für alle  $r, s, n \in \mathbb{N}$ ,  $r, s \geq 2$  die Ungleichung

$$\left| \frac{f(s)}{f(r)} - \frac{\log s}{\log r} \right| \leq \frac{1}{n}$$

folgt. Durch Grenzübergang  $n \rightarrow \infty$  schließt man  $\frac{f(s)}{\log s} = \frac{f(r)}{\log r} > 0$  für alle  $r, s \geq 2$ . Dies liefert die Existenz einer Konstanten  $c > 0$ , derart daß

$f(s) = c \log s$  für alle  $s \geq 2$ . Insgesamt wurde damit bewiesen, daß eine Konstante  $c > 0$  existiert mit

$$H\left(\frac{1}{m}, \dots, \frac{1}{m}\right) = c \cdot \log m. \quad (3.11)$$

für alle  $m \in \mathbb{N}$ .

Seien jetzt  $(p_1^*, \dots, p_m^*) \in \mathcal{P}_m$ ,  $p_i^* \in \mathbb{Q}$ ,  $i = 1, \dots, m$ , mit der Eigenschaft  $p_i^* = \frac{g_i}{g}$ ,  $\sum_{i=1}^m g_i = g$ ,  $g_i \in \mathbb{N}$ . Ausgangspunkt ist die zweidimensionale Verteilung  $p_{ij}$ ,  $i = 1, \dots, m$ ,  $j = 1, \dots, g$  aus, die durch folgende Tabelle beschrieben wird.

$$\begin{array}{cccc|c} \hline \underbrace{1/g \cdots 1/g}_{g_1} & \underbrace{0 \cdots 0}_{g_2} & \underbrace{0 \cdots 0 \cdots 0}_{g_m} & 0 \cdots 0 & g_1/g = p_1^* \\ 0 \cdots 0 & 1/g \cdots 1/g & 0 \cdots 0 \cdots 0 & 0 \cdots 0 & g_2/g = p_2^* \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 0 \cdots 0 & 0 \cdots 0 & 0 \cdots 0 \cdots 1/g \cdots 1/g & & g_m/g = p_m^* \\ \hline \end{array}$$

Hierbei gilt für alle  $i = 1, \dots, m$ ,  $j = 1, \dots, g$

$$p_{i|j} = \frac{p_{ij}}{p_i} = \begin{cases} 1/g_i, & \text{falls } j \text{ im } i\text{-ten Block,} \\ 0, & \text{sonst} \end{cases}.$$

Mit Eigenschaft (E2) folgt für diese Verteilung

$$\begin{aligned} & H\left(\underbrace{\frac{1}{g} \cdots \frac{1}{g}}_{g_1}, 0, \dots, 0, \underbrace{\frac{1}{g} \cdots \frac{1}{g}}_{g_2}, 0, \dots, 0, \underbrace{\frac{1}{g} \cdots \frac{1}{g}}_{g_m}\right) \\ &= H(p_1^*, \dots, p_m^*) + \sum_{i=1}^m p_i^* H\left(0, \dots, 0, \underbrace{\frac{1}{g_i} \cdots \frac{1}{g_i}}_{g_i}, 0, \dots, 0\right). \end{aligned}$$

Wendet man nun Eigenschaft (E3) und (3.11) auf diese Gleichung an, erhält man  $c \log g = H(p_1^*, \dots, p_m^*) + \sum_{i=1}^m p_i^* (c \log g_i)$ , so daß

$$\begin{aligned} H(p_1^*, \dots, p_m^*) &= c \left( \log g - \sum_{i=1}^m p_i^* \log g_i \right) \\ &= -c \sum_{i=1}^m p_i^* \log \frac{g_i}{g} = -c \sum_{i=1}^m p_i^* \log p_i^*. \end{aligned}$$

Damit ist die Behauptung für stochastische Vektoren mit Komponenten in den rationalen Zahlen  $\mathbb{Q}$  bewiesen.

Für beliebiges  $(p_1, \dots, p_m) \in \mathcal{P}_m$  existiert eine Folge  $(p_1^{(\ell)}, \dots, p_m^{(\ell)}) \in \mathcal{P}_m$ ,  $p_i^{(\ell)} \in \mathbb{Q}$ ,  $\ell \in \mathbb{N}$  mit  $(p_1^{(\ell)}, \dots, p_m^{(\ell)}) \rightarrow (p_1, \dots, p_m)$  ( $\ell \rightarrow \infty$ ). Wegen der Stetigkeit von  $H$  auf  $\mathcal{P}_m$  folgt

$$\begin{aligned} H(p_1, \dots, p_m) &= \lim_{\ell \rightarrow \infty} H(p_1^{(\ell)}, \dots, p_m^{(\ell)}) \\ &= \lim_{\ell \rightarrow \infty} \left( -c \sum_{i=1}^m p_i^{(\ell)} \log p_i^{(\ell)} \right) = -c \sum_{i=1}^m p_i \log p_i, \end{aligned}$$

womit die Behauptung für beliebiges  $(p_1, \dots, p_m) \in \mathcal{P}$  bewiesen ist.  $\blacksquare$

### 3.3 Übungsaufgaben

**Aufgabe 3.1** Man betrachte die Relation “ $\prec$ ” (majorisiert) auf  $\mathcal{P}_m$  (vgl. Def. 3.1). Zeigen Sie:

- a) “ $\prec$ ” definiert eine Präordnung auf  $\mathcal{P}_m$ , d.h. für alle  $p, q, r \in \mathcal{P}_m$  gilt
  - (i)  $p \prec q$ ,   (ii) aus  $p \prec q$  und  $q \prec r$  folgt  $p \prec r$ .
- b)  $q \in \mathcal{P}_m$  heißt größtes (kleinstes) Element von  $(\mathcal{P}_m, \prec)$ , falls  $p \prec q$  ( $q \prec p$ ) für alle  $p \in \mathcal{P}_m$  gilt. Man bestimme alle größten und kleinsten Elemente von  $(\mathcal{P}_m, \prec)$ .

**Aufgabe 3.2** Die binären Zufallsvariablen  $X$  und  $Y$  bezeichnen die Ein- bzw. Ausgabe bei einem binären symmetrischen Kanal mit Übertragungswahrscheinlichkeit  $\varepsilon \in [0, 1]$  (vgl. Beispiel 3.1). Zeigen Sie:  $H(X) \leq H(Y)$ . Wann tritt Gleichheit ein?

**Aufgabe 3.3** Bei einem Übertragungssystem bestehe Eingabe- und Ausgabealphabet aus den Buchstaben  $\{x_1, x_2, x_3\}$ . Jeder Buchstabe trete mit gleicher Wahrscheinlichkeit als Eingabe auf und werde mit Wahrscheinlichkeit  $1 - \varepsilon$ ,  $0 \leq \varepsilon \leq 1$ , nach Übertragung richtig ausgegeben oder je mit Wahrscheinlichkeit  $\varepsilon/2$  als einer der verbleibenden Buchstaben fälschlich ausgegeben. Bestimmen Sie ein stochastisches Modell, das dieses Übertragungssystem beschreibt, und bestimmen Sie die Entropie der Ausgabe, die Entropie der Eingabe, die bedingte Entropie der Ausgabe, gegeben die Eingabe, sowie die Transinformation von Ein- und Ausgabe. Skizzieren Sie diese Größen als Funktionen von  $\varepsilon$ ,  $0 \leq \varepsilon \leq 1$ . Interpretieren Sie diese Funktionen.

**Aufgabe 3.4**

- a) Man zeige:  $S = \sum_{n=2}^{\infty} n^{-1}(\log n)^{-2} < \infty$ .
- b) Auf  $(\mathbb{N}, \mathfrak{P}(\mathbb{N}))$  sei eine Wahrscheinlichkeitsverteilung  $P$  definiert durch

$$P(\{n\}) = p_n = \begin{cases} S^{-1} n^{-1}(\log n)^{-2}, & \text{falls } n \geq 2 \\ 0, & \text{falls } n = 1 \end{cases}$$

Man berechne  $-\sum_{n=1}^{\infty} p_n \log p_n$ .

- c) Auf  $(\mathbb{N}, \mathfrak{P}(\mathbb{N}))$  sei eine Wahrscheinlichkeitsverteilung  $P$  gegeben mit  $P(\{n\}) = p_n$ . Ferner gelte: (i)  $p_1 \geq p_2 \geq p_3 \geq \dots$  und (ii)  $\sum_{n=1}^{\infty} p_n \log n < \infty$ .  
Man zeige:  $-\sum_{n=1}^{\infty} p_n \log p_n < \infty$ .

Hinweis: Für monoton fallende Funktionen  $f$  gilt  $\sum_{k=2}^n f(k) \leq \int_1^n f(x) dx \leq \sum_{k=1}^{n-1} f(k)$ .

**Aufgabe 3.5**  $X : \Omega \rightarrow \mathcal{X}$ ,  $\mathcal{X} = \{x_1, \dots, x_m\}$  sei eine endlich diskrete Zufallsvariable auf einem Wahrscheinlichkeitsraum  $(\Omega, \mathcal{A}, P)$  und  $g : \mathcal{X} \rightarrow \mathcal{Y}$ ,  $\mathcal{Y} = \{y_1, \dots, y_n\}$  eine Abbildung. Man zeige, daß  $H(g(X)) \leq H(X)$ . Wann gilt Gleichheit?

**Aufgabe 3.6**  $(\Omega, \mathcal{A}, P)$  sei ein Wahrscheinlichkeitsraum und  $\{\mathcal{E}_n\}_{n \in \mathbb{N}}$  eine aufsteigende Folge von endlichen, meßbaren Partitionen von  $\Omega$ , d.h. für alle  $n \in \mathbb{N}$  gilt  $\mathcal{E}_n = \{E_{n,1}, E_{n,2}, \dots, E_{n,k(n)}\}$ ,  $k(n) \in \mathbb{N}$  mit  $E_{n,j} \in \mathcal{A}$ ,  $\Omega = \bigcup_{j=1}^{k(n)} E_{n,j}$ ,  $E_{n,i} \cap E_{n,j} = \emptyset$  für alle  $i \neq j$ ,  $i, j \in \{1, \dots, k(n)\}$ , und für alle  $j \in \{1, \dots, k(n)\}$  existieren Indizes  $T \subset \{1, \dots, k(n+1)\}$  mit  $E_{n,j} = \bigcup_{i \in T} E_{n+1,i}$ .  
Die Zufallsvariablen  $X_n$  seien definiert durch

$$X_n(\omega) = \sum_{i=1}^{k(n)} i \cdot \mathbb{I}_{E_{n,i}}(\omega), \quad n \in \mathbb{N},$$

wobei  $\mathbb{I}$  die Indikatorfunktion bezeichnet,  $\mathbb{I}_A(x) = 1$ , falls  $x \in A$ , und 0 sonst.  
Zeigen Sie:  $X_n$  ist für alle  $n \in \mathbb{N}$  eine endlich diskrete Zufallsvariable, und die Folge der Entropien  $\{H(X_n)\}_{n \in \mathbb{N}}$  ist schwach monoton wachsend. Ist  $\{H(X_n)\}_{n \in \mathbb{N}}$  stets nach oben beschränkt?

### Aufgabe 3.7

- a) Die Zufallsvariable  $X$  sei  $\text{Bin}(n, p)$ -verteilt, also  $P(X = k) = \binom{n}{k} p^k q^{n-k}$ ,  $0 \leq k \leq n$ , mit  $0 \leq p \leq 1$  und  $q = 1 - p$ .  
Man zeige:  $H(X) \leq -n(p \log p + q \log q)$ .
- b) In einer unabhängigen Versuchsserie mit den Ausgängen 0 und 1 mit Eintrittswahrscheinlichkeit  $1-p$  und  $p$ ,  $0 < p < 1$ , bezeichne  $X_n$  die Anzahl der Versuche unter den ersten  $n$ , die der ersten 1 vorausgehen.  
Berechnen Sie  $H(X_n)$ . Was ergibt sich hier für  $n \rightarrow \infty$ ?

**Aufgabe 3.8** Die Zufallsvariable  $X$  beschreibe die Eingabe, die Zufallsvariable  $Y$  die Ausgabe eines binären symmetrischen Kanals mit Fehlerwahrscheinlichkeit  $\varepsilon$ ,  $0 \leq \varepsilon \leq 1$  (s. Beispiel 3.1). Zeigen Sie, daß  $H(X) \leq H(Y)$ . Wann tritt Gleichheit ein?

**Aufgabe 3.9**

a)  $X, Y$  und  $Z$  seien endlich diskrete Zufallsvariablen. Man beweise

$$H((X, Y) | Z) \leq H(X | Z) + H(Y | Z).$$

b)  $\{X_n\}_{n \in \mathbb{N}}$  sei eine Folge von endlich diskreten Zufallsvariablen auf einem Wahrscheinlichkeitsraum  $(\Omega, \mathcal{A}, P)$ . Zeigen Sie: Die Folge der bedingten Entropien  $\{H(X_1 | (X_2, X_3, \dots, X_n))\}_{n=2}^{\infty}$  ist monoton fallend, dagegen ist die Folge der bedingten Entropien  $\{H((X_2, X_3, \dots, X_n) | X_1)\}_{n=2}^{\infty}$  monoton steigend.

**Aufgabe 3.10** Die Entfernung  $d(X, Y)$  zweier endlich diskreter Zufallsvariablen  $X, Y$  sei definiert durch  $d(X, Y) = H(X | Y) + H(Y | X)$ . Man zeige, daß für drei endlich diskrete Zufallsvariablen  $X, Y, Z$  die Dreiecksungleichung gilt:

$$d(X, Y) \leq d(X, Z) + d(Z, Y).$$

**Aufgabe 3.11** Man beweise: Für jede reelle Zahl  $b > 1$  ist die Funktion  $f = \log_b$  die einzige stetige Funktion  $f : (0, \infty) \rightarrow \mathbb{R}$ , die der Funktionalgleichung  $f(xy) = f(x) + f(y)$  und der Anfangsbedingung  $f(b) = 1$  genügt.

**Aufgabe 3.12**

a) Sei  $\mathcal{X}_\mu$  die Menge aller endlich diskreten Zufallsvariablen  $X$ , die höchstens die Werte  $v_1, \dots, v_n \in \mathbb{R}$  mit positiver Wahrscheinlichkeit annehmen und den Erwartungswert  $E(X) = \mu$  für ein festes  $\mu \in \mathbb{R}$  besitzen. Man zeige:

$$\max_{X \in \mathcal{X}_\mu} H(X) = H(X^*)$$

wird angenommen für eine Zufallsvariable  $X^*$  deren Verteilung gegeben ist durch  $p_i = P(X^* = v_i) = e^{\beta v_i} \left( \sum_{j=1}^n e^{\beta v_j} \right)^{-1}$ ,  $1 \leq i \leq n$ , mit einem eindeutigen "Verteilungsmodul"  $-\infty \leq \beta \leq \infty$  (Boltzmann-Verteilung).

b) Für  $n = 3$  und  $v_i = i$ ,  $1 \leq i \leq 3$ , bestimme man  $\beta$  explizit.

**Aufgabe 3.13**  $\mathcal{P}^* = \bigcup_{m \in \mathbb{N}} \mathcal{P}_m$  bezeichne die Menge aller Wahrscheinlichkeitsvektoren.  $H : \mathcal{P}^* \rightarrow \mathbb{R}$  sei eine Funktion mit folgenden Eigenschaften. Für alle  $m \in \mathbb{N}$ ,  $(p_1, \dots, p_m) \in \mathcal{P}_m$  gilt:



### 42 3 Information und Entropie

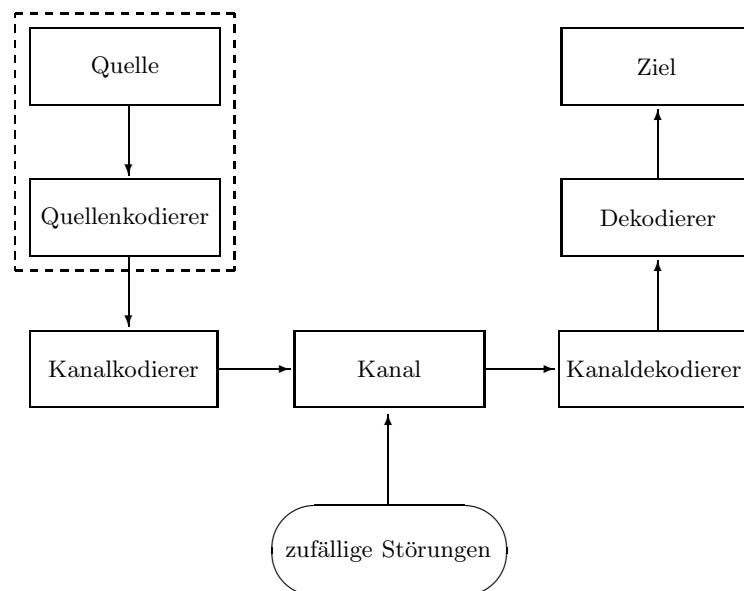
- 1.)  $H \neq 0$ ,
- 2.)  $H|_{\mathcal{P}_m}$  ist stetig und symmetrisch, d.h.  $H(p_1, \dots, p_m) = H(p_{\pi(1)}, \dots, p_{\pi(m)})$  für alle Permutationen  $\pi$ ,
- 3.)  $H(p_1, \dots, p_m) \leq H(\frac{1}{m}, \dots, \frac{1}{m})$ ,
- 4.)  $H(p_1, \dots, p_m, 0) = H(p_1, \dots, p_m)$ ,
- 5.)  $H(p_1, \dots, p_m) = H(p_1 + p_2, p_3, \dots, p_m) + (p_1 + p_2)H(\frac{p_1}{p_1+p_2}, \frac{p_2}{p_1+p_2})$ .

Zeigen Sie, daß dann eine Konstante  $b > 0$  existiert mit

$$H(p_1, \dots, p_m) = - \sum_{i=1}^m p_i \log_b p_i.$$

## 4 Kodierung diskreter Quellen

Bei dem in der Einleitung vorgestellten Standardmodell wird nun der Zusammenhang zwischen Quelle und Quellenkodierer untersucht, also der in der folgenden Graphik gestrichelt umrahmte Block.



Generelle Zielsetzung ist, die von der Quelle aus einem bestimmten Alphabet gebildeten Wörter über einem (in der Regel binären) Kodealphabet möglichst kompakt zu kodieren. Dieser Vorgang wird als fehlerfrei angenommen, so daß keine Redundanz zu einer eventuell nötigen Fehlerkorrektur vorgesehen werden muß. Die Hauptaufgabe des Quellenkodierers ist also Datenkompression und Übersetzung der eingehenden Quellsignale.

Die Quelle benutzt das  $m$ -elementige Alphabet  $\mathcal{X} = \{x_1, \dots, x_m\}$  und sendet Wörter  $(u_1, \dots, u_N) \in \mathcal{X}^* = \bigcup_{\ell=0}^{\infty} \mathcal{X}^\ell$ . Der Kodierer verwendet das Alphabet  $\mathcal{Y} = \{y_1, \dots, y_d\}$ , er kodiert  $(u_1, \dots, u_N)$  durch  $(w_1, \dots, w_L) \in$

$\mathcal{Y}^* = \bigcup_{\ell=0}^{\infty} \mathcal{Y}^\ell$ . Wenn mit Codes fester Länge gearbeitet wird, läßt sich die kleinste Wortlänge  $L$ , so daß jedes Quellwort der Länge  $N$  kodiert werden kann, aus folgender Ungleichung bestimmen.

$$|\mathcal{X}^N| = m^N \leq d^L = |\mathcal{Y}^L| \iff L \geq N \frac{\log m}{\log d}.$$

Bei natürlichsprachigen Quellen treten jedoch viele Quellwörter nur mit sehr kleiner Wahrscheinlichkeit auf. So gibt es  $26^N$  Wörter der Länge  $N$ , doch davon wird in einer natürlichen Sprache nur ein winziger Teil genutzt. Wieviele sinnvolle Wörter mit 5 Buchstaben sind Ihnen im Deutschen bekannt? Man kann aus 26 Buchstaben 11.881.376 verschiedene Wörter mit 5 Buchstaben bilden, einen Sinn haben die wenigsten davon. Es stellt sich die Frage, ob  $L$  nicht wesentlich kleiner gewählt werden kann als  $N \log m / \log d$ , und zwar so, daß die Wahrscheinlichkeit, für ein Quellwort kein Kodewort zur Verfügung zu haben, sehr klein ist. Dieses Problem wird in Abschnitt 4.1 über Codes fester Länge behandelt.

Werden hingegen Codes variabler Länge benutzt und ist die Verteilung des Auftretens einzelner Quellbuchstaben bekannt, wird man den Code so konstruieren, daß häufig auftretende Quellbuchstaben Kodewörter kurzer Länge und selten auftretende die längeren Kodewörter zugewiesen bekommen. Beabsichtigt wird hiermit, die Kodewortlänge im Mittel möglichst kurz zu machen. Dieses Ziel wird in Abschnitt 4.2 verfolgt.

Die beiden folgenden einführenden Beispiele verdeutlichen grundsätzliche Konzepte bei der Kodierung. Binäre Kodierung wird hierbei als ja- oder nein-Antwort auf Serien von entsprechenden Fragen interpretiert.

**Beispiel 4.1**  $\mathcal{X} = \{x_1, \dots, x_5\}$  sei eine fünfelementige Menge. Der stochastische Vektor  $\mathbf{p} = (0.3, 0.2, 0.2, 0.15, 0.15)$  beschreibe die Verteilung eines Zufallsexperiment mit Ergebnismenge  $\mathcal{X}$ . Wir stellen uns die Aufgabe, mit ja/nein-Fragen den Ausgang des Experiments von jemandem zu erfragen, der weiß, wie das Experiment ausgegangen ist. Die folgende naive Fragestrategie liegt direkt auf der Hand. Bei der  $i$ -ten Frage wird gefragt: ‘War der Ausgang  $x_i$ ?’,  $i = 1, \dots, 4$ . Spätestens nach 4 Fragen weiß man den Ausgang des Experiments, da viermal ‘nein’ auf den Ausgang  $x_5$  schließen läßt. Die Wahrscheinlichkeit, hierbei den Ausgang mit genau einer Frage zu erfahren, ist 0.3, mit genau zwei Fragen 0.5, u.s.w.. Eine Kodierung der Ergebnisse

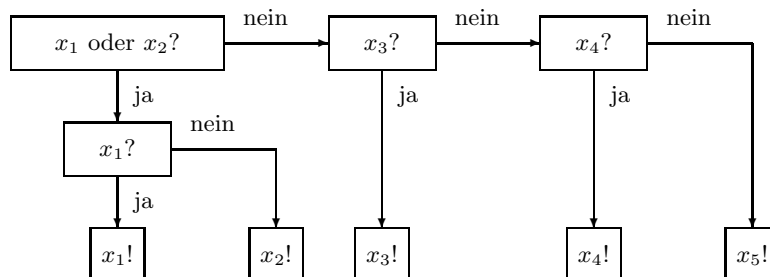
des Experiments  $(\mathcal{X}, \mathbf{p})$  kann unter dieser Strategie durch  $(1 \hat{=} \text{ja und } 0 \hat{=} \text{nein})$  mit den Kodewörtern

$$\begin{aligned} x_1 &\mapsto (1), & x_2 &\mapsto (0, 1), & x_3 &\mapsto (0, 0, 1), \\ x_4 &\mapsto (0, 0, 0, 1), & x_5 &\mapsto (0, 0, 0, 0, ) \end{aligned}$$

vorgenommen werden. Die erwartete Anzahl von Fragen ist

$$1 \cdot 0.3 + 2 \cdot 0.2 + 3 \cdot 0.2 + 4 \cdot 0.15 + 4 \cdot 0.15 = 2.5.$$

Es gibt jedoch bessere Fragestrategien, die im Mittel mit weniger Fragen auskommen. Die im folgenden Graphen dargestellte Strategie unterbietet die naive bezüglich der erwarteten Kodewortlänge.



Den Ausgang  $x_1$  des Experiments erfährt man nach zwei ja-Antworten,  $x_2$  nach einer ja- und einer nein-Antwort, u.s.w.. Eine Kodierung der Ergebnisse des Experiments erfolgt wie oben durch

$$x_1 \mapsto (1, 1), \quad x_2 \mapsto (1, 0), \quad x_3 \mapsto (0, 1), \quad x_4 \mapsto (0, 0, 1), \quad x_5 \mapsto (0, 0, 0).$$

Die Anzahl der Fragen oder äquivalent die Länge des Kodewortes bei dieser Strategie ist eine Zufallsvariable

$$\begin{aligned} Z : \mathcal{X} &\rightarrow \mathbb{N} : x_i \mapsto \text{Anzahl der Fragen, mit denen man } x_i \\ &\text{erfragen kann} \\ &= \text{Länge des zugehörigen Kodewortes.} \end{aligned}$$

Hier gilt  $Z(x_1) = Z(x_2) = Z(x_3) = 2$ ,  $Z(x_4) = Z(x_5) = 3$ , ferner  $P(Z = 2) = P(\{x_1, x_2, x_3\}) = 0.7$ ,  $P(Z = 3) = P(\{x_4, x_5\}) = 0.3$ . Die erwartete

Anzahl von Fragen, die in jedem Fall zum Ziel führen, oder äquivalent die erwartete Länge der verwendeten Kodewörter, beträgt

$$E(Z) = 0.7 \cdot 2 + 0.3 \cdot 3 = 2.3.$$

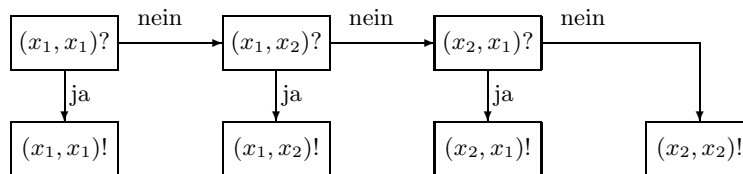
Die Strategie ist in der Tat besser als die naive Strategie von oben. Sie ist sogar die beste, wie später bewiesen wird. Die Entropie der Verteilung  $\mathbf{p}$  berechnet sich zu

$$\begin{aligned} H(p_1, \dots, p_5) &= -(0.3 \log_2 0.3 + 0.4 \log_2 0.2 + 0.3 \log_2 0.15) \\ &\approx 2.27095, \end{aligned}$$

sie ist etwas kleiner als  $E(Z)$ . Wir werden später sehen, daß dieser Sachverhalt allgemein gilt: Ist die Basis der verwendeten Logarithmen  $m$ , so ist die Entropie eines Zufallsexperiments stets kleiner als die erwartete Kodewortlänge eines Codes, der die zugehörigen Ergebnisse eindeutig kodiert. Die Entropie ist also eine untere Schranke für das Effizienzmaß "erwartete Kodewortlänge". ■

**Beispiel 4.2** (Blockfragestrategie)

Seien  $X_1, X_2$  stochastisch unabhängige Zufallsvariable, jeweils mit Träger  $\mathcal{X} = \{x_1, x_2\}$ .  $\mathbf{X} = (X_1, X_2)$  ist dann ein Zufallsvektor mit Träger  $\mathcal{X}^2$ .  $P(X_i = x_1) = 0.7$ ,  $P(X_i = x_2) = 0.3$ ,  $i = 1, 2$ , beschreibe die Verteilung von  $X_1$  bzw.  $X_2$ . Wegen der stochastischen Unabhängigkeit ist hiermit auch die gemeinsame Verteilung von  $(X_1, X_2)$  festgelegt. Die Aufgabe lautet, mit ja/nein-Fragen den Ausgang des Experiments zu erfragen, das durch  $(X_1, X_2)$  beschrieben wird. Eine zugehörige Blockfragestrategie wird analog zu Beispiel 4.1 graphisch dargestellt.



Wie im vorigen Beispiel lassen sich Antwortfolgen, die zu einem bestimmten Ergebnis führen, als Kodierung des Ergebnisses interpretieren, und zwar

$$\begin{aligned} (x_1, x_1) &\mapsto (1), & (x_1, x_2) &\mapsto (0, 1), \\ (x_2, x_1) &\mapsto (0, 0, 1), & (x_2, x_2) &\mapsto (0, 0, 0). \end{aligned}$$

Um die Effizienz der Fragestrategie (oder Kodierung) mit Fällen vergleichbar zu machen, bei denen lediglich ein Ergebnis erfragt wird, wählen wir als Maß die erwartete Anzahl der Fragen pro Komponente. Bezeichnet die Zufallsvariable  $Z$  wie oben die Anzahl der Fragen, so gilt im vorliegenden Beispiel

$$\frac{1}{2}E(Z) = \frac{1}{2}(1 \cdot 0.49 + 2 \cdot 0.21 + 3 \cdot 0.21 + 3 \cdot 0.09) = 0.905.$$

Für die Entropie von  $X_1$  bzw.  $X_2$  gilt

$$H(X_1) = H(X_2) = -0.7 \log_2 0.7 - 0.3 \log_2 0.3 = 0.8813,$$

sie ist wieder kleiner als das Effizienzmaß  $\frac{1}{2}E(Z)$ .

Wir werden später sehen, daß dies allgemein gilt. Mehr noch, bezeichnet  $E(Z_N)$  die erwartete Kodewortlänge bei Blockkodierung von Blöcken der Länge  $N$  und sind  $X_1, \dots, X_N$  stochastisch unabhängige, identisch verteilte Zufallsvariable ein Modell für die Quelle, so existieren Blockcodes mit  $\frac{1}{N} E(Z_N) \rightarrow H(X_1)$  ( $N \rightarrow \infty$ ). ■

## 4.1 Kodes fester Länge

Wie bereits festgestellt, benutzen natürliche Sprachen nur einen sehr kleinen Teil der Menge aller theoretisch möglichen Wörter bestimmter Länge über einem gegebenen Alphabet. Ein adäquates stochastisches Modell muß diesen Sachverhalt widerspiegeln und den meisten Wörtern eine nur sehr kleine Wahrscheinlichkeit geben. In der Tat läßt sich diese Situation schon im einfachen Modell von stochastisch unabhängigen Buchstaben beobachten. Hiermit beschäftigt sich der folgende Abschnitt.

Die Menge der Quellwörter läßt sich in *typische* und *untypische* einteilen. Die untypischen besitzen alle zusammen für große Wortlängen eine beliebig kleine Wahrscheinlichkeit. Werden jetzt nur für die typischen Quellwörter Kodewörter zur Verfügung gestellt, so ist die Wahrscheinlichkeit, für ein gegebenes Wort kein Kodewort zu haben, sehr klein. Wir beginnen mit der Definition eines einfachen Quellenmodells.

**Definition 4.1** (Diskrete gedächtnislose Quelle, DGQ)

$\{X_n\}_{n \in \mathbb{N}}$  sei eine Folge von stochastisch unabhängigen, identisch verteilten Zufallsvariablen mit Verteilung  $P^{X_n} = P^X$  für alle  $n \in \mathbb{N}$  und Träger  $\mathcal{X} = \{x_1, \dots, x_m\}$ .  $\{X_n\}_{n \in \mathbb{N}}$  (oder auch  $X$ ) heißt diskrete gedächtnislose Quelle (DGQ) (discrete memoryless source).  $\mathcal{X}$  heißt Alphabet der Quelle.

**Satz 4.1** (Asymptotische Gleichverteilungseigenschaft, Feinstein 1959)

$\{X_n\}_{n \in \mathbb{N}}$  sei eine DGQ mit Alphabet  $\mathcal{X} = \{x_1, \dots, x_m\}$  und Entropie  $H = H(X)$ , gebildet mit Logarithmen zur Basis 2. Dann gelten die folgenden Aussagen.

$\forall \varepsilon, \delta > 0 \exists N_0 \in \mathbb{N} \forall N > N_0 \exists T \subseteq \mathcal{X}^N :$

- a)  $2^{-N(H+\delta)} \leq P(X_1 = a_1, \dots, X_N = a_N) \leq 2^{-N(H-\delta)}$   
für alle  $(a_1, \dots, a_N) \in T$
- b)  $P((X_1, \dots, X_N) \in T^c) \leq \varepsilon$
- c)  $(1 - \varepsilon)2^{N(H-\delta)} \leq |T| \leq 2^{N(H+\delta)}$

**Beweis.** Bezeichne  $\mathbf{X}_N = (X_1, \dots, X_N)$  und  $\mathbf{a}_N = (a_1, \dots, a_N)$ . Der Informationsgehalt des Ereignisses  $\{\mathbf{X}_N = \mathbf{a}_N\}$ ,  $\mathbf{a}_N \in \mathcal{X}^N$ , ist wegen der stochastischen Unabhängigkeit von  $X_1, \dots, X_N$

$$\begin{aligned} I_{\mathbf{X}_N}(\mathbf{a}_N) &= -\log P(\mathbf{X}_N = \mathbf{a}_N) \\ &= -\log (P(X_1 = a_1) \cdots P(X_N = a_N)) \\ &= -\sum_{i=1}^N \log P(X_i = a_i) = \sum_{i=1}^N I_{X_i}(a_i). \end{aligned} \quad (4.1)$$

Die Zufallsvariablen

$$\begin{aligned} I_{X_k} &: \left( \mathcal{X}^\infty, \bigotimes_{i=1}^{\infty} \mathfrak{P}(\mathcal{X}), P^{\{X_n\}} \right) \rightarrow \mathbb{R} \\ &: \{a_n\}_{n \in \mathbb{N}} \mapsto -\log P(X_k = a_k), \quad k \in \mathbb{N}, \end{aligned}$$

bilden eine stochastisch unabhängige Folge  $\{I_{X_k}\}_{k \in \mathbb{N}}$ . Mit dem starken Gesetz großer Zahlen (vgl. Satz 2.3) folgt

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N I_{X_k} = \mathbb{E}(I_{X_1}) = H(X) \quad P^{\{X_n\}}\text{-fast sicher.}$$

Hieraus folgt wie in (2.4) stochastische Konvergenz, also für alle  $\varepsilon, \delta > 0$  existiert ein  $N_0$ , derart daß  $N > N_0$

$$P^{\{X_N\}} \left( \left| \frac{1}{N} \sum_{k=1}^N I_{X_k} - H(X) \right| > \delta \right) \leq \varepsilon.$$

Wegen (4.1) ist dies äquivalent zu

$$P^{(X_1, \dots, X_N)} \left( \left\{ (a_1, \dots, a_N) \in \mathcal{X}^N \mid \left| \frac{1}{N} I_{\mathbf{X}_N}(\mathbf{a}_N) - H(X) \right| > \delta \right\} \right) \leq \varepsilon. \quad (4.2)$$

Mit der Menge

$$T = \left\{ (a_1, \dots, a_N) \in \mathcal{X}^N \mid \left| \frac{1}{N} I_{\mathbf{X}_N}(\mathbf{a}_N) - H(X) \right| \leq \delta \right\}$$

wird nun die Menge der typischen Quellwörter bezeichnet. Wegen (4.2) gilt die Abschätzung  $P^{(X_1, \dots, X_N)}(T) \geq 1 - \varepsilon$  und  $P^{(X_1, \dots, X_N)}(T^c) \leq \varepsilon$ . Damit ist b) gezeigt.

Ferner gilt mit Logarithmen zur Basis 2 für alle  $(a_1, \dots, a_N) \in T$ , daß  $N(H - \delta) \leq I_{\mathbf{X}_N}(a_1, \dots, a_N) \leq N(H + \delta)$  genau dann, wenn  $2^{-N(H-\delta)} \geq P(X_1 = a_1, \dots, X_N = a_N) \geq 2^{-N(H+\delta)}$ , woraus a) folgt.

Die Mächtigkeit von  $T$  wird mit Hilfe der folgenden Ungleichungen abgeschätzt.

$$1 \geq P^{\mathbf{X}_N}(T) \geq |T| \cdot \min_{\mathbf{a}_N \in T} P(\mathbf{X}_N = \mathbf{a}_N) \geq |T| \cdot 2^{-N(H+\delta)},$$

so daß  $|T| \leq 2^{N(H+\delta)}$ . Ferner gilt

$$(1 - \varepsilon) \leq P^{\mathbf{X}_N}(T) \leq |T| \cdot \max_{\mathbf{a}_N \in T} P(\mathbf{X}_N = \mathbf{a}_N) \leq |T| \cdot 2^{-N(H-\delta)}.$$

Also ist  $|T| \geq (1 - \varepsilon) \cdot 2^{N(H-\delta)}$ . Hiermit ist c) gezeigt. ■

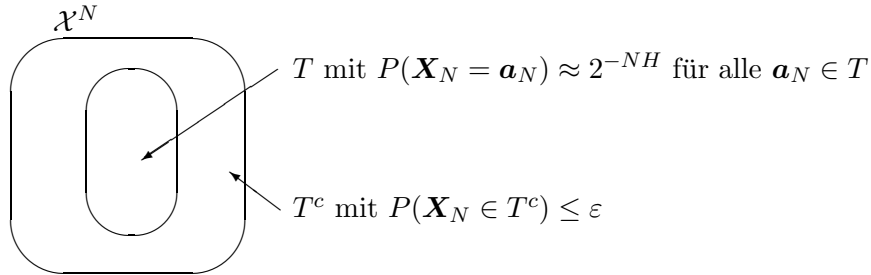
Satz 4.1 besagt, daß für beliebig kleines  $\varepsilon > 0$  für alle genügend großen Wortlängen  $N$  eine Partition von  $\mathcal{X}^N$  in eine Menge von typischen Quellwörtern  $T$  und eine von untypischen  $T^c$  existiert. Die untypischen Quellwörter



haben alle zusammen eine Wahrscheinlichkeit von höchstens  $\varepsilon$ . Für die typischen Quellwörter  $\mathbf{a} \in T$  gilt wegen Teil a)

$$\left| -\frac{1}{N} \log P(\mathbf{X}_N = \mathbf{a}_N) - H \right| \leq \delta.$$

Die mit  $1/N$  normierte logarithmierte Wahrscheinlichkeit für das Auftreten eines typischen Quellworts ist also approximativ konstant gleich der Entropie der Quelle. Können  $\delta$  und  $N$  so gewählt werden, daß für kleine  $\varepsilon$  auch  $N\delta$  nahe bei Null ist, so besitzt jedes  $\mathbf{a} \in T$  approximativ die Wahrscheinlichkeit  $2^{-NH}$ . In diesem Fall enthält  $T$  approximativ  $2^{NH}$  Elemente. Man sagt, die Quelle besitzt die asymptotische Gleichverteilungseigenschaft (asymptotic equipartition property). Diese Zusammenhänge können folgendermaßen skizziert werden.



Die Menge  $T$  in Satz 4.1 ergibt sich aus der Konvergenzaussage des Gesetzes großer Zahlen. Eine konstruktive Bestimmung von  $T$  ist hieraus jedoch nicht möglich. Bei Verwendung einer anderen Basis  $b$  der Logarithmen erhält man analoge Aussagen, indem die Zahl 2 durch die entsprechende Basis  $b$  ersetzt wird.

Satz 4.1 läßt sich ebenso für jede diskrete Quelle  $\{X_n\}_{n \in \mathbb{N}}$  beweisen, solange  $P^{X_N} = P^X$  und  $\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{k=1}^N I_{X_k} = E(X)$   $P^{\{X_n\}}$ -stochastisch gilt. Diese Eigenschaften besitzen etwa stationäre, ergodische Quellen (siehe z.B. [15], [10]). Stationäre Quellen werden in späteren Kapiteln behandelt.

Bei diskreten gedächtnislosen Quellen kann Satz 4.1 a) verschärft werden zu (s. Übungsaufgabe 4.3)  $\forall \varepsilon > 0 \exists \delta(\varepsilon), N_0 \in \mathbb{N} \forall N > N_0 \forall \mathbf{a}_N \in T$ :

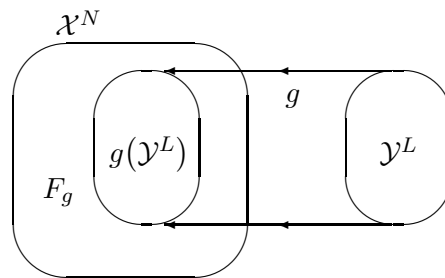
$$2^{-NH - \sqrt{N}\delta} \leq P(\mathbf{X}_N = \mathbf{a}_N) \leq 2^{-NH + \sqrt{N}\delta}.$$

Die Einteilung der Grundmenge in typische und untypische Quellwörter hat eine Reduktion des Aufwands bei Kodierungen fester Länge von Wörtern

über einem Quellalphabet  $\mathcal{X} = \{x_1, \dots, x_m\}$  durch Wörter über einem Kodealphabet  $\mathcal{Y} = \{y_1, \dots, y_d\}$ ,  $d \geq 2$  zum Ziel. Werden lediglich für die typischen Quellwörter eindeutige Kodewörter zur Verfügung gestellt, so treten wegen Satz 4.1 Dekodierfehler für große  $N$  nur mit sehr kleiner Wahrscheinlichkeit auf. Die Frage ist, wie “kurz” die Kodewortlänge dann sein darf.

Im folgenden werden Quellwörter der Länge  $N$  aus  $\mathcal{X}^N$  durch Kodewörter der Länge  $L$  aus  $\mathcal{Y}^L$  kodiert. Wir nehmen an, daß  $|\mathcal{Y}^L| \leq |\mathcal{X}^N|$ , d.h.  $L \leq N \frac{\log m}{\log d}$ . Ansonsten gibt es ja keine Probleme, da für jedes Quellwort ein Kodewort bereitsteht.

Zur Beschreibung der Menge von Quellwörtern, die ein Kodewort besitzen, wird im folgenden eine injektive Abbildung  $g : \mathcal{Y}^L \rightarrow \mathcal{X}^N$  verwendet. Die zugehörige Inverse  $g^{-1}$  auf dem Bildbereich  $g(\mathcal{Y}^L) \subseteq \mathcal{X}^N$  repräsentiert dann eine Abbildung, die jedem Quellwort aus  $g(\mathcal{Y}^L)$  ein Kodewort aus  $\mathcal{Y}^L$  zuordnet. Wir vereinbaren daher für den folgenden Satz die Sprechweisen:  $g$  heißt Kodierung,  $g(\mathcal{Y}^L)$  heißt Menge der Quellwörter mit Kodewort.  $F_g = \mathcal{X}^N \setminus g(\mathcal{Y}^L)$  bezeichnet die Menge der Quellwörter ohne Kodewort, die sogenannte Fehlmenge unter der Kodierung  $g$ . Graphisch läßt sich die Situation wie folgt veranschaulichen.



#### Satz 4.2

$\{X_n\}_{n \in \mathbb{N}}$  sei eine diskrete gedächtnislose Quelle mit Entropie  $H(X)$ . Dann gilt für alle  $\delta > 0$ :

- a) Falls  $L_N \geq N \frac{H(X) + \delta}{\log d}$ , existiert eine Folge von Kodierungen  $g_N : \mathcal{Y}^{L_N} \rightarrow \mathcal{X}^N$ ,  $N \in \mathbb{N}$ , mit

$$\lim_{N \rightarrow \infty} P(\mathbf{X}_N \in F_{g_N}) = 0.$$

b) Falls  $L_N \leq N \frac{H(X) - \delta}{\log d}$ , gilt für alle Folgen von Kodierungen  $g_N : \mathcal{Y}^{L_N} \rightarrow \mathcal{X}^N$ , daß

$$\lim_{N \rightarrow \infty} P(\mathbf{X}_N \in F_{g_N}) = 1.$$

**Beweis.** a) Für alle  $\delta > 0$ ,  $N > N_0$  gilt mit Satz 4.1 c), daß  $|T| \leq 2^{N(H+\delta)} \leq d^L$ , da  $L \geq N \frac{H(X) + \delta}{\log d}$ . Hieraus folgt die Existenz einer injektiven Abbildung  $g_N$  mit  $g_N(\mathcal{Y}^L) \supseteq T$ . Da  $F_{g_N} \subseteq T^C$ , gilt  $P(\mathbf{X}_N \in F_{g_N}) \leq P(\mathbf{X}_N \in T^C) \leq \varepsilon$  für alle  $\varepsilon > 0$  und alle genügend großen  $N$ , womit a) gezeigt ist.

b) Für alle  $\delta > 0$ , genügend großen  $N$  und  $(a_1, \dots, a_N) \in T$  gilt mit Satz 4.1 a)  $P(\mathbf{X}_N = \mathbf{a}_N) \leq 2^{-N(H-\delta)}$ . Nach Voraussetzung ist  $d^L \leq 2^{N(H-2\delta)}$ , und damit  $P(\mathbf{X}_N \in T \cap g_N(\mathcal{Y}^L)) \leq 2^{N(H-2\delta)} 2^{-N(H-\delta)} = 2^{-N\delta}$ . Ferner ist  $P(\mathbf{X}_N \in T^C \cap g_N(\mathcal{Y}^L)) \leq \varepsilon$  für alle  $\varepsilon > 0$  und alle genügend großen  $N$ . Insgesamt gilt  $P(\mathbf{X}_N \in g_N(\mathcal{Y}^L)) \leq \varepsilon + 2^{-N\delta}$ , woraus  $P(\mathbf{X}_N \in F_{g_N}) \rightarrow 1$  ( $N \rightarrow \infty$ ) für jede Kodierung  $g_N$  folgt. Dies zeigt b). ■

Interpretiert man  $\frac{L}{N}$  als Anzahl der Kodebuchstaben pro Quellbuchstabe, so macht Satz 4.2 die folgenden Aussagen. Ist  $\frac{L}{N}$  etwas größer als  $\frac{H(X)}{\log d}$ , so existiert ein Kode mit einer kleinen Fehlerwahrscheinlichkeit. Ist aber  $\frac{L}{N}$  nur etwas kleiner als  $\frac{H(X)}{\log d}$ , so hat jeder Kode bei großen Quellwortlängen eine Fehlerwahrscheinlichkeit nahe bei 1. Die normierte Entropie  $\frac{H(X)}{\log d}$  erweist sich als der kritische Wert für die Anzahl der Kodebuchstaben pro Quellbuchstabe, oberhalb dessen brauchbare Kodierungen existieren, unterhalb dessen aber keine funktionierende Kodierung konstruiert werden kann. Dies ist eine weitere Bestätigung dafür, daß die Entropie das richtige Maß für Unbestimmtheit ist.  $\log d$  ist lediglich eine Normierungskonstante, mit der die Mächtigkeit des Kodealphabets berücksichtigt wird. Sie verschwindet, wenn  $d$  als Basis des Logarithmus gewählt wird.

