

---

Skripte zur Methodenlehre, V

**Einführung in die statistische Analyse  
von Zustandsverläufen**

G. Rohwer

U. Pötter

Version 3

August 2008

---

**Vorbemerkung.** Bei den "Skripten zur Methodenlehre" handelt es sich um Texte, die als Leitfäden für Seminare zur sozialwissenschaftlichen Methodenlehre dienen sollen. Das vorliegende Skript beschäftigt sich mit statistischen Methoden zur Analyse von Längsschnittdaten. Dabei orientieren wir uns an Anwendungen dieser Methoden in der empirischen Sozialforschung, bei denen es in erster Linie um eine Untersuchung von Lebensverläufen geht und Daten dementsprechend in Gestalt von Zustandsverläufen gegeben sind. Darauf bezieht sich die im ersten Kapitel eingeführte Terminologie.

Der Text enthält zahlreiche Übungsaufgaben, die während der Bearbeitung des Stoffes gelöst werden sollten. Die meisten Aufgaben können mit Bleistift und Papier gelöst werden; für einige Aufgaben ist die Verwendung eines Taschenrechners hilfreich.

Für Anwendungen der Methoden in der empirischen Sozialforschung, bei denen man es meist mit größeren Datensätzen zu tun hat, muss man allerdings Computer und geeignete Statistikprogramme verwenden. Der Text enthält deshalb einen Anhang, anhand dessen man das Programm R kennenlernen kann, mit dem die meisten Fragestellungen der Verlaufsdatenanalyse bearbeitet werden können. Mit den Aufgaben dieses Anhangs kann man sich entweder parallel zur Behandlung des Haupttextes oder in einem sich anschließenden Workshop beschäftigen.

Über die hier behandelten statistischen Methoden gibt es eine sehr umfangreiche Literatur. Wer sein Wissen über die statistischen Aspekte der Methoden vertiefen möchte, sei auf Lawless (1982) und Cox und Oakes (1984) sowie Kalbfleisch und Prentice (2002) hingewiesen. Weiterführende Aspekte behandeln Andersen et al. (1993) und Martinussen und Scheike (2006). Für eine weiterführende Diskussion von Anwendungen in der empirischen Sozialstrukturforschung sei auf Blossfeld und Rohwer (1995) verwiesen. Therneau und Grambsch (2000) und Tableman und Kim (2004) geben weitere Anregungen für die Umsetzung mit R.

## Inhaltsverzeichnis

1	Einführung	1
1.1	Objekte und Lebensverläufe	1
1.2	Verhalten und Zustände	1
1.3	Der Zustandsraum	1
1.4	Biographieschema	2
1.5	Mehrdimensionale Zustandsräume	2
1.6	Die Zeitachse	3
1.7	Ereignisse als Zustandswechsel	3
1.8	Kalenderzeit und Prozeßzeit	3
1.9	Verlaufdiagramme	4
1.10	Kohorten	4
2	Statistische Beschreibungen	5
2.1	Statistische Variablen	5
2.2	Zustandsvariablen	6
2.3	Partielle Lebensverläufe	6
2.4	Statistische Verteilungen	7
2.5	Zustandsverteilungen	9
3	Verweildauerverteilungen	11
3.1	Episoden	11
3.2	Statistischer Begriffsrahmen	12
3.3	Ein möglicher Folgezustand	13
3.4	Mehrere mögliche Folgezustände	15
4	Zensierte Beobachtungen	17
4.1	Rechts zensierte Beobachtungen	17
4.2	Berechnung von Survivorfunktionen	18
4.3	Das Kaplan-Meier-Verfahren	18
4.4	Mehrere Folgezustände	20
4.5	Selbst-Konsistenz	21
5	Regressionsmodelle für Zustände	23
5.1	Der Modellansatz	23
5.2	Spekulation und Empirie	24
5.3	Modelle für zwei Zustände	25
5.4	Modelle mit Kovariablen	26
5.5	Binäre Logitmodelle	27
5.6	Maximum-Likelihood-Schätzung	28
6	Modelle für Verweildauern	33
6.1	Zeitkonstante Raten	33
6.2	Weibull-Verteilung	34