## 4.2 Kodes variabler Länge

Wie oben bezeichne  $\mathcal{X} = \{x_1, \dots, x_m\}$  das Quellalphabet und  $\mathcal{Y} = \{y_1, \dots, y_d\}$  das verwendete Kodealphabet. Im Unterschied zum vorhergehenden Abschnitt werden jetzt nicht Quellwörter der Länge  $N$  mit Kodewörtern fester Länge  $L$  kodiert, sondern jeder Buchstabe des Alphabets  $\mathcal{X}$

mit einem Kodewort variabler Länge über dem Alphabet  $\mathcal{Y}$ . Um die Anzahl zu übertragender Symbole klein zu halten, sollten häufig auftretende Buchstaben des Quellalphabets mit kurzen Kodewörtern und selten auftretende mit längeren Kodewörtern kodiert werden. Im Morsealphabet zum Beispiel wird der häufig anzutreffende Buchstabe **e** mit ‘·’ und das selten auftretende **q** mit ‘· · — —’ kodiert. Hiermit wird das Ziel verfolgt, die Kodewortlänge im Mittel möglichst klein zu halten.

In einem zweiten Schritt werden dann Blöcke aus Quellbuchstaben nach den gleichen Prinzipien mit Kodewörtern variabler Länge kodiert. Dies führt zu den sogenannten Blockcodes.

**Definition 4.2** (*Kode, Kodewort*)

Seien  $\mathcal{X} = \{x_1, \dots, x_m\}$  ein Quellalphabet,  $\mathcal{Y} = \{y_1, \dots, y_d\}$  ein Kodealphabet und

$$g : \mathcal{X} \rightarrow \bigcup_{\ell=1}^{\infty} \mathcal{Y}^{\ell} : x_j \mapsto (w_{j1}, \dots, w_{jn_j}),$$

wobei  $w_{jk} \in \mathcal{Y}$ ,  $j = 1, \dots, m$ ,  $k = 1, \dots, n_j$ ,  $n_j \in \mathbb{N}$ , eine injektive Abbildung.  $g(x_j) = (w_{j1}, \dots, w_{jn_j})$  heißt Kodewort des Buchstabens  $x_j \in \mathcal{X}$ .  $n_j$  heißt Länge des Kodewortes.  $g$  heißt Kode oder Kodierung mit Kodewortmenge  $\mathcal{K} = g(\mathcal{X}) = \{g(x_1), \dots, g(x_m)\}$ .

Durch den Kode  $g$  steht für jeden Quellbuchstaben ein Kodewort aus  $\bigcup_{\ell=1}^{\infty} \mathcal{Y}^{\ell}$  zur Verfügung. Zur Kodierung von ganzen Wörtern über dem Quellalphabet werden die Kodewörter der entsprechenden Buchstaben einfach aneinandergehängt. Hierbei ist wichtig, daß aus einer solchen Verkettung die einzelnen Kodewörter wieder identifiziert werden können. Mit Codes, die diese Eigenschaft besitzen, wird vermieden, daß bei der Übertragung von mehreren Kodewörtern hintereinander diese durch ein besonderes Symbol voneinander getrennt werden müssen.

**Definition 4.3** (*eindeutig dekodierbar*)

Ein Kode  $g$  heißt eindeutig dekodierbar (kurz: e.d., englisch: uniquely decodable, uniquely decipherable), wenn die Abbildung

$$G : \bigcup_{\ell=1}^{\infty} \mathcal{X}^{\ell} \rightarrow \bigcup_{\ell=1}^{\infty} \mathcal{Y}^{\ell} : (a_1, \dots, a_N) \mapsto (g(a_1), \dots, g(a_N))$$

injektiv ist.  $(g(a_1), \dots, g(a_N))$  ist dabei als Verkettung oder Konkatenation der Kodewörter  $g(a_1), \dots, g(a_N)$  zu lesen.

Im allgemeinen ist es schwierig festzustellen, ob ein Kode eindeutig dekodierbar ist. Nützlich für die Konstruktion solcher Kodes sind daher hinreichende Bedingungen für eindeutige Dekodierbarkeit. Die folgende Klasse der präfixfreien Kodes ist insbesondere unter Implementationsaspekten von besonderer Bedeutung.

**Definition 4.4** (Präfix, präfixfrei)

$\mathbf{b} = (b_1, \dots, b_r)$  und  $\mathbf{c} = (c_1, \dots, c_s) \in \bigcup_{\ell=1}^{\infty} \mathcal{Y}^\ell$ ,  $r \leq s$ , seien Kodewörter.  $\mathbf{b}$  heißt Präfix von  $\mathbf{c}$ , wenn ein  $\mathbf{d} = (d_1, \dots, d_{s-r}) \in \bigcup_{\ell=0}^{\infty} \mathcal{Y}^\ell$  existiert mit  $\mathbf{c} = (b_1, \dots, b_r, d_1, \dots, d_{s-r})$ . Ein Kode heißt präfixfrei (kurz: PF-Kode, englisch: prefix code, instantaneous code), wenn kein Kodewort aus  $\{g(x_1), \dots, g(x_m)\}$  Präfix eines anderen ist.

**Lemma 4.1** PF-Kodes sind eindeutig dekodierbar.

**Beweis.** Dem Beweis liegt die folgende Idee zugrunde. Suche die Folge der Kodebuchstaben von der ersten Stelle (d.h. von links) beginnend ab, bis das erste Kodewort identifiziert ist. Dieses ist aufgrund der Präfixfreiheit eindeutig. Fahre mit dem Rest der Folge genauso fort.

Zu zeigen ist, daß für alle  $N, M \in \mathbb{N}$ ,  $(a_1, \dots, a_N), (b_1, \dots, b_M) \in \bigcup_{\ell=1}^{\infty} \mathcal{X}^\ell$

$$\begin{aligned} (g(a_1), \dots, g(a_N)) &= (g(b_1), \dots, g(b_M)) \\ &\Rightarrow N = M \text{ und } (a_1, \dots, a_N) = (b_1, \dots, b_M), \end{aligned}$$

also die Injektivität von  $G$ . Angenommen, es gilt  $|g(a_1)| < |g(b_1)|$ , wobei  $|\cdot|$  die Länge eines Wortes bezeichnet. Dann ist  $g(a_1)$  Präfix von  $g(b_1)$ . Dies ist ein Widerspruch zur Präfixfreiheit von  $g$ . Analog ist  $|g(a_1)| \not> |g(b_1)|$ . Also gilt  $|g(a_1)| = |g(b_1)|$ , folglich  $g(a_1) = g(b_1)$  und wegen der Injektivität von  $g$  folgt  $a_1 = b_1$ . Mit Induktion ergibt sich  $g(a_\ell) = g(b_\ell)$  für alle  $\ell = 1, \dots, N$ , so daß  $a_\ell = b_\ell$  und  $N = M$ . ■

**Beispiel 4.3** Für das Alphabet  $\mathcal{X} = \{x_1, \dots, x_4\}$  betrachten wir die in folgender Tabelle definierten Abbildungen  $g_1, \dots, g_4 : \mathcal{X} \rightarrow \{0, 1\}^N$ .

	$g_1$	$g_2$	$g_3$	$g_4$
$x_1$	0	0	0	0
$x_2$	0	1	10	01
$x_3$	1	00	110	011
$x_4$	10	11	111	0111

Man sieht sofort, daß  $g_1$  nicht injektiv, also kein Kode ist.  $g_2$  ist ein binärer Kode, ist aber nicht eindeutig dekodierbar, denn sowohl  $(x_2, x_2) \rightarrow (1, 1)$  als auch  $(x_4) \rightarrow (1, 1)$ . Durch paarweise Vergleiche zeigt man, daß  $g_3$  ein eindeutig dekodierbarer PF-Kode ist.  $g_4$  ist nicht präfixfrei aber eindeutig dekodierbar. 0 dient bei diesem Kode als Trennzeichen. ■

Die Frage, wann zu einem Quellalphabet  $\mathcal{X} = \{x_1, \dots, x_m\}$ , Kodealphabet  $\mathcal{Y} = \{y_1, \dots, y_d\}$  und zu Kodewortlängen  $n_1, \dots, n_m \in \mathbb{N}$  ein eindeutig dekodierbarer Kode mit den Wortlängen  $n_1, \dots, n_m$  existiert, wird durch den folgenden Satz beantwortet.

**Satz 4.3** (a) McMillan (1959), b) Kraft (1949))

- a) Für alle eindeutig dekodierbaren Kodes mit Kodewortlängen  $n_1, \dots, n_m$  gilt  $\sum_{j=1}^m d^{-n_j} \leq 1$ .
- b) Gilt umgekehrt für Zahlen  $n_1, \dots, n_m \in \mathbb{N}$ , daß  $\sum_{j=1}^m d^{-n_j} \leq 1$ , so existiert ein PF-Kode, also ein e.d. Kode mit Kodewortlängen  $n_1, \dots, n_m$ .

**Beweis.** a) Sei  $g$  e.d. mit Kodewortlängen  $n_1, \dots, n_m \in \mathbb{N}$ . Es bezeichne  $r = \max_{1 \leq i \leq m} n_i$  und  $\beta_\ell = |\{j \mid n_j = \ell\}|$  die Anzahl der Kodewörter mit Länge  $\ell$ . Dann gilt für alle  $k \in \mathbb{N}$

$$\left( \sum_{j=1}^m d^{-n_j} \right)^k = \left( \sum_{\ell=1}^r \beta_\ell d^{-\ell} \right)^k = \sum_{\ell=k}^{k \cdot r} \gamma_\ell \cdot d^{-\ell},$$

wobei

$$\gamma_\ell = \sum_{\substack{1 \leq i_1, \dots, i_k \leq r \\ i_1 + \dots + i_k = \ell}} \beta_{i_1} \cdots \beta_{i_k}, \quad \ell = k, \dots, k \cdot r.$$

$\gamma_\ell$  ist die Anzahl der Quellwörter, die bei Verwendung von  $g$  Kodewortlänge  $\ell$  haben. Andererseits ist  $d^\ell$  die Anzahl aller Kodewörter der Länge  $\ell$ . Da

$g$  eindeutig dekodierbar also injektiv ist, besitzt jedes Kodewort höchstens ein Quellwort. Folglich ist  $\gamma_\ell \leq d^\ell$  für alle  $\ell = k, \dots, k \cdot r$ . Für alle  $k \in \mathbb{N}$  gilt also

$$\left( \sum_{j=1}^m d^{-n_j} \right)^k \leq \sum_{\ell=k}^{kr} d^\ell d^{-\ell} = kr - k + 1 \leq kr$$

und durch Grenzübergang  $k \rightarrow \infty$

$$\sum_{j=1}^m d^{-n_j} \leq (kr)^{1/k} \rightarrow 1 \quad (k \rightarrow \infty).$$

Hieraus folgt Behauptung a).

b) Ohne Einschränkung der Allgemeinheit können wir annehmen, daß  $n_1 \leq n_2 \leq \dots \leq n_m$ . Ansonsten können die Quellbuchstaben entsprechend umnummeriert werden. Gelte  $\sum_{j=1}^m d^{-n_j} \leq 1$ . Wir führen den Beweis konstruktiv und geben einen PF-Kode mit Wortlängen  $n_1, \dots, n_m$  an. Bezeichne hierzu

$$\mathcal{N}(\mathbf{a}) = \{\mathbf{b} \in \mathcal{Y}^{n_m} \mid \mathbf{a} \text{ ist Präfix von } \mathbf{b}\}.$$

$\mathcal{N}(\mathbf{a})$  ist die Menge aller Fortsetzungen von  $\mathbf{a}$  zu einem Wort der Länge  $n_m$ . Wähle nun  $\mathbf{z}_1 \in \mathcal{Y}^{n_1}$  beliebig. Offensichtlich gilt  $|\mathcal{N}(\mathbf{z}_1)| = d^{n_m - n_1}$ , und falls  $m \geq 2$ ,

$$d^{n_m - n_1} < d^{n_m}.$$

Also ist  $\mathcal{Y}^{n_m} \setminus \mathcal{N}(\mathbf{z}_1) \neq \emptyset$ .

Hieraus folgt, daß ein  $\mathbf{z}_2 \in \mathcal{Y}^{n_2}$  derart existiert, daß  $\mathbf{z}_1$  kein Präfix von  $\mathbf{z}_2$  ist. Es gilt  $|\mathcal{N}(\mathbf{z}_2)| = d^{n_m - n_2}$ . Falls also  $m \geq 3$ , hat man

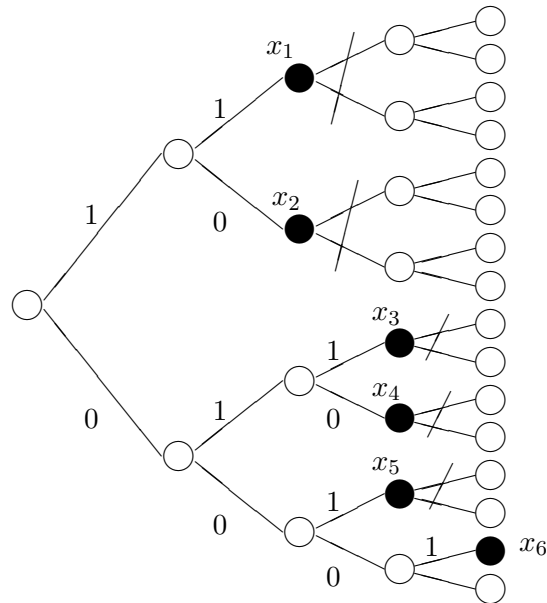
$$d^{n_m - n_1} + d^{n_m - n_2} < \sum_{j=1}^m d^{n_m - n_j} = d^{n_m} \sum_{j=1}^m d^{-n_j} \leq d^{n_m}$$

wegen der Voraussetzung  $\sum_{j=1}^m d^{-n_j} \leq 1$ . Also ist  $\mathcal{Y}^{n_m} \setminus (\mathcal{N}(\mathbf{z}_1) \cup \mathcal{N}(\mathbf{z}_2)) \neq \emptyset$ .

Damit existiert  $\mathbf{z}_3 \in \mathcal{Y}^{n_3}$  so, daß weder  $\mathbf{z}_1$  noch  $\mathbf{z}_2$  Präfix von  $\mathbf{z}_3$  sind. Die  $m$ -fache Anwendung dieses Verfahrens liefert eine präfixfreie Menge  $\mathcal{K} = \{\mathbf{z}_1, \dots, \mathbf{z}_m\}$ . Der gesuchte Kode  $g$  wird nun durch  $g(x_j) = \mathbf{z}_j$ ,  $j = 1, \dots, m$  festgelegt. Er ist präfixfrei, also eindeutig dekodierbar, womit b) bewiesen ist. ■

**Beispiel 4.4** (Konstruktion eines PF-Kodes mit gegebenen Wortlängen)

Für die Zahlenwerte  $m = 6$ ,  $d = 2$  und  $n_1 = n_2 = 2$ ,  $n_3 = n_4 = n_5 = 3$ ,  $n_6 = 4$  gilt  $\sum_{j=1}^6 2^{-n_j} = 2 \cdot \frac{1}{4} + 3 \cdot \frac{1}{8} + \frac{1}{16} = \frac{15}{16} < 1$ . Also existiert ein binärer PF-Kode, der wie im Beweis von Satz 4.3 b) konstruiert wird. Mit Hilfe eines binären Codebaums, dessen Tiefe der maximalen Kodewortlänge  $n_6 = 4$  entspricht, kann das Prinzip übersichtlich dargestellt werden.



Zunächst wird ein beliebiger Knoten der Tiefe 2 ausgewählt, der ein Kodewort für  $x_1$  repräsentiert. ‘Nach oben’ wird hierbei als 1 und ‘nach unten’ als 0 gelesen. In dem oben dargestellten Baum ergibt sich (1,1) als Kodewort für  $x_1$ . Alle Nachfolgerknoten würden den mit  $x_1$  markierten Knoten als Präfix haben, sie werden als mögliche Kodewörter ausgeschieden. Dies ist durch den schrägen Strich durch den entsprechenden Teilbaum angedeutet. Mit  $x_2, \dots, x_6$  verfährt man entsprechend. Die Existenz von verbleibenden Knoten wird durch den Beweis von Satz 4.3 b) sichergestellt. Nachdem alle 6 Buchstaben im Baum untergebracht sind, erhält man den in der folgenden Tabelle dargestellten Code.

$x_i$	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$
$g(x_i)$	11	10	011	010	001	0001



Dieser ist nicht eindeutig, er ist auch nicht optimal.  $x_6$  hätte durch den untersten Knoten der Tiefe 3, also  $(0, 0, 0)$ , kodiert werden können, ohne die Präfixfreiheit zu stören. In der Tat ist  $2 \cdot 2^{-2} + 4 \cdot 2^{-3} = 1$ , die hinreichende Bedingung für die Existenz eines PF-Kodes aus Satz 4.3 b) ist also auch für die Kodewortlängen  $n_1 = n_2 = 2$  und  $n_3 = n_4 = n_5 = n_6 = 3$  erfüllt. ■

Die Güte einer Kodierung wird im folgenden durch die erwartete Kodewortlänge gemessen. Sei  $X$  eine diskrete gedächtnislose Quelle und  $g$  ein Kode. Die Zufallsvariable

$$N_g : \mathcal{X} \rightarrow \mathbb{N} : x_j \mapsto n_j$$

bildet den Quellbuchstaben  $x_j$  auf die Länge des zugehörigen Kodeworts unter der Kodierung  $g$  ab. Die erwartete Kodewortlänge, bezeichnet mit  $\bar{n}(g)$  beträgt dann

$$\bar{n}(g) = \mathbb{E}(N_g) = \sum_{j=1}^m n_j P(X = x_j).$$

Ziel ist es, einen eindeutig dekodierbaren Kode  $g^*$  mit kürzester erwarteter Kodewortlänge zu finden, d.h. mit der Eigenschaft  $\bar{n}(g^*) \leq \bar{n}(g)$  für alle e.d. Kodes  $g$ .

Der folgende wichtige Satz gibt eine obere und untere Schranke für die erwartete Kodewortlänge eindeutig dekodierbarer Kodes an.

**Satz 4.4** (*Noiseless Coding Theorem, Shannon (1948)*)

Sei  $X$  eine diskrete gedächtnislose Quelle mit Alphabet  $\mathcal{X} = \{x_1, \dots, x_m\}$  und Entropie  $H(X) > 0$ .  $\mathcal{Y} = \{y_1, \dots, y_d\}$ ,  $d \geq 2$ , sei ein Kodealphabet.

a) Für alle e.d. Kodes  $g$  gilt  $\frac{H(X)}{\log d} \leq \bar{n}(g)$ .

b) Es existiert ein PF-Kode  $g$  mit  $\bar{n}(g) < \frac{H(X)}{\log d} + 1$ .

**Beweis.** a) Zur Abkürzung wird  $p_j = P(X = x_j)$  gesetzt,  $j = 1, \dots, m$ . Für alle e.d. Kodes  $g$  folgt dann mit Satz 4.3 a) und der Ungleichung  $\ln z \leq z - 1$ ,

falls  $z > 0$ , daß

$$\begin{aligned} H(X) - \bar{n}(g) \log d &= \sum_{j=1}^m p_j \log \frac{1}{p_j} - \sum_{j=1}^m p_j n_j \log d \\ &= \sum_{j=1}^m p_j \log \left( \frac{d^{-n_j}}{p_j} \right) \leq (\log e) \sum_{j:p_j>0} p_j \left( \frac{d^{-n_j}}{p_j} - 1 \right) \\ &\leq (\log e) \left( \sum_{j=1}^m d^{-n_j} - 1 \right) \leq 0. \end{aligned}$$

Durch Auflösen nach  $\bar{n}(g)$  folgt die Ungleichung a).

b) Wir weisen die Existenz eines PF-Kodes mit den gewünschten Eigenschaften nach. Wähle  $n_j$  so, daß  $d^{-n_j} \leq p_j < d^{-n_j+1}$ ,  $j = 1, \dots, m$ . Hierbei kann ohne Einschränkung der Allgemeinheit angenommen werden, daß  $p_j > 0$  für alle  $j$ . Es folgt  $\sum_{j=1}^m d^{-n_j} \leq \sum_{j=1}^m p_j = 1$ . Wegen Satz 4.3 b) existiert ein PF-Kode  $g$  mit Wortlängen  $n_1, \dots, n_m$ . Nach Konstruktion gilt für diesen Kode  $\log p_j < (-n_j + 1) \log d$ , so daß

$$\sum_{j=1}^m p_j \log p_j < (\log d) \sum_{j=1}^m p_j (-n_j + 1).$$

Dies bedeutet  $H(X) > (\log d)(\bar{n}(g) - 1)$ , und b) folgt nach Auflösen dieser Ungleichung. ■

Man beachte, daß bei Satz 4.4 die Gedächtnislosigkeit, also die stochastische Unabhängigkeit der Zufallsvariablen, keine Rolle spielt. Die Abhängigkeitsstruktur wird erst bei der Blockkodierung von Quellen relevant.

Ferner geht die Wahl der Basis des Logarithmus bei der Entropie in die Ungleichungen als multiplikative Konstante ein. Wird als Basis die Mächtigkeit des Kodealphabets  $d$  gewählt, gilt  $\log d = 1$ , so daß der normierende Faktor  $\log d$  in a) und b) verschwindet.

Satz 4.4 beantwortet auch die Frage aus Beispiel 4.1, wenn man ja/nein-Fragestrategien mit binären Kodierungen  $g$  identifiziert. Bei Verwendung von Logarithmen zur Basis 2 gilt für die erwartete Anzahl von Fragen  $E(Z)$ , daß  $H(X) \leq E(Z) = \bar{n}(g)$  für alle zum Ziel führenden, d.h. eindeutig dekodierbaren Fragestrategien.

Gegeben sei nun eine Quelle  $X$  mit Wahrscheinlichkeiten  $p_1, \dots, p_m$ , wobei  $p_j = P(X = x_j)$ . Wählt man  $n_j = \lceil -\log p_j / \log d \rceil$ ,  $j = 1, \dots, m$ , so gilt  $\sum_{j=1}^m d^{-n_j} \leq 1$ . Nach Satz 4.4 b) existiert ein PF-Kode  $g$  mit Wortlängen  $n_1, \dots, n_m$  und  $\bar{n}(g) < H(X)/\log d + 1$ , der höchstens um 1 (Bit) schlechter als der bestmögliche Wert  $H(X)/\log d$  ist. Die Konstruktion eines präfixfreien Codes mit  $\sum_{j=1}^m d^{-n_j} \leq 1$  wurde im Beweis von Satz 4.3 b) vorgeführt. Explizit bestimmt man den zugehörigen Code mit Hilfe eines Codebaums wie in Beispiel 4.4. Der nach diesem Verfahren zu gegebenen Wahrscheinlichkeiten  $p_1, \dots, p_m$  konstruierte Code  $\hat{g}$  heißt *Shannon-Fano-Kode*. Er hat die Eigenschaft  $\bar{n}(\hat{g}) \leq H(X)/\log d + 1$ .

Die bisher entwickelten Methoden lassen sich genauso zur Kodierung von Wörtern über dem Alphabet  $\mathcal{X}$  einsetzen. Blöcke  $(a_1, \dots, a_N) \in \mathcal{X}^N$  der Länge  $N$  werden dabei als Buchstaben des zusammengesetzten Alphabets  $\mathcal{X}^N$  aufgefaßt und mit Kodewörtern variabler Länge kodiert. Hierdurch entstehen sogenannte Blockcodes.

Zur formalen Beschreibung sei  $\{X_n\}_{n \in \mathbb{N}}$  eine diskrete gedächtnislose Quelle mit Entropie  $H(X)$  und zugehörigem Alphabet  $\mathcal{X} = \{x_1, \dots, x_m\}$ . Der Zufallsvektor  $\mathbf{X}_N = (X_1, \dots, X_N)$  ist dann eine endlich diskrete Zufallsvariable mit Träger  $\mathcal{X}^N$  und stochastisch unabhängigen Komponenten. Für die Entropie gilt mit Satz 3.1 d)

$$H(\mathbf{X}_N) = H(X_1) + \dots + H(X_N) = N \cdot H(X).$$

Ein Kode  $g^{(N)} : \mathcal{X}^N \rightarrow \bigcup_{\ell=1}^{\infty} \mathcal{Y}^{\ell}$  heißt in diesem speziellen Zusammenhang Blockcode. Mit  $n(a_1, \dots, a_N) = |g^{(N)}((a_1, \dots, a_N))|$  wird die Länge des Kodewortes von  $(a_1, \dots, a_N) \in \mathcal{X}^N$  bezeichnet und mit

$$\bar{n}(g^{(N)}) = \sum_{(a_1, \dots, a_N) \in \mathcal{X}^N} n(a_1, \dots, a_N) P(X_1 = a_1, \dots, X_N = a_N)$$

die erwartete Kodewortlänge des Blockcodes  $g^{(N)}$ .  $\bar{n}(g^{(N)})/N$  ist dann die erwartete Kodewortlänge pro Quellbuchstabe. Die Definitionen von eindeutiger Dekodierbarkeit und Präfixfreiheit werden entsprechend auf das zusammengesetzte Alphabet  $\mathcal{X}^N$  übertragen.

**Satz 4.5** (Blockkodierung)

$\{X_n\}_{n \in \mathbb{N}}$  sei eine diskrete gedächtnislose Quelle mit dem Alphabet  $\mathcal{X} = \{x_1, \dots, x_m\}$  und Entropie  $H(X)$ .  $\mathcal{Y} = \{y_1, \dots, y_d\}$ ,  $d \geq 2$ , sei ein Kodealphabet und  $N \in \mathbb{N}$ .

- a) Für alle e.d. Blockcodes  $g^{(N)}$  gilt  $\frac{H(X)}{\log d} \leq \frac{\bar{n}(g^{(N)})}{N}$ .
- b) Es existiert ein präfixfreier Blockcode  $g^{(N)}$  mit  $\frac{\bar{n}(g^{(N)})}{N} < \frac{H(X)}{\log d} + \frac{1}{N}$ .

**Beweis.** Zu beachten ist lediglich, daß Blockcodes Codes über dem Alphabet  $\mathcal{X}^N$  sind. Mit Satz 4.4 gilt für alle e.d. Codes  $g^{(N)}$ , daß  $\bar{n}(g^{(N)}) \geq \frac{H(X_1, \dots, X_N)}{\log d} = \frac{N \cdot H(X)}{\log d}$ . Nach Division durch  $N$  folgt Behauptung a).

Ferner existiert ein präfixfreier Code  $g^{(N)}$  mit  $\bar{n}(g^{(N)}) < \frac{H(X_1, \dots, X_N)}{\log d} + 1 = \frac{N \cdot H(X)}{\log d} + 1$ . Division durch  $N$  liefert Behauptung b). ■

Für jedes  $N \in \mathbb{N}$  existiert also ein eindeutig dekodierbarer Blockcode mit

$$\frac{H(X)}{\log d} \leq \frac{\bar{n}(g^{(N)})}{N} < \frac{H(X)}{\log d} + \frac{1}{N}.$$

Durch Limesbildung  $N \rightarrow \infty$  folgt sofort das folgende

**Korollar 4.1** Es existieren Folgen  $\{g^{(N)}\}_{N \in \mathbb{N}}$  von eindeutig dekodierbaren Blockcodes mit

$$\lim_{N \rightarrow \infty} \frac{\bar{n}(g^{(N)})}{N} = \frac{H(X)}{\log d}.$$

Für genügend große  $N$  existiert also ein eindeutig dekodierbarer Blockcode  $g^{(N)}$ , dessen Effizienz (im Sinn der erwarteten Kodewortlänge pro Quellbuchstabe) nahezu optimal ist. Durch Satz 4.5 und Korollar 4.1 werden die Behauptungen in Beispiel 4.2 für binäre Blockfragestrategien gezeigt.

Wieder erweist sich die (mit  $\log d$  normierte) Entropie als das richtige Maß für Unbestimmtheit. Sie ist eine untere Schranke für die erwartete Kodewortlänge jedes eindeutig dekodierbaren Codes, die zum Beispiel mit Shannon-Fano-Blockcodes asymptotisch erreicht wird.

Das Noiseless Coding Theorem, Satz 4.4, besagt  $H(X)/\log d \leq \bar{n}(g)$  für alle eindeutig dekodierbaren Codes  $g$ . Kein eindeutig dekodierbarer Code kann also eine kürzere erwartete Kodewortlänge als  $H(X)/\log d$  besitzen. Codes  $g^*$ , die diese Schranke erreichen, für die also  $\bar{n}(g^*) = H(X)/\log d$  gilt, heißen *absolut optimal*.

Eine notwendige und hinreichende Bedingung für Gleichheit ist  $p_j = d^{-n_j}$  für alle  $j \in \{1, \dots, m\}$  mit  $p_j > 0$ . Dies zeigt der Beweis von Satz 4.4 a).

Es folgt, daß kein absolut optimaler Kode existiert, falls ein  $p_j$  irrational ist. Dennoch existieren Codes kürzester erwarteter Kodewortlänge, deren Konstruktion unser nächstes Ziel ist.

**Definition 4.5** (*optimaler Kode*)

Ein eindeutig dekodierbarer Kode  $g^*$  mit Wortlängen  $n_1^*, \dots, n_m^*$  heißt optimal (englisch: compact), wenn

$$\bar{n}(g^*) = \sum_{i=1}^m p_i n_i^* \leq \sum_{i=1}^m p_i n_i = \bar{n}(g)$$

für alle e.d. Codes  $g$  mit Wortlängen  $n_1, \dots, n_m$ .

Aus Satz 4.3 folgt

$$\begin{aligned} & \{(n_1, \dots, n_m) \mid n_1, \dots, n_m \text{ sind Wortlängen eines e.d. Codes}\} \\ &= \{(n_1, \dots, n_m) \mid n_i \in \mathbb{N}, \sum_{j=1}^m d^{-n_j} \leq 1\}. \end{aligned}$$

Zur Bestimmung eines optimalen Codes für eine Quelle mit Wahrscheinlichkeiten  $p_1, \dots, p_m$  ist also das folgende Optimierungsproblem zu lösen.

$$\text{minimiere } \sum_{j=1}^m p_j n_j \text{ über alle } n_1, \dots, n_m \in \mathbb{N} \text{ mit } \sum_{j=1}^m d^{-n_j} \leq 1 \quad (4.3)$$

In einem ersten Lösungsansatz wird das relaxierte Problem mit stetigen Variablen  $z_1, \dots, z_m \in \mathbb{R}$  betrachten.

$$\text{minimiere } \sum_{j=1}^m p_j z_j \text{ über alle } z_1, \dots, z_m \in \mathbb{R} \text{ mit } \sum_{j=1}^m d^{-z_j} \leq 1 \quad (4.4)$$

Offensichtlich wird das Minimum bei  $\sum_{j=1}^m d^{-z_j} = 1$  angenommen. Ein Lagrange-Ansatz mit der Lagrangefunktion

$$L(z_1, \dots, z_m, \lambda) = \sum_{j=1}^m p_j z_j + \lambda \left( \sum_{j=1}^m d^{-z_j} - 1 \right)$$

und den zu Null gesetzten partiellen Ableitungen

$$\begin{aligned}\frac{\partial L}{\partial z_i} &= p_i - \lambda(\ln d)d^{-z_i} = 0, \quad i = 1, \dots, m, \\ \frac{\partial L}{\partial \lambda} &= \sum_{j=1}^m d^{-z_j} - 1 = 0\end{aligned}$$

hat als Lösung

$$\lambda = 1/\ln d \quad \text{und} \quad d^{-z_i} = p_i, \quad i = 1, \dots, m.$$

Falls  $p_i > 0$  für alle  $i = 1, \dots, m$ , folgt als notwendige Bedingung für eine Minimalstelle von (4.4)

$$z_i^* = -\ln p_i / \ln d, \quad i = 1, \dots, m.$$

Sind  $z_1, \dots, z_m$  ganzzahlig, liegt in der Tat eine Lösung von (4.3) vor, da die untere Schranke  $H(X)/\log d$  in Satz 4.3 a) erreicht wird.

In der Regel wird jedoch keine Ganzzahligkeit vorliegen. Dennoch besitzt das diskrete Optimierungsproblems (4.3) stets eine optimale Lösung  $n_1^*, \dots, n_m^*$  mit einem zugehörigen Kode  $g^*$ . Der Nachweis der Existenz wird als Übungsaufgabe empfohlen. Die zugehörige erwartete Kodewortlänge eines optimalen Kodes

$$\bar{n}(g^*) = \min\{\bar{n}(g) \mid g \text{ ist ein e.d. Kode}\}$$

heißt reale Entropie der Quelle. Mit Korollar 4.1 folgt für optimale Blockcodes  $g^{(N)*}$  bei diskreten gedächtnislosen Quellen, daß

$$\lim_{N \rightarrow \infty} \frac{1}{N} \bar{n}(g^{(N)*}) = \frac{H(X)}{\log d},$$

d.h. die normierte reale Entropie von optimalen Blockcodes konvergiert mit wachsender Blocklänge gegen die Entropie  $H(X)$ .

Im folgenden wird ein Verfahren zur expliziten Bestimmung optimaler Kodes hergeleitet, d.h. ein Algorithmus zur Lösung von (4.3). Mit Satz 4.3 ist klar, daß zu jedem eindeutig dekodierbaren Kode ein präfixfreier Kode mit

denselben Kodewortlängen existiert. Es reicht also, einen optimalen Kode in der Menge der präfixfreien zu suchen.

Wir beschränken uns im weiteren auf den Fall binärer Kodes, d.h.  $d = 2$  und  $\mathcal{Y} = \{0, 1\}$ . Für eine diskrete Quelle  $X$  mit Alphabet  $\mathcal{X} = \{x_1, \dots, x_m\}$  bezeichne  $p_i = P(X = x_i)$ ,  $i = 1, \dots, m$ . Im folgenden Lemma wird die Existenz von optimalen, binären präfixfreien Kodes nachgewiesen, bei denen die Wahrscheinlichkeiten  $p_i$  und die Kodewortlängen gegenläufig geordnet sind.

**Lemma 4.2**  *$X$  sei eine diskrete Quelle mit Alphabet  $\mathcal{X} = \{x_1, \dots, x_m\}$  und Wahrscheinlichkeiten  $p_1 \geq \dots \geq p_m > 0$ . Dann existiert ein optimaler, binärer präfixfreier Kode  $h$  mit Kodewortlängen  $\ell_1, \dots, \ell_m$ , derart daß*

- (i)  $\ell_1 \leq \dots \leq \ell_m$ ,
- (ii)  $\ell_{m-1} = \ell_m$ ,
- (iii)  $h(x_{m-1})$  und  $h(x_m)$  unterscheiden sich nur in der letzten Stelle.

**Beweis.**  $h$  sei ein optimaler, binärer PF-Kode mit Wortlängen  $\ell_1, \dots, \ell_m$ . Falls  $p_i > p_j$  für  $1 \leq i < j \leq m$ , gilt  $\ell_i \leq \ell_j$ . Ansonsten tauscht man die Kodewörter  $h(x_i)$  und  $h(x_j)$  aus und erhält hierdurch einen Kode  $h'$  mit

$$\bar{n}(h') - \bar{n}(h) = p_i \ell_j + p_j \ell_i - p_i \ell_i - p_j \ell_j = (p_i - p_j)(\ell_j - \ell_i) < 0,$$

im Widerspruch zur Optimalität von  $h$ . Falls  $p_i = p_j$  und  $\ell_i > \ell_j$ , erhält man durch Austauschen der Kodewörter  $h(x_i)$  und  $h(x_j)$  einen Kode  $h'$  mit  $\bar{n}(h) = \bar{n}(h')$ . Durch endlich viele solcher Austauschschritte ergibt sich ein optimaler PF-Kode, der (i) erfüllt.

Nach (i) existiert ein optimaler PF-Kode  $h$  mit  $\ell_1 \leq \dots \leq \ell_m$ . Gilt  $\ell_{m-1} < \ell_m$ , erhält man durch Streichen der letzten  $(\ell_m - \ell_{m-1})$  Komponenten des Kodeworts  $h(x_m)$  einen PF-Kode  $h'$  mit  $\bar{n}(h') < \bar{n}(h)$ , im Widerspruch zur Optimalität von  $h$ . Dies zeigt (ii).

Angenommen für einen optimalen PF-Kode  $h$  mit  $\ell_1 \leq \dots \leq \ell_{m-1} = \ell_m$  unterscheiden sich je zwei Kodewörter der Länge  $\ell_m$  bereits in den ersten  $\ell_m - 1$  Stellen. Dann erhält man wegen der Präfixfreiheit durch Streichen der letzten Komponente bei den Kodewörtern der Länge  $\ell_m$  einen Kode  $h'$  mit  $\bar{n}(h') < \bar{n}(h)$ . Hieraus folgt Eigenschaft (iii). ■

Das folgende Lemma zeigt, wie ein optimaler Kode durch Anhängen einer 0 bzw. 1 an das längste Kodewort zu einem optimalen Kode für eine Quelle

mit einem zusätzlichen Buchstaben aufgefaltet werden kann, vorausgesetzt die Wahrscheinlichkeiten der beiden Quellen sind geeignet verknüpft.

**Lemma 4.3** Sei  $X$  eine diskrete Quelle mit Alphabet  $\mathcal{X} = \{x_1, \dots, x_m\}$  und Wahrscheinlichkeiten  $p_1 \geq \dots \geq p_m > 0$ .  $X'$  sei eine diskrete Quelle mit  $\mathcal{X}' = \{x'_1, \dots, x'_{m-1}\}$ ,  $p'_j = p_j$ ,  $j = 1, \dots, m-2$  und  $p'_{m-1} = p_{m-1} + p_m$ .  $g'$  mit Kodewortmenge  $\mathcal{K}' = \{g'(x'_1), \dots, g'(x'_{m-1})\}$  sei ein optimaler binärer PF-Kode für  $X'$ . Dann ist der Kode  $g$  mit  $g(x_j) = g'(x'_j)$ , falls  $j = 1, \dots, m-2$ ,  $g(x_{m-1}) = (g'(x_{m-1}), 0)$  und  $g(x_m) = (g'(x_{m-1}), 1)$  ein optimaler PF-Kode für  $X$ .

**Beweis.** Bezeichne  $n_j$  und  $n'_j$  die Kodewortlängen von  $g$  bzw.  $g'$ . Dann gilt

$$\begin{aligned} \bar{n}(g) &= \sum_{j=1}^{m-2} p_j n'_j + (p_{m-1} + p_m)(n'_{m-1} + 1) \\ &= \sum_{j=1}^{m-2} p'_j n'_j + p'_{m-1}(n'_{m-1} + 1) \\ &= \sum_{j=1}^{m-1} p'_j n'_j + p_{m-1} + p_m = \bar{n}(g') + p_{m-1} + p_m. \end{aligned} \quad (4.5)$$

Angenommen  $g$  ist nicht optimal für  $X$ . Dann existiert ein optimaler PF-Kode  $h$  für  $X$ ,  $h(x_j) = (w_{j1}, \dots, w_{j\ell_j})$ , der den Bedingungen (i), (ii) und (iii) aus Lemma 4.2 genügt mit  $\bar{n}(h) < \bar{n}(g)$ . Setze  $h'(x'_j) = h(x_j)$ ,  $j = 1, \dots, m-2$  und durch Streichen der letzten Stelle  $h'(x'_{m-1}) = (w_{m-1,1}, \dots, w_{m-1,\ell_{m-1}})$ .  $h'$  ist ein binärer PF-Kode für  $X'$ , und es gilt analog zu (4.5)

$$\bar{n}(h') + p_{m-1} + p_m = \bar{n}(h) < \bar{n}(g) = \bar{n}(g') + p_{m-1} + p_m.$$

Es folgt  $\bar{n}(h') < \bar{n}(g')$ , im Widerspruch zur Optimalität von  $g'$ . ■

Einen optimalen Kode konstruiert man nun durch rekursive Anwendung von Lemma 4.3 nach folgendem Verfahren. Bestimme eine Folge von Quellen, indem jede aus der vorherigen durch Addition der beiden kleinsten Wahrscheinlichkeiten  $p_j$  verkürzt wird. Führe dies rekursiv fort, bis das Alphabet



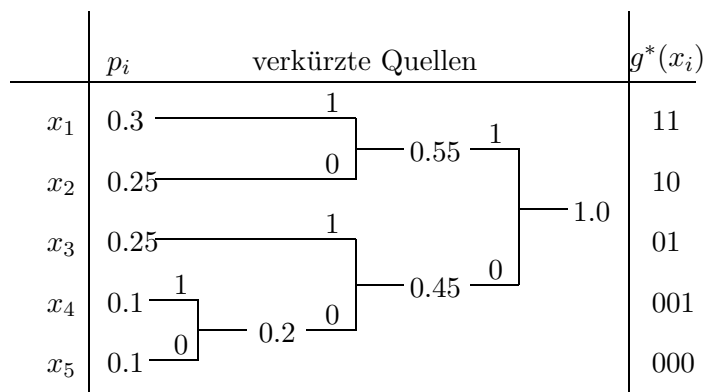


Abb. 4.1 Beispiel eines Huffman-Kodes mit binärem Baum.

der aktuellen Quelle aus genau zwei Buchstaben besteht. Kodiere diese optimal mit 0 und 1. Baue mit Lemma 4.3 daraus rekursiv einen optimalen Kode für die ursprüngliche Quelle auf. Dies ist das sogenannte *Huffman-Verfahren*, der zugehörige Kode heißt *Huffman-Kode*.

**Beispiel 4.5** (Huffman-Verfahren)

Für eine Quelle mit Alphabet  $\{x_1, \dots, x_5\}$  und Wahrscheinlichkeiten  $(0.3, 0.25, 0.25, 0.1, 0.1)$  wird ein Huffman-Kode berechnet. Es empfiehlt sich, die sukzessive verkürzten Quellen durch einen binären Baum darzustellen, dessen Knoten mit den jeweils addierten Wahrscheinlichkeiten markiert werden. Für obiges Beispiel ergibt sich die Graphik aus Abbildung 4.1, in der durch Rückwärtsgehen der Huffman-Kode abgelesen werden kann. Dieser ist in der rechten Spalte unter  $g^*(x_i)$  vermerkt. ■

Ein einfach zu implementierendes Verfahren, das allerdings nicht notwendig einen optimalen Kode liefert, ist die Fano-Kodierung. Die Idee hierzu stammt aus der folgenden Zerlegung der Entropie einer Quelle  $X$  mit Wahrscheinlichkeiten  $p_1, \dots, p_m$ .

Betrachtet wird ein PF-Kode  $g : \mathcal{X} = \{x_1, \dots, x_m\} \rightarrow \bigcup_{\ell=1}^{\infty} \{0, 1\}^{\ell}$  mit zugehörigen Kodewörtern  $g(x_1), \dots, g(x_m)$  und Wortlängen  $n_1, \dots, n_m \in \mathbb{N}$ . Es bezeichne  $r = \max_{1 \leq i \leq m} n_i$ . Um eine konstante Länge zu erreichen,

werden alle Kodewörter  $g(x_i)$  durch Auffüllen mit Nullen zu Kodewörtern der Länge  $r$  ergänzt, und zwar

$$\mathbf{y}_i = \left( \underbrace{g(x_i)}_{n_i}, \underbrace{0, \dots, 0}_{r-n_i} \right) \in \{0, 1\}^r, \quad i = 1, \dots, m.$$

$\tilde{g} : \mathcal{X} \rightarrow \{0, 1\}^r : x_i \mapsto \mathbf{y}_i$  ist dann ebenfalls ein präfixfreier Kode.

Der Zufallsvektor  $\mathbf{Y}$  wird nun durch  $\mathbf{Y} = (Y_1, \dots, Y_r) = \tilde{g}(X)$  definiert. Wegen der Injektivität von  $\tilde{g}$  gilt  $H(X) = H(\mathbf{Y})$  (vgl. Aufgabe 3.5). Mit Lemma 3.2 a) folgt dann die Zerlegung

$$\begin{aligned} H(X) &= H(Y_1, \dots, Y_r) \\ &= H(Y_1) + H(Y_2 | Y_1) + \dots + H(Y_r | Y_1, \dots, Y_{r-1}). \end{aligned} \quad (4.6)$$

Umgekehrt versucht man, bei der Konstruktion eines Kodes die Entropie der Quelle  $H(X)$  mit möglichst wenigen Komponenten  $Y_1, \dots, Y_r$  auszuschöpfen. Wir benötigen zunächst ein vorbereitendes Lemma.

**Lemma 4.4** Für  $(p_1, \dots, p_m) \in \mathcal{P}_m$  und  $T \subseteq \{1, \dots, m\}$  bezeichne  $\sigma_T = \sum_{i \in T} p_i$ . Dann wird

$$\max_{T \subseteq \{1, \dots, m\}} \left( -\sigma_T \log \sigma_T - (1 - \sigma_T) \log(1 - \sigma_T) \right)$$

angenommen für jedes  $T^* \subseteq \{1, \dots, m\}$  mit  $|\sum_{i \in T^*} p_i - \frac{1}{2}| \leq |\sum_{i \in T} p_i - \frac{1}{2}|$  für alle  $T \subseteq \{1, \dots, m\}$ .

**Beweis.**  $-x \log x - (1-x) \log(1-x)$ ,  $x \in [0, 1]$ , ist konkav, symmetrisch und maximal für  $x^* = \frac{1}{2}$ . Bei Maximierung dieser Funktion über eine endliche Menge wird das Maximum für das Argument angenommen, dessen Abstand zu  $\frac{1}{2}$  minimal ist. ■

$H(Y_1)$  aus (4.6) wird nun durch Wahl von  $g$  bzw.  $\tilde{g}$  maximiert. Wir setzen

$$q_0 = P(Y_1 = 0) = P\left(\left(\tilde{g}(X)\right)_1 = 0\right) = \sum_{i: (\tilde{g}(x_i))_1 = 0} P(X = x_i)$$

und  $P(Y_1 = 1) = 1 - q_0$ , wobei  $(\tilde{g}(X))_1$  die erste Komponente von  $\tilde{g}(X)$  bedeutet.

Dann ist  $H(Y_1) = -q_0 \log q_0 - (1 - q_0) \log(1 - q_0)$  maximal bezüglich  $q_0$ , falls  $|\sum_{i \in T_0} p_i - \frac{1}{2}|$  minimal ist über alle  $T \subseteq \{1, \dots, m\}$  und

$$(\tilde{g}(x_i))_1 = \begin{cases} 0, & \text{falls } i \in T_0 \\ 1, & \text{falls } i \notin T_0 \end{cases}.$$

Setze  $T_1 = \{1, \dots, m\} \setminus T_0$ .

Im nächsten Schritt wird  $H(Y_2 | Y_1)$  in (4.6) maximiert. Nach Definition 3.3 gilt

$$H(Y_2 | Y_1) = P(Y_1 = 0)H(Y_2 | Y_1 = 0) + P(Y_1 = 1)H(Y_2 | Y_1 = 1).$$

Zu maximieren sind  $H(Y_2 | Y_1 = 0)$  und  $H(Y_2 | Y_1 = 1)$ . Wir beginnen mit dem ersteren und setzen

$$\begin{aligned} q_{0|0} &= P(Y_2 = 0 | Y_1 = 0) = \frac{P(Y_2 = 0, Y_1 = 0)}{P(Y_1 = 0)} \\ &= \sum_{i: (\tilde{g}(x_i))_1=0, (\tilde{g}(x_i))_2=0} \frac{p_i}{q_0} \end{aligned}$$

sowie  $1 - q_{0|0} = P(Y_2 = 1 | Y_1 = 0)$ .  $H(Y_2 | Y_1 = 0)$  ist maximal bezüglich  $q_{0|0}$ , wenn  $|\sum_{i \in T_{0|0}} p_i/q_0 - \frac{1}{2}|$  minimal ist über alle  $T_{0|0} \subseteq T_0$  und

$$(\tilde{g}(x_i))_2 = \begin{cases} 0, & \text{falls } i \in T_{0|0} \\ 1, & \text{falls } i \in T_0 \setminus T_{0|0} = T_{1|0} \end{cases}.$$

Ganz analog ist  $H(Y_2 | Y_1 = 1)$  maximal bezüglich  $q_{0|1}$ , wenn  $|\sum_{i \in T_{0|1}} p_i/(1 - q_0) - \frac{1}{2}|$  minimal ist über alle  $T_{0|1} \subseteq T_1$  und

$$(\tilde{g}(x_i))_2 = \begin{cases} 0, & \text{falls } i \in T_{0|1} \\ 1, & \text{falls } i \in T_1 \setminus T_{0|1} = T_{1|1} \end{cases}.$$

Dieses Verfahren der sukzessiven Aufspaltung in Teilmengen mit möglichst gleicher Summe der verbleibenden Wahrscheinlichkeiten wird bei entsprechender Kodierung fortgesetzt bis erstmalig für  $k \in \mathbb{N}$

$$H(Y_k | Y_1 = y_1, \dots, Y_{k-1} = y_{k-1}) = 0, \quad y_1, \dots, y_{k-1} \in \{0, 1\}.$$

**Beispiel 4.6** Gegeben sei eine Quelle mit zehnelementigem Alphabet und den in der Tabelle angegebenen Wahrscheinlichkeiten. Berechnet werden ein Fano-Kode  $g$  und ein (optimaler) Huffman-Kode  $g^*$ .

$p_i$	Fano-Kode $g$		Huffman-Kode $g^*$				
0.3	0	0		1 0			
0.2	0	1		0 0			
0.1	1	0		0	0 1 0		
0.1	1	0		1	0 1 1		
0.05	1	0		1	1	1 1 0 0 0	
0.07	1	1		0	0	1 1 1 0	
0.05	1	1		0	1	1 1 0 0 1	
0.08	1	1		1	0	1 1 0 1	
0.03	1	1		1	1	0	1 1 1 1 0
0.02	1	1		1	1	1	1 1 1 1 1

Es gilt  $\bar{n}(g) = 2 \cdot 0.5 + 3 \cdot 0.1 + 4 \cdot 0.35 + 5 \cdot 0.05 = 2.95$  und  $\bar{n}(g^*) = 2.95$ . Der berechnete Fano-Kode ist also in diesem Fall optimal, da er die gleiche erwartete Kodewortlänge wie  $g^*$  erzielt. Bei Verwendung von Logarithmen zur Basis 2 gilt  $H(X) = 2.8728 \leq 2.95$ . ■

### 4.3 Binäre Suchbäume

Mit denselben Prinzipien, die bei der Konstruktion optimaler Codes verwendet werden, können auch effiziente binäre Suchbäume bestimmt werden. Die Elemente einer endlichen, total geordneten Menge  $\mathcal{X}$  werden hierbei so als Knoten eines binären Baumes angeordnet, daß ein vorgegebenes Element von  $\mathcal{X}$  durch einen festen Algorithmus in möglichst kurzer Zeit gefunden wird.

Die im vorigen Abschnitt bestimmten binären Codes lassen sich mit den Knoten eines binären Baums identifizieren, wenn man in der Wurzel eines vollständigen binären Baumes genügend großer Tiefe startet und jede 1 mit der Anweisung “besuche den rechten Nachfolger” und jede 0 mit “besuche den linken Nachfolger” identifiziert. Das  $j$ -te Kodewort  $(w_{j1}, \dots, w_{jn_j})$ ,

$w_{ji} \in \{0, 1\}$ , wird mit dem Knoten identifiziert, in dem man durch sukzessives Abarbeiten des Kodewortes  $(w_{j1}, \dots, w_{jn_j})$  nach obiger Vorschrift endet. Bei präfixfreien Codes ist kein Knoten eines binären Baumes Nachfolger eines anderen (vergleiche Beispiel 4.4). Wäre nämlich Knoten  $\ell$  Nachfolger von Knoten  $k$ , so müßte das zu  $k$  gehörige Kodewort Präfix des zu  $\ell$  gehörigen sein. Dieses "Nachfolerverbot" wird bei der jetzt folgenden Konstruktion von Suchbäumen aufgegeben und verhindert eine direkte Anwendung der im vorigen Abschnitt erhaltenen Ergebnisse. Verwandte Methoden werden jedoch auch für die neuen Probleme zum Ziel führen.

Zunächst wird der Begriff des binären Suchbaums definiert, wobei die Kenntnis binärer Bäume und ihre Implementation in rekursiver PASCAL-Notation vorausgesetzt wird. Wenn man die Knoten mit einer Teilmenge der natürlichen Zahlen identifiziert, ist ein geeigneter Datentyp

```

TYPE node = RECORD no:   INTEGER;
                    left,right: ↑node
                    END.

```

(4.7)

**Definition 4.6** (*binärer Suchbaum*)

$\mathcal{X} = \{x_1, \dots, x_m\}$ ,  $m \in \mathbb{N}$ ,  $m \geq 2$ , sei eine bezüglich " $<$ " vollständig geordnete endliche Menge mit  $x_1 < x_2 < \dots < x_m$ . Ein binärer Baum  $T$  mit Knoten  $\{x_1, \dots, x_m\}$  heißt binärer Suchbaum über  $\mathcal{X}$ , wenn für alle Knoten  $x_j \in \mathcal{X}$ , für alle Knoten  $x_i$  im linken und alle Knoten  $x_k$  im rechten Teilbaum mit Wurzel  $x_j$  gilt, daß  $x_i < x_j < x_k$ .

Die in Abbildung 4.2 dargestellten Bäume  $T_1$  und  $T_2$  sind binäre Suchbäume über der Menge  $\mathcal{X} = \{0, 1, 2, \dots, 9\} \subset \mathbb{N}_0$  mit der Ordnung  $0 < 1 < 2 < \dots < 9$ .

Der folgende Algorithmus findet Elemente  $x \in \mathcal{X}$  in einem binären Suchbaum  $T$  über  $\mathcal{X}$  wieder. Beginne mit der Wurzel des Baumes  $x^*$ . Ist  $x < x^*$ , setze  $x^* =$  Wurzel des linken Teilbaums, sonst  $x^* =$  Wurzel des rechten Teilbaums, bis  $x = x^*$  oder ein Nachfolgerbaum leer ist.

Dieser Algorithmus benötigt bis zum Zugriff auf Knoten  $x_i$  im Baum  $T$  genau  $Z_T(x_i) = t_i + 1$  Vergleiche, wobei  $t_i$  die Tiefe (= Anzahl der Vorgänger) des Knotens  $x_i$  ist. Die sogenannte Zugriffszeit  $Z_T(x_i)$  des Knotens  $x_i$  ist eine Abbildung von  $\mathcal{X}$ , der Knotenmenge des Baumes, in die natürlichen Zahlen,  $Z_T : \mathcal{X} \rightarrow \mathbb{N}$ , und hängt natürlich von der speziellen Gestalt des Baums  $T$  ab.

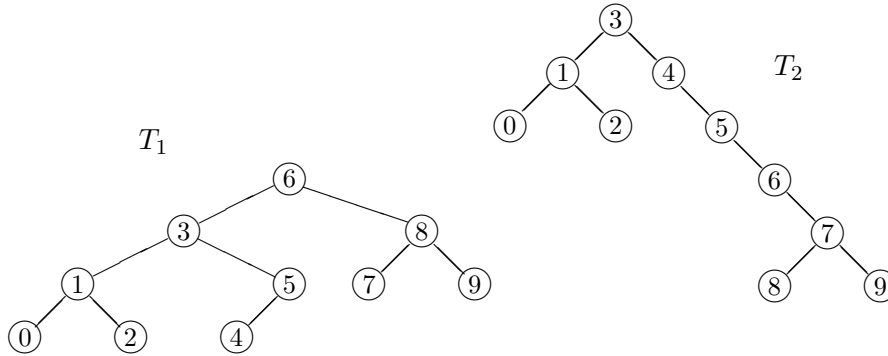


Abb. 4.2 Binäre Suchbäume für  $\{0, 1, 2, \dots, 9\}$ .

Wir nehmen jetzt an, daß Vorkenntnisse über die Häufigkeiten, mit der auf ein Element von  $\mathcal{X}$  zugegriffen wird, in Form einer Wahrscheinlichkeitsverteilung  $\mathbf{p} = (p_1, \dots, p_m)$  vorliegen. Diese Verteilung heißt Zugriffsverteilung. Die Zugriffszeit ist unter der Zugriffsverteilung eine Zufallsvariable, ihr Erwartungswert  $E(Z_T)$  ein Maß für die Qualität des unterliegenden binären Suchbaums  $T$ ,

$$E(Z_T) = \sum_{i=1}^m (t_i + 1)P(Z_T = t_i + 1) = \sum_{i=1}^m p_i(t_i + 1).$$

Für die in Abbildung 4.2 betrachteten Bäume  $T_1$  bzw.  $T_2$  beträgt bei Vorliegen einer Gleichverteilung  $p_i = 0.1, i = 0, 1, \dots, 9$ ,

$$E(Z_{T_1}) = 2.9 \quad \text{und} \quad E(Z_{T_2}) = 3.5.$$

Der erste Suchbaum ist daher effizienter als der zweite. Unter der Zugriffsverteilung

$$\mathbf{p} = (0.2, 0.2, 0.2, 0.1, 0.05, 0.05, 0.05, 0.05, 0.05, 0.05)$$

gilt jedoch

$$E(Z_{T_1}) = 3.2 \quad \text{und} \quad E(Z_{T_2}) = 2.55.$$

In diesem Fall hat der zweite Suchbaum die höhere Effizienz, die damit offensichtlich von der speziellen Gestalt der Zugriffsverteilung abhängt. Bei

fester Zugriffsverteilung sind solche Suchbäume günstig, die den Knoten mit geringer Wahrscheinlichkeit große Zugriffszeiten geben und denen mit großer Wahrscheinlichkeit eine kleine Zugriffszeit.

Analog zu Satz 4.4 werden nun mit Hilfe der Entropie der Zugriffsverteilung obere und untere Schranken für die erwartete Zugriffszeit bestimmt. Auf dem Weg dorthin wird das folgende Ergebnis benötigt.

**Lemma 4.5**  *$T$  sei ein binärer Suchbaum über der Knotenmenge  $\mathcal{X} = \{x_1, \dots, x_m\}$ ,  $m \geq 2$ ,  $t_i$  bezeichne die Tiefe des Knotens  $x_i$ . Dann gilt für alle  $0 \leq c \leq 1$  die Ungleichung*

$$c \sum_{i=1}^m ((1-c)/2)^{t_i} \leq 1.$$

**Beweis.** Für  $m(k) = |\{x_i | t_i = k\}|$ , die Anzahl der Knoten mit Tiefe  $k$ ,  $k \in \mathbb{N}_0$ , gilt  $m(k) \leq 2^k$ , da  $T$  ein binärer Baum ist. Es folgt

$$\begin{aligned} c \sum_{i=1}^m \left(\frac{1-c}{2}\right)^{t_i} &= c \sum_{k=0}^m m(k) \left(\frac{1-c}{2}\right)^k \\ &\leq c \sum_{k=0}^m (1-c)^k = 1 - (1-c)^{m+1} \leq 1. \end{aligned}$$

■

**Satz 4.6**  *$H = H(p_1, \dots, p_m)$  sei die Entropie der Zugriffsverteilung  $\mathbf{p} = (p_1, \dots, p_m)$ . Dann gilt für alle binären Suchbäume  $T$  über der Knotenmenge  $\mathcal{X} = \{x_1, \dots, x_m\}$*

$$\max_{y \in \mathbb{R}} \left\{ \frac{H - y}{\log(2 + b^{-y})} \right\} \leq E(Z_T), \quad (4.8)$$

wobei  $b > 1$  die Basis des Logarithmus ist, die auch bei der Berechnung von  $H$  verwendet wird.

**Beweis.** Für  $0 \leq c < 1$  bezeichne  $c' = (1 - c)/2$  und  $q_i = c \left(\frac{1-c}{2}\right)^{t_i}$ ,  $i = 1, \dots, m$ , wobei  $t_i$  wieder die Tiefe des Knotens  $x_i$  bedeutet. Dann gilt  $t_i + 1 = (\log q_i - \log c)/\log c' + 1$  und

$$\mathbb{E}(Z_T) = \sum_{i=1}^m p_i(t_i + 1) = 1 - \frac{\log c}{\log c'} + \frac{1}{\log c'} \sum_{i=1}^m p_i \log q_i \quad (4.9)$$

Wegen  $\log z \leq (\log e)(z - 1)$ ,  $z > 0$ , und Lemma 4.5 folgt die Abschätzung

$$\begin{aligned} H(p_1, \dots, p_m) + \sum_{i=1}^m p_i \log q_i &= \sum_{i=1}^m p_i \log \frac{q_i}{p_i} \\ &= \log e \sum_{\substack{i=1 \\ p_i \neq 0}}^m p_i \ln \frac{q_i}{p_i} = \log e \left( \sum_{\substack{i=1 \\ p_i \neq 0}}^m p_i \frac{q_i}{p_i} - 1 \right) \leq 0, \end{aligned}$$

die mit (4.9), da  $\log c' \leq 0$ ,

$$\mathbb{E}(Z_T) \geq 1 - \frac{\log c}{\log c'} - \frac{H}{\log c'} = \frac{1}{\log c'} \left( \log \frac{c'}{c} - H \right) = \frac{H - \log(c'/c)}{\log(1/c')}$$

liefert.  $y = \log \frac{c'}{c} = \log \frac{1-c}{2c}$  durchläuft mit  $0 < c < 1$  alle reellen Zahlen. Bei Basis  $b > 1$  gilt also  $b^{-y} = c/c'$ , bzw.  $2 + b^{-y} = 1/c'$ . Insgesamt erhält man  $\mathbb{E}(Z_T) \geq \frac{H-y}{\log(2+b^{-y})}$  für alle  $y \in \mathbb{R}$ , woraus die Behauptung folgt. ■

Aus Ungleichung (4.8) folgt speziell für  $y = 0$  die untere Schranke

$$H/\log 3 \leq \mathbb{E}(Z_T) \quad (4.10)$$

für alle binären Suchbäume  $T$ . Eine bessere, implizite Schranke läßt sich aus Satz 4.6 ableiten, indem man dort  $y = \log(\mathbb{E}(Z_T)/2)$  einsetzt und die entstehende Ungleichung nach  $H$  auflöst.

**Korollar 4.2** *Unter den Bezeichnungen von Satz 4.6 gilt für alle binären Suchbaume  $T$ , daß  $H \leq \mathbb{E}(Z_T) + \log \mathbb{E}(Z_T) + \log e - 1$ .*

Ein binärer Suchbaum  $T$ , der die untere Schranke in (4.8) erreicht, ist bezüglich der erwarteten Zugriffszeit absolut optimal im Sinn von Satz 4.4. Im allgemeinen existiert ein solcher Suchbaum jedoch nicht. Analog zu Satz 4.4 b) wird jetzt ein binärer Suchbaum  $T^*$  konstruiert, für den  $\mathbb{E}(Z_{T^*})$  nahe an der unteren Schranke (4.8) liegt.



**Satz 4.7**  $H = H(p_1, \dots, p_m)$  sei die Entropie der Zugriffsverteilung. Dann existiert ein binärer Suchbaum  $T^*$  für  $\mathcal{X} = \{x_1, \dots, x_m\}$  mit

$$E(Z_{T^*}) \leq \frac{H}{\log 2} + 1.$$

Die Konstruktionsidee eines solchen Baums  $T^*$  beruht darauf, die Summe der Wahrscheinlichkeiten in jedem rechten und linken Teilbaum jedes Knotens möglichst gleich zu machen. Sind  $x_\ell, \dots, x_r \in S$ ,  $1 \leq \ell < r \leq m$ , als Knoten eines binären Teilbaumes angeordnet, so wähle  $x_k$  als Wurzel dieses Teilbaumes, wenn

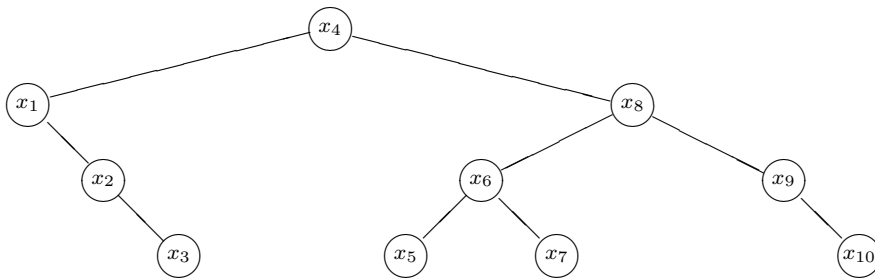
$$\sum_{j=\ell}^{k-1} p_j < \frac{1}{2} \sum_{j=\ell}^r p_j \quad \text{und} \quad \sum_{j=\ell}^k p_j \geq \frac{1}{2} \sum_{j=\ell}^r p_j. \quad (4.11)$$

Verfahre mit den Knoten  $x_\ell, \dots, x_{k-1}$  bzw.  $x_{k+1}, \dots, x_r$  genauso bis  $\ell > k-1$  bzw.  $r < k+1$ .

**Beispiel 4.7** Sei  $\mathcal{X} = \{x_1, \dots, x_{10}\}$ ,  $x_1 < x_2 < \dots < x_{10}$  mit Zugriffsverteilung

$$\mathbf{p} = (0.2, 0.1, 0.05, 0.3, 0.05, 0.03, 0.08, 0.03, 0.1, 0.06).$$

Dann ist  $\sum_{j=1}^4 p_j \geq 1/2$  erstmalig,  $x_4$  bildet daher die Wurzel des Baumes. Im linken Teilbaum verbleiben die Knoten  $x_1, x_2, x_3$ . Es gilt  $p_1 = 0.2 \geq \frac{1}{2} \sum_{j=1}^3 p_j = 0.35/2$ , so daß  $x_1$  Wurzel des linken Teilbaums ist. Im rechten Teilbaum gilt  $\sum_{j=5}^8 p_j = 0.18 \geq \frac{1}{2} \sum_{j=5}^{10} p_j = 0.35/2$  erstmalig, so daß  $x_8$  Wurzel des rechten Teilbaums wird. Führt man das Verfahren vollständig durch, entsteht der folgende binäre Suchbaum.



■

Der Algorithmus zur allgemeinen Konstruktion solcher ausgewogenen binären Suchbäume wird durch die folgende rekursive Prozedur in PASCAL-Notation mit dem Datentyp `node` aus (4.7) beschrieben.

```

FUNCTION bintree(l,r: INTEGER): ↑node;
  VAR z: ↑node;
      m: INTEGER;
BEGIN
  { bestimme Index m zwischen l und r mit
    ausgeglichenen Wahrscheinlichkeiten, d.h.
     $\sum_{j=l}^{m-1} p_j < \frac{1}{2} \sum_{j=l}^r p_j, \sum_{j=l}^m p_j \geq \frac{1}{2} \sum_{j=l}^r p_j$  };
  new(z); z↑.no:=m;
  IF l<=m-1 THEN z↑.left:=bintree(l,m-1)
    ELSE z↑.left:=NUL;
  IF m+1<=r THEN z↑.right:=bintree(m+1,r)
    ELSE z↑.right:=NUL;
  bintree:=z;
END;

```

(4.12)

Zum Beweis von Satz 4.7 wird noch die folgende Aussage benötigt.

**Lemma 4.6** *Für jeden binären Teilbaum mit Wurzel  $x_k$  und Knotenmenge  $\{x_\ell, \dots, x_k, \dots, x_r\}$  eines ausgewogenen Suchbaums  $T^*$  aus (4.12) gilt  $\sum_{j=\ell}^r p_j \leq 2^{-t_k^*}$ , insbesondere also  $p_k \leq 2^{-t_k^*}$ , wobei  $t_k^*$  die Tiefe des Knotens  $x_k$  bezeichnet.*

**Beweis.** Der Beweis wird durch Induktion über die Tiefe der Wurzeln geführt. Ist  $x_{k_0}$  Wurzel von  $T^*$ , so gilt  $t_{k_0}^* = 0$ , also  $\sum_{j=1}^m p_j \leq 2^{-t_{k_0}^*} = 1$ .  $x_k$  sei Wurzel des Teilbaums mit Knoten  $x_\ell, \dots, x_k, \dots, x_r$ , und es gelte  $\sum_{j=\ell}^r p_j \leq 2^{-t_k^*}$ . Für den linken und rechten Teilbaum gilt dann wegen (4.11)  $\sum_{j=\ell}^{k-1} p_j \leq 2^{-t_k^*-1}$  und  $\sum_{j=k+1}^r p_j \leq 2^{-t_k^*-1}$ , wobei Summen mit leerem Indexbereich zu 0 gesetzt werden. Die Tiefe der Wurzel des linken bzw. rechten Teilbaums beträgt gerade  $t_k+1$ , falls dieser nichtleer ist, woraus die Behauptung folgt. ■

Nach diesen Vorbereitungen läßt sich der Beweis von Satz 4.7 wie folgt führen. Mit Lemma 4.6 gilt für einen nach (4.12) konstruierten Baum  $T^*$ , daß  $p_j \leq 2^{-t_j^*}$  für alle  $j = 1, \dots, m$ , also  $\log p_j \leq -t_j^* \log 2$ . Es folgt

$$\begin{aligned} E(Z_{T^*}) &= \sum_{j=1}^m p_j (t_j^* + 1) \leq -\frac{1}{\log 2} \sum_{j=1}^m p_j \log p_j + 1 \\ &= \frac{H(p_1, \dots, p_m)}{\log 2} + 1. \end{aligned}$$

Durch Algorithmus (4.12) wird also ein binärer Suchbaum  $T^*$  berechnet, für dessen erwartete Zugriffszeit wegen (4.10) und Satz 4.7

$$\frac{H}{\log 3} \leq E(Z_{T^*}) \leq \frac{H}{\log 2} + 1 \quad (4.13)$$

gilt. Mit Ungleichung (4.13) stehen wie bei der Kodierung obere und untere Schranken für die Güte binärer Suchbäume zur Verfügung. Die Konstruktion optimaler Suchbäume wird in Mehlhorn [25] beschrieben. Diese Algorithmen sind das Äquivalent des Huffman-Verfahrens, mit dem optimale Codes konstruiert werden.

## 4.4 Stationäre Quellen, Markoff-Quellen

Bei der buchstabenweisen Kodierung einer diskreten gedächtnislosen Quelle  $\{X_n\}_{n \in \mathbb{N}}$  spielt die stochastische Unabhängigkeit der Komponenten  $X_n$  natürlich keine Rolle. Es interessiert lediglich die identische Randverteilung  $P^{X_n} = P^X$ . Bei Blockkodierung in Satz 4.5 und Korollar 4.1 wurde allerdings wesentlich benutzt, daß aufeinanderfolgende Symbole der Quelle unabhängig voneinander generiert werden,  $\{X_n\}_{n \in \mathbb{N}}$  also eine stochastisch unabhängige Folge ist.

Die Annahme der stochastischen Unabhängigkeit einzelner Buchstaben ist jedoch zur Beschreibung realer Quellen in vielen Fällen zu restriktiv. In natürlichen Sprachen bestehen sehr wohl Abhängigkeiten zwischen aufeinanderfolgenden Buchstaben. Solche Abhängigkeiten zu modellieren und einen Entropiebegriff zu entwickeln, der die Abhängigkeiten berücksichtigt, ist Inhalt dieses Abschnitts. Für eine sehr allgemeine Abhängigkeitsstruktur können bei Blockkodierung dieselben Aussagen wie in Satz 4.5 erzielt werden, und zwar für diskrete stationäre Quellen.

**Definition 4.7** (Diskrete stationäre Quelle, DSQ)

Eine stationäre Folge von Zufallsvariablen  $\{X_n\}_{n \in \mathbb{N}}$  (vgl. Definition 2.4),  $X_n$  jeweils diskret mit endlichem Träger  $\mathcal{X} = \{x_1, \dots, x_m\}$ , heißt diskrete stationäre Quelle (kurz: DSQ).

Stationäre Folgen von Zufallsvariablen besitzen insbesondere identische Randverteilungen  $P^{X_n} = P^X$ , wie man aus Definition 2.4 mit  $n = 1$  sieht. Es gilt also  $H(X_n) = H(X_1)$  für alle  $n \in \mathbb{N}$ .

Im folgenden bezeichne  $\mathbf{X}_N = (X_1, \dots, X_N)$  den Zufallsvektor aus den ersten  $N$  Gliedern von  $\{X_n\}_{n \in \mathbb{N}}$ .  $\frac{1}{N}H(\mathbf{X}_N)$  ist dann die durchschnittliche Unbestimmtheit pro Quellbuchstabe der ersten  $N$  Symbole. Bei der Definition der Entropie pro Quellbuchstabe müssen Abhängigkeiten zwischen den Zufallsvariablen  $X_n$  berücksichtigt werden, die gegebenenfalls unendlich weit reichen. Die entscheidenden Grundlagen für die Definition stellt der folgende Satz bereit.

**Satz 4.8**  $\{X_n\}_{n \in \mathbb{N}}$  sei eine diskrete stationäre Quelle. Dann gilt

- a)  $H(X_N | X_1, \dots, X_{N-1})$  ist monoton fallend in  $N$ .
- b)  $H(X_N | X_1, \dots, X_{N-1}) \leq \frac{1}{N}H(\mathbf{X}_N)$  für alle  $N \in \mathbb{N}$ .
- c)  $\frac{1}{N}H(\mathbf{X}_N)$  ist monoton fallend in  $N$ .
- d)  $\lim_{N \rightarrow \infty} H(X_N | X_1, \dots, X_{N-1}) = \lim_{N \rightarrow \infty} \frac{1}{N}H(\mathbf{X}_N)$ , wobei wegen a) und c) beide Grenzwerte existieren.

**Beweis.** a) Wegen Satz 3.1 e) und der Stationarität von  $\{X_n\}_{n \in \mathbb{N}}$  folgt

$$\begin{aligned} H(X_N | X_1, \dots, X_{N-1}) &\leq H(X_N | X_2, \dots, X_{N-1}) \\ &= H(X_{N-1} | X_1, \dots, X_{N-2}) \end{aligned}$$

für alle  $N \in \mathbb{N}$ .  $H(X_N | X_1, \dots, X_{N-1})$  ist also monoton fallend, und der erste Limes in d) existiert, da eine monoton fallende und durch 0 nach unten beschränkte Folge vorliegt.

b) Lemma 3.2 a) liefert für alle  $N \in \mathbb{N}$

$$\begin{aligned} \frac{1}{N}H(\mathbf{X}_N) &= \frac{1}{N} \left( H(X_1) + H(X_2 | X_1) + \dots \right. \\ &\quad \left. + H(X_N | X_1, \dots, X_{N-1}) \right) \\ &\geq H(X_N | X_1, \dots, X_{N-1}). \end{aligned}$$

Die letzte Ungleichung folgt hierbei aus Teil a).

c) Mit b) folgt

$$\begin{aligned} H(\mathbf{X}_N) &= H(X_1, \dots, X_{N-1}) + H(X_N | X_1, \dots, X_{N-1}) \\ &\leq H(\mathbf{X}_{N-1}) + \frac{1}{N}H(\mathbf{X}_N). \end{aligned}$$

Folglich ist  $\frac{N-1}{N}H(\mathbf{X}_N) \leq H(\mathbf{X}_{N-1})$ , bzw.  $\frac{1}{N}H(\mathbf{X}_N) \leq \frac{1}{N-1}H(\mathbf{X}_{N-1})$  für alle  $N \in \mathbb{N}$ ,  $N \geq 2$ . Dies zeigt die Monotonie der Folge  $\frac{1}{N}H(\mathbf{X}_N)$  und damit die Existenz des zweiten Limes in d).

d) Für alle  $k \in \mathbb{N}$  gilt mit a)

$$\begin{aligned} \frac{1}{N+k}H(\mathbf{X}_{N+k}) &= \frac{1}{N+k} \left( H(X_{N+k} | X_1, \dots, X_{N+k-1}) \right. \\ &\quad \left. + \dots + H(X_{N+1} | X_1, \dots, X_N) \right. \\ &\quad \left. + H(X_N | X_1, \dots, X_{N-1}) + H(X_1, \dots, X_{N-1}) \right) \\ &\leq \frac{1}{N+k}H(X_1, \dots, X_{N-1}) + \frac{k+1}{k+N}H(X_N | X_1, \dots, X_{N-1}). \end{aligned}$$

Der erste Summand der letzten Zeile konvergiert mit  $k \rightarrow \infty$  gegen 0, der zweite gegen  $H(X_N | X_1, \dots, X_{N-1})$ . Setzt man jetzt  $L = N + k$ , so folgt  $\lim_{L \rightarrow \infty} \frac{1}{L}H(\mathbf{X}_L) \leq H(X_N | X_1, \dots, X_{N-1})$  für alle  $N \in \mathbb{N}$ . Mit Hilfe von b) schließen wir für die Grenzwerte

$$\begin{aligned} \lim_{N \rightarrow \infty} H(X_N | X_1, \dots, X_{N-1}) &\leq \lim_{N \rightarrow \infty} \frac{1}{N}H(\mathbf{X}_N) \\ &\leq \lim_{N \rightarrow \infty} H(X_N | X_1, \dots, X_{N-1}) \end{aligned}$$

woraus schließlich d) folgt. ■

Die Grenzwerte in Teil d) des obigen Satzes sind die richtigen Begriffe für die Entropie pro Quellbuchstabe unter Berücksichtigung von Abhängigkeiten.

**Definition 4.8** (*Entropie einer diskreten stationären Quelle*)

$\{X_n\}_{n \in \mathbb{N}}$  sei eine diskrete stationäre Quelle. Dann heißt

$$H_\infty(X) = \lim_{N \rightarrow \infty} \frac{1}{N}H(\mathbf{X}_N) = \lim_{N \rightarrow \infty} H(X_N | X_1, \dots, X_{N-1})$$

Entropie der Quelle.

Wegen Satz 4.8 c) und Satz 3.1 d) gilt für alle  $N \in \mathbb{N}$

$$H_\infty(X) \leq \frac{1}{N} H(\mathbf{X}_N) \leq \frac{1}{N} [H(X_1) + \dots + H(X_N)] = H(X_1).$$

Damit ist die Entropie einer DSQ stets kleiner als die Entropie der (identischen) Randverteilung. Dies sollte man aus eventuell bestehenden Abhängigkeiten in der Folge  $\{X_n\}_{n \in \mathbb{N}}$  auch erwarten. Im Spezialfall einer diskreten gedächtnislosen Quelle gilt  $\frac{1}{N} H(\mathbf{X}_N) = H(X_1) = H_\infty(X)$ .

Als direktes Analogon zu Korollar 4.1 erhalten wir nun

**Korollar 4.3**  $\{X_n\}_{n \in \mathbb{N}}$  sei eine diskrete stationäre Quelle. Dann existiert eine Folge von e.d. Blockcodes  $g^{(N)}$ ,  $N \in \mathbb{N}$ , mit

$$\lim_{N \rightarrow \infty} \frac{1}{N} \bar{n}(g^{(N)}) = \frac{H_\infty(X)}{\log d}.$$

**Beweis.** Nach Satz 4.4 existiert für alle  $N \in \mathbb{N}$  ein eindeutig dekodierbarer Blockcode  $g^{(N)}$  mit

$$\frac{H(\mathbf{X}_N)}{\log d} \leq \bar{n}(g^{(N)}) < \frac{H(\mathbf{X}_N)}{\log d} + 1.$$

Division durch  $N$  liefert

$$\frac{\frac{1}{N} H(\mathbf{X}_N)}{\log d} \leq \frac{\bar{n}(g^{(N)})}{N} < \frac{\frac{1}{N} H(\mathbf{X}_N)}{\log d} + \frac{1}{N}.$$

Die linke Seite ist größer gleich  $H_\infty(X)/\log d$ , die rechte konvergiert mit  $N \rightarrow \infty$  gegen  $H_\infty(X)/\log d$ . Dies zeigt die Behauptung. ■

**Beispiel 4.8** Sei  $\mathcal{X} = \{x_1, x_2, x_3\} \subset \mathbb{R}$ ,  $U$  eine binomialverteilte Zufallsvariable mit  $P(U = 0) = P(U = 1) = \frac{1}{2}$ .  $V$  sei einpunktverteilt in  $x_1$ , d.h.  $P(V = x_1) = 1$ . Weiterhin sei  $\{Z_n\}_{n \in \mathbb{N}}$  eine Folge von stochastisch unabhängigen, identisch verteilten Zufallsvariablen mit  $P(Z_n = x_2) = P(Z_n = x_3) = \frac{1}{2}$ .  $U, V$  und  $\{Z_n\}$  seien gemeinsam stochastisch unabhängig. Die Zufallsvariablen  $X_n = VU + Z_n(1 - U)$ ,  $n \in \mathbb{N}$ , bilden dann eine diskrete stationäre Quelle, was in Aufgabe 4.11 nachzuweisen ist.

Das Verhalten der Quelle kann folgendermaßen interpretiert werden. Mit Wahrscheinlichkeit  $\frac{1}{2}$  sendet die Quelle gleichverteilt die Buchstaben  $x_2$  oder

$x_3$ , oder aber mit der Wahrscheinlichkeit  $\frac{1}{2}$  die konstante Folge  $(x_1, x_1, \dots)$ , etwa wenn sie defekt ist. In beiden möglichen Zuständen (defekt oder intakt) befindet sich die Quelle mit der Wahrscheinlichkeit  $\frac{1}{2}$ .

Eine optimale binäre Blockkodierung erfolgt nun mit dem Huffman-Verfahren. Der zugehörige Kodebaum läßt sich leicht überblicken. Die resultierenden Kodewortlängen eines optimalen Kodes sind in der dritten Spalte der folgenden Tabelle angegeben.

$\mathbf{a}_N$	$P(\mathbf{X}_N = \mathbf{a}_N)$	$n_j, j = 1, \dots, 2^N + 1$
$\underbrace{(x_1, \dots, x_1)}_N$	$\frac{1}{2}$	1
$(a_1, \dots, a_N) \in \{x_2, x_3\}^N$	$\frac{1}{2} \cdot \frac{1}{2^N} = \frac{1}{2^{N+1}}$	$N + 1$

Für die erwartete Kodewortlänge bei optimaler Blockkodierung von Blöcken der Länge  $N$  ergibt sich

$$\bar{n}(g^{(N)*}) = 1 \cdot \frac{1}{2} + (N + 1) \frac{2^N}{2^{N+1}} = \frac{N + 2}{2}.$$

Mit Korollar 4.3 folgt bei Verwendung von Logarithmen zur Basis 2, daß

$$H_\infty(X) = \lim_{n \rightarrow \infty} \frac{1}{N} \bar{n}(g^{(N)*}) = \lim_{n \rightarrow \infty} \frac{N + 2}{2N} = \frac{1}{2}. \quad \blacksquare$$

Empirische Untersuchungen haben gezeigt, daß das Auftreten einzelner Buchstaben in natürlichen Sprachen sehr gut durch ein stochastisches Modell beschrieben werden kann, bei dem stochastische Abhängigkeiten nur auf wenige Vorgängerbuchstaben bestehen. Mit Hilfe eines Lexikons kann man sich zum Beispiel davon überzeugen, daß in einem fehlerfreien deutschen Text auf die drei ersten Buchstaben **ins** eines Wortes mit positiver Wahrscheinlichkeit nur ein **p** oder ein **t** folgen kann.

Die folgende Methode eignet sich dazu, simulativ Buchstaben zu generieren, die auf einen vorgegebenen Buchstabenblock der Länge  $k$  folgen. Nach Wahl einer Anfangssequenz der Länge  $k$  wird in einem repräsentativen Text die nächste Stelle gesucht, an der diese Buchstabenfolge auftritt. Das dort erscheinende Folgezeichen wird als Nachfolgesymbol eingesetzt. Mit den letzten  $k$  Zeichen der jetzt vorliegenden  $k + 1$  Buchstaben fährt man genauso fort und iteriert das Verfahren.

Der Text des ersten Kapitels aus “Stochastik für Informatiker” [24] wurde für jeweils  $k = 3$  Vorgängerbuchstaben so behandelt. Als Symbole wurden nur die 26 Buchstaben des lateinischen Alphabets und das Leerzeichen zugelassen, wobei nicht zwischen Groß- und Kleinschreibung unterschieden wurde. Die resultierende Zeichenkette war

dinn arlenkest nurbott gibl einsatz ...

Hierdurch werden die zufälligen Übergänge von drei gegebenen Buchstaben auf den nachfolgenden vierten aufgrund der Datenbasis eines gegebenen Textes simuliert. Es ist nicht zu erwarten, daß sich ein sinnvoller Text ergibt. Dennoch ist das Ergebnis eine Zeichenfolge, die von der sprachlichen Struktur sehr nahe am Deutschen liegt. Für informationstheoretische Zwecke ist dies der entscheidende Aspekt. Der syntaktische Inhalt einer Zeichenfolge spielt für Kodierungs- und Übertragungsverfahren keine Rolle.

Beschreibt die Zufallsvariable  $X_n$  den Buchstaben an der  $n$ -ten Stelle, so wird bei diesem Verfahren eine Abhängigkeit von  $X_n$  auf jeweils  $k$  Vorgänger  $X_{n-k}, \dots, X_{n-1}$  angenommen. Diese kann durch Markoff-Ketten (vgl. Definition 2.8) beschrieben werden. Für stationäre Markoff-Ketten berechnet sich die Entropie dann mit den bekannten Formeln für stationäre diskrete Quellen.

In der folgenden allgemeinen Definition wird das Auftreten von Buchstaben in einem Text durch eine unterliegende Markoff-Kette mit Hilfe einer Funktion  $f$  gesteuert.  $f$  modelliert speziell die Rückwirkung auf mehr als nur den unmittelbaren Vorgängerbuchstaben. Die im folgenden benötigten Begriffe zu Markoff-Ketten sind in Abschnitt 2.2 zusammengestellt.

**Definition 4.9** (*diskrete Markoff-Quelle*)

$\{Z_n\}_{n \in \mathbb{N}_0}$  sei eine homogene Markoff-Kette mit endlichem Zustandsraum  $\mathcal{S} = \{s_1, \dots, s_r\}$ ,  $r \in \mathbb{N}$ .  $\mathcal{X} = \{x_1, \dots, x_m\}$ ,  $m \in \mathbb{N}$ , sei ein Alphabet und  $f: \mathcal{S} \rightarrow \mathcal{X}$  eine Abbildung. Die Folge  $\{X_n\}_{n \in \mathbb{N}_0}$  der Zufallsvariablen  $X_n = f(Z_n)$ ,  $n \in \mathbb{N}_0$ , heißt *diskrete Markoff-Quelle* (kurz: *DMQ*) mit Alphabet  $\mathcal{X} = \{x_1, \dots, x_m\}$ .  $f$  heißt *assoziierte Funktion*.

Im Spezialfall  $\mathcal{S} = \{x_1, \dots, x_m\}$  und  $f = \text{Identität}$  bildet die Markoff-Kette  $X_n = Z_n$ ,  $n \in \mathbb{N}_0$ , bereits selbst eine diskrete Markoff-Quelle. Mehrstufige Abhängigkeiten können durch eine geeignete assoziierte Funktion modelliert werden, wie folgendes Beispiel zeigt.



**Beispiel 4.9**  $\{X_n\}_{n \in \mathbb{N}_0}$  sei eine diskrete Markoff-Quelle mit Alphabet  $\mathcal{X} = \{0, 1\}$ . Stochastische Abhängigkeiten sollen über  $k = 3$  sukzessive Buchstaben bestehen. Man wähle nun  $\mathcal{S} = \{0, 1\}^3 = \{(0, 0, 0), \dots, (1, 1, 1)\} = \{s_1, \dots, s_8\}$ .  $\{Z_n\}_{n \in \mathbb{N}_0}$  sei eine homogene Markoff-Kette mit Zustandsraum  $\mathcal{S}$  und der Übergangsmatrix

$$\mathbf{H} = \begin{array}{c|cccccccc} & 000 & 001 & 010 & 011 & 100 & 101 & 110 & 111 \\ \hline 000 & p_{11} & p_{12} & 0 & 0 & 0 & 0 & 0 & 0 \\ 001 & 0 & 0 & p_{23} & p_{24} & 0 & 0 & 0 & 0 \\ 010 & 0 & 0 & 0 & 0 & p_{35} & p_{36} & 0 & 0 \\ 011 & 0 & 0 & 0 & 0 & 0 & 0 & p_{47} & p_{48} \\ 100 & p_{51} & p_{52} & 0 & 0 & 0 & 0 & 0 & 0 \\ 101 & 0 & 0 & p_{63} & p_{64} & 0 & 0 & 0 & 0 \\ 110 & 0 & 0 & 0 & 0 & p_{75} & p_{76} & 0 & 0 \\ 111 & 0 & 0 & 0 & 0 & 0 & 0 & p_{87} & p_{88} \end{array}$$

Die Übergangswahrscheinlichkeiten  $p_{ij}$  sind eventuell positiv, alle anderen Einträge haben den Wert 0. Mit Hilfe der assoziierten Funktion  $f$  wird jeweils das letzte Symbol des internen Zustands aus  $\mathcal{S}$  ausgegeben. Dies wird erreicht, wenn  $f$  die Projektion auf die 3. Komponente ist, also  $f : \mathcal{S} \rightarrow \mathcal{X} : (a_1, a_2, a_3) \mapsto a_3$ . ■

Ist die Anfangsverteilung  $\mathbf{p}(0)$  der Markoff-Kette  $\{Z_n\}_{n \in \mathbb{N}_0}$  stationär, so ist nach Lemma 2.3 die ganze Folge  $\{Z_n\}_{n \in \mathbb{N}_0}$  stationär. Wie man leicht mit Definition 2.4 überprüft, ist dann auch  $X_n = f(Z_n)$ ,  $n \in \mathbb{N}_0$ , eine stationäre Folge. Mit Satz 4.8 folgt, daß  $H_\infty(Z)$  und  $H_\infty(X)$  existieren.

Die Berechnung der Entropie ist für stationäre Markoff-Quellen im allgemeinen ein schwieriges Problem. Für gewisse Typen assoziierter Funktionen läßt sich die Berechnung aber explizit durchführen. Hierauf zielt die Begriffsbildung der folgenden Definition.

**Definition 4.10** (*unifilar*)

$\{X_n\}_{n \in \mathbb{N}_0}$  sei eine diskrete Markoff-Quelle mit zugehöriger Markoff-Kette  $\{Z_n\}_{n \in \mathbb{N}_0}$  und assoziierter Funktion  $f$ . Für  $s_k \in \mathcal{S}$  bezeichne  $R(s_k) = \{s_i \in \mathcal{S} \mid p_{ki} > 0\}$  die Menge der Zustände, die von  $s_k$  aus in einem Schritt

mit positiver Wahrscheinlichkeit erreichbar sind.  $\{X_n\}_{n \in \mathbb{N}_0}$  heißt unifilar (einfädig), wenn

$$f(s_i) \neq f(s_j) \text{ für alle } s_i \neq s_j \in R(s_k) \text{ und alle } s_k \in \mathcal{S},$$

wenn also  $f|_{R(s_k)}$  für alle  $s_k \in \mathcal{S}$  injektiv ist.

Gegeben seien  $s_k$  als Zustand zur Zeit  $n$  und  $f(s_i)$  mit  $s_i \in R(s_k)$  als Buchstabe zur Zeit  $n+1$ . Dann ist für eine DMQ mit unifilarem  $f$  der Zustand zur Zeit  $n+1$  eindeutig bestimmt durch  $s_i = (f|_{R(s_k)})^{-1}(f(s_i))$ . Der Zustand zur Zeit  $n$  und der Buchstabe zur Zeit  $n+1$  bestimmen also eindeutig den Zustand zur Zeit  $n+1$ . Dies wird bei der Berechnung der Entropie unifilarer Markoff-Quellen benutzt.

**Satz 4.9**  $\{X_n\}_{n \in \mathbb{N}_0}$  sei eine unifilare, diskrete Markoff-Quelle mit unterliegender stationärer Markoff-Kette  $\{Z_n\}_{n \in \mathbb{N}_0}$ . Dann gilt

$$H_\infty(X) = \sum_{j=1}^r p_j^* H(p_{j1}, \dots, p_{jr}),$$

wobei  $\mathbf{p}^* = (p_1^*, \dots, p_r^*)$  die stationäre Verteilung von  $\{Z_n\}_{n \in \mathbb{N}_0}$  bezeichnet und  $(p_{j1}, \dots, p_{jr})$  den  $j$ -ten Zeilenvektor der Übergangsmatrix  $\mathbf{II} = (p_{ij})_{1 \leq i, j \leq r}$  von  $\{Z_n\}_{n \in \mathbb{N}_0}$ .

**Beweis.** Da  $\{X_n\}_{n \in \mathbb{N}_0}$  stationär ist, existiert  $H_\infty(X)$  nach Satz 4.8. Seien nun  $z_0, \dots, z_n \in \mathcal{S}$  mit  $0 < P(Z_0 = z_0, \dots, Z_n = z_n)$ . Wegen

$$\begin{aligned} P(Z_0 = z_0, \dots, Z_n = z_n) \\ &= P(Z_0 = z_0) \cdot P(Z_1 = z_1 \mid Z_0 = z_0) \cdot P(Z_2 = z_2 \mid Z_1 = z_1) \\ &\quad \cdots P(Z_n = z_n \mid Z_{n-1} = z_{n-1}). \end{aligned}$$

ist  $z_j \in R(z_{j-1})$  für alle  $j = 1, \dots, n$ . Da  $f|_{R(z_{j-1})}$  für alle  $j = 1, \dots, n$  injektiv ist, gilt

$$\begin{aligned} P(Z_n = z_n, \dots, Z_0 = z_0) \\ &= P(Z_n = z_n, \dots, Z_1 = z_1 \mid Z_0 = z_0) \cdot P(Z_0 = z_0) \\ &= P(X_n = f(z_n), \dots, X_1 = f(z_1) \mid Z_0 = z_0) \cdot P(Z_0 = z_0) \\ &= P(X_n = u_n, \dots, X_1 = u_1 \mid Z_0 = z_0) \cdot P(Z_0 = z_0), \end{aligned}$$

wobei  $u_j = f(z_j)$  gesetzt wurde. Hiermit folgt

$$\begin{aligned}
& \frac{1}{n} H(Z_0, \dots, Z_n) \\
&= -\frac{1}{n} \sum_{z_0, \dots, z_n \in \mathcal{S}} P(Z_0 = z_0, \dots, Z_n = z_n) \\
&\quad \cdot \log P(Z_0 = z_0, \dots, Z_n = z_n) \\
&= -\frac{1}{n} \sum_{z_0 \in \mathcal{S}, u_1, \dots, u_n \in \mathcal{X}} P(X_n = u_n, \dots, X_1 = u_1 \mid Z_0 = z_0) \\
&\quad P(Z_0 = z_0) \cdot (\log P(X_n = u_n, \dots, X_1 = u_1 \mid Z_0 = z_0) \\
&\quad + \log P(Z_0 = z_0)) \\
&= \frac{1}{n} (H(Z_0) + H(X_1, \dots, X_n \mid Z_0)).
\end{aligned}$$

Mit Lemma 3.2 erhalten wir

$$\begin{aligned}
H(X_1, \dots, X_n \mid Z_0) &= H(Z_0, X_1, \dots, X_n) - H(Z_0) \\
&= H(X_1, \dots, X_n) + H(Z_0 \mid X_1, \dots, X_n) - H(Z_0).
\end{aligned}$$

Die beiden letzten Gleichungen zusammen ergeben

$$\begin{aligned}
H_\infty(Z) &= \lim_{n \rightarrow \infty} \frac{1}{n} H(Z_0, \dots, Z_n) \\
&= \lim_{n \rightarrow \infty} \frac{1}{n} (H(X_1, \dots, X_n) + H(Z_0 \mid X_1, \dots, X_n)) \\
&= \lim_{n \rightarrow \infty} \frac{1}{n} H(X_1, \dots, X_n) = H_\infty(X).
\end{aligned}$$

Für die Entropie der unterliegenden homogenen Markoff-Kette gilt mit Satz 4.8 und der Markoff-Eigenschaft, daß

$$\begin{aligned}
H_\infty(Z) &= \lim_{n \rightarrow \infty} H(Z_n \mid Z_0, \dots, Z_{n-1}) = \lim_{n \rightarrow \infty} H(Z_n \mid Z_{n-1}) \\
&= \lim_{n \rightarrow \infty} H(Z_1 \mid Z_0) = H(Z_1 \mid Z_0) \\
&= \sum_{j=1}^r P(Z_0 = s_j) H(Z_1 \mid Z_0 = s_j) \\
&= \sum_{j=1}^r p_j^* H(p_{j1}, \dots, p_{jr}).
\end{aligned}$$

■

## 4.5 Übungsaufgaben

**Aufgabe 4.1** Von 12 äußerlich gleichen Kugeln besitzen 11 gleiches Gewicht. Eine der Kugeln hat abweichendes Gewicht, jedoch ist nicht bekannt, ob sie leichter oder schwerer als die übrigen ist. Mit Hilfe einer Balkenwaage soll durch Vergleichswägungen herausgefunden werden, welche der Kugeln abweichendes Gewicht besitzt, und gleichzeitig, ob diese leichter oder schwerer ist.

Zeigen Sie: Mit drei Wägungen kann man obige Aufgabe lösen, mit weniger Wägungen jedoch nicht.

**Aufgabe 4.2** Das Quellalphabet  $\mathcal{X} = \{x_1, x_2, x_3\}$  werde durch den binären Kode  $g$  kodiert, wobei  $g(x_1) = (0)$ ,  $g(x_2) = (1, 0)$ ,  $g(x_3) = (1, 1)$ . Die Abbildung

$$G: \bigcup_{j=1}^{\infty} \mathcal{X}^j \rightarrow \bigcup_{j=1}^{\infty} \{0, 1\}^j : (x_1, \dots, x_k) \mapsto (g(x_1), \dots, g(x_k))$$

beschreibe die Kodierung von endlichen Wörtern über dem Quellalphabet.

Bestimmen Sie eine allgemeine Formel, die für gegebenes  $\ell \in \mathbb{N}$  die Mächtigkeit der Menge

$$\mathcal{M} = \{(x_1, \dots, x_k) \in \bigcup_{j=1}^{\infty} \mathcal{X}^j \mid G(x_1, \dots, x_k) \in \mathcal{Y}^{\ell}, k \in \mathbb{N}\}$$

angibt.  $\mathcal{M}$  ist die Menge der Nachrichten, die bei Verwendung von  $g$  mit Kodewörtern der Länge  $\ell$  dargestellt werden können. Welche Anzahlen ergeben sich für  $\ell = 1, \dots, 5$ ?

**Aufgabe 4.3** Man beweise die folgende Verschärfung des Satzes von Feinstein (Satz 4.1):

Sei  $\{X_n\}_{n \in \mathbb{N}}$  eine diskrete, gedächtnislose Quelle mit Entropie  $H = H(X)$  und Alphabet  $\mathcal{X} = \{x_1, \dots, x_m\}$ . Dann gilt  $\forall \varepsilon > 0 \exists N_0 \in \mathbb{N} \exists A > 0 \forall N \geq N_0 \exists T \subset \mathcal{X}^N$ :

- $P((X_1, \dots, X_N) \in T^c) \leq \varepsilon$
- $\forall (a_1, \dots, a_N) \in T: 2^{-NH-A\sqrt{N}} \leq P(X_1 = a_1, \dots, X_N = a_N) \leq 2^{-NH+A\sqrt{N}}$
- $(1 - \varepsilon) 2^{NH-A\sqrt{N}} \leq |T| \leq 2^{NH+A\sqrt{N}}$

Hinweise:

1.) Man definiere für eine geeignete Konstante  $K > 0$  und  $p_i = P(X = x_i)$ ,  $q_i = 1 - p_i$ :

$$T = \left\{ (a_1, \dots, a_N) \in \mathcal{X}^N \mid \left| \frac{\sum_{1 \leq j \leq N, a_j = x_i} 1 - Np_i}{\sqrt{N}p_iq_i} \right| \leq K, 1 \leq i \leq m \right\}.$$

2.) Man benutze die Tschebyscheff-Ungleichung:  $\forall \varepsilon > 0$  gilt

$$P(|Y - \mu| > \varepsilon) \leq \frac{1}{\varepsilon^2} E((Y - \mu)^2)$$

für jede (endlich diskrete) Zufallsvariable  $Y$  mit Erwartungswert  $E(Y) = \mu$ .

**Aufgabe 4.4**  $\mathcal{X} = \{x_1, \dots, x_m\}$  sei das Quellalphabet einer diskreten, gedächtnislosen Quelle  $X$ ,  $\mathcal{Y} = \{y_1, \dots, y_d\}$ ,  $d \geq 2$ , ein Kodealphabet und  $g : \mathcal{X} \rightarrow \bigcup_{\ell=1}^{\infty} \mathcal{Y}^{\ell}$  ein eindeutig dekodierbarer Kode. Zeigen Sie:

$\sum_{j=1}^m n_j P(X = x_j) = H(X)/\log d$  gilt genau dann, wenn  $P(X = x_j) = d^{-n_j}$  für alle  $j = 1, \dots, m$  mit  $P(X = x_j) > 0$ .  $n_j$ ,  $j = 1, \dots, m$ , bezeichnet hierbei die Länge des Kodeworts  $g(x_j)$ .

**Aufgabe 4.5** Es sei  $\mathcal{X} = \{x_1, \dots, x_m\}$  das Quellalphabet einer diskreten, gedächtnislosen Quelle.  $\mathcal{Y} = \{y_1, \dots, y_d\}$ ,  $d \geq 2$ , sei ein Kodealphabet und  $g'$  ein Kode mit Kodewortlängen  $n'_1, \dots, n'_m$ , der den Buchstaben  $y_1$  als Trennzeichen benutzt, d.h. jedes Kodewort beginnt mit  $y_1$  und  $y_1$  wird nur als erster Buchstabe verwendet. Zeigen Sie:

- $g'$  ist eindeutig dekodierbar.
- Es existiert ein eindeutig dekodierbarer Kode  $g$  mit Kodewortlängen  $n_1, \dots, n_m$  mit  $n_i \leq n'_i$  für alle  $i = 1, \dots, m$  und  $n_j < n'_j$  für mindestens ein  $j \in \{1, \dots, m\}$ .

**Aufgabe 4.6** Sei  $m \in \mathbb{N}$ ,  $m \geq 2$ , und  $\mathbf{p} = (p_1, \dots, p_m) \in \mathcal{P}_m$  ein stochastischer Vektor mit  $p_i > 0$  für alle  $i = 1, \dots, m$ .  $\mathbf{n}^* = (n_1^*, \dots, n_m^*)$ ,  $n_1^* \leq \dots \leq n_m^*$ , sei eine Lösung von

$$\min \sum_{j=1}^m p_j n_j \quad \text{über} \quad n_1, \dots, n_m \in \mathbb{N} \quad \text{mit} \quad \sum_{j=1}^m 2^{-n_j} \leq 1.$$

Zeigen Sie, daß a)  $\sum_{j=1}^m 2^{-n_j^*} = 1$ , b)  $n_{m-1}^* = n_m^*$ .

**Aufgabe 4.7**  $X$  sei eine diskrete, gedächtnislose Quelle mit Alphabet  $\mathcal{X} = \{x_1, \dots, x_m\}$ . Man zeige: Es existiert stets ein optimaler präfixfreier Kode für  $X$ .

**Aufgabe 4.8** In einer Datenverarbeitungsanlage sollen die Inhalte von verschiedenen Magnetbändern zusammengestellt werden. Die Magnetbänder enthalten jeweils eine geordnete Menge von Daten, die aus einer gewissen Anzahl von Posten besteht. Es stehen 3 Magnetbandeinheiten zur Verfügung. Die Anlage kann in einer einzigen Operation nur 2 Bänder zusammenfügen. Dies geschieht in der folgenden Weise: In 2 Magnetbandeinheiten werden die 2 Bänder, die zusammengefügt werden sollen, und die etwa  $n_1$  und  $n_2$  Posten enthalten, eingelegt. In die letzte Station legt man

ein leeres Band. Die Datenverarbeitungsanlage kann nun die Bänder so bearbeiten, daß sich zum Schluß die gesamte aus  $n_1 + n_2$  bestehende Datenmenge der 2 Bänder, in der richtigen Weise geordnet, auf dem ursprünglich leeren Band befindet. Die Zeit, die eine solche Zusammenfügung erfordert, ist proportional zu  $n_1 + n_2$ . Insgesamt sind nun 20 Bänder mit 200, 190, 185, 170, 168, 150, 150, 150, 145, 120, 120, 120, 110, 90, 80, 60, 20, 10, 4 und 2 Posten gegeben.

Man gebe eine optimale Strategie an, unter der die zum Zusammenfügen aller 20 Bänder benötigte Zeit so klein wie möglich wird.

**Aufgabe 4.9**  $X$  sei eine diskrete, gedächtnislose Quelle mit Alphabet  $\mathcal{X} = \{x_1, \dots, x_m\}$ ,  $m \geq 2$ , und zugehöriger Verteilung  $\mathbf{p} = (\frac{1}{m}, \dots, \frac{1}{m}) \in \mathcal{P}_m$ . Zeigen Sie: Für jeden optimalen binären Kode  $g$  gilt

$$\bar{n}(g) = \lfloor \log_2 m \rfloor + \frac{2}{m} \left( m - 2^{\lfloor \log_2 m \rfloor} \right),$$

wobei  $\lfloor x \rfloor$  die größte ganze Zahl kleiner oder gleich  $x \in \mathbb{R}$  bedeutet.

**Aufgabe 4.10**  $\{X_n\}_{n \in \mathbb{N}}$  sei eine diskrete gedächtnislose Quelle mit Alphabet  $\mathcal{X} = \{x_1, \dots, x_5\}$  und Verteilung  $\mathbf{p} = (0.3, 0.2, 0.2, 0.15, 0.15)$ .

Bestimmen Sie einen optimalen Blockkode für Blöcke der Länge  $N = 2$ . Berechnen Sie die erwartete Kodewortlänge pro Quellbuchstabe und vergleichen Sie mit  $H(X_1)$  und der in Beispiel 4.1 berechneten erwarteten Kodewortlänge bei buchstabenweiser Kodierung.

**Aufgabe 4.11** Zeigen Sie, daß die in Beispiel 4.8 definierte Quelle stationär ist.

**Aufgabe 4.12** Bestimmen Sie die Entropie einer stationären diskreten Markoff-Quelle mit zweielementigem Alphabet, assoziierter Funktion  $f = \text{Identität}$  und Übergangsmatrix ( $0 < \alpha, \beta < 1$ )

$$\mathbf{H} = \begin{pmatrix} \alpha & 1 - \alpha \\ 1 - \beta & \beta \end{pmatrix}.$$

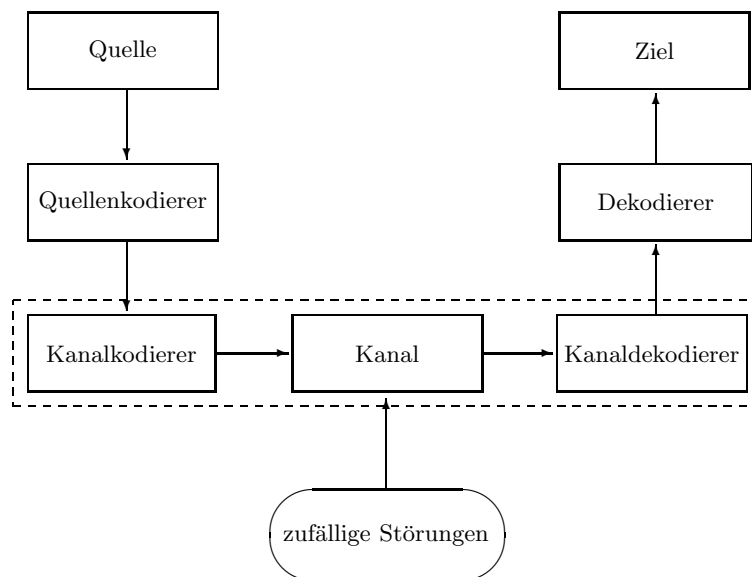
**Aufgabe 4.13**  $\{X_n\}_{n \in \mathbb{N}}$  sei eine Markoff-Quelle mit Alphabet  $\mathcal{X} = \{a, b\}$ , wobei  $X_n = f(Z_n)$  für alle  $n \in \mathbb{N}$  gilt und  $\{Z_n\}_{n \in \mathbb{N}}$  eine homogene Markoff-Kette mit Zustandsraum  $\mathcal{S} = \{s_1, s_2, s_3, s_4\}$ , Übergangsmatrix

$$\mathbf{H} = \begin{pmatrix} 0.2 & 0.8 & 0 & 0 \\ 0 & 0 & 0.1 & 0.9 \\ 0 & 0 & 0.2 & 0.8 \\ 0.7 & 0.3 & 0 & 0 \end{pmatrix}$$

und assoziierter Funktion  $f: \mathcal{S} \rightarrow \mathcal{X}$  mit  $f(s_1) = f(s_4) = a$ ,  $f(s_2) = f(s_3) = b$  ist. Die Anfangsverteilung der Markoff-Kette sei stationär. Berechnen Sie die Entropie  $H_\infty(X)$  und zeigen Sie, daß  $P(Z_n = z_n \mid X_n = x_n, X_{n-1} = x_{n-1}) \in \{0, 1\}$  für alle  $n \geq 2$ ,  $z_n \in \mathcal{S}$ ,  $x_n, x_{n-1} \in \mathcal{X}$ .

## 5 Diskrete gedächtnislose Kanäle

Mit den im vorigen Kapitel eingeführten Kodierungsverfahren soll eine hohe Datenkompression im Sinn von kurzen erwarteten Kodewortlängen erreicht werden. Bei der Übertragung der kodierten Information in einem gestörten Kanal können jedoch Fehler auftreten, zu deren Kompensation wieder gezielt Redundanz hinzugefügt werden muß. Diese Arbeit erledigt der Quellenkodierer. Er wird zur Übertragung solche Kodewörter aussuchen, die in einem noch zu spezifizierenden Sinn weit auseinanderliegen und auch nach Störung noch korrekt unterscheidbar sind. Auf der Empfängerseite müssen im Kanaldekodierer Algorithmen implementiert werden, die aus den empfangenen, gestörten Wörtern wieder die ursprünglich gesendeten mit kleiner Fehlerwahrscheinlichkeit rekonstruieren.



Im vorliegenden Kapitel werden stochastische Modelle zur Beschreibung der zufälligen Störungen im Kanal sowie des Kanalkodierers und -dekodierers

untersucht. Im oben dargestellten Standardmodell der Informationsübertragung ist der jetzt interessierende Block gestrichelt umrahmt.

## 5.1 Kanalkapazität

Im folgenden sei  $\mathcal{X} = \{x_1, \dots, x_m\}$  das Eingabe- und  $\mathcal{Y} = \{y_1, \dots, y_d\}$  das Ausgabealphabet des Kanals, die im allgemeinen verschieden sein können. Die Wirkungsweise des Kanals bei der Übertragung eines Buchstabens wird beschrieben durch Übertragungswahrscheinlichkeiten  $p(y_j | x_i)$ ,  $1 \leq i \leq m$ ,  $1 \leq j \leq d$ .  $p(y_j | x_i)$  bedeutet die Wahrscheinlichkeit,  $y_j$  zu empfangen, wenn  $x_i$  gesendet wurde.

Bei der Übertragung von Wörtern der Länge  $N$  beschreibt  $p_N(\mathbf{b}_N | \mathbf{a}_N)$  analog die Wahrscheinlichkeit,  $\mathbf{b}_N \in \mathcal{Y}^N$  zu empfangen, wenn  $\mathbf{a}_N \in \mathcal{X}^N$  gesendet wurde. Als gedächtnislos bezeichnen wir einen Kanal, bei dem jeder Ausgabebuchstabe nur vom entsprechenden Eingabebuchstaben, nicht von Vorgängern oder Nachfolgern abhängt und die Übertragung einzelner Symbole unabhängig geschieht. Dies ist die Aussage der folgenden Definition.

**Definition 5.1** (*diskreter gedächtnisloser Kanal*)

$\mathcal{X} = \{x_1, \dots, x_m\}$  und  $\mathcal{Y} = \{y_1, \dots, y_d\}$  seien Eingabe- bzw. Ausgabealphabet. Ein System von bedingten Zehldichten

$$\left\{ p_N(\cdot | \mathbf{a}_N) \mid p_N(\cdot | \mathbf{a}_N) : \mathcal{Y}^N \rightarrow [0, 1], \right. \\ \left. \sum_{\mathbf{b}_N \in \mathcal{Y}^N} p_N(\mathbf{b}_N | \mathbf{a}_N) = 1, \mathbf{a}_N \in \mathcal{X}^N, N \in \mathbb{N} \right\}$$

heißt *diskreter Kanal*. Ein diskreter Kanal heißt *gedächtnislos* (kurz: DMC, englisch: *discrete memoryless channel*), wenn

$$p_N(\mathbf{b}_N | \mathbf{a}_N) = \prod_{i=1}^N p_1(b_i | a_i)$$

für alle  $N \in \mathbb{N}$ ,  $\mathbf{a}_N = (a_1, \dots, a_N) \in \mathcal{X}^N$ ,  $\mathbf{b}_N = (b_1, \dots, b_N) \in \mathcal{Y}^N$ .



Man beachte hierbei, daß ein diskreter gedächtnisloser Kanal bereits eindeutig durch die stochastische Matrix

$$\mathbf{\Pi} = (p_1(y_j | x_i))_{1 \leq i \leq m, 1 \leq j \leq d}$$

beschrieben wird, die sogenannte *Kanalmatrix*.

Oft ist es bequemer, diskrete Kanäle durch Zufallsvariablen statt durch ein System bedingter Zähldichten zu beschreiben. Sei hierzu  $\{(X_n, Y_n)\}_{n \in \mathbb{N}}$  eine Folge von Zufallsvariablen,  $(X_n, Y_n)$  jeweils mit Träger  $\mathcal{X} \times \mathcal{Y}$ . Gilt für alle  $N \in \mathbb{N}$ ,  $(a_1, \dots, a_N) \in \mathcal{X}^N$ ,  $(b_1, \dots, b_N) \in \mathcal{Y}^N$ , daß

$$\begin{aligned} P(Y_1 = b_1, \dots, Y_N = b_N | X_1 = a_1, \dots, X_N = a_N) \\ = \prod_{i=1}^N P(Y_i = b_i | X_i = a_i), \end{aligned} \quad (5.1)$$

so heißt die Folge  $\{(X_n, Y_n)\}_{n \in \mathbb{N}}$  ebenfalls diskreter gedächtnisloser Kanal. Identifiziert man in Definition 5.1  $p_N(\mathbf{b}_N | \mathbf{a}_N) = P(\mathbf{Y}_N = \mathbf{b}_N | \mathbf{X}_N = \mathbf{a}_N)$ , wobei  $\mathbf{X}_N = (X_1, \dots, X_N)$  und  $\mathbf{Y}_N = (Y_1, \dots, Y_N)$ , so ist die Äquivalenz der beiden Definitionen offensichtlich.

Eine zur definierenden Eigenschaft eines diskreten gedächtnislosen Kanals äquivalente Bedingung enthält das folgende Lemma.

**Lemma 5.1**  $\{(X_n, Y_n)\}_{n \in \mathbb{N}}$  ist genau dann ein diskreter gedächtnisloser Kanal, wenn für alle  $1 \leq \ell \leq N \in \mathbb{N}$ ,  $(a_1, \dots, a_N) \in \mathcal{X}^N$ ,  $(b_1, \dots, b_N) \in \mathcal{Y}^N$  gilt

$$\begin{aligned} P(Y_\ell = b_\ell | X_1 = a_1, \dots, X_N = a_N, Y_1 = b_1, \dots, Y_{\ell-1} = b_{\ell-1}) \\ = P(Y_1 = b_\ell | X_1 = a_\ell), \end{aligned}$$

d.h. der Kanal überträgt ohne Gedächtnis, ohne Vorgriff und unabhängig von der Position  $\ell$ .

**Beweis.** Wir zeigen zunächst, daß aus der angegebenen Bedingung die DMC-Eigenschaft folgt. Für alle  $N \in \mathbb{N}$ ,  $\mathbf{a}_N = (a_1, \dots, a_N) \in \mathcal{X}^N$ ,  $\mathbf{b}_N =$

$(b_1, \dots, b_N) \in \mathcal{Y}^N$  gilt

$$\begin{aligned}
& P(\mathbf{Y}_N = \mathbf{b}_N \mid \mathbf{X}_N = \mathbf{a}_N) \\
&= P(Y_N = b_N \mid \mathbf{X}_N = \mathbf{a}_N, \mathbf{Y}_{N-1} = \mathbf{b}_{N-1}) \\
&\quad \cdot \frac{P(\mathbf{Y}_{N-1} = \mathbf{b}_{N-1}, \mathbf{X}_N = \mathbf{a}_N)}{P(\mathbf{X}_N = \mathbf{a}_N)} \\
&= P(Y_1 = b_N \mid X_1 = a_N) P(\mathbf{Y}_{N-1} = \mathbf{b}_{N-1} \mid \mathbf{X}_N = \mathbf{a}_N) \\
&= P(Y_1 = b_N \mid X_1 = a_N) P(Y_1 = b_{N-1} \mid X_1 = a_{N-1}) \\
&\quad \cdot P(\mathbf{Y}_{N-2} = \mathbf{b}_{N-2} \mid \mathbf{X}_N = \mathbf{a}_N) \\
&= \dots = \prod_{i=1}^N P(Y_1 = b_i \mid X_1 = a_i).
\end{aligned}$$

Umgekehrt erhält man aus der DMC-Eigenschaft, daß für alle  $\ell \leq N$ ,  $\mathbf{a}_N \in \mathcal{X}^N$ ,  $\mathbf{b}_N \in \mathcal{Y}^N$

$$\begin{aligned}
& P(Y_\ell = b_\ell \mid \mathbf{X}_N = \mathbf{a}_N, \mathbf{Y}_{\ell-1} = \mathbf{b}_{\ell-1}) \\
&= \frac{P(\mathbf{Y}_\ell = \mathbf{b}_\ell \mid \mathbf{X}_N = \mathbf{a}_N)}{P(\mathbf{Y}_{\ell-1} = \mathbf{b}_{\ell-1} \mid \mathbf{X}_N = \mathbf{a}_N)} \\
&= \frac{\sum_{b_{\ell+1}, \dots, b_N \in \mathcal{Y}} P(\mathbf{Y}_N = \mathbf{b}_N \mid \mathbf{X}_N = \mathbf{a}_N)}{\sum_{b_\ell, \dots, b_N \in \mathcal{Y}} P(\mathbf{Y}_N = \mathbf{b}_N \mid \mathbf{X}_N = \mathbf{a}_N)} \\
&= \frac{\prod_{i=1}^{\ell} P(Y_1 = b_i \mid X_1 = a_i)}{\prod_{i=1}^{\ell-1} P(Y_1 = b_i \mid X_1 = a_i)} \\
&\quad \cdot \frac{\sum_{b_{\ell+1}, \dots, b_N \in \mathcal{Y}} \prod_{i=\ell+1}^N P(Y_1 = b_i \mid X_1 = a_i)}{\sum_{b_\ell, \dots, b_N \in \mathcal{Y}} \prod_{i=\ell}^N P(Y_1 = b_i \mid X_1 = a_i)} \\
&= P(Y_1 = b_\ell \mid X_1 = a_\ell).
\end{aligned}$$

Die letzte Gleichheit folgt, da beide Summen im Zähler und Nenner den Wert 1 haben. Damit ist die behauptete Äquivalenz gezeigt. ■

Offensichtlich erfüllen stochastisch unabhängige, identisch verteilte Zufallsvektoren  $(X_n, Y_n)$ ,  $n \in \mathbb{N}$ , die Bedingungen aus Lemma 5.1, bilden also einen diskreten gedächtnislosen Kanal.

**Beispiel 5.1** (binärer symmetrischer Kanal, BSC)

Anknüpfend an Beispiel 3.1 seien Ein- und Ausgabealphabet beide binär,

$\mathcal{X} = \mathcal{Y} = \{0, 1\}$ . Für die Übertragungswahrscheinlichkeiten einzelner Symbole gelte

$$p_1(0|0) = 1 - \varepsilon = p_1(1|1) \quad \text{und} \quad p_1(1|0) = \varepsilon = p_1(0|1), \quad 0 \leq \varepsilon \leq 1.$$

Im Fall eines gedächtnislosen Kanals folgt für die Übertragungswahrscheinlichkeiten von Symbolpaaren

$$p_2(ij|k\ell) = \begin{cases} (1 - \varepsilon)^2, & \text{falls } (i, j) = (k, \ell) \\ \varepsilon^2, & \text{falls } i \neq k \text{ und } j \neq \ell. \\ \varepsilon(1 - \varepsilon), & \text{sonst} \end{cases}$$

Bei der Modellierung mit stochastisch unabhängigen, identisch verteilten zweidimensionalen Zufallsvariablen  $(X_n, Y_n)$ ,  $n \in \mathbb{N}$ , führt die folgende gemeinsame Verteilung von  $(X_1, Y_1)$  zu dem gleichen Kanalmodell. Mit den Abkürzungen  $p_i = P(X_1 = i)$ ,  $i = 0, 1$ , also  $p_1 = 1 - p_0$ , wird gesetzt

$$\begin{aligned} P(X_1 = 0, Y_1 = 0) &= (1 - \varepsilon)p_0, & P(X_1 = 0, Y_1 = 1) &= \varepsilon p_0, \\ P(X_1 = 1, Y_1 = 1) &= (1 - \varepsilon)p_1, & P(X_1 = 1, Y_1 = 0) &= \varepsilon p_1. \end{aligned}$$

Als Übertragungswahrscheinlichkeiten ergeben sich

$$P(Y_1 = 0 \mid X_1 = 0) = \frac{(1 - \varepsilon)p_0}{(1 - \varepsilon)p_0 + \varepsilon p_0} = 1 - \varepsilon,$$

und analog

$$\begin{aligned} P(Y_1 = 1 \mid X_1 = 1) &= 1 - \varepsilon, \\ P(Y_1 = 1 \mid X_1 = 0) &= P(Y_1 = 0 \mid X_1 = 1) = \varepsilon. \end{aligned}$$

Man sieht, daß verschiedene gemeinsame Verteilungen von  $(X_1, Y_1)$  zu dem gleichen Kanalmodell führen können. ■

Gegeben sei nun ein diskreter gedächtnisloser Kanal  $\{(X_n, Y_n)\}_{n \in \mathbb{N}}$  mit Eingabealphabet  $\mathcal{X} = \{x_1, \dots, x_m\}$  und Ausgabealphabet  $\mathcal{Y} = \{y_1, \dots, y_d\}$ . Wegen (5.1) ist nur das Paar  $(X_1, Y_1)$  von Interesse. Zur Vereinfachung werden die Indizes weggelassen, also  $(X, Y) = (X_1, Y_1)$ . Abkürzend bezeichne  $p(j|i) = p_1(y_j \mid x_i) = P(Y = y_j \mid X = x_i)$ ,  $i = 1, \dots, m$ ,  $j = 1, \dots, d$ .

Die gemeinsame Verteilung von  $(X, Y)$  ist eindeutig bestimmt, wenn eine Inputverteilung  $P^X$  mit  $p_i = P(X = x_i)$  festliegt. Denn für alle  $i, j$  gilt  $P(X = x_i, Y = y_j) = P(Y = y_j | X = x_i)P(X = x_i) = p(j|i) p_i$ .

Die Transinformation oder Synentropie von  $X$  und  $Y$  hat gemäß Definition 3.4 die folgende Gestalt.

$$\begin{aligned} I(X, Y) &= H(Y) - H(Y|X) \\ &= \sum_{i=1}^m \sum_{j=1}^d P(X = x_i, Y = y_j) \log \frac{P(Y = y_j | X = x_i)}{P(Y = y_j)} \\ &= \sum_{i=1}^m \sum_{j=1}^d p_i p(j|i) \log \frac{p(j|i)}{P(Y = y_j)} \\ &= \sum_{i=1}^m \sum_{j=1}^d p_i p(j|i) \log \frac{p(j|i)}{\sum_{\ell=1}^m p_\ell p(j|\ell)} \end{aligned}$$

Die Transinformation ist eine Maßzahl, um wieviel die Unbestimmtheit von  $Y$  im Mittel sinkt, wenn man die Eingabe  $X$  kennt. Ein gegebener Kanal wird nun optimal ausgenutzt, wenn die Inputverteilung so gewählt wird, daß die Transinformation maximal ist. Dies führt zum Begriff der Kanalkapazität. Man beachte, daß in der folgenden Definition die gemeinsame Verteilung von  $(\mathbf{X}_N, \mathbf{Y}_N)$  keine Rolle spielt. Kanalkapazität kann in dieser Bedeutung für beliebige diskrete zweidimensionale Zufallsvariable  $(X, Y)$  definiert werden, sofern die bedingte Verteilung  $P^{Y|X}$  vorgegeben ist.

**Definition 5.2** (Kanalkapazität)

$(X, Y)$  beschreibe einen diskreten gedächtnislosen Kanal mit Kanalmatrix  $\mathbf{\Pi} = (p(j|i))_{1 \leq i \leq m, 1 \leq j \leq d}$ . Die maximale Transinformation von  $X$  und  $Y$  über alle Inputverteilungen  $P^X = (p_1, \dots, p_m) \in \mathcal{P}_m$  heißt Kanalkapazität  $C$ , also

$$\begin{aligned} C &= \max_{(p_1, \dots, p_m) \in \mathcal{P}_m} I(X, Y) \\ &= \max_{(p_1, \dots, p_m) \in \mathcal{P}_m} \sum_{i,j} p_i p(j|i) \log \frac{p(j|i)}{\sum_{\ell} p_\ell p(j|\ell)}. \end{aligned}$$

Auch hier wird die Konvention  $0 \cdot * = 0$  benutzt. Undefinierte Ausdrücke ‘\*’ treten auf, falls  $P(X = x_i) = 0$  oder  $P(Y = y_j) = 0$  und  $P(X = x_i) > 0$ .

Im ersten Fall liegt ein Ausdruck der Form  $0 \cdot \log 0$  vor, im zweiten liefert die Multiplikation mit  $p(j|i) = 0$  den Wert Null. Die kritischen Fälle werden also durch die Konvention erfaßt. Hierdurch ist die Transinformation bei gegebener Kanalmatrix eine stetige Funktion auf der kompakten Menge  $\mathcal{P}_m$ . Das Maximum in der Definition der Kanalkapazität wird also angenommen.

**Beispiel 5.2** (binärer symmetrischer Kanal, BSC)

Die Übertragungswahrscheinlichkeiten seien wie in Beispiel 3.1 und 5.1, also  $p(0|0) = p(1|1) = 1 - \varepsilon$  und  $p(1|0) = p(0|1) = \varepsilon$ ,  $0 \leq \varepsilon \leq 1$ . Es bezeichne  $p_0 = P(X = 0)$  und  $p_1 = P(X = 1)$  die Verteilung von  $X$ .

Zur Berechnung der Transinformation von  $X$  und  $Y$  wird die Identität  $I(X, Y) = H(Y) - H(Y | X)$  benutzt. Offensichtlich ist

$$\begin{aligned} H(Y | X = x_i) &= - \sum_{j=1}^d P(Y = y_j | X = x_i) \log P(Y = y_j | X = x_i) \\ &= -\varepsilon \log \varepsilon - (1 - \varepsilon) \log(1 - \varepsilon) = a \end{aligned}$$

unabhängig von  $x_i$ . Folglich ist die bedingte Entropie

$$H(Y | X) = \sum_{i=1}^m P(X = x_i) H(Y | X = x_i) = a$$

unabhängig von  $P^X$ . Bei der Maximierung von  $I(X, Y)$  braucht also nur  $H(Y)$  berücksichtigt zu werden, so daß zur Bestimmung der Kanalkapazität lediglich das Problem  $\max_{(p_1, p_2) \in \mathcal{P}_2} H(Y)$  zu lösen ist. Nach Satz 3.1 a) ist  $H(Y)$  maximal, wenn  $Y$  gleichverteilt ist, also  $P(Y = 0) = P(Y = 1) = \frac{1}{2}$ . Aus der in Beispiel 5.1 bestimmten gemeinsamen Verteilung von  $X$  und  $Y$  kann leicht die Randverteilung von  $Y$  in Abhängigkeit von  $p_0$ ,  $p_1$  und  $\varepsilon$  berechnet werden. Es gilt

$$P(Y = 0) = (1 - \varepsilon)p_0 + \varepsilon p_1 \quad \text{und} \quad P(Y = 1) = \varepsilon p_0 + (1 - \varepsilon)p_1.$$

$I(X, Y)$  ist also maximal, wenn

$$(1 - \varepsilon)p_0 + \varepsilon p_1 = \frac{1}{2} \quad \text{und} \quad \varepsilon p_0 + (1 - \varepsilon)p_1 = \frac{1}{2}.$$

Eine Lösung dieser Gleichungen ist gegeben durch  $p_0^* = p_1^* = \frac{1}{2}$ . Die Kanalkapazität des BSC wird also für gleichverteiltes  $X$  angenommen. Ihr Wert beträgt nach Beispiel 3.1 mit Logarithmen zur Basis 2

$$C = 1 + (1 - \varepsilon) \log_2(1 - \varepsilon) + \varepsilon \log_2 \varepsilon = 1 - H(\varepsilon, 1 - \varepsilon). \quad (5.2)$$

■

Im folgenden wird benutzt, daß für einen DMC  $\{(X_n, Y_n)\}_{n \in \mathbb{N}}$  die Übertragungswahrscheinlichkeiten  $P(Y_\ell = y_j | X_\ell = x_i)$  für alle  $x_i \in \mathcal{X}$ ,  $y_j \in \mathcal{Y}$  unabhängig von  $\ell$  den gleichen Wert haben. Der Nachweis hiervon ist Inhalt von Übungsaufgabe 5.2.

**Satz 5.1**  $\{(X_n, Y_n)\}_{n \in \mathbb{N}}$  sei ein DMC mit Kanalkapazität  $C$ . Dann gilt für alle  $N \in \mathbb{N}$

$$I(\mathbf{X}_N, \mathbf{Y}_N) \leq \sum_{\ell=1}^N I(X_\ell, Y_\ell) \leq N \cdot C,$$

wobei Gleichheit links gilt, wenn  $X_1, \dots, X_N$  stochastisch unabhängig sind, und Gleichheit rechts, wenn  $P^{X_\ell}$  für alle  $\ell = 1, \dots, N$  eine die Transinformation  $I(X_\ell, Y_\ell)$  maximierende Verteilung ist.

**Beweis.** Mit Lemma 3.2 b) gilt  $I(\mathbf{X}_N, \mathbf{Y}_N) = H(\mathbf{Y}_N) - H(\mathbf{Y}_N | \mathbf{X}_N)$ . Ferner

$$\begin{aligned} & H(\mathbf{Y}_N | \mathbf{X}_N) \\ &= \sum_{\mathbf{a}_N \in \mathcal{X}^N, \mathbf{b}_N \in \mathcal{Y}^N} P(\mathbf{X}_N = \mathbf{a}_N) p_N(\mathbf{b}_N | \mathbf{a}_N) \log \frac{1}{p_N(\mathbf{b}_N | \mathbf{a}_N)} \\ &= \mathbb{E} \left( \log \frac{1}{p_N(\mathbf{Y}_N | \mathbf{X}_N)} \right) = \mathbb{E} \left( \sum_{\ell=1}^N \log \frac{1}{p_1(Y_\ell | X_\ell)} \right) \\ &= \sum_{\ell=1}^N \mathbb{E} \left( \log \frac{1}{p_1(Y_\ell | X_\ell)} \right) = \sum_{\ell=1}^N H(Y_\ell | X_\ell). \end{aligned}$$

Mit Satz 3.1 d) folgt

$$\begin{aligned}
 I(\mathbf{X}_N, \mathbf{Y}_N) &= H(\mathbf{Y}_N) - \sum_{\ell=1}^N H(Y_\ell | X_\ell) \\
 &\leq \sum_{\ell=1}^N (H(Y_\ell) - H(Y_\ell | X_\ell)) \\
 &= \sum_{\ell=1}^N I(X_\ell, Y_\ell) \leq N \cdot C.
 \end{aligned}$$

Gleichheit gilt hierbei genau dann, wenn  $Y_1, \dots, Y_N$  stochastisch unabhängig sind. Hinreichend hierfür ist die stochastische Unabhängigkeit der Zufallsvariablen  $X_1, \dots, X_N$ . In diesem Fall gilt nämlich

$$\begin{aligned}
 P(\mathbf{X}_N = \mathbf{a}_N, \mathbf{Y}_N = \mathbf{b}_N) &= P(\mathbf{Y}_N = \mathbf{b}_N | \mathbf{X}_N = \mathbf{a}_N)P(\mathbf{X}_N = \mathbf{a}_N) \\
 &= \prod_{\ell=1}^N P(Y_\ell = b_\ell | X_\ell = a_\ell)P(X_\ell = a_\ell) \\
 &= \prod_{\ell=1}^N P(X_\ell = a_\ell, Y_\ell = b_\ell)
 \end{aligned}$$

für alle  $a_\ell \in \mathcal{X}$ ,  $b_\ell \in \mathcal{Y}$ , und weiterhin

$$P(\mathbf{Y}_N = \mathbf{b}_N) = \sum_{\mathbf{a}_N \in \mathcal{X}^N} \prod_{\ell=1}^N P(X_\ell = a_\ell, Y_\ell = b_\ell) = \prod_{\ell=1}^N P(Y_\ell = b_\ell)$$

für alle  $b_1, \dots, b_N \in \mathcal{Y}$ .  $Y_1, \dots, Y_N$  sind damit stochastisch unabhängig, woraus Gleichheit in der linken Ungleichung folgt.

Die zweite Bedingung für Gleichheit folgt aus der Definition der Kanalkapazität, wenn man noch beachtet, daß  $P(Y_\ell = y_j | X_\ell = x_i)$  unabhängig von  $\ell \in \mathbb{N}$  ist. ■

## 5.2 Kanaldekodierung

Symbole werden über einen diskreten gedächtnislosen Kanal übertragen, der Eingabealphabet  $\mathcal{X} = \{x_1, \dots, x_m\}$ , Ausgabealphabet  $\mathcal{Y} = \{y_1, \dots, y_d\}$  und Übertragungswahrscheinlichkeiten

$$\begin{aligned} p_N(\mathbf{b}_N | \mathbf{a}_N) &= P(\mathbf{Y}_N = \mathbf{b}_N | \mathbf{X}_N = \mathbf{a}_N) \\ &= \prod_{i=1}^N p_1(b_i | a_i) = \prod_{i=1}^N P(Y_1 = b_i | X_1 = a_i) \end{aligned}$$

für Wörter der Länge  $N$  hat, wobei  $\mathbf{a}_N = (a_1, \dots, a_N) \in \mathcal{X}^N$  und  $\mathbf{b}_N = (b_1, \dots, b_N) \in \mathcal{Y}^N$ .

Der Kanalkodierer hat  $M$  Kodewörter der Länge  $N$  zur Verfügung, die durch den Kanal übertragen und wieder richtig dekodiert werden sollen. Die Menge dieser Eingabewörter wird mit

$$\mathcal{C} = \{\mathbf{c}_1, \dots, \mathbf{c}_M\} \subseteq \mathcal{X}^N$$

bezeichnet. Wird nun  $\mathbf{c}_j \in \mathcal{C}$  in den Kanal eingegeben und übertragen, so erhält man auf der Ausgabeseite ein Wort  $\mathbf{b}_N \in \mathcal{Y}^N$ , das wegen der zufälligen Störungen im allgemeinen nicht mehr mit  $\mathbf{c}_j$  übereinstimmt.

Das Problem lautet dann, das empfangene  $\mathbf{b}_N$  richtig zu dekodieren, d.h. aus der Kenntnis von  $\mathbf{b}_N$  zu entscheiden, welches Kodewort aus  $\mathcal{C}$  gesendet wurde. Vernünftigerweise wird man dasjenige  $\mathbf{c}_j$  dekodieren, welches bei empfangenem  $\mathbf{b}_N$  die größte Wahrscheinlichkeit als Eingabewort aufweist. Dieses Konzept wird in verschiedener Hinsicht konkretisiert. Zunächst formalisieren wir den Begriff einer Dekodierregel.

**Definition 5.3** (*Dekodierregel*)

Eine Partition  $\mathcal{R} = \{R_1, \dots, R_M\}$  von  $\mathcal{Y}^N$  heißt Dekodierregel (englisch: *decoding rule, decision scheme*).

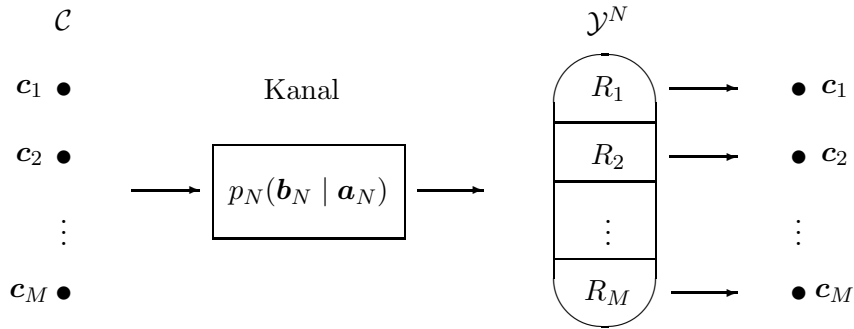
Dabei heißt  $R_1, \dots, R_M \subseteq \mathcal{Y}^N$  Partition von  $\mathcal{Y}^N$ , wenn  $R_1, \dots, R_M$  paarweise disjunkt sind und  $\bigcup_{i=1}^M R_i = \mathcal{Y}^N$ .

Jede Dekodierregel  $\mathcal{R}$  kann äquivalent durch eine Funktion  $h : \mathcal{Y}^N \rightarrow \mathcal{C}$  beschrieben werden, indem

$$h(\mathbf{b}_N) = \mathbf{c}_j, \text{ falls } \mathbf{b}_N \in R_j, j = 1, \dots, m \quad (5.3)$$



gesetzt wird. Graphisch ergibt sich das folgende Bild.



Dekodierregeln werden wie folgt angewendet. Ist  $\mathbf{b}_N \in R_j$  das empfangene Wort, so dekodiert man  $\mathbf{c}_j \in \mathcal{C}$ , entscheidet also, daß  $\mathbf{c}_j \in \mathcal{C}$  gesendet wurde. Gute Dekodierregeln sind offensichtlich solche, die kleine Fehlerwahrscheinlichkeiten besitzen.

**Definition 5.4** (ME-Dekodierung)  
 Eine Dekodierregel  $R_1, \dots, R_M$  mit

$$\mathbf{b}_N \in R_j \Rightarrow P(\mathbf{X}_N = \mathbf{c}_j | \mathbf{Y}_N = \mathbf{b}_N) \geq P(\mathbf{X}_N = \mathbf{c}_i | \mathbf{Y}_N = \mathbf{b}_N)$$

für alle  $i = 1, \dots, M$  und alle  $\mathbf{b}_N \in \mathcal{Y}^N$  mit  $P(\mathbf{Y}_N = \mathbf{b}_N) > 0$  heißt ME-Dekodierung (englisch: minimum error rule, ideal observer).

Bei ME-Dekodierung wird zugunsten eines Kodeworts  $\mathbf{c}_j$  entschieden, das maximale Wahrscheinlichkeit hat, gesendet worden zu sein, wenn  $\mathbf{b}_N$  empfangen wurde.  $R_j$  enthält also nur solche Wörter  $\mathbf{b}_N$ , für die obige Wahrscheinlichkeit maximal ist.

Für eine gegebene Kodewortmenge  $\mathcal{C}$ , Dekodierregel  $\mathcal{R} = \{R_1, \dots, R_M\}$  und Inputverteilung  $\mathbf{p} = (p_1, \dots, p_m)$  mit  $p_j = P(\mathbf{X}_N = \mathbf{c}_j)$ ,  $j = 1, \dots, m$ ,  $\sum_{j=1}^m p_j = 1$  bezeichne

$$e_j = P(\mathbf{Y}_N \notin R_j | \mathbf{X}_N = \mathbf{c}_j) \tag{5.4}$$

die Wahrscheinlichkeit für eine Fehldekodierung, wenn  $\mathbf{c}_j$  gesendet wurde, und

$$e = \sum_{j=1}^M e_j p_j \quad (5.5)$$

die Wahrscheinlichkeit für einen Dekodierfehler. ME-Dekodierungen sind in folgendem Sinn optimal.

**Satz 5.2** *Bei fester Inputverteilung  $\mathbf{p} = (p_1, \dots, p_m)$  minimieren ME-Dekodierungen unter allen Dekodierregeln die Fehlerwahrscheinlichkeit  $e$  aus (5.5).*

**Beweis.** Sei  $\mathcal{R} = \{R_1, \dots, R_M\}$  eine ME-Dekodierung. Dann gilt

$$\begin{aligned} e &= \sum_{j=1}^M P(\mathbf{Y}_N \notin R_j \mid \mathbf{X}_N = \mathbf{c}_j) P(\mathbf{X}_N = \mathbf{c}_j) \\ &= \sum_{j=1}^M P(\mathbf{Y}_N \notin R_j, \mathbf{X}_N = \mathbf{c}_j) \\ &= \sum_{j=1}^M \left( P(\mathbf{X}_N = \mathbf{c}_j) - P(\mathbf{Y}_N \in R_j, \mathbf{X}_N = \mathbf{c}_j) \right) \\ &= 1 - \sum_{j=1}^M \sum_{\mathbf{b}_N \in R_j} P(\mathbf{Y}_N = \mathbf{b}_N, \mathbf{X}_N = \mathbf{c}_j) \\ &= 1 - \sum_{j=1}^M \sum_{\mathbf{b}_N \in R_j} P(\mathbf{X}_N = \mathbf{c}_j \mid \mathbf{Y}_N = \mathbf{b}_N) P(\mathbf{Y}_N = \mathbf{b}_N) \\ &\leq 1 - \sum_{j=1}^M \sum_{\mathbf{b}_N \in S_j} P(\mathbf{X}_N = \mathbf{c}_j \mid \mathbf{Y}_N = \mathbf{b}_N) P(\mathbf{Y}_N = \mathbf{b}_N) \end{aligned}$$

für alle Dekodierregeln  $\mathcal{S} = \{S_1, \dots, S_M\}$ . Obige Ungleichung gilt, da  $P(\mathbf{X}_N = \mathbf{c}_j \mid \mathbf{Y}_N = \mathbf{b}_N) = \max_{i=1, \dots, M} P(\mathbf{X}_N = \mathbf{c}_i \mid \mathbf{Y}_N = \mathbf{b}_N)$  für alle  $\mathbf{b}_N$  mit  $P(\mathbf{Y}_N = \mathbf{b}_N) > 0$ . Durch Rückwärtsumformen mit  $S_j$  anstelle von  $R_j$  entlang obiger Gleichungen erhält man den Dekodierfehler bezüglich der Dekodierregel  $\mathcal{S}$ . ■

Ein Nachteil der ME-Dekodierregel ist, daß sie von der Inputverteilung  $\mathbf{p}$  abhängt, die aber oft nicht bekannt ist. Selbst wenn sie bekannt ist, hängt bei Anwendung der ME-Dekodierung die Kanaldekodierung von stochastischen Eigenschaften der Quelle ab, was bei der Verwendung eines Kanals zur Übertragung von Signalen aus verschiedenen Quellen nicht erwünscht ist. Diese Nachteile werden durch die folgende Dekodierregel vermieden.

**Definition 5.5** (*ML-Dekodierung*)

Eine Dekodierregel  $\mathcal{R} = \{R_1, \dots, R_M\}$  mit

$$\mathbf{b}_N \in R_j \Rightarrow P(\mathbf{Y}_N = \mathbf{b}_N | \mathbf{X}_N = \mathbf{c}_j) \geq P(\mathbf{Y}_N = \mathbf{b}_N | \mathbf{X}_N = \mathbf{c}_i)$$

für alle  $i = 1, \dots, M$ , heißt *ML-Dekodierung* (englisch: *maximum-likelihood decoding rule*).

ML-Dekodierungen entscheiden also bei Empfang von  $\mathbf{b}_N$  zugunsten des  $\mathbf{c}_j$ , bei dessen Sendung der Empfang von  $\mathbf{b}_N$  maximale bedingte Wahrscheinlichkeit hat.

Man beachte, daß im allgemeinen weder ME- noch ML-Dekodierungen eindeutig sind. Wörter  $\mathbf{b}_N$ , bei denen die bedingten Wahrscheinlichkeiten in Definition 5.4 und 5.5 übereinstimmen, können verschiedenen Mengen  $R_j$  zugeordnet werden, ohne die Optimalitätseigenschaften zu verletzen.

Im Fall einer auf  $\mathcal{C}$  gleichverteilten Eingabe  $\mathbf{X}_N$  sind beide Dekodierregeln äquivalent, wie das folgende Lemma zeigt.

**Lemma 5.2** Wenn  $P(\mathbf{X}_N = \mathbf{c}_j) = \frac{1}{M}$  für alle  $j = 1, \dots, M$ , so ist jede ME-Dekodierung eine ML-Dekodierung und umgekehrt, d.h. ME- und ML-Dekodierung sind in diesem Fall äquivalent.

**Beweis.** Für alle  $\mathbf{b}_N \in \mathcal{Y}^N$  mit  $P(\mathbf{Y}_N = \mathbf{b}_N) > 0$  gilt unter der Gleichverteilungsannahme

$$\begin{aligned} P(\mathbf{X}_N = \mathbf{c}_j | \mathbf{Y}_N = \mathbf{b}_N) &= \frac{P(\mathbf{Y}_N = \mathbf{b}_N | \mathbf{X}_N = \mathbf{c}_j)}{P(\mathbf{Y}_N = \mathbf{b}_N) P(\mathbf{X}_N = \mathbf{c}_j)} \\ &= \frac{1}{M P(\mathbf{Y}_N = \mathbf{b}_N)} P(\mathbf{Y}_N = \mathbf{b}_N | \mathbf{X}_N = \mathbf{c}_j). \end{aligned}$$

Hieraus folgt, daß

$$P(\mathbf{X}_N = \mathbf{c}_j | \mathbf{Y}_N = \mathbf{b}_N) \geq P(\mathbf{X}_N = \mathbf{c}_i | \mathbf{Y}_N = \mathbf{b}_N)$$

für alle  $i = 1, \dots, M$  genau dann, wenn

$$P(\mathbf{Y}_N = \mathbf{b}_N \mid \mathbf{X}_N = \mathbf{c}_j) \geq P(\mathbf{Y}_N = \mathbf{b}_N \mid \mathbf{X}_N = \mathbf{c}_i)$$

für alle  $i = 1, \dots, M$ . Die Behauptung folgt jetzt aus den Definitionen von ME- und ML-Dekodierung. ■

Im weiteren wird die sogenannte Hamming-Distanz als Abstands begriff zwischen Wörtern benötigt.

**Definition 5.6** (*Hamming-Distanz*)

Seien  $\mathbf{b}_N = (b_1, \dots, b_N)$ ,  $\mathbf{b}'_N = (b'_1, \dots, b'_N) \in \mathcal{Y}^N$ . Die Anzahl verschiedener Komponenten von  $\mathbf{b}_N$  und  $\mathbf{b}'_N$ ,

$$d(\mathbf{b}_N, \mathbf{b}'_N) = |\{i \mid b_i \neq b'_i, i = 1, \dots, N\}|,$$

heißt Hamming-Distanz zwischen  $\mathbf{b}_N$  und  $\mathbf{b}'_N$ .

Die Hamming-Distanz ist eine Metrik auf  $\mathcal{Y}^N$ , denn sie genügt für alle  $\mathbf{a}_N, \mathbf{b}_N, \mathbf{c}_N \in \mathcal{Y}^N$  den Forderungen (i)  $d(\mathbf{a}_N, \mathbf{b}_N) = 0 \Leftrightarrow \mathbf{a}_N = \mathbf{b}_N$ , (ii)  $d(\mathbf{a}_N, \mathbf{b}_N) = d(\mathbf{b}_N, \mathbf{a}_N)$  und (iii)  $d(\mathbf{a}_N, \mathbf{c}_N) \leq d(\mathbf{a}_N, \mathbf{b}_N) + d(\mathbf{b}_N, \mathbf{c}_N)$ . Stimmen nun Eingabe- und Ausgabealphabet eines Kanals überein, kann mit Hilfe dieses Abstandsmaßes auf der Menge der Wörter der Länge  $N$  die folgende einfache Dekodierregel implementiert werden.

**Definition 5.7** (*MD-Dekodierung*)

Ist  $\mathcal{X} = \mathcal{Y}$ , so heißt eine Dekodierregel  $\mathcal{R} = \{R_1, \dots, R_M\}$  MD-Dekodierung, wenn

$$\mathbf{b}_N \in R_j \Rightarrow d(\mathbf{b}_N, \mathbf{c}_j) \leq d(\mathbf{b}_N, \mathbf{c}_i) \text{ für alle } i = 1, \dots, M.$$

Obwohl MD-Dekodierung nicht von den Übertragungswahrscheinlichkeiten des Kanals abhängt, können bei einem binären symmetrischen Kanal ML-Dekodierregeln als MD-Dekodierer implementiert werden.

**Satz 5.3** Für einen binären symmetrischen Kanal mit  $0 < \varepsilon \leq \frac{1}{2}$  sind ML-Dekodierung und MD-Dekodierung äquivalent.

**Beweis.**  $\mathcal{C} = \{\mathbf{c}_1, \dots, \mathbf{c}_M\}$  sei der Eingabekode, wobei  $\mathbf{c}_j = (c_{j1}, \dots, c_{jN})$ ,  $j = 1, \dots, M$ . Es gilt

$$\begin{aligned} P(\mathbf{Y}_N = \mathbf{b}_N \mid \mathbf{X}_N = \mathbf{c}_j) &= \prod_{i=1}^N P(Y_i = b_i \mid X_i = c_{ji}) \\ &= \varepsilon^{d(\mathbf{b}_N, \mathbf{c}_j)} (1 - \varepsilon)^{N - d(\mathbf{b}_N, \mathbf{c}_j)}. \end{aligned}$$

$f_\varepsilon(x) = \varepsilon^x (1 - \varepsilon)^{N-x}$  ist monoton fallend in  $x \in [0, N]$ , falls  $0 < \varepsilon \leq \frac{1}{2}$ , da  $f'_\varepsilon(x) = \varepsilon^x (1 - \varepsilon)^{N-x} (\ln \varepsilon - \ln(1 - \varepsilon)) \leq 0$ . Es folgt  $d(\mathbf{b}_N, \mathbf{c}_j) \leq d(\mathbf{b}_N, \mathbf{c}_i)$  für alle  $i = 1, \dots, M$  genau dann, wenn  $P(\mathbf{Y}_N = \mathbf{b}_N \mid \mathbf{X}_N = \mathbf{c}_j) \geq P(\mathbf{Y}_N = \mathbf{b}_N \mid \mathbf{X}_N = \mathbf{c}_i)$  für alle  $i = 1, \dots, M$ . Hieraus folgt die Behauptung. ■

### 5.3 Der Shannonsche Fundamentalsatz

Es ist einsichtig, daß man auch bei einem gestörten Kanal mit kleiner Fehlerwahrscheinlichkeit übertragen kann, wenn ein hoher Aufwand zur Fehlerkorrektur betrieben wird. Man könnte etwa jedes Symbol  $n$ -mal hintereinander wiederholen, und der Empfänger entscheidet sich für das Symbol, das er am häufigsten dekodiert hat. Offensichtlich wird bei diesem naiven Verfahren der Durchsatz mit sinkender Fehlerrate immer geringer. Generell könnte dies Anlaß zu der Vermutung geben, daß bei einem verrauschten Kanal die Zuverlässigkeit der Übertragung mit steigendem Aufwand für fehlerkorrigierende Symbole erkauft werden muß.

Daß dies nicht so ist, zeigt der in diesem Abschnitt behandelte Shannonsche Fundamentalsatz. In der englischsprachigen Literatur findet er sich unter dem Stichwort “Fundamental Theorem of Information Theory” oder “Noisy Coding Theorem”. Entgegen der obigen Vermutung zeigt dieses wichtige Ergebnis, daß eine Übertragung mit beliebig kleiner Fehlerwahrscheinlichkeit ohne Reduzierung der Datenrate möglich ist, sofern die Datenrate kleiner als die Kanalkapazität (s. Definition 5.2) ist. Um kleine Fehlerwahrscheinlichkeiten zu erreichen, werden Blockcodes eingesetzt, deren Länge allerdings mit fallender Fehlerwahrscheinlichkeit wächst. Die hierdurch bedingte aufwendige Kodierung und Dekodierung sind der Preis, den man für die zuverlässige Übertragung zahlen muß.

Andererseits kann diese Schranke nicht überschritten werden. Die Umkehrung des Fundamentalsatzes besagt, daß bei Datenraten oberhalb der Kapazität und einer Gleichverteilung als Inputverteilung stets mit positiver Wahrscheinlichkeit Fehler auftreten, egal wie ausgeklügelt die Kodier- und Dekodierregeln sind.

Wie bisher sei  $\mathcal{C} = \{\mathbf{c}_1, \dots, \mathbf{c}_M\} \subseteq \mathcal{X}^N$  ein Eingabekode aus  $M$  Wörtern der Länge  $N$ . Zugehörige Dekodierregeln werden mit  $\mathcal{R} = \{R_1, \dots, R_M\} \subseteq \mathcal{Y}^N$  bezeichnet, die in der Regel von dem verwendeten Eingabekode abhängen, also  $R_j = R_j(\mathcal{C}) = R_j(\mathbf{c}_1, \dots, \mathbf{c}_M)$ .

Zur Beurteilung der Güte eines Kodes werden die Fehlerwahrscheinlichkeiten (5.4) und (5.5) herangezogen. Mit den folgenden Notationen wird die Abhängigkeit von  $\mathcal{C}$  betont. Es bezeichne

$$e_j(\mathcal{C}) = P(\mathbf{Y}_N \notin R_j(\mathcal{C}) \mid \mathbf{X}_N = \mathbf{c}_j) \quad (5.6)$$

die Wahrscheinlichkeit für eine Fehldekodierung, wenn  $\mathbf{c}_j$  gesendet wurde,

$$e(\mathcal{C}) = \sum_{j=1}^M e_j(\mathcal{C}) P(\mathbf{X}_N = \mathbf{c}_j) \quad (5.7)$$

die Wahrscheinlichkeit für eine Fehldekodierung und

$$\bar{e}(\mathcal{C}) = \frac{1}{M} \sum_{j=1}^M e_j(\mathcal{C}) \quad (5.8)$$

die mittlere Wahrscheinlichkeit für eine Fehldekodierung. Diese stimmt mit (5.7) überein, wenn die Inputverteilung eine Gleichverteilung auf  $\mathcal{C}$  ist. Schließlich bezeichnet

$$\hat{e}(\mathcal{C}) = \max_{1 \leq j \leq M} e_j(\mathcal{C}) \quad (5.9)$$

die maximale Wahrscheinlichkeit für Fehldekodierung.

Offensichtlich gilt  $e(\mathcal{C}) \leq \hat{e}(\mathcal{C})$  für jede Inputverteilung  $P^{\mathbf{X}_N}$  auf  $\mathcal{C}$ , so daß mit jeder oberen Schranke für  $\hat{e}(\mathcal{C})$  auch eine obere Schranke für  $e(\mathcal{C})$  unabhängig von der Inputverteilung gewonnen wird.

Um das Verständnis des folgenden allgemeinen Fundamentalsatzes zu erleichtern, formulieren wir vorab eine spezielle Version für den binären symmetrischen Kanal und interpretieren die Aussage in einem nachfolgenden Beispiel. Verwendet werden hierbei Logarithmen zur Basis 2. Man beachte, daß dann für die Kanalkapazität  $C \leq 1$  gilt (s. Beispiel 5.2). Es werden Eingabekodes mit höchstens  $2^{CN}$  Wörtern der Länge  $N$  verwendet, wobei  $2^{CN} \leq 2^N$ , der Anzahl aller Eingabewörter der Länge  $N$ .

**Satz 5.4** *Gegeben sei ein binärer symmetrischer Kanal mit Kapazität  $C$ . Seien  $R$  eine Konstante mit  $0 < R < C$  und  $\varepsilon > 0$ . Für  $M_N = \lfloor 2^{RN} \rfloor$  existiert eine Folge von Kodes  $\mathcal{C}_N \subseteq \{0, 1\}^N$  mit  $M_N$  Kodewörtern der Länge  $N$  derart, daß bei ML-Dekodierung  $\hat{e}(\mathcal{C}_N) \leq \varepsilon$  für alle genügend großen  $N$ .*

Die ML-Dekodierung darf hierbei durch eine beliebige andere Dekodierregel ersetzt werden, die einen kleineren maximalen Fehler  $\hat{e}(\mathcal{C}_N)$  liefert.

**Beispiel 5.3** Betrachtet wird ein BSC mit einer Wahrscheinlichkeit für Fehlübertragung von  $\varepsilon = 0.03$ . Die Kanalkapazität beträgt dann nach Gleichung (5.2) in Beispiel 5.2

$$C = 1 - H(\varepsilon, 1 - \varepsilon) = 1 + 0.03 \log_2 0.03 + 0.97 \log_2 0.97 = 0.8056.$$

Wegen Satz 5.4 existieren dann für  $R = 0.75 < C$  für genügend große  $N$   $\lfloor 2^{0.75N} \rfloor$  Kodewörter der Länge  $N$  aus  $\mathcal{X}^N$  mit  $\hat{e}(\mathcal{C}_N) \leq \varepsilon$ . Konkret ergeben sich die in der folgenden Tabelle angegebenen Werte. In der letzten Zeile ist der prozentuale Anteil der für den Kode  $\mathcal{C}_n$  benutzten Kodewörter aus der Menge aller Kodewörter der Länge  $N$  angegeben. Man sieht, daß dieser rapide fällt.

$N$	10	20	30
$ \mathcal{X}^N  = 2^N$	1024	1048576	$1.0737 \cdot 10^9$
$\lfloor 2^{0.75N} \rfloor$	181	32768	$5.931 \cdot 10^6$
$100 \cdot 2^{0.75N} / 2^N$	17.7 %	3.1 %	0.552 %

Ist der Eingabekode  $\mathcal{C}_N$  bekannt, wird Satz 5.4 wie folgt umgesetzt. Der Nachrichtenstrom aus Nullen und Einsen wird in Blöcke der Länge  $K =$

$0.75N$  geteilt, wobei im folgenden zur Vereinfachung  $0.75N \in \mathbb{N}$  angenommen wird. Für jeden dieser Blöcke steht dann eines der  $2^K = 2^{0.75N}$  Kodewörter der Länge  $N$  zur Verfügung, die bei genügend großem  $N$  mit beliebig kleiner Fehlerwahrscheinlichkeit übertragen werden können.

Diese Aussage kann auch in zeitbezogenen Einheiten interpretiert werden. Die Zeit wird durch die Forderung normiert, daß der Kanal 1 Symbol pro Zeiteinheit (ZE) überträgt. Die angeschlossene Quelle, bzw. der zugehörige Quellenkodierer produziere 0.75 Symbole pro ZE. Wartet man  $K = 0.75N$  Symbole der Quelle ab und kodiert diesen Block durch ein Kodewort der Länge  $N$ , so steht für jedes der  $2^{0.75N}$  möglichen Wörter der Quelle ein Kodewort der Länge  $N$  aus  $\mathcal{C}_N$  zur Verfügung, das mit kleiner Fehlerwahrscheinlichkeit übertragen werden kann. Quelle und Kanal können also synchron arbeiten.

$R$  kann maximal in der Größenordnung von 0.8 gewählt werden. Mit der zeit-synchronen Kopplung von Quelle und Kanal bedeutet dies, daß ein Übertragungsfehler von 3 % eine Reduktion der Datenrate auf 80 % erforderlich macht, um zuverlässige Übertragung erzielen zu können. Eine analoge Rechnung zeigt, daß bei einem Übertragungsfehler von 1 % nur 91.9 % der maximalen Datenrate zur Verfügung steht.

Der Preis für die kleine Fehlerwahrscheinlichkeit ist das für großes  $N$  komplizierte und zeitaufwendige Kodier- bzw. Dekodierverfahren. Wie die zugehörigen Codes  $\mathcal{C}_N$  konkret aussehen, ist nicht bekannt, da bisher vorliegende Beweise des Fundamentalsatzes nicht konstruktiv sind. Die wesentliche Aussage des Fundamentalsatzes ist, daß es solche Codes gibt. ■

Im folgenden wird ein diskreter Kanal  $\{(X_n, Y_n)\}_{n \in \mathbb{N}}$  mit Übertragungswahrscheinlichkeiten

$$P(\mathbf{Y}_N = \mathbf{b}_N \mid \mathbf{X}_N = \mathbf{a}_N) = p_N(\mathbf{b}_N \mid \mathbf{a}_N), \quad \mathbf{a}_N \in \mathcal{X}^N, \quad \mathbf{b}_N \in \mathcal{Y}^N \quad (5.10)$$

und Ein- bzw. Ausgabealphabet  $\mathcal{X} = \{x_1, \dots, x_m\}$  bzw.  $\mathcal{Y} = \{y_1, \dots, y_d\}$  betrachtet. Der Beweis des Fundamentalsatzes beruht auf der Herleitung einer exponentiellen Schranke für die maximale Fehlerwahrscheinlichkeit. Als beweistechnisches Hilfsmittel werden Zufallskodes benutzt.

**Definition 5.8** *Stochastisch unabhängige Zufallsvariable  $C_1, \dots, C_M$ , identisch verteilt jeweils mit Träger  $\mathcal{X}^N$ , heißen  $(M, N)$ -Zufallskode.*

Im folgenden wird stets angenommen, daß die Zufallskodes stochastisch unabhängig von den Zufallsvariablen sind, die den Kanal definieren. Durch eine



Realisation  $\mathbf{c}_1, \dots, \mathbf{c}_M$  von  $\mathbf{C}_1, \dots, \mathbf{C}_M$  wird zufällig ein Kode bestimmt, wobei einzelne Kodewörter eventuell mehrfach vorkommen können, wenn auch für großes  $M$  und  $N$  mit kleiner Wahrscheinlichkeit.

ML-Dekodierung bei Zufallskodes  $\mathbf{C}_1, \dots, \mathbf{C}_M: (\Omega, \mathcal{A}, P) \rightarrow (\mathcal{X}^N, \mathfrak{P}(\mathcal{X}^N))$  ist punktweise für  $\omega \in \Omega$  zu verstehen. Die Partition

$$R_1(\mathbf{C}_1(\omega), \dots, \mathbf{C}_M(\omega)), \dots, R_M(\mathbf{C}_1(\omega), \dots, \mathbf{C}_M(\omega))$$

ist wie in Definition 5.5 bei festem  $\omega \in \Omega$  definiert. Meßbarkeitsfragen interessieren in diesem Zusammenhang nicht, die in späteren Beweisen benötigten Wahrscheinlichkeiten sind alle wohldefiniert.

**Lemma 5.3**  $\mathbf{C}_1, \dots, \mathbf{C}_M$  sei ein  $(M, N)$ -Zufallskode. Dann gilt bei ML-Dekodierung für alle  $j = 1, \dots, M$  und  $0 \leq \rho \leq 1$

$$\begin{aligned} \mathbb{E}(e_j(\mathbf{C}_1, \dots, \mathbf{C}_M)) \\ \leq (M-1)^\rho \sum_{\mathbf{b} \in \mathcal{Y}^N} \left( \sum_{\mathbf{a} \in \mathcal{X}^N} P(\mathbf{C}_1 = \mathbf{a}) (p_N(\mathbf{b}|\mathbf{a}))^{1/(1+\rho)} \right)^{1+\rho}. \end{aligned}$$

**Beweis.** Für alle  $\mathbf{c}_1, \dots, \mathbf{c}_M \in \mathcal{X}^N$ ,  $j = 1, \dots, M$  gilt

$$\begin{aligned} e_j(\mathbf{c}_1, \dots, \mathbf{c}_M) &= P(\mathbf{Y}_N \notin R_j(\mathbf{c}_1, \dots, \mathbf{c}_M) \mid \mathbf{X}_N = \mathbf{c}_j) \\ &= \sum_{\mathbf{b} \notin R_j(\mathbf{c}_1, \dots, \mathbf{c}_M)} p_N(\mathbf{b}|\mathbf{c}_j) \\ &= \sum_{\mathbf{b} \in \mathcal{Y}^N} p_N(\mathbf{b}|\mathbf{c}_j) \mathbb{I}_{R_j^c(\mathbf{c}_1, \dots, \mathbf{c}_M)}(\mathbf{b}). \end{aligned}$$

Setzt man nun einen Zufallskode  $\mathbf{C}_1, \dots, \mathbf{C}_M$  ein, folgt unter Verwendung elementarer Rechenregeln für bedingte Erwartungswerte, daß für jedes  $j = 1, \dots, M$

$$\begin{aligned} \mathbb{E}(e_j(\mathbf{C}_1, \dots, \mathbf{C}_M)) \\ &= \sum_{\mathbf{a} \in \mathcal{X}^N} \mathbb{E}(e_j(\mathbf{C}_1, \dots, \mathbf{C}_M) \mid \mathbf{C}_j = \mathbf{a}) P(\mathbf{C}_j = \mathbf{a}) \\ &= \sum_{\mathbf{a} \in \mathcal{X}^N} \mathbb{E}(e_j(\mathbf{C}_1, \dots, \mathbf{a}, \dots, \mathbf{C}_M)) P(\mathbf{C}_j = \mathbf{a}) \end{aligned}$$

$$\begin{aligned}
&= \sum_{\mathbf{a} \in \mathcal{X}^N} \mathbb{E} \left( \sum_{\mathbf{b} \in \mathcal{Y}^N} p_N(\mathbf{b}|\mathbf{a}) \mathbb{I}_{R_j^c(\mathbf{C}_1, \dots, \mathbf{a}, \dots, \mathbf{C}_M)}(\mathbf{b}) \right) P(\mathbf{C}_j = \mathbf{a}) \\
&= \sum_{\mathbf{a} \in \mathcal{X}^N} P(\mathbf{C}_j = \mathbf{a}) \sum_{\mathbf{b} \in \mathcal{Y}^N} p_N(\mathbf{b}|\mathbf{a}) P(R_j(\mathbf{C}_1, \dots, \mathbf{a}, \dots, \mathbf{C}_M) \not\equiv \mathbf{b}),
\end{aligned}$$

wobei  $\mathbf{a}$  jeweils an der  $j$ -ten Stelle steht. Bei Anwendung von ML-Dekodierung gilt punktweise für jedes Argument  $\omega$  des Zufallskodes und alle  $\mathbf{a} \in \mathcal{X}^N$

$$\mathbf{b} \notin R_j(\mathbf{C}_1, \dots, \mathbf{a}, \dots, \mathbf{C}_M) \Rightarrow \exists i \neq j \text{ mit } p_N(\mathbf{b}|\mathbf{a}) \leq p_N(\mathbf{b}|\mathbf{C}_i).$$

Hiermit erhalten wir für alle  $0 \leq \rho \leq 1$  und  $s > 0$  die folgenden Abschätzungen

$$\begin{aligned}
&P(\mathbf{b} \notin R_j(\mathbf{C}_1, \dots, \mathbf{a}, \dots, \mathbf{C}_M)) \\
&\leq P\left(\bigcup_{i \neq j} \{p_N(\mathbf{b}|\mathbf{a}) \leq p_N(\mathbf{b}|\mathbf{C}_i)\}\right) \\
&\leq \left(\sum_{i \neq j} P(p_N(\mathbf{b}|\mathbf{a}) \leq p_N(\mathbf{b}|\mathbf{C}_i))\right)^\rho \\
&\leq \left(\sum_{i \neq j} \sum_{\mathbf{c} \in \mathcal{X}^N, p_N(\mathbf{b}|\mathbf{a}) \leq p_N(\mathbf{b}|\mathbf{c})} P(\mathbf{C}_i = \mathbf{c})\right)^\rho \\
&\leq \left(\sum_{i \neq j} \sum_{\mathbf{c} \in \mathcal{X}^N} P(\mathbf{C}_1 = \mathbf{c}) \left(\frac{p_N(\mathbf{b}|\mathbf{c})}{p_N(\mathbf{b}|\mathbf{a})}\right)^s\right)^\rho \\
&= \left((M-1) \sum_{\mathbf{c} \in \mathcal{X}^N} P(\mathbf{C}_1 = \mathbf{c}) \frac{p_N(\mathbf{b}|\mathbf{c})^s}{p_N(\mathbf{b}|\mathbf{a})^s}\right)^\rho.
\end{aligned}$$

Insgesamt folgt für den Erwartungswert von  $e_j$

$$\begin{aligned}
&\mathbb{E}(e_j(\mathbf{C}_1, \dots, \mathbf{C}_M)) \\
&\leq \sum_{\mathbf{a} \in \mathcal{X}^N} P(\mathbf{C}_j = \mathbf{a}) \sum_{\mathbf{b} \in \mathcal{Y}^N} p_N(\mathbf{b}|\mathbf{a}) \\
&\quad \cdot \left((M-1) \sum_{\mathbf{c} \in \mathcal{X}^N} P(\mathbf{C}_1 = \mathbf{c}) \frac{p_N(\mathbf{b}|\mathbf{c})^s}{p_N(\mathbf{b}|\mathbf{a})^s}\right)^\rho \\
&= (M-1)^\rho \sum_{\mathbf{b} \in \mathcal{Y}^N} \left(\sum_{\mathbf{a} \in \mathcal{X}^N} P(\mathbf{C}_1 = \mathbf{a}) p_N(\mathbf{b}|\mathbf{a})^{1-s\rho}\right) \\
&\quad \cdot \left(\sum_{\mathbf{c} \in \mathcal{X}^N} P(\mathbf{C}_1 = \mathbf{c}) p_N(\mathbf{b}|\mathbf{c})^s\right)^\rho.
\end{aligned}$$

Setzt man nun  $s = 1/(1 + \rho)$ , äquivalent  $1 - s\rho = 1/(1 + \rho)$ , ergibt sich schließlich

$$\begin{aligned} & \mathbb{E}(e_j(\mathbf{C}_1, \dots, \mathbf{C}_M)) \\ & \leq (M - 1)^\rho \sum_{\mathbf{b} \in \mathcal{Y}^N} \left( \sum_{\mathbf{a} \in \mathcal{X}^N} P(\mathbf{C}_1 = \mathbf{a}) (p_N(\mathbf{b}|\mathbf{a}))^{1/(1+\rho)} \right)^{1+\rho}. \end{aligned}$$

■

Im Beweis von Lemma 5.3 wird die Gedächtnislosigkeit des Kanals nicht gebraucht. Die Aussage gilt also für einen beliebigen diskreten Kanal, dessen Übertragungswahrscheinlichkeiten durch ein System bedingter Zähldichten (5.10) gegeben ist. Wird zusätzlich die Gedächtnislosigkeit eingesetzt, erhält man die folgende Abschätzung.

**Lemma 5.4**  $\{(X_n, Y_n)\}_{n \in \mathbb{N}}$  sei ein diskreter gedächtnisloser Kanal und  $(\mathbf{C}_1, \dots, \mathbf{C}_M)$  ein Zufallskode,  $\mathbf{C}_j = (C_{j1}, \dots, C_{jN})$  mit stochastisch unabhängigen, identisch verteilten  $C_{j\ell}$ ,  $j = 1, \dots, M$ ,  $\ell = 1, \dots, N$ . Die Verteilung von  $C_{j\ell}$  auf  $\mathcal{X}$  werde durch den stochastischen Vektor  $\mathbf{q} = (q_1, \dots, q_m) \in \mathcal{P}_m$  beschrieben. Bei ML-Dekodierung gilt dann für alle  $0 \leq \rho \leq 1$ ,  $j = 1, \dots, M$ , daß

$$\begin{aligned} & \mathbb{E}(e_j(\mathbf{C}_1, \dots, \mathbf{C}_M)) \\ & \leq (M - 1)^\rho \left( \sum_{j=1}^d \left( \sum_{i=1}^m q_i p_1(y_j|x_i)^{1/(1+\rho)} \right)^{1+\rho} \right)^N. \end{aligned} \quad (5.11)$$

**Beweis.** Einsetzen in Lemma 5.3 liefert

$$\begin{aligned} & \mathbb{E}(e_j(\mathbf{C}_1, \dots, \mathbf{C}_M)) \\ & \leq (M - 1)^\rho \cdot \\ & \quad \sum_{b_1, \dots, b_N \in \mathcal{Y}} \left( \sum_{a_1, \dots, a_N \in \mathcal{X}} \prod_{\ell=1}^N P(C_{1\ell} = a_\ell) p_1(b_\ell|a_\ell)^{1/(1+\rho)} \right)^{1+\rho} \\ & = (M - 1)^\rho \sum_{b_1, \dots, b_N \in \mathcal{Y}} \prod_{\ell=1}^N \left( \sum_{i=1}^m q_i p_1(b_\ell|x_i)^{1/(1+\rho)} \right)^{1+\rho} \\ & = (M - 1)^\rho \prod_{\ell=1}^N \sum_{j=1}^d \left( \sum_{i=1}^m q_i p_1(y_j|x_i)^{1/(1+\rho)} \right)^{1+\rho}. \end{aligned}$$

Die Summe der betrachteten Wahrscheinlichkeiten ist unabhängig von  $\ell$ , so daß die Behauptung folgt. ■

Die Doppelsumme in (5.11) wird im folgenden mit

$$F(\rho, \mathbf{q}) = \sum_{j=1}^d \left( \sum_{i=1}^m q_i p_1(y_j|x_i)^{1/(1+\rho)} \right)^{1+\rho}$$

abgekürzt. Dann ist

$$\begin{aligned} \mathbb{E}(e_j(\mathbf{C}_1, \dots, \mathbf{C}_M)) &\leq (M-1)^\rho F(\rho, \mathbf{q})^N \\ &= \exp(\rho \ln(M-1) + N \ln F(\rho, \mathbf{q})) \\ &\leq \exp(\rho \ln M + N \ln F(\rho, \mathbf{q})). \end{aligned}$$

Setzt man nun noch

$$G(\rho, \mathbf{q}) = -\ln F(\rho, \mathbf{q}) \quad \text{und} \quad R = (\ln M)/N, \quad (5.12)$$

so folgt für alle  $M, N \in \mathbb{N}$  und  $0 \leq \rho \leq 1$ , daß

$$\mathbb{E}(e_j(\mathbf{C}_1, \dots, \mathbf{C}_M)) \leq \exp(-N(G(\rho, \mathbf{q}) - \rho R)). \quad (5.13)$$

$R$  ist die Quellenrate. Sie ergibt sich in Satz 5.4 durch Auflösen der Gleichung  $M = 2^{RN}$  nach  $R$  mit Logarithmen zur Basis 2.

Um eine möglichst scharfe obere Schranke für den erwarteten Dekodierfehler zu erhalten, wählen wir  $\rho \in [0, 1]$  und  $\mathbf{q} = (q_1, \dots, q_m) \in \mathcal{P}_m$  so, daß die rechte Seite von (5.13) möglichst klein wird. Zu bestimmen ist dann

$$\max_{0 \leq \rho \leq 1} \max_{\mathbf{q} \in \mathcal{P}_m} (G(\rho, \mathbf{q}) - \rho R) = G^*(R). \quad (5.14)$$

Obiges Maximum wird angenommen, da eine stetige Funktion mit kompaktem Definitionsbereich vorliegt.

**Satz 5.5**  $\{(X_n, Y_n)\}_{n \in \mathbb{N}}$  sei ein diskreter gedächtnisloser Kanal. Bei ML-Dekodierung existieren für alle  $N \in \mathbb{N}$  Kodewörter  $\mathbf{c}_1, \dots, \mathbf{c}_M \in \mathcal{X}^N$  so, daß

$$\hat{e}(\mathbf{c}_1, \dots, \mathbf{c}_M) \leq 4e^{-NG^*(R)}.$$

Gilt  $4e^{-NG^*(R)} < 1$ , so existieren paarweise verschiedene  $\mathbf{c}_1, \dots, \mathbf{c}_M$ .

**Beweis.** Wir ersetzen in Lemma 5.4  $M$  durch  $2M$  und benutzen für  $C_{j\ell}$  die (5.14) maximierende Verteilung  $\mathbf{q}^*$ . Für den mittleren erwarteten Dekodierfehler der  $2M$  Wörter gilt dann

$$\frac{1}{2M} \sum_{j=1}^{2M} \mathbb{E}(e_j(\mathbf{C}_1, \dots, \mathbf{C}_{2M})) \leq e^{-NG^*(R)} = e^{-NG^*(\ln 2M/N)}$$

Dann existieren ein Argument  $\omega$  des zugrundeliegenden Wahrscheinlichkeitsraums  $\Omega$  und  $2M$  Kodewörter  $\mathbf{c}_1 = \mathbf{C}_1(\omega), \dots, \mathbf{c}_{2M} = \mathbf{C}_{2M}(\omega) \in \mathcal{X}^N$  mit

$$\frac{1}{2M} \sum_{j=1}^{2M} e_j(\mathbf{c}_1, \dots, \mathbf{c}_{2M}) \leq e^{-NG^*(\ln 2M/N)}. \quad (5.15)$$

Die Existenz von  $\omega \in \Omega$  mit der angegebenen Eigenschaft ist klar. (“Es können nicht alle mehr verdienen als der Durchschnitt.”) Allerdings ist diese Stelle im Beweis hochgradig nicht konstruktiv. Da lediglich die Verteilung des Zufallskodes, nicht aber das punktweise Verhalten spezifiziert ist, besteht keine Möglichkeit, die Realisation  $\mathbf{c}_1, \dots, \mathbf{c}_{2M}$  konstruktiv zu bestimmen. Wir müssen uns hier mit der Existenzaussage begnügen.

Nun entfernt man aus  $\mathbf{c}_1, \dots, \mathbf{c}_{2M}$  genau  $M$  Kodewörter, insbesondere alle diejenigen mit  $e_k(\mathbf{c}_1, \dots, \mathbf{c}_{2M}) \geq 2e^{-NG^*(\ln 2M/N)}$ . Höchstens  $M$  Stück erfüllen diese Bedingung, sonst würde ein Widerspruch zu (5.15) bestehen. Die verbleibenden  $M$  Kodewörter bilden einen Teilkode  $\mathbf{c}_{i_1}, \dots, \mathbf{c}_{i_M}$  mit

$$\begin{aligned} e_j(\mathbf{c}_{i_1}, \dots, \mathbf{c}_{i_M}) &\leq 2 \exp(-NG^*(\ln 2M/N)) \\ &= 2 \exp(-N \max_{0 \leq \rho \leq 1} \max_{\mathbf{q} \in \mathcal{P}_m} \{G(\rho, \mathbf{q}) - \rho \ln 2M/N\}) \\ &= 2 \exp(-N \max_{0 \leq \rho \leq 1} \max_{\mathbf{q} \in \mathcal{P}_m} \{G(\rho, \mathbf{q}) - \rho \ln M/N - \rho \ln 2/N\}) \\ &\leq 2 \exp(-NG^*(\ln M/N) + \ln 2) = 4 \exp(-NG^*(R)) \end{aligned}$$

für alle  $j \in \{i_1, \dots, i_M\}$ .

Die Behauptung, daß paarweise verschiedene Kodewörter existieren, sieht man so ein. Angenommen zwei Kodewörter sind gleich, etwa  $\mathbf{c}_i = \mathbf{c}_j$ . Dann existiert eine ML-Delodierung mit  $R_j = \emptyset$ , und es gilt  $e_j(\mathbf{c}_1, \dots, \mathbf{c}_M) = P(\mathbf{Y}_n \notin R_j \mid \mathbf{X}_N = \mathbf{c}_j) = 1$ , im Widerspruch zu der Annahme, daß die obere Schranke für den maximalen Fehler kleiner als 1 ist. ■

Satz 5.5 liefert eine in  $N$  exponentiell schnell fallende Fehlerschranke, wenn  $G^*(R) > 0$  ist. Wir werden jetzt untersuchen, wann dies der Fall ist. Benötigt werden zwei vorbereitende Lemmata.

**Lemma 5.5** Sei  $g : [0, 1] \rightarrow \mathbb{R}$  eine stetig differenzierbare Funktion mit  $g(0) = 0$ . Dann folgt aus  $R < g'(0)$ , daß  $\max_{0 \leq \rho \leq 1} \{g(\rho) - \rho R\} > 0$ .  $g'(0)$  bedeutet hierbei die rechtsseitige Ableitung im Punkt 0.

**Beweis.** Das Maximum der stetigen Funktion  $g(\rho) - \rho R$  wird auf dem Intervall  $[0, 1]$  angenommen. Angenommen, für alle  $\rho \in [0, 1]$  wäre  $g(\rho) \leq \rho R$ . Sei  $\rho_n \in [0, 1]$ ,  $n \in \mathbb{N}$ , eine Folge mit  $\lim_{n \rightarrow \infty} \rho_n = 0$ . Der Mittelwertsatz liefert die Existenz von Zwischenpunkten  $\xi_n \in [0, \rho_n]$  mit  $g'(\xi_n) = g(\rho_n)/\rho_n \leq R$  für alle  $n \in \mathbb{N}$ . Also ist  $g'(\xi_n) \leq R$ . Aufgrund der Stetigkeit folgt  $g'(0) \leq R$ . Dies ist ein Widerspruch zur Voraussetzung und zeigt die Behauptung. ■

**Lemma 5.6** Für alle  $0 \leq \rho \leq 1$ ,  $\mathbf{q} \in \mathcal{P}_m$  gilt

$$0 \leq \left. \frac{d}{d\rho} G(\rho, \mathbf{q}) \right|_{\rho=0+} \leq C,$$

wobei Gleichheit gilt, wenn  $\mathbf{q}^*$  eine Verteilung ist, die die Transinformation des diskreten gedächtnislosen Kanals maximiert.

**Beweis.** Es gilt

$$\begin{aligned} \left. \frac{d}{d\rho} G(\rho, \mathbf{q}) \right|_{\rho=0+} &= \sum_{i=1}^m \sum_{j=1}^d q_i p_1(y_j|x_i) \ln \frac{p_1(y_j|x_i)}{\sum_{k=1}^m q_k p_1(y_j|x_i)} \\ &= I(X_1, Y_1). \end{aligned}$$

Dies ist die Transinformation des Kanals bei Inputverteilung  $(q_1, \dots, q_m)$ . Nach Definition der Kanalkapazität ist  $\left. \frac{d}{d\rho} G(\rho, \mathbf{q}) \right|_{\rho=0+} \leq C$ , wobei Gleichheit gilt, wenn  $\mathbf{q} = \mathbf{q}^*$ . ■

Insgesamt kann nun wie folgt geschlossen werden. Aus  $R = \ln M/N < C = \left. \frac{d}{d\rho} G(\rho, \mathbf{q}^*) \right|_{\rho=0+}$  folgt mit Lemma 5.5

$$\max_{0 \leq \rho \leq 1} \{G(\rho, \mathbf{q}^*) - \rho R\} > 0,$$

so daß

$$G^*(R) = \max_{\mathbf{q} \in \mathcal{P}_m} \max_{0 \leq \rho \leq 1} \{G(\rho, \mathbf{q}) - \rho R\} \geq \max_{0 \leq \rho \leq 1} \{G(\rho, \mathbf{q}^*) - \rho R\} > 0.$$

Wendet man noch Satz 5.5 an, ergibt sich die Aussage des nun folgenden Fundamentalsatzes.

**Satz 5.6** (*Shannonscher Fundamentalsatz*)

Gegeben sei ein diskreter gedächtnisloser Kanal mit Kapazität  $C$  und eine Konstante  $R$  mit  $0 < R < C$ .  $M_N \in \mathbb{N}$  sei eine Folge mit  $\frac{\log M_N}{N} < R$ . Dann existieren eine Folge von Codes  $\mathcal{C}_N$  mit  $M_N$  Kodewörtern der Länge  $N$  und eine Konstante  $A > 0$  mit

$$\hat{e}(\mathcal{C}_N) \leq 4e^{-NA} \rightarrow 0 \quad (N \rightarrow \infty).$$

Man könnte versuchen, die obige Beweismethode mit Zufallskodes in die Praxis umzusetzen, indem  $2M$  Kodewörter zufällig erzeugt werden und aus diesen  $M$  Kodewörter mit kleiner Fehlerwahrscheinlichkeit ausgesucht werden. Der Beweismethode von Lemma 5.5 entsprechend kann man hoffen, daß hierdurch in der Mehrzahl der Fälle brauchbare Codes entstehen. Dieses Verfahren ist jedoch nicht praktikabel, da zufällig erzeugte Codes in der Regel keine besondere Struktur haben. Diese ist jedoch nötig, um effiziente Kodier- und Dekodieralgorithmen zu entwerfen. Aus diesem Grund sind fehlerkorrigierende Codes insbesondere unter algorithmischen und algebraischen Aspekten entwickelt worden.

Es ist wichtig, alle auftretenden Logarithmen im Fundamentalsatz zur gleichen Basis zu wählen, auch die in der Transformation  $I(X_1, Y_1)$  zur Bestimmung der Kanalkapazität  $C$  auftretenden. Es erweist sich als bequem, Logarithmen zur Basis  $m$  zu verwenden. Die Bedingung  $\log M_N/N < R$  kann dann äquivalent umgeformt werden in  $M_N < m^{NR} < m^N$ , da  $R < C \leq 1$ .  $m^N$  ist die Anzahl aller möglichen Wörter der Länge  $N$  über  $\mathcal{X} = \{x_1, \dots, x_m\}$ . Der Kanal überträgt also höchstens  $m^{NC}$  Kodewörter mit kleiner Fehlerwahrscheinlichkeit bei großen  $N$ . Wie diese Kodewörter zu wählen sind, ist wegen obiger Bemerkungen nicht klar.

Alle Kanaleigenschaften und die Mächtigkeiten der verwendeten Alphabete gehen in einer einzigen Konstanten ein, nämlich der Kanalkapazität  $C$ . Die Konstante  $A$  kann gewählt werden als  $A = G^*(R) = \max_{\rho} \max_{\mathbf{q}} \{G(\rho, \mathbf{q}) - \rho R\}$ . Am Beweis sieht man, daß bereits  $A = \max_{\rho} \{G(\rho, \mathbf{q}^*) - \rho R\}$  ausreicht.

Die Anwendung und Interpretation von Satz 5.6 erfolgt wie in Beispiel 5.3. Die Quelle benutze ein Alphabet der Mächtigkeit  $k$ . Werden dann Blöcke der Länge  $L \leq NR/\log_m k$  mit  $M_N = \lfloor m^{NR} \rfloor$  Kodewörtern der Länge  $N$  kodiert, so existieren Codes mit exponentiell fallender maximaler Fehlerwahrscheinlichkeit  $\hat{e}$ . Zeitsynchrones Übertragen ist möglich, wenn der Kanal ein Symbol pro Zeiteinheit überträgt und die Quelle in  $N$  Zeiteinheiten höchstens  $\lfloor m^{NR} \rfloor$  verschiedene Wörter produzieren kann.

Der Fundamentalsatz 5.6 kann nicht verschärft werden in folgendem Sinn. Gilt  $R > C$ , so existiert keine Folge von Codes mit  $\lfloor m^{NR} \rfloor$  Kodewörtern der Länge  $N$  und der Eigenschaft  $\hat{e}(\mathcal{C}_N) \rightarrow 0$  ( $N \rightarrow \infty$ ), egal welche Dekodierregel eingesetzt wird. Diese Aussage wird im nächsten Abschnitt behandelt.

Bezeichnet  $\mathbf{X}_N$  die Kanaleingabe und  $\mathbf{V}_N$  die Ausgabe des Kanaldekodierers, so tritt ein Fehler auf, wenn  $\mathbf{X}_N \neq \mathbf{V}_N$ . Die zugehörige Fehlerwahrscheinlichkeit  $P(\mathbf{X}_N \neq \mathbf{V}_N)$  entspricht der Fehlerwahrscheinlichkeit (5.7), die explizit die Dekodierregel und den Code  $\mathcal{C}$  einbezieht. Es gilt nämlich

$$\begin{aligned} e(\mathcal{C}) &= \sum_{j=1}^M P(\mathbf{Y}_N \notin R_j(\mathcal{C}) \mid \mathbf{X}_N = \mathbf{c}_j) P(\mathbf{X}_N = \mathbf{c}_j) \\ &= \sum_{j=1}^M P(\mathbf{V}_N \neq \mathbf{c}_j, \mathbf{X}_N = \mathbf{c}_j) = P(\mathbf{V}_N \neq \mathbf{X}_N). \end{aligned}$$

Diese Fehlerwahrscheinlichkeit läßt sich allgemein für beliebige diskrete Zufallsvariable  $\mathbf{U}$  und  $\mathbf{V}$  ohne ein spezielles Kanal- oder Kodiermodell formulieren. In der folgenden Fano-Ungleichung wird die bedingte Entropie mit ihrer Hilfe abgeschätzt.

**Lemma 5.7** (*Fano-Ungleichung*)

Seien  $\mathbf{U}, \mathbf{V}$  diskrete Zufallsvariablen mit Träger  $\{\mathbf{c}_1, \dots, \mathbf{c}_M\}$ ,  $M > 2$ , und

$$p_e = P(\mathbf{U} \neq \mathbf{V}) = \sum_{\mathbf{u}, \mathbf{v} \in \{\mathbf{c}_1, \dots, \mathbf{c}_M\}, \mathbf{u} \neq \mathbf{v}} P(\mathbf{U} = \mathbf{u}, \mathbf{V} = \mathbf{v}).$$

Dann gilt  $H(\mathbf{U} \mid \mathbf{V}) \leq H(p_e, 1 - p_e) + p_e \log(M - 1)$ , wobei  $H(p_e, 1 - p_e)$  die Entropie des stochastischen Vektors  $(p_e, 1 - p_e) \in \mathcal{P}_2$  bezeichnet.



**Beweis.** Es gilt

$$H(\mathbf{U} | \mathbf{V}) = \sum_{\mathbf{u} \neq \mathbf{v}} P(\mathbf{U} = \mathbf{u}, \mathbf{V} = \mathbf{v}) \log \frac{1}{P(\mathbf{U} = \mathbf{u} | \mathbf{V} = \mathbf{v})} \\ + \sum_{\mathbf{u}} P(\mathbf{U} = \mathbf{u}, \mathbf{V} = \mathbf{u}) \log \frac{1}{P(\mathbf{U} = \mathbf{u} | \mathbf{V} = \mathbf{u})}.$$

Mit den Identitäten

$$p_e \log(M-1) = \sum_{\mathbf{u} \neq \mathbf{v}} \log(M-1) P(\mathbf{U} = \mathbf{u}, \mathbf{V} = \mathbf{v}) \quad \text{und} \\ H(p_e, 1-p_e) = -p_e \log p_e - (1-p_e) \log(1-p_e)$$

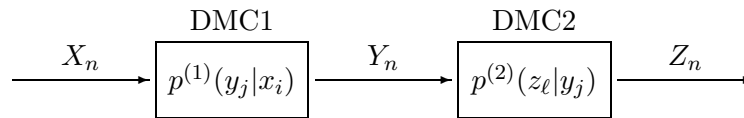
folgt

$$H(\mathbf{U} | \mathbf{V}) - p_e \log(M-1) - H(p_e, 1-p_e) \\ = \sum_{\mathbf{u} \neq \mathbf{v}} P(\mathbf{U} = \mathbf{u}, \mathbf{V} = \mathbf{v}) \log \frac{p_e}{P(\mathbf{U} = \mathbf{u} | \mathbf{V} = \mathbf{v})(M-1)} \\ + \sum_{\mathbf{u}} P(\mathbf{U} = \mathbf{u}, \mathbf{V} = \mathbf{u}) \log \frac{1-p_e}{P(\mathbf{U} = \mathbf{u} | \mathbf{V} = \mathbf{u})} \\ \leq \log e \left( \sum_{\mathbf{u} \neq \mathbf{v}} P(\mathbf{U} = \mathbf{u}, \mathbf{V} = \mathbf{v}) \left( \frac{p_e}{(M-1)P(\mathbf{U} = \mathbf{u} | \mathbf{V} = \mathbf{v})} - 1 \right) \right. \\ \left. + \sum_{\mathbf{u}} P(\mathbf{U} = \mathbf{u}, \mathbf{V} = \mathbf{u}) \left( \frac{1-p_e}{P(\mathbf{U} = \mathbf{u} | \mathbf{V} = \mathbf{u})} - 1 \right) \right) \\ = \log e \left( \frac{p_e}{M-1} \sum_{\mathbf{u} \neq \mathbf{v}} P(\mathbf{V} = \mathbf{v}) - \sum_{\mathbf{u} \neq \mathbf{v}} P(\mathbf{U} = \mathbf{u}, \mathbf{V} = \mathbf{v}) \right. \\ \left. + (1-p_e) \sum_{\mathbf{u}} P(\mathbf{V} = \mathbf{u}) - \sum_{\mathbf{u}} P(\mathbf{U} = \mathbf{u}, \mathbf{V} = \mathbf{u}) \right) \\ = \log e (p_e - p_e + (1-p_e) - (1-p_e)) = 0.$$

Dies zeigt die Behauptung. ■

## 5.4 Kaskadenkanäle und Umkehrung des Fundamentalsatzes

Stimmt das Ausgabealphabet eines diskreten gedächtnislosen Kanals mit dem Eingabealphabet eines zweiten überein, können die beiden Kanäle zu einem *Kaskadenkanal* hintereinandergeschaltet werden. Es ergibt sich das folgende Bild.



Die einzelnen diskreten gedächtnislosen Kanäle werden mit DMC1 und DMC2 bezeichnet, die spezifizierenden Größen sind in der folgenden Übersicht zusammengefaßt.

DMC1:	Eingabealphabet	$\mathcal{X} = \{x_1, \dots, x_m\}$
	Ausgabealphabet	$\mathcal{Y} = \{y_1, \dots, y_d\}$
	Kanalmatrix	$\mathbf{\Pi}^{(1)} = (p^{(1)}(y_j x_i))_{1 \leq i \leq m, 1 \leq j \leq d}$
DMC2:	Eingabealphabet	$\mathcal{Y} = \{y_1, \dots, y_d\}$
	Ausgabealphabet	$\mathcal{Z} = \{z_1, \dots, z_k\}$
	Kanalmatrix	$\mathbf{\Pi}^{(2)} = (p^{(2)}(z_\ell y_j))_{1 \leq j \leq d, 1 \leq \ell \leq k}$

Die Ausgabe von DMC1 ist also die Eingabe in DMC2. Es wird angenommen, daß die Verteilung der Ausgabe von DMC2 nur von der Eingabe in DMC2 (der Ausgabe von DMC1) abhängt, nicht aber von der Eingabe in DMC1. Dieses Verhalten wird durch die folgende Definition beschrieben.

**Definition 5.9** (*Kaskadenkanal*)

Eine Folge von Zufallsvariablen  $\{(X_n, Z_n)\}_{n \in \mathbb{N}}$  heißt *diskreter gedächtnisloser Kaskadenkanal* (kurz: *Kaskaden-DMC*), wenn eine Folge von Zufallsvariablen  $\{Y_n\}_{n \in \mathbb{N}}$  existiert, mit

- a)  $\{(X_n, Y_n)\}_{n \in \mathbb{N}}$  ist ein DMC mit Kanalmatrix  $\mathbf{\Pi}^{(1)}$  und Alphabeten  $\mathcal{X}, \mathcal{Y}$ ,

- b)  $\{(Y_n, Z_n)\}_{n \in \mathbb{N}}$  ist ein DMC mit Kanalmatrix  $\mathbf{\Pi}^{(2)}$  und Alphabeten  $\mathcal{Y}, \mathcal{Z}$ ,  
c) für alle  $N \in \mathbb{N}$ ,  $\mathbf{a}_N \in \mathcal{X}^N$ ,  $\mathbf{b}_N \in \mathcal{Y}^N$ ,  $\mathbf{c}_N \in \mathcal{Z}^N$  mit  $P(\mathbf{X}_N = \mathbf{a}_N, \mathbf{Y}_N = \mathbf{b}_N) > 0$  gilt

$$P(\mathbf{Z}_N = \mathbf{c}_N \mid \mathbf{Y}_N = \mathbf{b}_N, \mathbf{X}_N = \mathbf{a}_N) = P(\mathbf{Z}_N = \mathbf{c}_N \mid \mathbf{Y}_N = \mathbf{b}_N).$$

$\mathbf{Z}_N$  und  $\mathbf{X}_N$  sind bei einem Kaskadenkanal also bedingt stochastisch unabhängig, gegeben  $\mathbf{Y}_N$ .

Es ist noch zu klären, ob ein Kaskaden-DMC wirklich gedächtnislos ist. Insbesondere interessiert dann die Gestalt der gemeinsamen Kanalmatrix  $\mathbf{\Pi}$ .

**Lemma 5.8**  $\{(X_n, Z_n)\}_{n \in \mathbb{N}}$  sei ein Kaskaden-DMC. Für alle Wörter  $\mathbf{a}_N = (a_1, \dots, a_N) \in \mathcal{X}^N$ ,  $\mathbf{c}_N = (c_1, \dots, c_N) \in \mathcal{Z}^N$  gilt dann

$$P(\mathbf{Z}_N = \mathbf{c}_N \mid \mathbf{X}_N = \mathbf{a}_N) = \prod_{i=1}^N P(Z_i = c_i \mid X_i = a_i).$$

Für die Kanalmatrix  $\mathbf{\Pi}$  des Gesamtkanals gilt also  $\mathbf{\Pi} = \mathbf{\Pi}^{(1)} \cdot \mathbf{\Pi}^{(2)}$ , bzw. elementweise  $p(z_\ell | x_i) = \sum_{j=1}^d p^{(1)}(y_j | x_i) \cdot p^{(2)}(z_\ell | y_j)$ ,  $i = 1, \dots, m$ ,  $\ell = 1, \dots, k$ .

**Beweis.** Für  $\mathbf{a}_N \in \mathcal{X}^N$ ,  $\mathbf{c}_N \in \mathcal{Z}^N$  berechnen sich die Übertragungswahrscheinlichkeiten zu

$$\begin{aligned} & P(\mathbf{Z}_N = \mathbf{c}_N \mid \mathbf{X}_N = \mathbf{a}_N) \\ &= \sum_{\mathbf{b}_N \in \mathcal{Y}^N} \frac{P(\mathbf{X}_N = \mathbf{a}_N, \mathbf{Y}_N = \mathbf{b}_N, \mathbf{Z}_N = \mathbf{c}_N)}{P(\mathbf{X}_N = \mathbf{a}_N)} \\ & \quad \frac{P(\mathbf{X}_N = \mathbf{a}_N, \mathbf{Y}_N = \mathbf{b}_N)}{P(\mathbf{X}_N = \mathbf{a}_N, \mathbf{Y}_N = \mathbf{b}_N)} \\ &= \sum_{\mathbf{b}_N \in \mathcal{Y}^N} P(\mathbf{Z}_N = \mathbf{c}_N \mid \mathbf{X}_N = \mathbf{a}_N, \mathbf{Y}_N = \mathbf{b}_N) \cdot \\ & \quad P(\mathbf{Y}_N = \mathbf{b}_N \mid \mathbf{X}_N = \mathbf{a}_N) \\ &= \sum_{b_1, \dots, b_N \in \mathcal{Y}} \prod_{i=1}^N P(Z_i = c_i \mid Y_i = b_i) P(Y_i = b_i \mid X_i = a_i) \end{aligned}$$

$$\begin{aligned}
&= \prod_{i=1}^N \sum_{b \in \mathcal{Y}} P(Z_1 = c_i | Y_1 = b) P(Y_1 = b | X_1 = a_i) & (5.16) \\
&= \prod_{i=1}^N \sum_{b \in \mathcal{Y}} \frac{P(Z_1 = c_i, Y_1 = b, X_1 = a_i)}{P(X_1 = a_i)} \\
&= \prod_{i=1}^N P(Z_1 = c_i | X_1 = a_i).
\end{aligned}$$

Die dritte Gleichheit folgt hierbei aus (iii) in Definition 5.9 und die vierte, da  $\{(X_n, Y_n)\}_{n \in \mathbb{N}}$  und  $\{(Y_n, Z_n)\}_{n \in \mathbb{N}}$  beides diskrete gedächtnislose Kanäle sind.

Die Produktdarstellung der Kanalmatrix  $\Pi$  folgt aus (5.16) mit  $N = 1$ . ■

**Satz 5.7** (*Data Processing Theorem*)

Sei  $\{(X_n, Z_n)\}_{n \in \mathbb{N}}$  ein Kaskaden-DMC. Dann gilt:

- a)  $I((\mathbf{X}_N, \mathbf{Z}_N) | \mathbf{Y}_N) = 0$ ,
- b)  $I(\mathbf{X}_N, \mathbf{Z}_N) \leq \min \{I(\mathbf{X}_N, \mathbf{Y}_N), I(\mathbf{Y}_N, \mathbf{Z}_N)\}$ .

**Beweis.** a)  $\mathbf{Z}_N$  und  $\mathbf{X}_N$  sind bedingt stochastisch unabhängig, gegeben  $\mathbf{Y}_N$ , und daher ist  $H(\mathbf{Z}_N | (\mathbf{X}_N, \mathbf{Y}_N)) = H(\mathbf{Z}_N | \mathbf{Y}_N)$ . Mit Lemma 3.2 b) folgt die Behauptung a) aus

$$I((\mathbf{X}_N, \mathbf{Z}_N) | \mathbf{Y}_N) = H(\mathbf{Z}_N | \mathbf{Y}_N) - H(\mathbf{Z}_N | (\mathbf{X}_N, \mathbf{Y}_N)) = 0.$$

b) Nach Satz 3.1 e) gilt  $H(\mathbf{Z}_N | (\mathbf{X}_N, \mathbf{Y}_N)) \leq H(\mathbf{Z}_N | \mathbf{X}_N)$ . Dies zeigt, daß

$$\begin{aligned}
I(\mathbf{Y}_N, \mathbf{Z}_N) &= H(\mathbf{Z}_N) - H(\mathbf{Z}_N | \mathbf{Y}_N) \\
&= H(\mathbf{Z}_N) - H(\mathbf{Z}_N | (\mathbf{X}_N, \mathbf{Y}_N)) \\
&\geq H(\mathbf{Z}_N) - H(\mathbf{Z}_N | \mathbf{X}_N) = I(\mathbf{X}_N, \mathbf{Z}_N).
\end{aligned}$$

Mit a) und obigem folgt weiterhin

$$\begin{aligned}
&I(\mathbf{X}_N, (\mathbf{Y}_N, \mathbf{Z}_N)) \\
&= H(\mathbf{X}_N) - H(\mathbf{X}_N | (\mathbf{Y}_N, \mathbf{Z}_N)) - H(\mathbf{X}_N | \mathbf{Y}_N) + H(\mathbf{X}_N | \mathbf{Y}_N) \\
&= I(\mathbf{X}_N, \mathbf{Y}_N) + I((\mathbf{X}_N, \mathbf{Z}_N) | \mathbf{Y}_N) = I(\mathbf{X}_N, \mathbf{Y}_N)
\end{aligned}$$

und völlig analog

$$\begin{aligned} I(\mathbf{X}_N, (\mathbf{Y}_N, \mathbf{Z}_N)) &= H(\mathbf{X}_N) - H((\mathbf{X}_N | (\mathbf{Y}_N, \mathbf{Z}_N))) - H(\mathbf{X}_N | \mathbf{Z}_N) + H(\mathbf{X}_N | \mathbf{Z}_N) \\ &= I(\mathbf{X}_N, \mathbf{Z}_N) + I((\mathbf{X}_N, \mathbf{Y}_N) | \mathbf{Z}_N). \end{aligned}$$

Insgesamt erhalten wir

$$I(\mathbf{X}_N, \mathbf{Y}_N) = I(\mathbf{X}_N, \mathbf{Z}_N) + I((\mathbf{X}_N, \mathbf{Y}_N) | \mathbf{Z}_N) \geq I(\mathbf{X}_N, \mathbf{Z}_N).$$

Beide Ungleichungen zusammen verifizieren b). ■

Aussage b) von Satz 5.7 kann wie folgt interpretiert werden. Beim Hintereinanderschalten von zwei DMCs ist die gesamte Transinformation kleiner als jede der beiden einzelnen Transformationen. Verlorene Äquivokation kann also nicht wiederhergestellt werden.

Man beachte, daß im Beweis von Satz 5.7 lediglich Bedingung c) aus Definition 5.9 verwendet wurde. Die Aussagen des Data Processing Theorems gelten also allgemein für endlich diskrete Zufallsvariable  $X, Y, Z$  mit  $P^{Z|(X,Y)} = P^{Z|Y}$ , für die also  $X$  und  $Z$  bedingt stochastisch unabhängig sind, gegeben  $Y$ .

**Korollar 5.1**  $\{(X_n, Z_n)\}_{n \in \mathbb{N}}$  sei ein Kaskaden-DMC.  $C_i$  bezeichne die Kapazität der Teilkanäle mit Übertragungswahrscheinlichkeiten  $p^{(i)}(\cdot | \cdot)$ ,  $i = 1, 2$ , und  $C$  die Gesamtkapazität. Dann gilt  $C \leq \min\{C_1, C_2\}$ .

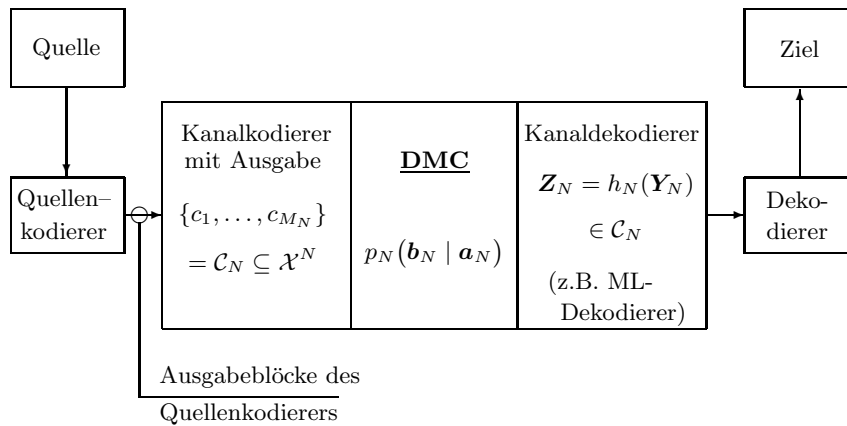
**Beweis.** Da nach Satz 5.7 für alle Verteilungen  $\mathbf{p} \in \mathcal{P}_m$  von  $X_1$ , die zu  $I(X_1, Z_1)$  führen,

$$\begin{aligned} I(X_1, Z_1) \leq I(X_1, Y_1) &\leq \max_{\mathbf{q} \in \mathcal{P}_m, X_1 \sim \mathbf{q}} I(X_1, Y_1) = C_1 \quad \text{und} \\ I(X_1, Z_1) \leq I(Y_1, Z_1) &\leq \max_{\mathbf{q} \in \mathcal{P}_d, Y_1 \sim \mathbf{q}} I(Y_1, Z_1) = C_2 \end{aligned}$$

gilt, folgt  $C = \max_{\mathbf{p} \in \mathcal{P}_m, X_1 \sim \mathbf{p}} I(X_1, Z_1) \leq \min\{C_1, C_2\}$ . ■

Das folgende Bild gibt nun eine detaillierte Darstellung der Übertragung in einem gestörten Kanal. Es werden Blöcke des Ausgabestroms des Quellencodierers mit den Kodewörtern  $\mathcal{C}_N = \{\mathbf{c}_1, \dots, \mathbf{c}_{M_N}\} \subseteq \mathcal{X}^N$  kodiert und

übertragen. Auf der Ausgabeseite des Kanals wird die Kanaldekodierung mit Hilfe einer Funktion  $h_N$  wie in (5.3) dargestellt. Mit dem Index  $N$  wird die Abhängigkeit von  $N$  betont. Mit wachsendem  $N$  steigt im allgemeinen auch die Anzahl der verwendeten Kodewörter. Weiterhin wird im folgenden angenommen, daß Eingabe und Ausgabealphabet des Kanals übereinstimmen.



Typischerweise hat der Quellenkodierer eine Datenkompression vorgenommen, die zumindest approximativ auf der Menge der Ausgabeblöcke bestimmter Länge eine Gleichverteilung liefert. Es ist daher vernünftig, für die Eingabevariable  $\mathbf{X}_N$  des Kanals eine Gleichverteilung auf  $\mathcal{C}_N$  anzunehmen, d.h.  $P(\mathbf{X}_N = \mathbf{c}_j) = 1/M_N$  für alle  $j = 1, \dots, M_N$ . Dieser Fall ist für die Fehlerwahrscheinlichkeit auch der ungünstigste, wie in Lemma 5.9 gezeigt wird.

$\mathbf{Z}_N = h_N(\mathbf{Y}_N) \in \mathcal{C}_N$  bezeichnet die Ausgabe des Kanals nach der Kanaldekodierung bei Verwendung der Dekodierregel  $h_N : \mathcal{Y}^N \rightarrow \mathcal{C}_N$ . Für die Fehlerwahrscheinlichkeiten (5.6), (5.7) und (5.8) gilt

$$\bar{e}(\mathcal{C}_N) = \frac{1}{M_N} \sum_{j=1}^{M_N} e_j(\mathcal{C}_N) \leq \hat{e}(\mathcal{C}_N) = \max_{1 \leq j \leq M_N} e_j(\mathcal{C}_N),$$

wobei  $e_j(\mathcal{C}_N) = P(\mathbf{Z}_N \neq \mathbf{c}_j \mid \mathbf{X}_N = \mathbf{c}_j)$ ,  $j = 1, \dots, M_N$ , die Wahrscheinlichkeit für Fehldekodierung ist, wenn  $\mathbf{c}_j$  gesendet wurde.

**Satz 5.8** (Schwache Umkehrung des Fundamentalsatzes)

$\{(X_n, Y_n)\}_{n \in \mathbb{N}}$  sei ein diskreter gedächtnisloser Kanal mit Kapazität  $C$  und  $R > C$ .  $\mathcal{C}_N$  sei eine Folge von Codes mit  $M_N$  Kodewörtern der Länge  $N$ ,  $M_N \geq m^{NR}$ , und  $\mathbf{X}_N$  sei gleichverteilt auf  $\mathcal{C}_N$  für alle  $N \in \mathbb{N}$ . Dann gilt

$$\liminf_{N \rightarrow \infty} \bar{e}(\mathcal{C}_N) > 0.$$

**Beweis.** Sei  $R = C + \varepsilon$  für ein  $\varepsilon > 0$ . Für die Funktion  $h_N$  gilt

$$\mathbf{Z}_N = h_N(\mathbf{Y}_N) \text{ mit } h_N(\mathbf{b}_N) = \mathbf{c}_j, \text{ falls } \mathbf{b}_N \in R_j, j = 1, \dots, M_N.$$

$R_j$  sind hierbei die Elemente der zur verwendeten Dekodierregel gehörenden Partition.  $\mathbf{X}_N$  und  $\mathbf{Z}_N$  sind offensichtlich bedingt stochastisch unabhängig, gegeben  $\mathbf{Y}_N$ . Da  $\mathbf{X}_N$  auf  $\mathcal{C}_N$  gleichverteilt ist, folgt mit den Sätzen 5.1 und 5.7 b)

$$\begin{aligned} NC &\geq I(\mathbf{X}_N, \mathbf{Y}_N) \geq I(\mathbf{X}_N, \mathbf{Z}_N) \\ &= H(\mathbf{X}_N) - H(\mathbf{X}_N | \mathbf{Z}_N) = \log(M_N) - H(\mathbf{X}_N | \mathbf{Z}_N). \end{aligned}$$

Wegen Lemma 5.7 (Fano-Ungleichung) gilt

$$\begin{aligned} H(\mathbf{X}_N | \mathbf{Z}_N) &\leq H(p_e, 1 - p_e) + p_e \log(M_N - 1) \\ &\leq \log 2 + p_e \log M_N, \end{aligned}$$

wobei  $p_e = P(\mathbf{X}_N \neq \mathbf{Z}_N)$ . Durch Zusammensetzen der beiden Ungleichungen folgt

$$\log M_N \leq NC + H(\mathbf{X}_N | \mathbf{Z}_N) \leq NC + \log 2 + p_e \log M_N.$$

Auflösen nach  $p_e$  liefert

$$p_e \geq 1 - \frac{NC + \log 2}{\log M_N} \geq 1 - \frac{NC + \log 2}{N(C + \varepsilon)} \xrightarrow{(N \rightarrow \infty)} 1 - \frac{C}{C + \varepsilon} > 0.$$

Schließlich gilt wegen der Gleichverteilung von  $\mathbf{X}$

$$\begin{aligned} p_e &= P(\mathbf{X}_N \neq \mathbf{Z}_N) = \sum_i \sum_{j \neq i} P(\mathbf{Z}_N = \mathbf{c}_j, \mathbf{X}_N = \mathbf{c}_i) \\ &= \sum_i P(\mathbf{X}_N = \mathbf{c}_i) \sum_{j \neq i} P(\mathbf{Z}_N = \mathbf{c}_j | \mathbf{X}_N = \mathbf{c}_i) \\ &= \sum_i \frac{1}{M_N} P(\mathbf{Z}_N \neq \mathbf{c}_i | \mathbf{X}_N = \mathbf{c}_i) = \bar{e}(\mathcal{C}_N), \end{aligned}$$

woraus die Behauptung folgt.  $\blacksquare$

Werden also für ein  $\varepsilon > 0$  mehr als  $m^{N(C+\varepsilon)}$  Kodewörter zur Übertragung im Kanal verwendet, so bleibt für beliebig große  $N$  immer noch eine positive mittlere Fehlerwahrscheinlichkeit  $\bar{e}$ .

Unter der Annahme, daß der Kanal ein Symbol pro Zeiteinheit überträgt, hat Satz 5.8 in zeitbezogenen Einheiten folgende Interpretation. Der Quellenkodierer produziere  $\tau_Q$  Symbole pro Zeiteinheit aus einem Alphabet der Mächtigkeit  $k$ . Gilt für ein  $\varepsilon > 0$ , daß die Anzahl der kodierten Wörter der Quelle nach  $N$  Zeiteinheiten  $k^{N\tau_Q}$  größer als  $m^{N(C+\varepsilon)}$ , d.h.  $\tau_Q > \frac{C+\varepsilon}{\log_m k}$ , so bleibt für alle  $N \in \mathbb{N}$  eine positive mittlere Fehlerwahrscheinlichkeit, egal welche Dekodierregel verwendet wird, vorausgesetzt jedes kodierte Quellwort tritt mit derselben Wahrscheinlichkeit auf. Im Fall  $m = k$  reduziert sich die obige Ungleichung auf  $\tau_Q > C + \varepsilon$ .

Man vergleiche dies mit der Aussage des Shannonschen Fundamentalsatzes. Gilt für ein  $\varepsilon > 0$ , daß  $k^{N\tau_Q} < m^{N(C-\varepsilon)}$ , d.h.  $\tau_Q < \frac{C-\varepsilon}{\log_m k}$ , so existiert für genügend große  $N$  ein Kode mit beliebig kleiner maximaler Fehlerwahrscheinlichkeit, unabhängig von der Eingabeverteilung  $P^{\mathbf{X}_N}$ .

Wolfowitz [37] hat 1961 gezeigt, daß unter den Voraussetzungen von Satz 5.8 die stärkere Aussage  $\lim_{N \rightarrow \infty} \bar{e}(\mathcal{C}_N) \rightarrow 1$  gilt. Der hier bewiesene Satz 5.8 wird deshalb als "schwache" Umkehrung des Fundamentalsatzes bezeichnet.

Zum Abschluß des Kapitels wird noch gezeigt, daß eine Gleichverteilung auf  $\mathcal{C}_N$  die ungünstigste Inputverteilung für die Wahrscheinlichkeit einer Fehldekodierung ist. Diese hängt bei fester Kodewortmenge  $\mathcal{C}_N = \{\mathbf{c}_1, \dots, \mathbf{c}_M\}$  von der Inputverteilung  $\mathbf{p} = (p_1, \dots, p_M) \in \mathcal{P}_M$  ab, was durch die Notation

$$e(\mathbf{p}) = \sum_{j=1}^M p_j e_j = P(\mathbf{Z}_N \neq \mathbf{X}_N)$$

ausgedrückt wird, wobei  $e_j = e_j(\mathcal{C}_N)$  in (5.4) definiert ist. Der stochastische Vektor, der aus  $\mathbf{p}$  durch eine Permutation  $\sigma$  der Komponenten entsteht, wird mit

$$\mathbf{p}_\sigma = (p_{\sigma(1)}, \dots, p_{\sigma(M)}) \in \mathcal{P}_M$$

notiert. Zur Abkürzung bezeichne  $\tilde{\mathbf{p}} = (1/M, \dots, 1/M)$  die Gleichverteilung auf  $\mathcal{C}_N$ . Zu jeder Inputverteilung findet man nun eine Permutation der Kodewörter, derart daß die Wahrscheinlichkeit für eine Fehldekodierung kleiner ist als bei einer Gleichverteilung als Inputverteilung.



**Lemma 5.9** Für alle  $\mathbf{p} \in \mathcal{P}_M$  existiert eine Permutation  $\sigma$  mit  $e(\mathbf{p}_\sigma) \leq e(\tilde{\mathbf{p}})$ .

**Beweis.** Für alle  $\mathbf{p} \in \mathcal{P}_M$  gilt

$$\begin{aligned} \min_{\sigma} e(\mathbf{p}_\sigma) &= \min_{\sigma} \sum_{j=1}^M p_{\sigma(j)} e_j = \min_{\sigma} \sum_{j=1}^M p_j e_{\sigma(j)} \\ &\leq \frac{1}{M!} \sum_{\sigma} \sum_{j=1}^M p_j e_{\sigma(j)} = \sum_{j=1}^M p_j \frac{1}{M!} \sum_{\sigma} e_{\sigma(j)} \\ &= \sum_{j=1}^M p_j \frac{1}{M!} (M-1)! \sum_{i=1}^M e_i = \frac{1}{M} \sum_{i=1}^M e_i = e(\tilde{\mathbf{p}}). \end{aligned}$$

■

Man beachte, daß  $\min_{\sigma} e(\mathbf{p}_\sigma)$  für jede Permutation  $\sigma$  angenommen wird, bei der die bedingten Wahrscheinlichkeiten für eine Fehldekodierung  $e_1, \dots, e_M$  und die Komponenten der Inputverteilung  $\mathbf{p}_\sigma$  gegenläufig geordnet sind, etwa  $e_1 \geq \dots \geq e_M$  und  $p_{\sigma(1)} \leq \dots \leq p_{\sigma(M)}$ .

## 5.5 Übungsaufgaben

**Aufgabe 5.1** Zeigen Sie, daß die in Definition 5.6 eingeführte Hamming-Distanz eine Metrik auf  $\mathcal{Y}^N$  ist.

**Aufgabe 5.2**  $\{(X_n, Y_n)\}_{n \in \mathbb{N}}$  sei eine Folge von endlich diskreten Zufallsvariablen jeweils mit Träger  $\mathcal{X} \times \mathcal{Y}$ . Es bezeichne  $p(y|x) = P(Y_1 = y | X_1 = x)$ ,  $x \in \mathcal{X}$ ,  $y \in \mathcal{Y}$ . Zeigen Sie:

Bildet  $\{(X_n, Y_n)\}_{n \in \mathbb{N}}$  einen diskreten, gedächtnislosen Kanal, d.h.

$$P(Y_n = y_n, \dots, Y_1 = y_1 | X_n = x_n, \dots, X_1 = x_1) = \prod_{\ell=1}^n p(y_\ell | x_\ell)$$

für alle  $n \in \mathbb{N}$ ,  $x_1, \dots, x_n \in \mathcal{X}$  und  $y_1, \dots, y_n \in \mathcal{Y}$ , so gilt  $P(Y_n = y | X_n = x) = p(y|x)$  für alle  $n \in \mathbb{N}$ ,  $x \in \mathcal{X}$  und  $y \in \mathcal{Y}$ .

**Aufgabe 5.3** Ein diskreter, gedächtnisloser Kanal mit den Übertragungswahrscheinlichkeiten  $p(j|k)$ ,  $j = 0, 1, \dots, J-1$ ,  $k = 0, 1, \dots, K-1$ , heißt symmetrisch, wenn für alle  $k = 1, \dots, K-1$  eine Permutation  $\pi_k$  auf  $\{0, \dots, J-1\}$  existiert mit  $p(\pi_k(j)|k) = p(j|0)$  und für alle  $j = 1, \dots, J-1$  eine Permutation  $\rho_j$  auf  $\{0, \dots, K-1\}$  existiert mit  $p(j|\rho_j(k)) = p(0|k)$ . Zeigen Sie: Für die Kanalkapazität  $C$  eines symmetrischen Kanals gilt

$$C = \log J + \sum_{j=0}^{J-1} p(j|0) \cdot \log p(j|0).$$

**Aufgabe 5.4**  $\mathbf{II}_i$ ,  $i = 1, 2$ , seien  $K_i \times J_i$ -Matrizen aus Übertragungswahrscheinlichkeiten von zwei diskreten, gedächtnislosen Kanälen mit Kapazität  $C_1$  bzw.  $C_2$ . Man betrachte den Summenkanal aus beiden Kanälen, das ist der diskrete, gedächtnislose Kanal, der durch die  $(K_1 + K_2) \times (J_1 + J_2)$ -Matrix

$$\mathbf{II} = \begin{pmatrix} \mathbf{II}_1 & 0 \\ 0 & \mathbf{II}_2 \end{pmatrix}$$

von Übertragungswahrscheinlichkeiten mit entsprechend zusammengesetzten Alphabeten bestimmt wird. Zeigen Sie:

Für die Kapazität  $C$  des Summenkanals gilt

$$C = \log(a^{C_1} + a^{C_2}),$$

wobei alle auftretenden Logarithmen zur Basis  $a > 1$  gebildet werden.

**Aufgabe 5.5** Die Matrix  $\mathbf{II} = (p(j|k))_{k,j=1,2} = \begin{pmatrix} 1-\varepsilon & \varepsilon \\ \varepsilon & 1-\varepsilon \end{pmatrix}$ ,  $0 \leq \varepsilon < 1/2$ , beschreibe einen diskreten, gedächtnislosen Kanal mit binärem Ein- und Ausgabealphabet  $\mathcal{X} = \mathcal{Y} = \{0, 1\}$ . Bestimmen Sie für  $N = 3$  eine Maximum-Likelihood-Dekodierung für die Wörter  $(0, 0, 0)$ ,  $(0, 1, 0)$ ,  $(1, 1, 1)$ .

**Aufgabe 5.6** Es seien zwei diskrete, gedächtnislose Kanäle  $K_1$  und  $K_2$  mit den Kapazitäten  $C_1$  und  $C_2$  gegeben. Der Produktkanal  $K$  aus beiden Kanälen ist der diskrete, gedächtnislose Kanal, dessen In- und Outputs geordnete Paare  $(x_i, x'_j)$  und  $(y_k, y'_\ell)$  sind, wobei die erste Koordinate zum entsprechenden Alphabet von  $K_1$  und die zweite Koordinate zum Alphabet von  $K_2$  gehört. Die Kanalmatrix des Produktkanals wird dann durch

$$p((y_k, y'_\ell)|(x_i, x'_j)) = p(y_k|x_i) p(y'_\ell|x'_j)$$

festgelegt. Zeigen Sie:

Für die Kapazität  $C$  des Produktkanals  $K$  gilt  $C = C_1 + C_2$ .

**Aufgabe 5.7** Gegeben sei ein diskreter, gedächtnisloser Kanal mit Eingabe  $X$ , Ausgabe  $Y$  und Übertragungswahrscheinlichkeiten  $p(j|k)$ ,  $k = 0, 1, \dots, K-1$ ,  $j = 0, 1, \dots, J-1$ .  $\mathcal{P}_K = \{(q_0, \dots, q_{K-1}) \mid q_k \geq 0, \sum_{k=0}^{K-1} q_k = 1\}$  bezeichne die Menge aller Eingabeverteilungen. Zeigen Sie:

Die Transinformation  $I(X, Y) : \mathcal{P}_K \rightarrow \mathbb{R}$ , definiert durch

$$I(X, Y) = \sum_{j,k} q_k p(j|k) \log \frac{p(j|k)}{\sum_{\ell} q_{\ell} p(j|\ell)},$$

ist eine auf  $\mathcal{P}_K$  konkave Funktion.

**Aufgabe 5.8** Welche Kapazitäten besitzen die durch die folgenden Kanalmatrizen beschriebenen diskreten, gedächtnislosen Kanäle ?

$$\mathbf{II} = \begin{pmatrix} \frac{1}{3} & \frac{1}{3} & \frac{1}{6} & \frac{1}{6} \\ \frac{1}{6} & \frac{1}{6} & \frac{1}{3} & \frac{1}{3} \end{pmatrix}$$

$$\mathbf{II}' = \begin{pmatrix} cae & caf & c\bar{a}e & c\bar{a}f & \bar{c}ae & \bar{c}af & \bar{c}\bar{a}e & \bar{c}\bar{a}f \\ caf & cae & c\bar{a}f & c\bar{a}e & \bar{c}af & \bar{c}ae & \bar{c}\bar{a}f & \bar{c}\bar{a}e \\ cbe & cbf & c\bar{b}e & c\bar{b}f & \bar{c}be & \bar{c}bf & \bar{c}\bar{b}e & \bar{c}\bar{b}f \\ cbf & cbe & c\bar{b}f & c\bar{b}e & \bar{c}bf & \bar{c}be & \bar{c}\bar{b}f & \bar{c}\bar{b}e \\ dae & daf & d\bar{a}e & d\bar{a}f & \bar{d}ae & \bar{d}af & \bar{d}\bar{a}e & \bar{d}\bar{a}f \\ daf & dae & d\bar{a}f & d\bar{a}e & \bar{d}af & \bar{d}ae & \bar{d}\bar{a}f & \bar{d}\bar{a}e \\ dbe & dbf & d\bar{b}e & d\bar{b}f & \bar{d}be & \bar{d}bf & \bar{d}\bar{b}e & \bar{d}\bar{b}f \\ dbf & dbe & d\bar{b}f & d\bar{b}e & \bar{d}bf & \bar{d}be & \bar{d}\bar{b}f & \bar{d}\bar{b}e \end{pmatrix}$$

mit  $a, b, c, d, e, f \in (0, 1)$  mit  $e+f = 1$ ,  $a \neq b$ ,  $c \neq d$  und  $\bar{a} = 1-a$ ,  $\bar{b} = 1-b$ ,  $\bar{c} = 1-c$  sowie  $\bar{d} = 1-d$ .

**Aufgabe 5.9** Gegeben sei ein DMC mit Eingabealphabet  $\mathcal{X} = \{x_1, \dots, x_m\}$ , Ausgabealphabet  $\mathcal{Y} = \{y_1, \dots, y_d\}$ ,  $d \geq 2$ , und Kanalmatrix

$$\mathbf{II} = (p_1(y_j|x_i))_{1 \leq i \leq m, 1 \leq j \leq d}.$$

Es sei  $(\mathbf{C}_1, \dots, \mathbf{C}_M)$  ein  $(M, N)$ -Zufallskode,  $\mathbf{C}_j = (C_{j1}, \dots, C_{jN})$  mit stochastisch unabhängigen, identisch verteilten  $C_{j\ell}$ ,  $j = 1, \dots, M$ ,  $\ell = 1, \dots, N$ . Die Verteilung von  $C_{j\ell}$  werde durch den stochastischen Vektor  $\mathbf{q} = (q_1, \dots, q_m) \in \mathcal{P}_m$  beschrieben. Es wird eine ML-Dekodierung unter der (irrtümlichen) Annahme, daß die Kanalmatrix  $\mathbf{II}^* = (p_1(y_j|x_i))_{1 \leq i \leq m, 1 \leq j \leq d}$  vorliegt, benutzt. Zeigen Sie:

a) Für alle  $j \in \{1, \dots, M\}$  und alle  $0 \leq \rho \leq 1$  gilt

$$\begin{aligned} & \mathbb{E}(e_j(\mathbf{C}_1, \dots, \mathbf{C}_M)) \\ & \leq (M-1)^\rho \sum_{\mathbf{b}_N \in \mathcal{Y}^N} \left( \sum_{\mathbf{a}_N \in \mathcal{X}^N} P(\mathbf{C}_1 = \mathbf{a}_N) \frac{p_N(\mathbf{b}_N | \mathbf{a}_N)}{q_N(\mathbf{b}_N | \mathbf{a}_N)^{\frac{\rho}{1+\rho}}} \right) \cdot \\ & \quad \left( \sum_{\mathbf{c}_N \in \mathcal{X}^N} P(\mathbf{C}_1 = \mathbf{c}_N) q_N(\mathbf{b}_N | \mathbf{c}_N)^{\frac{1}{1+\rho}} \right)^\rho. \end{aligned}$$

b) Gibt es Konstanten  $K_1, K_2 > 0$ , so daß für alle  $\mathbf{a}_N \in \mathcal{X}^N$  und alle  $\mathbf{b}_N \in \mathcal{Y}^N$

$$K_1 p_N(\mathbf{b}_N | \mathbf{a}_N) \leq q_N(\mathbf{b}_N | \mathbf{a}_N) \leq K_2 p_N(\mathbf{b}_N | \mathbf{a}_N)$$

gilt, so existiert eine (nur von  $K_1$  und  $K_2$  abhängende) Konstante  $K$  mit

$$\mathbb{E}(e_j(\mathbf{C}_1, \dots, \mathbf{C}_M)) \leq K \exp\{-N(G(\rho, q) - \rho R)\},$$

wobei  $G(\rho, q)$  und  $R$  wie in (5.12) definiert sind.

c) Bestimmen Sie in der Situation von b) eine geeignete Konstante  $K$ .

**Aufgabe 5.10** Gegeben sei ein binärer symmetrischer Kanal mit Fehlerwahrscheinlichkeit  $\varepsilon$ . Zeichnen Sie die maximale Datenrate einer binären Quelle, bei der noch zuverlässige Übertragung möglich ist, als Funktion von  $\varepsilon$ . Wie verhält sich die Datenrate, wenn die Übertragungswahrscheinlichkeiten  $p(0|1) = \frac{1}{2}p(1|0)$  erfüllen?

## 6 Fehlerkorrigierende Codes

Der Shannonsche Fundamentalsatz sichert die Existenz von Blockcodes mit kleiner Fehlerwahrscheinlichkeit, sagt jedoch nichts über deren Konstruktion aus. Die zum Beweis verwendeten Zufallskodes könnten zwar in der Praxis simulativ erzeugt werden, haben aber den Nachteil, in der Regel keine Struktur zu besitzen, die eine effiziente Kodierung oder Dekodierung erlaubt. In diesem Kapitel werden fehlerrobuste Blockcodes konstruiert, die diese Eigenschaft besitzen. Ziel ist es, eine Menge von Wörtern der Länge  $N$  über einem bestimmten Eingabealphabet derart auszuwählen, daß möglichst viele bei der Übertragung entstandene Fehler wieder korrigiert werden können.

Natürlich können in diesem Kapitel nur einige grundlegende Aspekte der Kodierungstheorie behandelt werden. Weiterführende Darstellungen finden sich zum Beispiel in [2], [17] oder [32]. Hier werden lediglich nach einigen allgemeinen Untersuchungen zur Hamming-Distanz zwei wichtige Typen von Codes näher beleuchtet, zum einen lineare Codes als typische Vertreter eines algebraisch motivierten Zugangs, zum anderen Faltungskodes und der Viterbi-Algorithmus als Beispiel eines algorithmisch begründeten Typs.

Im folgenden wird unter den üblichen Notationen ein Kanal mit identischem Eingabe- und Ausgabealphabet  $\mathcal{X} = \mathcal{Y} = \{x_1, \dots, x_m\}$  betrachtet. Zugehörige Eingabekodes werden mit  $\mathcal{C} = \{\mathbf{c}_1, \dots, \mathbf{c}_M\} \subseteq \mathcal{X}^N$  bezeichnet.

### 6.1 Blockcodes und Hamming-Distanz

Für  $\mathbf{a}, \mathbf{b} \in \mathcal{X}^N$  bezeichnet  $d(\mathbf{a}, \mathbf{b})$  die Hammingdistanz (siehe Definition 5.6). Weiterhin heißt

$$d(\mathcal{C}) = \min_{1 \leq i < j \leq M} d(\mathbf{c}_i, \mathbf{c}_j)$$

minimaler Abstand des Codes  $\mathcal{C}$ . Mit

$$S_e(\mathbf{c}_i) = \{\mathbf{a} \in \mathcal{X}^N \mid d(\mathbf{a}, \mathbf{c}_i) \leq e\}$$

wird die Kugel mit Radius  $e \in \mathbb{N}$  und Mittelpunkt  $\mathbf{c}_i$  bezüglich der Metrik  $d$  bezeichnet.

**Satz 6.1**  $\mathcal{C} = \{\mathbf{c}_1, \dots, \mathbf{c}_M\}$  sei ein Kode mit  $d(\mathcal{C}) = d$ . Dann werden bei MD-Dekodierung bis zu  $\lfloor \frac{1}{2}(d-1) \rfloor$  Fehler korrigiert.

**Beweis.** Sei  $e = \lfloor \frac{1}{2}(d-1) \rfloor$ . Dann gilt  $S_e(\mathbf{c}_i) \cap S_e(\mathbf{c}_j) = \emptyset$  für alle  $\mathbf{c}_i \neq \mathbf{c}_j \in \mathcal{C}$ . Denn angenommen, es existieren Indizes  $i, j$  und  $\mathbf{a} \in \mathcal{X}^N$  mit  $\mathbf{a} \in S_e(\mathbf{c}_i) \cap S_e(\mathbf{c}_j)$ . Dann wäre  $d(\mathbf{c}_i, \mathbf{a}) \leq e$  und  $d(\mathbf{c}_j, \mathbf{a}) \leq e$ , woraus  $d(\mathbf{c}_i, \mathbf{c}_j) \leq d(\mathbf{c}_i, \mathbf{a}) + d(\mathbf{c}_j, \mathbf{a}) \leq 2e = 2\lfloor \frac{1}{2}(d-1) \rfloor \leq d-1$  folgen würde. Dies ist ein Widerspruch zu  $d(\mathcal{C}) = d$ .

$R_1, \dots, R_M$  sei eine MD-Dekodierung. Es gilt  $S_e(\mathbf{c}_i) \subseteq R_i$  für alle  $i = 1, \dots, M$ . Dies zeigt die Behauptung. ■

Ein Kode  $\mathcal{C} = \{\mathbf{c}_1, \dots, \mathbf{c}_M\} \subseteq \mathcal{X}^N$  mit  $d(\mathcal{C}) = d$  heißt  $(N, M, d)$ -Kode. Man beachte, daß  $M$  und  $d$  sich bei festem  $N$  gegenläufig verhalten, d.h. große  $M$  führen zu kleinen  $d$  und umgekehrt. Mit

$$A_m(N, d) = \max \{M \in \mathbb{N} \mid \text{es existiert ein } (N, M, d)\text{-Kode}\}$$

wird die maximale Anzahl von Kodewörtern aus  $\mathcal{X}^N$  bezeichnet, die paarweise den Mindestabstand  $d$  haben. Offensichtlich gilt  $A_m(N, 1) = m^N = |\mathcal{X}^N|$ .

Im allgemeinen ist die Bestimmung von  $A_m(N, d)$  ein schwieriges Problem. Zum Beispiel ist lediglich bekannt, daß  $72 \leq A_2(10, 3) \leq 79$  oder  $144 \leq A_2(11, 3) \leq 158$ .

Eine untere Schranke für  $A_2(3, 2)$  ist 4, denn wie man leicht nachprüft ist die Menge  $\{(0,0,0), (1,1,0), (1,0,1), (0,1,1)\}$  ein  $(3,4,2)$ -Kode. Allgemeine Schranken für  $A_m(N, d)$  können aus den im folgenden Satz angegebenen Ungleichungen durch Auflösen nach  $A_m(N, d)$  hergeleitet werden. Die durch die linke Ungleichung induzierte Schranke heißt *Hamming-Schranke*, die rechte *Gilbert-Varshamov-Schranke*.

**Satz 6.2** (*Hamming- und Gilbert-Varshamov-Schranke*)

Sei  $d = 2e + 1$ ,  $e \in \mathbb{N}$ . Dann gilt

$$A_m(N, d) \sum_{i=0}^e \binom{N}{i} (m-1)^i \leq m^N \leq A_m(N, d) \sum_{i=0}^{d-1} \binom{N}{i} (m-1)^i.$$

**Beweis.** Zum Beweis der linken Ungleichung sei  $\mathcal{C} = \{\mathbf{c}_1, \dots, \mathbf{c}_M\}$  ein Kode mit  $d(\mathcal{C}) = 2e + 1$ . Dann sind  $S_e(\mathbf{c}_1), \dots, S_e(\mathbf{c}_M)$  paarweise disjunkt und es gilt

$$|S_e(\mathbf{c}_j)| = \sum_{i=0}^e \binom{N}{i} (m-1)^i \text{ für alle } j = 1, \dots, M.$$

Folglich ist

$$\left| \bigcup_{j=1}^M S_e(\mathbf{c}_j) \right| = M \sum_{i=0}^e \binom{N}{i} (m-1)^i \leq m^N$$

für alle  $\mathcal{C}$  mit  $d(\mathcal{C}) = 2e + 1$ . Dies zeigt  $A_m(N, d) \sum_{i=0}^e \binom{N}{i} (m-1)^i \leq m^N$  und damit die linke Ungleichung.

Zum Beweis der rechten Ungleichung sei  $\mathcal{C}$  ein  $(N, M, d)$ -Kode mit maximaler Anzahl von Kodewörtern  $M = A_m(N, d)$ . Dann existiert kein  $\mathbf{a} \in \mathcal{X}^N$  so, daß  $d(\mathbf{a}, \mathbf{c}_i) \geq d$  für alle  $i \in \{1, \dots, M\}$  gilt, d.h. für alle  $\mathbf{a} \in \mathcal{X}^N$  existiert ein  $i \in \{1, \dots, M\}$  mit  $d(\mathbf{a}, \mathbf{c}_i) \leq d - 1$ . Es folgt

$$m^N \leq \left| \bigcup_{i=1}^M S_{d-1}(\mathbf{c}_i) \right| \leq A_m(N, d) \sum_{i=0}^{d-1} \binom{N}{i} (m-1)^i. \quad \blacksquare$$

Ein günstiger Fall liegt vor, wenn die gesamte Wortmenge  $\mathcal{X}^N$  mit disjunkten Kugeln gleichen Radius'  $t$  bezüglich der Metrik  $d(\cdot, \cdot)$  ausgeschöpft werden kann. Die zugehörigen Mittelpunkte bilden dann einen Kode mit Minimaldistanz  $2t + 1$ , und die Kugeln bestimmen gleichzeitig die Partition eine MD-Dekodierung. Solche Codes heißen perfekt.

**Definition 6.1** (*perfekter Kode*)

Ein Kode  $\mathcal{C} = \{\mathbf{c}_1, \dots, \mathbf{c}_M\} \subseteq \mathcal{X}^N$  heißt *perfekt*, wenn ein  $t > 0$  existiert, so daß  $S_t(\mathbf{c}_1), \dots, S_t(\mathbf{c}_M)$  eine Partition von  $\mathcal{X}^N$  bildet.

Ein kartesisches Produkt mit Kugeln auszuschöpfen, ist unter der üblichen Euklidischen Metrik unserer Raumvorstellung nicht möglich. Man beachte, daß Hilfsvorstellungen zur Hamming-Distanz auf der Basis von zweidimensionalen Zeichnungen nicht in allen Punkten mit den entsprechenden mathematischen Eigenschaften übereinstimmen.

Aus den bisherigen Untersuchungen lassen sich die nachstehenden Eigenschaften folgern.

**Lemma 6.1**

- a) *Perfekte Kodes korrigieren bis zu  $t$  Fehler pro Wort, nicht aber  $t + 1$  Fehler.*
- b) *Für perfekte  $(N, M, d)$ -Kodes ist  $d$  notwendig ungerade.*
- c) *Ein  $(N, M, d)$ -Kode ist perfekt genau dann, wenn*  

$$M \sum_{i=0}^{(d-1)/2} \binom{N}{i} (m-1)^i = m^N.$$

**Beweis.** a) ist eine direkte Konsequenz aus Satz 6.1.

b)  $\mathcal{C} = \{\mathbf{c}_1, \dots, \mathbf{c}_M\}$  sei ein perfekter  $(N, M, d)$ -Kode mit umgebenden Kugeln  $S_t(\mathbf{c}_1), \dots, S_t(\mathbf{c}_M)$ . Seien  $i, j \in \mathbb{N}$  so, daß  $d(\mathbf{c}_i, \mathbf{c}_j) = \min_{u,v} d(\mathbf{c}_u, \mathbf{c}_v)$ . Angenommen  $d(\mathbf{c}_i, \mathbf{c}_j) = 2k$  ist gerade. Dann unterscheiden sich  $\mathbf{c}_i$  und  $\mathbf{c}_j$  an genau  $2k$  Stellen, etwa bei den Indizes  $\ell_1, \dots, \ell_k$ . Setze  $a_m = c_{im}$ , falls  $m \neq \ell_1, \dots, \ell_k$ , und  $a_m = c_{jm}$ , sonst. Das so konstruierte  $\mathbf{a} = (a_1, \dots, a_N)$  erfüllt  $\mathbf{a} \in S_k(\mathbf{c}_i) \cap S_k(\mathbf{c}_j)$ , so daß  $t < k$  gilt. Da  $\mathcal{C}$  perfekt ist, existiert  $\mathbf{c}_\ell$  mit  $\mathbf{a} \in S_t(\mathbf{c}_\ell)$ , d.h.  $d(\mathbf{a}, \mathbf{c}_\ell) < k$ . Insgesamt folgt nun  $d(\mathbf{c}_i, \mathbf{c}_\ell) \leq d(\mathbf{c}_i, \mathbf{a}) + d(\mathbf{a}, \mathbf{c}_\ell) < 2k$ , im Widerspruch zu  $d(\mathbf{c}_i, \mathbf{c}_j) = 2k$ .

c) ergibt sich durch Anzahlüberlegungen wie im Beweis von Satz 6.2. ■

Die Dekodierung von  $(N, M, d)$ -Kodes ist im allgemeinen sehr aufwendig. Da über den Abstandsbegriff hinaus keine weitere Struktur vorliegt, muß im wesentlichen in einer großen Tabelle die entsprechende Dekodierung abgefragt werden. Für eine spezielle Klasse von Kodes kann dieses Problem leichter gelöst werden, für lineare Kodes nämlich.

## 6.2 Lineare Kodes

In diesem Abschnitt wird allgemein vorausgesetzt, daß das Alphabet  $\mathcal{X} = \{x_1, \dots, x_m\}$  mit den Operationen '+' und '·' einen Körper bildet. Im Fall, daß  $m = p$  eine Primzahl ist, kann  $\mathcal{X}$  mit den Zahlen  $\{0, \dots, m-1\}$  identifiziert werden; '+' und '·' werden dann durch die modulo- $p$ -Arithmetik realisiert. Ist  $m = p^\ell$  eine Primzahlpotenz, existieren Operationen '+' und '·', die  $\mathcal{X}$  zu einem endlichen Körper machen. Konstruktiv kann man diesen Körper mit Hilfe von Polynomrestklassen bezüglich eines irreduziblen Polynoms bestimmen (s. Anhang).



**Beispiel 6.1** Unter den durch die folgenden Tabellen definierten Operationen ist  $\mathcal{X} = \{x_0, x_1, x_2, x_3\}$  ein Körper ( $m = 4 = 2^2$ ).

+	$x_0$	$x_1$	$x_2$	$x_3$
$x_0$	$x_0$	$x_1$	$x_2$	$x_3$
$x_1$	$x_1$	$x_0$	$x_3$	$x_2$
$x_2$	$x_2$	$x_3$	$x_0$	$x_1$
$x_3$	$x_3$	$x_2$	$x_1$	$x_0$

·	$x_0$	$x_1$	$x_2$	$x_3$
$x_0$	$x_0$	$x_0$	$x_0$	$x_0$
$x_1$	$x_0$	$x_1$	$x_2$	$x_3$
$x_2$	$x_0$	$x_2$	$x_3$	$x_1$
$x_3$	$x_0$	$x_3$	$x_1$	$x_2$

$x_0$  spielt die Rolle des Nullelements in der additiven Gruppe, und  $x_1$  ist das Einselement in der entsprechenden multiplikativen Gruppe ohne  $x_0$ . ■

Bildet  $\mathcal{X}$  unter ‘+’ und ‘·’ einen endlichen Körper, so ist  $\mathcal{X}^N = V_N(m)$  ein Vektorraum der Dimension  $N$  über  $\mathcal{X}$ . Seine Elemente werden dargestellt als Zeilenvektoren über  $\mathcal{X}$

$$V_N(m) = \{\mathbf{a} = (a_1, \dots, a_N) \mid a_1, \dots, a_N \in \mathcal{X}\}.$$

Alle nachfolgend durchgeführten arithmetischen Operationen werden bezüglich der entsprechenden Körperarithmetik bzw. Vektorraumarithmetik durchgeführt.

**Definition 6.2** (*linearer Kode*)

$\mathcal{C} \subseteq V_N(m)$  heißt *linearer Kode*, wenn  $\mathcal{C}$  ein  $k$ -dimensionaler Unterraum in  $V_N(m)$  ist. Bezeichnung:  $[N, k]$ -Kode bzw.  $[N, k, d]$ -Kode, falls  $d = d(\mathcal{C})$ .

Ist  $F$  ein  $k$ -dimensionaler Unterraum in  $V_N(m)$ , so gilt bekanntlich  $|F| = m^k$ . Folglich ist jeder  $[N, k, d]$ -Kode  $\mathcal{C}$  ein  $(N, m^k, d)$ -Kode.

**Definition 6.3** (*Generatormatrix*)

Eine  $(k \times N)$ -Matrix  $\mathbf{G}$  heißt *Generatormatrix* des  $[N, k]$ -Kodes  $\mathcal{C}$ , wenn die Zeilen von  $\mathbf{G}$  linear unabhängige Vektoren in  $\mathcal{C}$  sind, also den Unterraum  $\mathcal{C}$  aufspannen.

Durch elementare Zeilenumformungen und Permutation der Spalten kann jede Generatormatrix in die folgende Standardform transformiert werden.

$$\mathbf{G} = (\mathbf{I}_k, \mathbf{A}), \text{ wobei } \mathbf{A} \text{ eine } (k \times (N - k))\text{-Matrix ist.}$$

Elementare Zeilenumformungen lassen sich durch Linksmultiplikation mit einer regulären  $(k \times k)$ -Matrix  $\mathbf{T}$  und Spaltenvertauschungen durch Rechtsmultiplikation mit einer Permutationsmatrix  $\mathbf{II}$  beschreiben.  $\mathbf{G}_0$  sei nun die Generatormatrix eines linearen Kodes  $\mathcal{C}_0$  und

$$\mathbf{G} = (\mathbf{I}_k, \mathbf{A}) = \mathbf{T}\mathbf{G}_0\mathbf{II} \quad (6.1)$$

die zugehörige Generatormatrix in Standardform. Mit Hilfe dieser Darstellung sieht man, daß der zu  $\mathbf{G}$  gehörige Kode  $\mathcal{C}$  ein linearer Kode ist, dessen Wörter durch Permutation der Komponenten der Wörter von  $\mathcal{C}_0$  entstehen. Die Kodes  $\mathcal{C}$  und  $\mathcal{C}_0$  heißen dann äquivalent.

Für  $\mathbf{a} = (a_1, \dots, a_N)$  bezeichne  $w(\mathbf{a})$  die Anzahl der von Null verschiedenen Komponenten.

$$w(\mathbf{a}) = |\{a_i \mid a_i \neq 0\}| = d(\mathbf{a}, \mathbf{0}) \quad (6.2)$$

heißt *Gewicht* von  $\mathbf{a}$ .

**Satz 6.3**  $\mathcal{C}$  sei ein linearer  $[N, k, d]$ -Kode. Dann gilt

$$d = d(\mathcal{C}) = \min\{w(\mathbf{c}) \mid \mathbf{c} \in \mathcal{C}, \mathbf{c} \neq \mathbf{0}\}.$$

**Beweis.** Seien  $\mathbf{a}, \mathbf{b} \in \mathcal{C}$  mit  $d(\mathbf{a}, \mathbf{b}) = d = d(\mathcal{C})$ . Da  $\mathcal{C}$  ein linearer Unterraum ist, gilt  $\mathbf{a} - \mathbf{b} \in \mathcal{C}$  und  $d(\mathbf{a}, \mathbf{b}) = d(\mathbf{a} - \mathbf{b}, \mathbf{0}) = w(\mathbf{a} - \mathbf{b})$ . Also ist  $\min\{w(\mathbf{c}) \mid \mathbf{c} \in \mathcal{C}, \mathbf{c} \neq \mathbf{0}\} \leq d$ .

Angenommen,  $\min\{w(\mathbf{c}) \mid \mathbf{c} \in \mathcal{C}, \mathbf{c} \neq \mathbf{0}\} < d$ . Dann existiert ein  $\mathbf{z} \in \mathcal{C}$ ,  $\mathbf{z} \neq \mathbf{0}$  mit  $w(\mathbf{z}) = d(\mathbf{z}, \mathbf{0}) < d$ . Dies ist ein Widerspruch, da  $\mathbf{0} \in \mathcal{C}$ . ■

$\mathcal{C}$  sei nun ein  $[N, k]$ -Kode mit Generatormatrix  $\mathbf{G} = (\mathbf{I}_k, \mathbf{A})$ . Jedes  $k$ -stellige Wort  $(a_1, \dots, a_k) \in \mathcal{X}^k$  wird bei linearen Kodes durch

$$(c_1, \dots, c_N) = (a_1, \dots, a_k)(\mathbf{I}_k, \mathbf{A})$$

kodiert. Die  $((N - k) \times N)$ -Matrix  $\mathbf{H} = (-\mathbf{A}', \mathbf{I}_{N-k})$  heißt *Kontrollmatrix* (englisch: parity check matrix).  $\mathbf{A}'$  bezeichnet hierbei die Transponierte der Matrix  $\mathbf{A}$ . Die Bedeutung der Kontrollmatrix erklärt sich aus folgendem Zusammenhang.

**Lemma 6.2** Sei  $\mathbf{c} \in V_N(m)$ . Unter obigen Bezeichnungen gilt  $\mathbf{c} \in \mathcal{C}$  genau dann, wenn  $\mathbf{c}\mathbf{H}' = \mathbf{o}$ .

**Beweis.** Ist  $\mathbf{c} = (c_1, \dots, c_N) \in \mathcal{C}$ , so existiert  $(s_1, \dots, s_k) \in V_k(m)$  mit der Darstellung  $(c_1, \dots, c_N) = (s_1, \dots, s_k)(\mathbf{I}_k, \mathbf{A})$ . Es folgt  $(c_1, \dots, c_k) = (s_1, \dots, s_k)$  und  $(c_{k+1}, \dots, c_N) = (s_1, \dots, s_k)\mathbf{A} = (c_1, \dots, c_k)\mathbf{A}$ . Die letzte Gleichung besagt in Matrixform  $(c_1, \dots, c_N) \begin{pmatrix} - \\ \mathbf{A} \end{pmatrix}$

$\mathbf{I}_{N-k} = \mathbf{o}$ , also  $\mathbf{c}\mathbf{H}' = \mathbf{o}$ .

Gilt umgekehrt  $(c_{k+1}, \dots, c_N) = (c_1, \dots, c_k)\mathbf{A}$ , so folgt die Gleichheit  $(c_1, \dots, c_N) = (c_1, \dots, c_k)(\mathbf{I}_k, \mathbf{A})$ , also  $(c_1, \dots, c_N) \in \mathcal{C}$ . ■

Insgesamt gilt für einen  $[N, k]$ -Code  $\mathcal{C}$  mit Generatormatrix  $\mathbf{G}$  und Kontrollmatrix  $\mathbf{H}$

$$\mathcal{C} = \{\mathbf{s}\mathbf{G} \mid \mathbf{s} = (s_1, \dots, s_k) \in V_k(m)\} = \{\mathbf{c} \in V_N(m) \mid \mathbf{c}\mathbf{H}' = \mathbf{o}\}.$$

Dies gibt die Möglichkeit, nachzuprüfen, ob ein empfangenes Wort  $\mathbf{b} = (b_1, \dots, b_N)$  ein zulässiges Kodewort aus  $\mathcal{C}$  ist. Falls  $\mathbf{b}\mathbf{H}' \neq \mathbf{o}$  gilt, ist bei der Übertragung ein Fehler passiert. Andererseits zeigt  $\mathbf{b}\mathbf{H}' = \mathbf{o}$ , daß zumindest ein zulässiges Kodewort empfangen wurde. Hieraus erklärt sich der Name Kontrollmatrix für  $\mathbf{H}$ .

Mit Hilfe der Kontrollmatrix können Aussagen über den Minimalabstand der Kodewörter des zugehörigen linearen Codes gemacht werden.

**Satz 6.4**  $\mathcal{C} \subset V_N(m)$  sei ein linearer Code mit Kontrollmatrix  $\mathbf{H}$ . Dann gilt  $d(\mathcal{C}) \geq d$  genau dann, wenn je  $d - 1$  Spalten von  $\mathbf{H}$  linear unabhängig sind. Besitzt  $\mathbf{H}$  zusätzlich  $d$  linear abhängige Spalten, so gilt  $d(\mathcal{C}) = d$ .

**Beweis.** (durch Kontraposition) Wenn es  $d - 1$  linear abhängige Spalten gibt, existiert  $\mathbf{c} \in V_N(m) \setminus \{\mathbf{o}\}$  mit höchstens  $d - 1$  von Null verschiedenen Komponenten, derart daß  $\mathbf{c}\mathbf{H}' = \mathbf{o}$ . Wegen Lemma 6.2 gilt  $\mathbf{c} \in \mathcal{C}$ , und mit Satz 6.3 folgt  $d(\mathcal{C}) \leq w(\mathbf{c}) \leq d - 1$ .

Gilt umgekehrt  $1 \leq w(\mathbf{c}) \leq d - 1$  für ein  $\mathbf{c} \in \mathcal{C}$ , also mit  $\mathbf{c}\mathbf{H}' = \mathbf{o}$ , so existieren  $d - 1$  linear unabhängige Spalten, da  $\mathbf{o}$  aus höchstens  $d - 1$  Spalten nicht-trivial linear kombiniert werden kann.

Existieren weiterhin  $d$  linear abhängige Spalten, so existiert  $\mathbf{c} \in V_N(m)$  mit  $w(\mathbf{c}) \leq d$  und  $\mathbf{c}\mathbf{H}' = \mathbf{o}$ . Dann gilt  $\mathbf{c} \in \mathcal{C}$  und  $d(\mathcal{C}) \leq d$ , insgesamt also  $d(\mathcal{C}) = d$ . ■

**Beispiel 6.2** Der binäre lineare Kode  $\mathcal{C}$  mit Generatormatrix

$$\mathbf{G} = \begin{pmatrix} 1 & 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 1 & 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 & 0 & 1 & 1 \end{pmatrix}$$

und Kontrollmatrix

$$\mathbf{H} = \begin{pmatrix} 1 & 1 & 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 & 1 & 0 \\ 1 & 1 & 0 & 1 & 0 & 0 & 1 \end{pmatrix}$$

heißt (7,4)-Hamming-Kode. Allgemein heißt ein linearer Kode  $\mathcal{C}$  mit Kodewörtern der Länge  $N = 2^k - 1$  binärer  $(2^k - 1, 2^k - 1 - k)$ -Hamming-Kode, wenn die zugehörige Kontrollmatrix  $\mathbf{H}$  alle von  $\mathbf{o}$  verschiedenen  $2^k - 1$  Vektoren aus  $V_k(2)$  als Spalten enthält.

Je 2 Spalten von  $\mathbf{H}$  sind linear unabhängig, die vierte, sechste und siebte sind jedoch gemeinsam linear abhängig. Nach Lemma 6.4 hat der Kode den Minimalabstand  $d(\mathcal{C}) = 3$ . Wegen Satz 6.1 ist er 1-fehlerkorrigierend.  $\mathcal{C}$  enthält  $M = 2^4$  Kodewörter der Länge 7. Die paarweise disjunkten Kugeln mit Radius 1 um jedes Kodewort enthalten genau  $7 + 1 = 2^3$  Kodewörter. Weil  $2^7 = 2^4 \cdot 2^3$ , ist der Kode perfekt (vgl. auch Lemma 6.1 c)). ■

Wir kommen nun zur MD-Dekodierung linearer Kodes. Die durch den linearen Unterraum  $\mathcal{C}$  definierten Nebenklassen (englisch: cosets) besitzen eine Darstellung  $\mathbf{a} + \mathcal{C}$ ,  $\mathbf{a} \in V_N(m)$ .  $\mathbf{z} = \mathbf{a} + \mathcal{C}$  heißt *Anführer* (englisch: coset leader) der entsprechenden Nebenklasse, wenn das Gewicht  $w(\mathbf{z})$  minimal ist für alle Elemente der Nebenklasse.

**Lemma 6.3**  $\mathcal{C}$  sei ein  $[N, k]$ -Kode mit Kontrollmatrix  $\mathbf{H}$ .  $\mathbf{x}, \mathbf{y} \in V_N(m)$  liegen genau dann in derselben Nebenklasse, wenn  $\mathbf{x}\mathbf{H}' = \mathbf{y}\mathbf{H}'$ .

**Beweis.** Es gilt  $\mathbf{x}, \mathbf{y} \in \mathbf{a} + \mathcal{C}$  für ein  $\mathbf{a} \in V_N(m)$  genau dann, wenn  $\mathbf{x} - \mathbf{y} \in \mathcal{C}$ , d.h. mit Lemma 6.2  $(\mathbf{x} - \mathbf{y})\mathbf{H}' = \mathbf{o}$ . Dies ist äquivalent zu  $\mathbf{x}\mathbf{H}' = \mathbf{y}\mathbf{H}'$ . ■

Für jede Nebenklasse  $\mathbf{a} + \mathcal{C}$  ist  $\mathbf{a}\mathbf{H}'$  wegen Lemma 6.3 unabhängig von der Wahl des speziellen Repräsentanten.  $\mathbf{a}\mathbf{H}'$  heißt Syndrom der Nebenklasse

$\mathbf{a} + \mathcal{C}$ . Im folgenden wird angenommen, daß für alle  $m^{N-k}$  verschiedenen Nebenklassen ihr Syndrom  $\mathbf{a}\mathbf{H}'$  und Anführer  $\mathbf{z}$  bekannt sind.

MD-Dekodierung für ein empfangenes  $\mathbf{y} \in V_N(m)$  bedeutet, ein Element  $\mathbf{c} \in \mathcal{C}$  mit minimaler Hamming-Distanz  $d(\mathbf{y}, \mathbf{c})$  zu bestimmen. Da  $d(\mathbf{y}, \mathbf{c}) = d(\mathbf{y} - \mathbf{c}, \mathbf{0}) = w(\mathbf{y} - \mathbf{c})$ , ist ein  $\mathbf{c} \in \mathcal{C}$  zu berechnen, für das  $w(\mathbf{y} - \mathbf{c})$  minimal ist. Offensichtlich gilt

$$\{\mathbf{y} - \mathbf{c} \mid \mathbf{c} \in \mathcal{C}\} = \{\mathbf{y} + \mathbf{c} \mid \mathbf{c} \in \mathcal{C}\} = \mathbf{y} + \mathcal{C}.$$

Bei MD-Dekodierung ist also  $\min\{w(\mathbf{e}) \mid \mathbf{e} \in \mathbf{y} + \mathcal{C}\}$  zu bestimmen. Eine Lösung hiervon ist durch den Anführer  $\mathbf{z}$  der Nebenklasse  $\mathbf{y} + \mathcal{C}$  gegeben.

Dieses Dekodierverfahren wird folgendermaßen realisiert. In einer Tabelle werden die Syndrome und zugehörigen Anführer aller Nebenklassen gespeichert. Bekannt ist ferner die Kontrollmatrix  $\mathbf{H}$ . Wird nun  $\mathbf{y} \in V_N(m)$  empfangen, verfährt man in drei Schritten.

1. Berechne  $\mathbf{y}\mathbf{H}'$ . Dies liefert das Syndrom der Restklasse  $\mathbf{y} + \mathcal{C}$ .
2. Bestimme aus der Tabelle zu  $\mathbf{y}\mathbf{H}'$  den Anführer  $\mathbf{z}$ .
3. Dekodiere  $\mathbf{y}$  durch  $\mathbf{c} = \mathbf{y} - \mathbf{z}$ .

Dann ist  $\mathbf{c} \in \mathcal{C}$  ein Element minimalen Abstands zu  $\mathbf{y} \in V_N(m)$ .

**Beispiel 6.3** Sei  $m = 2$ ,  $\mathcal{X} = \{0, 1\}$  und  $\mathcal{C}$  ein linearer Kode mit Generatormatrix  $\mathbf{G}$  und Kontrollmatrix  $\mathbf{H}$ , wobei  $k = 2$ ,  $N = 4$ ,

$$\mathbf{G} = \begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 1 & 1 \end{pmatrix} \quad \text{und} \quad \mathbf{H} = \begin{pmatrix} 1 & 1 & 1 & 0 \\ 0 & 1 & 0 & 1 \end{pmatrix}.$$

Der zugehörige  $[4, 2]$ -Kode ist

$$\mathcal{C} = \{(0000), (1010), (0111), (1101)\}.$$

Syndrome und Anführer sind in folgender Tabelle angegeben.

Syndrom	Anführer
(00)	(0000)
(10)	(0010)
(01)	(0001)
(11)	(0100)

Wird zum Beispiel  $\mathbf{y} = (1111)$  empfangen, so führt  $\mathbf{y}\mathbf{H}' = (10)$  zum Anführer  $\mathbf{z} = (0010)$ . Dekodiert wird also  $\mathbf{y}$  durch  $\mathbf{y} - \mathbf{z} = (1101)$ .

Das MD-Dekodierverfahren ist allerdings nicht notwendig eindeutig. Alternativ ist  $(1000)$  in der zweiten Zeile obiger Tabelle ebenfalls Anführer der zugehörigen Restklasse. Dies würde zur gleichwertigen Dekodierung  $\mathbf{y} - \mathbf{z} = (0111)$  führen. ■

### 6.3 Faltungskodes und der Viterbi-Algorithmus

Faltungskodes sind als Verallgemeinerung der linearen Blockcodes entstanden. Faltungskodierung ist eine Methode, die nahe am Schaltkreisentwurf entwickelt werden kann und daher in vielen Fällen effizient implementierbar ist. Mit dem Viterbi-Algorithmus steht auf der Seite des Kanaldekodierers ein Verfahren zur Verfügung, mit dem Faltungskodes auch wieder effizient dekodiert werden können.

Im folgenden wird nur das binäre Alphabet  $\mathcal{X} = \{0, 1\}$  mit der entsprechenden modulo-2-Arithmetik als Körper  $\text{GF}(2)$  betrachtet. Faltungskodierer kann man durch ihr Schaltdiagramm aus Schieberegistern und modulo-2-Addierern beschreiben. Verarbeitet werden  $k$  Eingabebits zu einem Ausgabeblock der Länge  $N$  durch sogenannte  $(k, N)$ -Faltungskodierer.

**Beispiel 6.4** In Abbildung 6.1 wird ein  $(1,3)$ -Faltungskodierer dargestellt. Jeweils  $k = 1$  Eingabebit wird mit der *Rückgrifftiefe* (englisch: constraint length)  $\nu = 2$  zu  $N = 3$  Ausgabebits kodiert. Bei Eingabe eines Bits  $u_t$  zur Zeit  $t \in \mathbb{N}$  wird die Ausgabe  $(a_{1t}, a_{2t}, a_{3t})$  erzeugt. Anschließend wird der Registerinhalt  $u_{t-2}$  durch  $u_{t-1}$  und das Vorgängerregister mit  $u_t$  überschrieben. Der Kodierer ist dann zur Eingabe des nächsten Bits bereit.  $\oplus$  bedeutet modulo-2-Addition entsprechend der XOR-Verknüpfung.

Zu Beginn des Kodiervorgangs für einen Eingabeblock sind alle Registerinhalte  $= 0$ . Zum Ende des Blocks werden die Register durch Eingabe von Nullen wieder zurückgesetzt. Kodiert wird zum Beispiel

$$(10011|00) \mapsto (111|011|001|111|100|010|001).$$

Die letzten beiden Nullen der Eingabe bewirken das Rücksetzen der Register und verursachen die letzten beiden Ausgabeblocke. ■

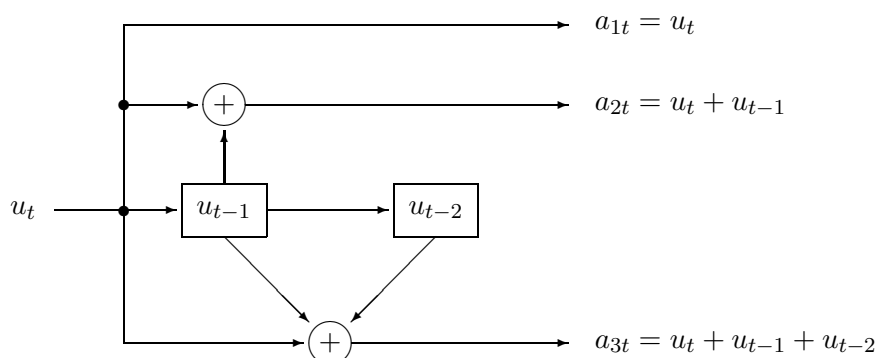


Abb. 6.1 Das Schaltdiagramm eines (1,3)-Faltungskodierers.

**Beispiel 6.5** In Abbildung 6.2 wird ein (2,3)-Faltungskodierer der Rückgriffiefe  $\nu = 1$  dargestellt. Die Bits  $u_{1t}, u_{2t}$  werden direkt ausgegeben,  $a_{3t}$  spielt die Rolle eines Kontrollbits. ■

In obigen Beispielen kann die Wirkungsweise der Faltungskodierer durch Matrixmultiplikation beschrieben werden, wobei alle Operationen in der modulo-2-Arithmetik durchgeführt werden. Bei Beispiel 6.4 etwa werden die drei Ausgabebits  $(a_{1t}, a_{2t}, a_{3t})$  zur Zeit  $t$  bestimmt durch

$$(a_{1t}, a_{2t}, a_{3t}) = (u_{t-2}, u_{t-1}, u_t) \begin{pmatrix} 0 & 0 & 1 \\ 0 & 1 & 1 \\ 1 & 1 & 1 \end{pmatrix}, \quad t \in \mathbb{N}. \quad (6.3)$$

Man beachte, daß hierbei die Anfangsregisterinhalte  $u_{-1} = u_0 = 0$  gesetzt werden. In Beispiel 6.4 entsteht jede Ausgabesequenz durch Rechtsmultiplikation mit der Matrix

$$\mathbf{G} = \begin{pmatrix} 1 & 1 & 1 & 0 & 1 & 1 & 0 & 0 & 1 & \mathbf{0} \\ & & & 1 & 1 & 1 & 0 & 1 & 1 & 0 & 0 & 1 \\ & & & & & & 1 & 1 & 1 & 0 & 1 & 1 \\ & & & & \mathbf{0} & & & & & 1 & 1 & 1 \\ & & & & & & & & & & & \ddots \end{pmatrix}$$

entsprechender Dimension, wobei die freigelassenen Einträge dieser Matrix alle gleich 0 sind.  $\mathbf{G}$  heißt *Generatormatrix*.

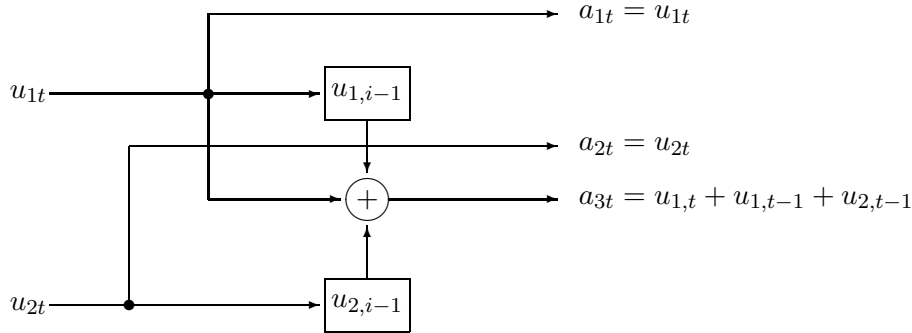


Abb. 6.2 Das Schaltdiagramm eines (2,3)-Faltungskodierers.

Die folgende Darstellung von Faltungskodierern betont den dynamischen Charakter des Systems und legt eine Repräsentation durch Zustandsübergangsdiagramme nahe. Aus Gründen der Übersichtlichkeit wird im folgenden nur der Fall  $k = 1$  betrachtet. Es bedeuten

$$\begin{array}{ll}
 u_t \in \{0, 1\} & \text{das Eingabebit zur Zeit } t, \\
 \mathbf{s}_t \in \{0, 1\}^\nu & \text{den Zustand zur Zeit } t, \\
 \mathbf{a}_t \in \{0, 1\}^N & \text{die Ausgabe zur Zeit } t.
 \end{array}$$

Ferner sind

$$\mathbf{A} \in \{0, 1\}^{N \times \nu}, \quad \mathbf{b} \in \{0, 1\}^N, \quad \mathbf{C} \in \{0, 1\}^{\nu \times \nu}, \quad \mathbf{d} \in \{0, 1\}^\nu$$

festen binäre Matrizen bzw. Vektoren.

Das Fortschreiten der internen Zustände (der Registerinhalte) und die Ausgabe des Faltungskodierers werden beschrieben durch die folgenden Gleichungen.

$$\mathbf{a}_t = \mathbf{A}\mathbf{s}_{t-1} + u_t \mathbf{b} \tag{6.4}$$

$$\mathbf{s}_t = \mathbf{C}\mathbf{s}_{t-1} + u_t \mathbf{d} \tag{6.5}$$

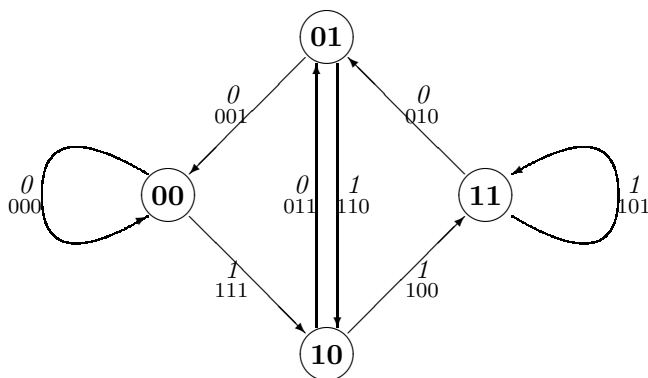
In Beispiel 6.4 ergibt sich konkret,  $t \in \mathbb{N}$ ,

$$\mathbf{a}_t = \begin{pmatrix} 0 & 0 \\ 1 & 0 \\ 1 & 1 \end{pmatrix} \mathbf{s}_{t-1} + u_t \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}, \quad \mathbf{s}_t = \begin{pmatrix} 0 & 0 \\ 1 & 0 \end{pmatrix} \mathbf{s}_{t-1} + u_t \begin{pmatrix} 1 \\ 0 \end{pmatrix}.$$



Auf der Basis dieser Gleichungen können Faltungskodierer durch ihr Zustandsänderungsdiagramm beschrieben werden. Dies ist ein gerichteter Graph, dessen Knoten die Menge der Registerzustände repräsentieren. Die gerichteten Kanten beschreiben den Zustandsübergang bei Eingabebit  $u_t$  gemäß Gleichung (6.5). An den Kanten wird die Eingabe  $u_t$  und die Ausgabe  $a_t$  aus Gleichung (6.5) notiert.

Bei Beispiel 6.4 erhält man den folgenden Graphen.



Hiermit läßt sich die Kodierung einzelner Bits genau verfolgen. Ist der Registerzustand etwa  $(1, 0)$  (entsprechend dem unteren Knoten) und wird eine 1 eingegeben, so geht der interne Zustand in  $(1, 1)$  über (rechter Knoten) und es wird der Block  $(100)$  ausgegeben.

Unbequem ist dieses Diagramm, wenn man das Verhalten bei ganzen Eingabesequenzen verfolgen will. Man muß dann in dem Zustandsänderungsdiagramm umherwandern, Vorgängerzustände und Ausgaben auf dem Weg einer bestimmten Kodierung werden jedoch vergessen, wenn sie nicht separat notiert wurden. Hier hilft die Darstellung mit Hilfe des *Trellis-Diagramms* oder *Spalierdiagramms*.

Dies ist ein gerichteter bipartiter Graph, dessen Knoten die Menge der Registerzustände zur Zeit  $t$  und  $t+1$  darstellen. Die Kanten beschreiben wie beim Zustandsänderungsdiagramm die internen Zustandsübergänge bei Eingabebit  $u_t \in \{0, 1\}$ . An den Kanten werden jeweils die Ausgabeblöcke vermerkt. Für den Faltungskodierer aus Beispiel 6.4 ergibt sich der Graph aus Abbildung 6.3. Ein Pfeil nach oben bedeutet Eingabe 0, ein Pfeil nach unten Eingabe 1. Seinen Namen bezieht das Trellis-Diagramm aus der Form des

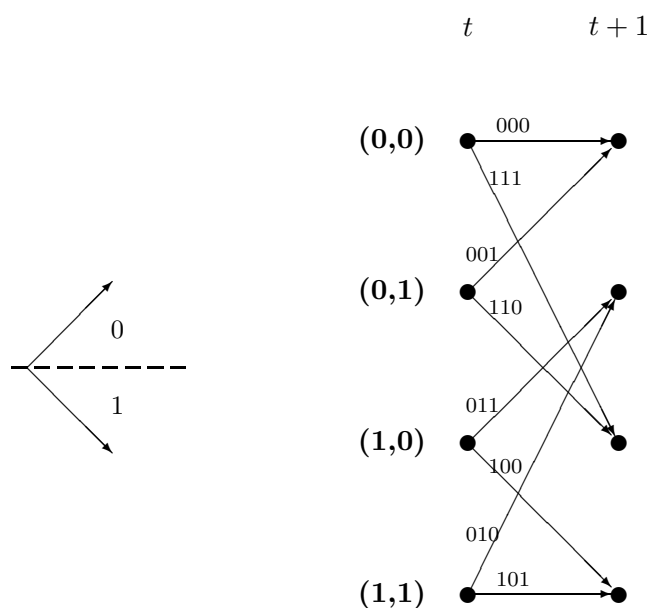


Abb. 6.3 Das Trellis-Diagramm zu Beispiel 6.4

resultierenden Graphen, der einem Spalier ähnelt. (Spalier = trellis im Englischen)

Der zugehörige Kode kann zu einem Kodebaum aufgefaltet werden. Dies ist ein binärer Baum mit einer Verzweigung nach oben bei Eingabe 0 und einer Verzweigung nach unten bei Eingabe 1. An den Kanten sind wieder die Ausgabeblocks vermerkt. Für obiges Beispiel erhält man den in Abbildung 6.4 dargestellten Baum. Das Eingabewort (10011) wird durch Verfolgen der entsprechenden Kanten zu der Sequenz (111|011|001|111|100) kodiert. MD-Dekodierung in einem Kodebaum ist einfach aber ineffizient. Auf jeder Stufe berechnet man den Hamming-Abstand des empfangenen Blocks mit allen Blöcken im Kodebaum. An den Knoten werden die aufsummierten Hamming-Distanzen vermerkt. Der in dem Endknoten mit kleinster kumulierter Hamming-Distanz endende Weg repräsentiert das Ergebnis der MD-Dekodierung. In Abbildung 6.4 ist dies für das empfangene Wort (110|011|101|111|110) durchgeführt; dekodiert wird (10011).

Zur Dekodierung von Faltungskodes hat Viterbi 1979 einen effizienten Al-

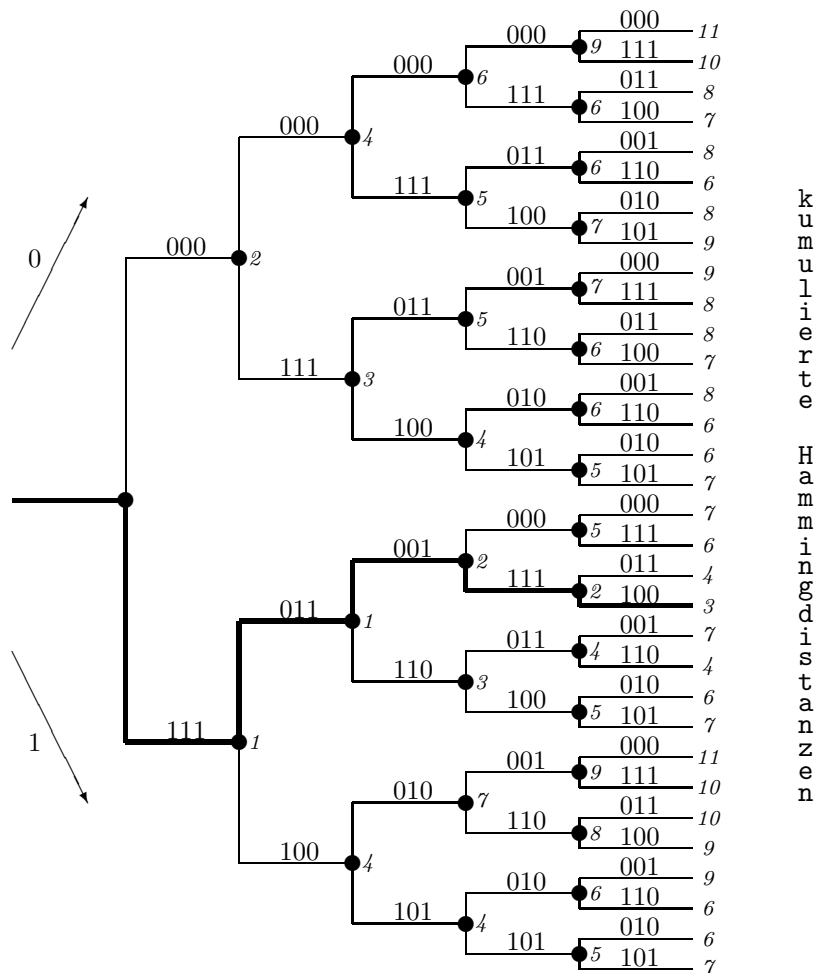


Abb. 6.4 Der binäre Kodebaum zu Beispiel 6.4

gorithmus entworfen. Er basiert auf dem Finden eines kürzesten Weges in einem iterierten Trellis-Diagramm und eignet sich daher allgemein für Codes, die eine Trellis-Darstellung mit gewichteten Kanten haben.

Den Ausgangspunkt bildet ein binärer diskreter Kanal, bei dem das Wort  $\mathbf{b} = (b_1, \dots, b_{NL})$  empfangen wurde. Gesucht wird nun ein zulässiges Codewort  $\mathbf{c} = (c_1, \dots, c_{NL})$  kleinsten Hamming-Abstands zu  $\mathbf{b}$ . Diese MD-Dekodierung ist bei einem binären symmetrischen Kanal sogar äquivalent

zu ML-Dekodierung, wie Satz 5.3 zeigt. Das Dekodierschema beruht 1.) auf einer geschickten Darstellung des Kodierers, 2.) auf der Identifikation des Dekodierungsproblems mit dem Finden eines kürzesten Wegs, und 3.) einem effizienten Algorithmus zur Berechnung des kürzesten Weges. Die einzelnen Stufen werden im folgenden dargestellt.

1.) Darstellung des Kodierers für Nachrichten der Länge  $L$ .

- Die Trellisdiagramme der einstufigen Übergänge zu den Zeiten  $t = 0, \dots, L$  werden konkateniert. Die Knoten des entstehenden gerichteten Graphen werden mit  $\{0, 1\}^v \times \{0, \dots, L\}$ , den Zuständen zur Zeit  $t$ , identifiziert.
- Es werden alle Wege weggelassen, die keinen Ursprung im Startzustand  $(0, \dots, 0)$  zur Zeit  $t = 0$  haben.
- Jede gerichtete Kante wird mit dem Eingabebit  $u_t$  und dem Ausgabeblock  $\mathbf{a}_t$  markiert,  $t = 1, \dots, L$ .

Für den in Beispiel 6.4 eingeführten Faltungskodierer entsteht bei Eingabeblöcken der Länge 5 das in Abbildung 6.5 dargestellte Trellis-Diagramm. Der hierin mit Kreisen markierte Weg entspricht der Bitfolge (10011).

2.) Dekodierung des empfangenen Worts  $\mathbf{b} = (b_1, \dots, b_{NL})$ .

- Jede Kante des  $L$ -stufigen Trellis-Diagramms wird mit der Hamming-Distanz zwischen dem empfangenen Block und dem zugehörigen Ausgabeblock gewichtet.
- Alle Endknoten werden mit einem imaginären Knoten  $E$  durch gerichtete Kanten vom Gewicht 0 verbunden.
- Die Aufgabe, ein Kodewort kleinster Hamming-Distanz zu dem empfangenen zu bestimmen, besteht in der Bestimmung eines kürzesten Wegs in dem oben beschriebenen Graphen (im Sinn eines minimalen kumulierten Kantengewichts).

3.) Ein Algorithmus zur Bestimmung eines kürzesten Wegs.

Für alle  $t = 1, \dots, L$  und alle Zustände  $\mathbf{s}$ :

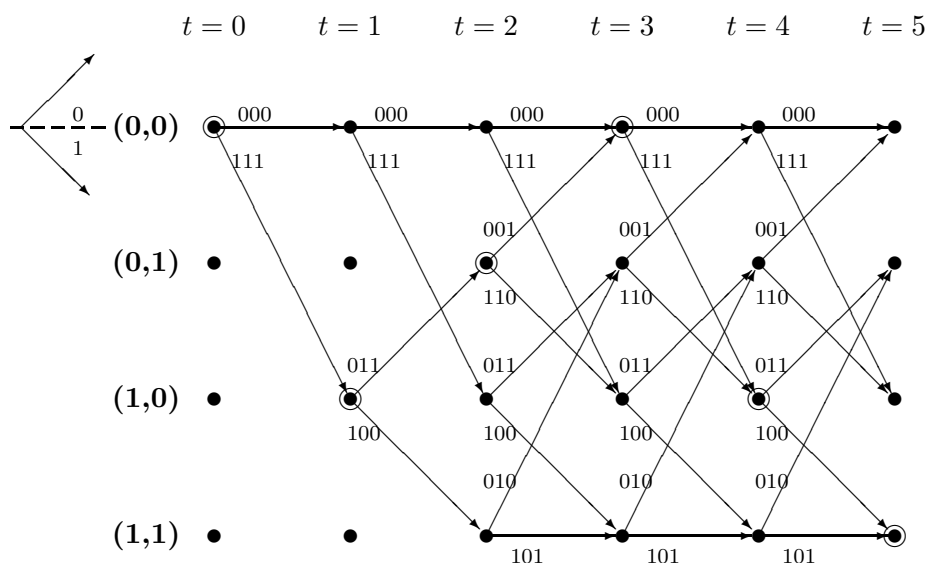


Abb. 6.5 Das 5-stufige Trellis-Diagramm zu Beispiel 6.4

- Bestimme auf der Stufe  $t$  die kumulierte Hamming-Distanz über alle eingehenden Wege durch Summation der am Vorgängerknoten notierten minimalen kumulierten Hamming-Distanz der Stufe  $t - 1$  mit dem jeweiligen Kantengewicht.
- Notiere das Minimum der ermittelten kumulierten Kantengewichte am jeweiligen Knoten.
- Dekodiere auf dem Weg, der für  $t = L$  minimale kumulierte Hammingdistanz hat.

Dies ist gerade der Algorithmus von Dykstra für den speziellen Fall der Trellis-Diagramme.

In Abbildung 6.6 wird dieses Verfahren für den Faltungskodierer aus Beispiel 6.4 dargestellt. Empfangen wurde das Wort  $\mathbf{b} = (110|011|101|111|110)$ , dekodiert wird  $(10011)$ . Es ist also anzunehmen, daß in  $\mathbf{b}$  das 3., 7. und 14. Bit gekippt waren. Die in Klammern an den Knoten notierten Werte sind die nicht minimalen kumulierten Kantengewichte. Das minimale Gewicht steht über dem Knoten, wenn es zu einer von oben einlaufenden Kante gehört, ansonsten unter dem Knoten.

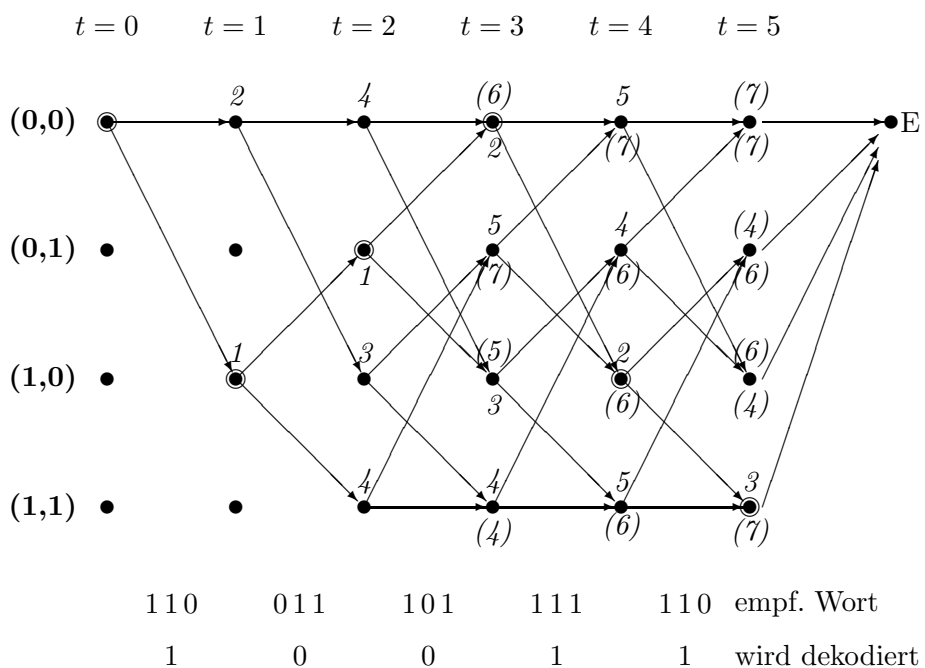


Abb. 6.6 Kumulierte Hamming-Distanzen und Dekodierung im 5-stufigen Trellis-Diagramm zu Beispiel 6.4

Der durch obigen Algorithmus bestimmte Weg ist minimal, da nicht betrachtete Wege schon bei Zwischenknoten größeres Gewicht hatten. Dieses Vorgehen ist ein Grundprinzip der dynamische Optimierung. Es läßt sich anschaulich an folgendem Beispiel begreifen. Ist von zwei Wegen von Aachen nach München über Stuttgart der erste schon in Stuttgart der längere, kann dieser insgesamt nicht der kürzeste sein, und man kann ihn schon in Stuttgart ausscheiden.

Allgemeine ML-Dekodierung bei beliebigen binären, diskreten gedächtnislosen Kanälen kann nach demselben Prinzip durchgeführt werden. Ist  $\mathbf{b} = (b_1, \dots, b_{NL})$  das empfangene Kodewort, so stellt sich die Aufgabe

$$\text{maximiere } \prod_{i=1}^{NL} p_1(b_i | c_i)$$

über alle zulässigen Kodewörter  $(c_1, \dots, c_{NL})$ . Dies ist äquivalent zu

$$\text{minimiere } \sum_{i=1}^{NL} (-\ln p_1(b_i|c_i)) \quad (6.6)$$

über alle zulässigen Kodewörter  $(c_1, \dots, c_{NL})$ . Um eine Lösung von (6.6) zu bestimmen, sind in obigem Trellis-Diagramm lediglich die Kantengewichte zu verändern. Die Kanten von  $t = \ell$  nach  $t = \ell + 1$  erhalten die Gewichte

$$-\sum_{j=1}^N \ln p_1(b_{\ell N+j}|c_{\ell N+j}).$$

Der Algorithmus selbst bleibt unverändert und liefert eine Maximum-Likelihood-Dekodierung des empfangenen Worts  $\mathbf{b}$ .

## 6.4 Übungsaufgaben

**Aufgabe 6.1**  $\mathcal{C}_1$  sei ein  $(N, M_1, d_1)$ -Kode und  $\mathcal{C}_2$  sei ein  $(N, M_2, d_2)$ -Kode über dem Alphabet  $\mathcal{X}$ . Der Kode  $\mathcal{C}_3$  sei definiert durch

$$\mathcal{C}_3 = \{(\mathbf{a}, \mathbf{a} + \mathbf{b}) \mid \mathbf{a} \in \mathcal{C}_1, \mathbf{b} \in \mathcal{C}_2\}.$$

Zeigen Sie, daß  $\mathcal{C}_3$  ein  $(2N, M_1 M_2, d_3)$ -Kode ist mit  $d_3 = \min\{2d_1, d_2\}$ .

**Aufgabe 6.2** Zeigen Sie, daß  $A_m(N, d) \leq m^{N-d+1}$ . Die Bezeichnungen sind hierbei wie in Satz 6.2 gewählt.

**Aufgabe 6.3**  $\mathcal{C} \subseteq V_N(m)$  sei ein linearer Kode. Der Kode  $\mathcal{C}'$  entstehe aus  $\mathcal{C}$ , indem das letzte Symbol jedes Kodeworts aus  $\mathcal{C}$  gestrichen wird. Zeigen Sie, daß  $\mathcal{C}' \subseteq V_{N-1}(m)$  auch ein linearer Kode ist.  $\mathcal{C}'$  heißt gestutzter Kode (englisch: truncated).

**Aufgabe 6.4** Man zeige, daß  $\mathbf{G} = (\mathbf{I}_{10}, \mathbf{B}')$  mit

$$\mathbf{B} = \begin{pmatrix} 0 & 0 & 2 & 2 & 2 & 2 & 2 & 2 & 2 & 2 \\ 2 & 2 & 0 & 0 & 2 & 2 & 2 & 1 & 1 & 1 \\ 2 & 1 & 2 & 1 & 0 & 2 & 1 & 0 & 2 & 1 \end{pmatrix}$$

die Generatormatrix eines perfekten linearen Kodes  $\mathcal{C}$  über  $\text{GF}(3)$  mit  $d(\mathcal{C}) = 3$  ist.

**Aufgabe 6.5** Beim Fußballtoto bedeutet 0 “unentschieden”, 1 “Heimsieg” und 2 “Heimniederlage”. Wieviel Tippreihen aus 13-mal 0, 1 oder 2 müssen Sie abgeben, um bei 13 Spielen garantiert einmal mindestens 12 Richtige zu haben. Beschreiben Sie die Tips durch die Kontrollmatrix eines linearen Codes.

Hinweis: Aufgabe 6.4

**Aufgabe 6.6** Man benutze den  $(7, 4)$ -Hamming-Code aus Beispiel 6.2 zum Kodieren der Wörter  $(1100)$ ,  $(1010)$ ,  $(1001)$ ,  $(0110)$ ,  $(0101)$ ,  $(0011)$ . Man dekodiere die empfangenen Wörter  $(1000000)$ ,  $(0001000)$ ,  $(0000001)$  und  $(1111111)$ .

**Aufgabe 6.7** Zeigen Sie, daß für  $k \geq 2$  jeder binäre  $(2^k - 1, 2^k - 1 - k)$ -Hamming-Code  $\mathcal{C}$  perfekt ist und  $d(\mathcal{C}) = 3$  erfüllt.

**Aufgabe 6.8** Man betrachte den Faltungskodierer, dessen Funktion wie in (6.3) durch die Matrix  $\begin{pmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 0 & 1 \end{pmatrix}$  beschrieben wird. Skizzieren Sie ein zugehöriges Schaltdiagramm. Dekodieren Sie das empfangene Wort  $(000\ 110\ 110\ 110\ 001)$  mit Hilfe des Viterbi-Algorithmus.

**Aufgabe 6.9** Bestimmen Sie einen binären Faltungskodierer, der ein Wort unendlicher Länge  $\mathbf{u} = (u_1, u_2, \dots)$  mit unendlichem Gewicht  $w(\mathbf{u}) = \infty$  in ein Kodewort  $\mathbf{a} = (a_1, a_2, \dots)$  mit endlichem Gewicht  $w(\mathbf{a}) < \infty$  kodiert.  $w$  bedeutet hierbei die Hamming-Distanz aus (6.2). Faltungskodierer mit dieser Eigenschaft heißen katastrophal.

**Aufgabe 6.10**  $\mathbf{a}(\mathbf{u})$  bezeichne das durch einen Faltungskodierer erzeugte Ausgabewort  $\mathbf{a} \in \{0, 1\}^\infty$  bei Eingabe eines unendlichen Worts  $\mathbf{u} = (u_1, u_2, \dots)$ . Zur Beurteilung der Güte eines Faltungskodierers wird der mit Hilfe des Hamming-Gewichts (6.2) durch

$$d_{\text{free}} = \min \{w(\mathbf{a}(\mathbf{u}_1) + \mathbf{a}(\mathbf{u}_2)) \mid \mathbf{u}_1 \neq \mathbf{u}_2 \in \{0, 1\}^\infty\}$$

definierte freie Abstand herangezogen.

Zeigen Sie, daß  $d_{\text{free}} = \min \{w(\mathbf{a}(\mathbf{u})) \mid \mathbf{u} \in \{0, 1\}^\infty, \mathbf{u} \neq \mathbf{o}\}$ . Man bestimme  $d_{\text{free}}$  für den Faltungskodierer aus Beispiel 6.4.



## 7 Anhang: endliche Körper

Endliche Körper und Vektorräume sind ein zentrales Hilfsmittel zur Konstruktion linearer Codes in Kapitel 6.2. Die grundlegenden Begriffe aus der Theorie endlicher Körper werden in einem kurzen Überblick zusammengestellt.

**Definition 7.1** Eine Menge  $\mathcal{K}$  mit  $|\mathcal{K}| \geq 2$ , die unter den Operationen ‘+’ und ‘·’ abgeschlossen ist (d.h.  $a+b \in \mathcal{K}$  und  $a \cdot b \in \mathcal{K}$  für alle  $a, b \in \mathcal{K}$ ), heißt Körper (englisch: field), wenn die folgenden Bedingungen erfüllt sind.

- 1)  $\mathcal{K}$  ist bezüglich + eine Abelsche Gruppe (d.h. es existiert ein neutrales Element 0 und inverse Elemente, es gilt das Assoziativgesetz  $(a+b)+c = a+(b+c)$  für alle  $a, b, c \in \mathcal{K}$  und Kommutativität  $a+b = b+a$  für alle  $a, b \in \mathcal{K}$ ).
- 2)  $\mathcal{K} \setminus \{0\}$  ist bezüglich · eine Abelsche Gruppe mit neutralem Element 1.
- 3) Es gilt das Distributivgesetz  $(a+b) \cdot c = a \cdot c + b \cdot c$  für alle  $a, b, c \in \mathcal{K}$ .

Aus obigen Axiomen leiten sich die üblichen Rechenregeln wie beim Umgang mit reellen Zahlen her. Bei der Konstruktion von Codes liegt das Hauptinteresse auf dem Fall, daß  $\mathcal{K}$  endlich ist.  $\mathcal{K}$  heißt dann endlicher Körper oder Galois-Feld der Ordnung  $m$ , wenn  $|\mathcal{K}| = m$ , kurz:  $\text{GF}(m)$ .

Bekannte Beispiele für endliche Körper sind die Restklassenkörper  $\mathbb{Z}_p$  in den ganzen Zahlen, wenn  $m = p$  eine Primzahl ist. Addition und Multiplikation werden durch den Rest bei Division durch  $p$  erklärt, also  $a+b = a+b \bmod p$  und  $a \cdot b = a \cdot b \bmod p$ , wobei der Bequemlichkeit halber für ‘+’ und ‘·’ auf beiden Seiten dieselben Symbole verwendet werden. Auf der rechten Seite ist die übliche Addition und Multiplikation in den ganzen Zahlen gemeint.  $\text{GF}(2)$  zum Beispiel hat die Rechenregeln  $0+0 = 1+1 = 0$ ,  $0+1 = 1+0 = 1$ ,  $0 \cdot 0 = 0 \cdot 1 = 1 \cdot 0 = 0$  und  $1 \cdot 1 = 1$ .

Die Frage ist nun, zu welchen  $m$  es überhaupt endliche Körper gibt und wie diese konstruiert werden. Hierbei helfen Polynome.

**Definition 7.2** Ein Ausdruck der Form  $a(D) = a_n D^n + a_{n-1} D^{n-1} + \dots + a_1 D + a_0$  mit  $a_0, a_1, \dots, a_n \in \text{GF}(m)$ ,  $a_n \neq 0$ , heißt Polynom vom Grad  $n$  über  $\text{GF}(m)$ .

Das Symbol  $D$  in Definition 7.2 ist nicht als Variable im Sinn einer Darstellung des Polynoms als Funktion zu betrachten. Im folgenden interessiert lediglich die Menge der Polynome als algebraische Objekte, für die Addition und Multiplikation im weiteren erklärt werden. Hierzu ist es bequem, ein Polynom  $a(D)$  vom Grad  $n$  als  $a(D) = \sum_{i=0}^{\infty} a_i D^i$  zu schreiben, wobei die Koeffizienten  $a_{n+1} = a_{n+2} = \dots = 0$  gesetzt werden. Der Grad eines Polynoms ist dann der größte Index der von Null verschiedenen Koeffizienten. Summe und Produkt von Polynomen  $a(D) = \sum_{i=0}^{\infty} a_i D^i$  und  $b(D) = \sum_{i=0}^{\infty} b_i D^i$  werden nun definiert als

$$\begin{aligned} a(D) + b(D) &= \sum_{i=0}^{\infty} (a_i + b_i) D^i, \\ a(D) \cdot b(D) &= \sum_{i=0}^{\infty} \left( \sum_{j=0}^i a_j b_{j-i} \right) D^i. \end{aligned}$$

Über  $\text{GF}(2)$  gilt zum Beispiel

$$\begin{aligned} (D^3 + D^2 + D + 1) + (D^2 + 1) &= D^3 + D, \\ (D^3 + D^2 + D + 1) \cdot (D^2 + 1) &= D^5 + D^4 + D + 1. \end{aligned}$$

Nimmt man noch das Polynom vom Grad 0, das alle Koeffizienten  $a_n = 0$  hat, als neutrales Element hinzu, so kann leicht überprüft werden, daß die Menge der Polynome über  $\text{GF}(m)$  bezüglich  $+$  eine Abelsche Gruppe bildet. Außerdem ist die Multiplikation kommutativ, assoziativ und distributiv mit der Addition. Allerdings fehlen im allgemeinen die Inversen bezüglich der Multiplikation, so daß kein Körper gemäß Definition 7.1 vorliegt.

Für je zwei Polynome  $a(D)$  und  $b(D)$  existieren nach dem Euklidischen Algorithmus eindeutige Polynome  $c(D)$  und  $r(D)$  mit  $\text{Grad } r(D) \leq \text{Grad } b(D)$ , derart daß

$$a(D) = b(D) \cdot c(D) + r(D).$$

$r(D)$  heißt Rest bei Division von  $a(D)$  durch  $b(D)$ . Er wird in Analogie zum Modulo-Operator der Division in  $\mathbb{Z}$  mit

$$r(D) = a(D) \pmod{b(D)}$$

bezeichnet.

$c(D)$  und  $r(D)$  können mit dem üblichen Divisionsalgorithmus berechnet werden. Zum Beispiel gilt über  $\text{GF}(2)$  für  $a(D) = D^5 + D^4 + D^2 + 1$  und  $b(D) = D^2 + 1$ , daß

$$\begin{array}{r} (D^5 + D^4 + D^2 + 1) : (D^2 + 1) = (D^3 + D^2 + D) \\ \underline{D^5 + D^3} \\ D^4 + D^3 \\ \underline{D^4 + D^2} \\ D^3 \\ \underline{D^3 + D} \\ D + 1 \end{array}$$

Es folgt also  $(D^5 + D^4 + D^2 + 1) = (D^2 + 1) \cdot (D^3 + D^2 + D) + (D + 1)$  mit  $c(D) = D^3 + D^2 + D$  und  $r(D) = D + 1$ .

Man sagt, das Polynom  $b(D)$  *teilt*  $a(D)$ , wenn ein Polynom  $c(D)$  existiert mit  $a(D) = b(D) \cdot c(D)$ . Ein Polynom heißt *irreduzibel*, wenn keine Polynome  $b(D)$  und  $c(D)$  vom Grad größer gleich 1 existieren mit  $a(D) = b(D) \cdot c(D)$ .

Der Grad einer Summe von Polynomen ist kleiner oder gleich dem Maximum der Grade der Summanden. Der Grad eines Produkts ist jedoch größer als dieses Maximum. Damit ist die Menge der Polynome vom Grad kleiner oder gleich  $n$  nicht abgeschlossen bezüglich der Multiplikation. Die folgende Polynommultiplikation garantiert die Abgeschlossenheit; außerdem sind unter ihr die Körperaxiome erfüllt.

**Satz 7.1** Sei  $q(D)$  ein irreduzibles Polynom vom Grad  $n$  über  $\text{GF}(M)$ . Die Menge der Polynome vom Grad kleiner oder gleich  $n - 1$  über  $\text{GF}(M)$  bildet unter Polynomaddition und der Multiplikation  $*$ , definiert durch

$$a(D) * b(D) = a(D) \cdot b(D) \pmod{q(D)},$$

einen Körper.

Die Menge der Polynome vom Grad kleiner gleich  $n - 1$  über  $\text{GF}(M)$  enthält offensichtlich  $m^n$  Elemente. Ist  $m = p$  eine Primzahl, kann mit Satz 7.1 ein Körper  $\text{GF}(p^n)$  der Mächtigkeit  $p^n$  als Polynomreste über  $\mathbb{Z}_p$  konstruiert werden. Dies sind bis auf Isomorphie auch die einzigen endlichen Körper. Es kann gezeigt werden, daß endliche Körper der Mächtigkeit  $m$  nur dann existieren, wenn  $m = p^n$  eine Primzahlpotenz ist, und daß endliche Körper gleicher Mächtigkeit bis auf Ummumerierung der Elemente übereinstimmen.

**Beispiel 7.1** Die Polynome vom Grad  $\leq 2$  über  $\text{GF}(2)$  sind

$$\{a_0, a_1, \dots, a_7\} = \{0, 1, D, D+1, D^2, D^2+1, D^2+D, D^2+D+1\}.$$

$D^3 + D + 1$  ist ein irreduzibles Polynom über  $\text{GF}(2)$ . Addition und Multiplikation in  $\text{GF}(8)$  werden gemäß Satz 7.1 durch die folgenden Tafeln definiert.

+	$a_0$	$a_1$	$a_2$	$a_3$	$a_4$	$a_5$	$a_6$	$a_7$
$a_0$	$a_0$	$a_1$	$a_2$	$a_3$	$a_4$	$a_5$	$a_6$	$a_7$
$a_1$	$a_1$	$a_0$	$a_3$	$a_2$	$a_5$	$a_4$	$a_7$	$a_6$
$a_2$	$a_2$	$a_3$	$a_0$	$a_1$	$a_6$	$a_7$	$a_4$	$a_5$
$a_3$	$a_3$	$a_2$	$a_1$	$a_0$	$a_7$	$a_6$	$a_5$	$a_4$
$a_4$	$a_4$	$a_5$	$a_6$	$a_7$	$a_0$	$a_1$	$a_2$	$a_3$
$a_5$	$a_5$	$a_4$	$a_7$	$a_6$	$a_1$	$a_0$	$a_3$	$a_2$
$a_6$	$a_6$	$a_7$	$a_4$	$a_5$	$a_2$	$a_3$	$a_0$	$a_1$
$a_7$	$a_7$	$a_6$	$a_5$	$a_4$	$a_3$	$a_2$	$a_1$	$a_0$

·	$a_0$	$a_1$	$a_2$	$a_3$	$a_4$	$a_5$	$a_6$	$a_7$
$a_0$	$a_0$	$a_0$	$a_0$	$a_0$	$a_0$	$a_0$	$a_0$	$a_0$
$a_1$	$a_0$	$a_1$	$a_2$	$a_3$	$a_4$	$a_5$	$a_6$	$a_7$
$a_2$	$a_0$	$a_2$	$a_4$	$a_6$	$a_3$	$a_1$	$a_7$	$a_5$
$a_3$	$a_0$	$a_3$	$a_6$	$a_5$	$a_7$	$a_4$	$a_1$	$a_2$
$a_4$	$a_0$	$a_4$	$a_3$	$a_7$	$a_6$	$a_2$	$a_5$	$a_1$
$a_5$	$a_0$	$a_5$	$a_1$	$a_4$	$a_2$	$a_7$	$a_3$	$a_6$
$a_6$	$a_0$	$a_6$	$a_7$	$a_1$	$a_5$	$a_3$	$a_2$	$a_4$
$a_7$	$a_0$	$a_7$	$a_5$	$a_2$	$a_1$	$a_6$	$a_4$	$a_3$

■



## Literaturverzeichnis

- [1] N. Abramson. *Information Theory and Coding*. McGraw-Hill, New York, 1963.
- [2] J.B. Anderson and S. Mohan. *Source and Channel Coding*. Kluwer Academic Publishers, Boston, 1991.
- [3] R. Ash. *Information Theory*. Interscience Publishers (Wiley) und Dover Publications (corrected republication 1990), New York, 1965.
- [4] P. Billingsley. *Ergodic Theory and Information*. Wiley, New York, 1965.
- [5] R.E. Blahut. *Principles and Practice of Information Theory*. Addison-Wesley, Reading, Mass., 1987.
- [6] K.W. Cattermole. *Statistische Analyse und Struktur von Information*. VCH Verlagsgesellschaft, Weinheim, 1990.
- [7] G.J. Chaitin. *Algorithmic Information Theory*. Cambridge University Press, Cambridge, 1990.
- [8] T.M. Cover and J.A. Thomas. *Elements of Information Theory*. Wiley, New York, 1991.
- [9] I. Csizár and J. Körner. *Coding Theorems for Discrete Memoryless Systems*. Academic Press, New York, 1981.
- [10] R.G. Gallager. *Information Theory and Reliable Communication*. Wiley, New-York, 1968.
- [11] C.M. Goldie and R.G. Pinch. *Communication Theory*. Cambridge University Press, Cambridge, 1991.
- [12] T. Grams. *Codierungsverfahren*. BI-Verlag, Mannheim, 1986.
- [13] R.M. Gray. *Entropy and Information Theory*. Springer, New York, 1990.
- [14] H. Grell. *Arbeiten zur Informationstheorie, Bd I-IV*. VEB-Verlag, Berlin, 1960-1963.

- [15] S. Guiasu. *Information Theory with Applications*. McGraw-Hill, New-York, 1977.
- [16] R.W. Hamming. *Coding and Information Theory*. Prentice-Hall, Englewood Cliffs, 1980.
- [17] W. Heise and P. Quattrocchi. *Informations- und Codierungstheorie*. Springer-Verlag, Berlin, 1983.
- [18] E. Henze and H.H. Homuth. *Einführung in die Informationstheorie*. Vieweg, Braunschweig, 1974.
- [19] R. Hill. *A First Course in Coding Theory*. Clarendon Press, Oxford, 1986.
- [20] A.M. Jaglom and I.M. Jaglom. *Wahrscheinlichkeit und Information*. Verlag Harri Deutsch, Thun, 1984.
- [21] T. Kameda and K. Weihrauch. *Einführung in die Codierungstheorie*. BI-Verlag, Mannheim, 1973.
- [22] S. Kullback, J.C. Keegel, and J.H. Kullback. *Topics in Statistical Information Theory*. Springer-Verlag, Berlin, 1987.
- [23] A.W. Marshall and I. Olkin. *Inequalities: Theory of Majorization and its Applications*. Academic Press, New York, 1979.
- [24] R. Mathar and D. Pfeifer. *Stochastik für Informatiker*. Teubner-Verlag, Stuttgart, 1990.
- [25] K. Mehlhorn. *Datenstrukturen und effiziente Algorithmen*. Teubner-Verlag, Stuttgart, 1988.
- [26] O. Mildenerger. *Informationstheorie und Codierung*. Vieweg-Verlag, Wiesbaden, 1992.
- [27] G. Raisbeck. *Informationstheorie – Eine Einführung für Naturwissenschaftler und Ingenieure*. Oldenbourg-Verlag, München, 1970.
- [28] A. Renyi. *Wahrscheinlichkeitsrechnung mit einem Anhang über Informationstheorie*. VEB-Verlag, Berlin, 1979.
- [29] F.M. Reza. *An Introduction to Information Theory*. McGraw-Hill, New York, 1961.
- [30] S. Roman. *Coding and Information Theory*. Springer-Verlag, New York, 1992.

- [31] E. Schultze. *Einführung in die mathematischen Grundlagen der Informationstheorie*. Springer, Berlin, 1969.
- [32] R.H. Schulz. *Codierungstheorie, eine Einführung*. Vieweg-Verlag, Braunschweig, 1991.
- [33] C.E. Shannon. A mathematical theory of communication. *Bell Systems Technical Journal*, 27:379–423, 623–656, 1948.
- [34] C.E. Shannon and W. Weaver. *The Mathematical Theory of Communication*. University of Illinois Press, Urbana, 1962.
- [35] F. Topsøe. *Informationstheorie*. Teubner, Stuttgart, 1974.
- [36] D. Welsh. *Codes and Cryptography*. Oxford University Press, Oxford, 1989.
- [37] I. Wolfowitz. *Coding Theorems of Information Theory*. Springer-Verlag, Berlin, 1961.
- [38] P.M. Woodward. *Probability and Information Theory with Applications to Radar*. McGraw-Hill, New York, 1953.



# Sachverzeichnis

- ( $N, M, d$ )-Kode, 128
- Abelsche Gruppe, 147
- abgeschlossen, 147
- AEP, 50
- äquivalente Codes, 132
- Algorithmus von Dykstra, 143
- Anfangsverteilung, 17
- Anführer, 134
- aperiodisch, 19
- Assoziativgesetz, 147
- assoziierte Funktion, 81
- asymptotic equipartition property, 50
- asymptotische Gleichverteilungseigenschaft, 48, 50
- Ausgabealphabet, 90
  
- Baum, binärer, 140
- bedingt stochastisch unabhängig, 117
- bedingte Verteilung, 15
- bedingter Erwartungswert, 15
- binärer Suchbaum, 70
- binärer symmetrischer Kanal, 29, 39, 105
- binary symmetric channel, 29
- Binomialverteilung, 79
- Blockfragestrategie, 46
- Blockkode, 53, 60
- Blockkodierung, 60
- Boltzmann-Verteilung, 41
- BSC, 29, 92, 95, 105
  - Fundamentalsatz für, 105
  
- constraint length, 136
- coset leader, 134
- cosets, 134
  
- Data Processing Theorem, 118
- Datenkompression, 43
  
- Datenrate, 106, 126
- Dekodierfehler, 100
- Dekodierregel, 98
  - äquivalente, 101
  - Eindeutigkeit, 101
- Dekodierung linearer Codes, 134
- DGQ, 48
- diskrete Zufallsvariable, 12
- Distributivgesetz, 147
- Divisionsalgorithmus, 149
- DMC, 90
- DMQ, 81
- doppelt stochastisch, 25
- DSQ, 77
- dynamische Optimierung, 144
  
- e.d., 53
- eindeutig dekodierbar, 53
- Eingabealphabet, 90
- Einpunktverteilung, 79
- Entropie, 25
  - als Erwartungswert, 26
  - axiomatische Charakterisierung, 36
  - bedingte, 27, 31
  - einer diskreten stationären Quelle, 78
  - einer stationären Markoff-Quelle, 83
  - eines Zufallsvektors, 27
  - pro Quellbuchstabe, 78
  - reale, 63
- Entropieungleichungen, 32
- Erwartungswert, 14
- Euklidischer Algorithmus, 148
- Eulersche  $\varphi$ -Funktion, 20
  
- Faltungskode, 127

- Faltungskodierer, 8, 136
- Fano-Kode, 66
- Fano-Ungleichung, 114
- fehlerkorrigierende Codes, 127
- field, 147
- Folgen von Zufallsvariablen, 12
- Fragestrategie, 45
- freier Abstand, 146
- Fundamental Theorem of Information Theory, 103
- Fundamentalsatz
  - für BSC, 105
  - schwache Umkehrung, 121
  - von Shannon, 113
- Galois-Feld, 147
- gemeinsame Verteilung, 12
- Generatormatrix, 131, 137
- Gesetz großer Zahlen
  - schwaches, 16
  - starkes, 16
- Gewicht, 132
- $GF(m)$ , 147
- Gilbert-Varshamov-Schranke, 128
- Hamming-Distanz, 102, 142
- Hamming-Kode, 134
- Hamming-Schranke, 128
- homogen, 17
- Huffman-Kode, 66
- Huffman-Verfahren, 66
- ideal observer, 99
- Indikatorfunktion, 40
- Informationsgehalt, 26
- Informationsgewinn, 22
- Inputverteilung, 94, 99
  - ungünstigste, 122
- inverses Element, 147
- irreduzibel, 17, 149
- Kanal
  - binärer symmetrischer, 92
  - diskreter gedächtnisloser, 90
- Kanal, binärer symmetrischer, 95
- Kanalkapazität, 94
- Kanalmatrix, 91, 94
  - eines Kaskaden-DMC, 117
- Kaskaden-DMC, 116
- Kaskadenkanal, 116
- katastrophal, 146
- Kode, 53
  - absolut optimaler, 61
  - gestutzter, 145
  - linearer, 131
  - minimaler Abstand, 127
  - optimaler, 62
  - perfekter, 129
- Kodealphabet, 51, 52
- Kodebaum, 60, 140
  - binärer, 57
- Kodewort, 53
- Kodewortlänge, 47
  - erwartete, 58
  - mittlere, 44
- Kodewortmenge, 53
- Kodierung, 51, 53
  - bei linearen Codes, 132
- Kodierungstheorie, 127
- Körper, 147
  - endlicher, 147
- Kommutativgesetz, 147
- konkav, 26, 125
- Kontrollmatrix, 132
- Lagrange-Ansatz, 62
- Lagrangefunktion, 62
- linearer Kode, 127
- Lorenzkurve, 23
- majorisiert, 23
- Majorisierung, 23, 25
- Markoff-Kette, 17
  - Limesverteilung, 19
- maximum-likelihood
  - decoding rule, 101
  - Dekodierung, 101
- MD-Dekodierung, 102, 140

- ME-Dekodierung, 99
- Meßbarkeit, 11
- Metrik, 102
- minimaler Abstand, 127
- minimum distance decoding, 102
- minimum error rule, 99
- ML-Dekodierung, 101, 144
- modulo-2-Arithmetik, 136
  
- natürlicher Logarithmus, 26
- Nebenklassen, 134
- neutrales Element, 147
- Noiseless Coding Theorem, 58
- Noisy Coding Theorem, 103
  
- parity check matrix, 132
- partielle Ordnung, 23
- Partition, 98
- perfekter Kode, 129
- Permutation, 23
- Permutationsmatrix, 132
- PF-Kode, 54
- Polynom, 148
  - Summe und Produkt, 148
- Präfix, 54
- präfixfrei, 54
- Präordnung, 23
- Produktkanal, 124
  
- Quellalphabet, 51, 52
- Quelle
  - diskrete gedächtnislose, 48
  - diskrete Markoff-, 81
  - diskrete stationäre, 77
  - ergodische, 50
  - stationäre, 50
- Quellenrate, 110
- Quellwörter, typische, 49
  
- Randverteilung, 18
- Rückgriftiefe, 136
  
- Satz von Feinstein, 48
- Satz von Kraft, 55
- Satz von McMillan, 55
  
- Schaltdiagramm, 136
- Schieberegister, 136
- schwache Umkehrung des Fundamentalsatzes, 121
- Shannon-Fano-Kode, 60
- Shannonsche Ungleichung, 33
- Shannonscher Fundamentalsatz, 103, 113
- Spalierdiagramm, 139
- Standardform einer Generatormatrix, 131
- Standardmodell der Informationsverarbeitung, 7, 43
- stationär, 14
- stationäre Markoff-Kette, 19
- stationäre Verteilung, 18
- stochastisch unabhängig, 13
  - bedingt, 117
- stochastischer Vektor, 18
- Suchbaum
  - ausgewogener, 75
  - binärer, 70
- Summenkanal, 124
- Syndrom, 134
- Synentropie, 28, 94
  
- teilt, 149
- total abhängig, 33
- Totalordnung, 24
- Transinformation, 28, 31, 94
  - als Erwartungswert, 29
  - bedingte, 28
- Trellis-Diagramm, 139
- Trennzeichen, 55
- truncated, 145
  
- Übergangsmatrix, 17
- Übergangswahrscheinlichkeit, 17
- Übertragungsrate, 106
- Übertragungswahrscheinlichkeiten, 90
- Unbestimmtheit, 22
- Ungleichungen, 32
- unifilar, 82
  
- Verteilung, 12

Viterbi-Algorithmus, 127, 141

Wahrscheinlichkeit  
für Dekodierfehler, 100  
für Fehldekodierung, 104  
bedingte, 104  
maximale, 104  
mittlere, 104

Wahrscheinlichkeitsverteilung, 12  
wechselseitige Information, 29

Zeilenvektor, 18  
Zufallskode, 106  
Zufallsvariable, 11  
Zufallsvektor, 12  
Zugriffsverteilung, 71  
Zugriffszeit, 70  
Zustandsänderungsdiagramm, 139  
Zustandsraum, 17