ii	INHALTSVERZEICHNIS	0
6.3	Loglogistische Verteilung . . . . .	35
6.4	Lognormal-Verteilung . . . . .	36
6.5	Mehrere Zielzustände . . . . .	37
6.6	Mischungen . . . . .	37
7	Ratenmodelle mit Kovariablen	39
7.1	Das Exponentialmodell . . . . .	39
7.2	Parameterschätzungen . . . . .	40
7.3	Ein allgemeiner Modellansatz . . . . .	42
7.4	Mehrere Folgezustände . . . . .	42
7.5	Pseudo-Residuen . . . . .	44
8	Zeitveränderliche Kovariablen	47
8.1	Konditionale Survivorfunktionen . . . . .	47
8.2	Reformulierte Likelihoodfunktion . . . . .	47
8.3	Zeitveränderliche Indikatorvariablen . . . . .	48
8.4	Episodensplitting . . . . .	49
A	Übungen mit R	51
	Literatur	58

# 1 Einführung

In diesem Kapitel besprechen wir Grundzüge des begrifflichen Rahmens, der in den nachfolgenden Kapiteln vorausgesetzt wird.

## 1.1 Objekte und Lebensverläufe

Wir beziehen uns zunächst ganz allgemein auf Objekte. Jedes Objekt existiert in der Form eines Lebensverlaufs: Es wird geboren, dann macht es einen gewissen Entwicklungsprozeß durch, und schließlich stirbt es. Unser Ziel ist es, uns mit einigen statistischen Begriffen und Modellen zu beschäftigen, die vorgeschlagen worden sind, um Lebensverläufe beschreiben und über ihre Entwicklung nachdenken zu können.

Wie wir sehen werden, sind diese Begriffe und Modelle sehr allgemein. Bei ihrer Verwendung in der empirischen Sozialforschung ist darauf zu achten, daß wir es dann meistens mit spezifischen Objekten zu tun haben, nämlich sozialen Akteuren (sowohl individuelle als auch korporative Akteure), die selbst Anteil daran nehmen, wie sich ihre Lebensverläufe entwickeln. Wir werden die Objekte, mit denen wir uns beschäftigen, in allgemeiner Weise als *Individuen* bezeichnen.

## 1.2 Verhalten und Zustände

Man kann Individuen unter zwei komplementären Betrachtungsweisen vergegenständlichen: als Objekte, die sich *verhalten* können, und als Objekte, die sich *in wechselnden Zuständen befinden* können. Der in diesem Text behandelte Ansatz geht von der zweiten Betrachtungsweise aus: Lebensverläufe von Individuen werden als Folgen von Zuständen konzipiert.

## 1.3 Der Zustandsraum

Ausgangspunkt ist also die Konzeption eines Zustandsraum. Wir setzen voraus, daß es stets nur eine endliche Menge möglicher Zustände gibt und bezeichnen den Zustandsraum mit dem Symbol  $Y$ . Der Lebensverlauf eines Individuums besteht dann in einer Folge von Zuständen aus dem vorgegebenen Zustandsraum. Die Aufenthaltsdauer in den Zuständen ist unbestimmt, und es wird auch nicht vorausgesetzt, daß alle Zustände durchlaufen werden müssen. Aus diesem Ansatz folgt, daß der hier verwendete Begriff des Lebensverlaufs wesentlich davon abhängt, welcher Zustandsraum vorausgesetzt wird.

Man beachte, daß ein Zustandsraum *eindeutig* sein muß. Damit ist gemeint, daß sich die zu betrachtenden Individuen zu jedem Zeitpunkt in genau einem der möglichen Zustände befinden müssen.

Wir sprechen von einem *vollständigen Zustandsraum*, wenn der Zustandsraum insbesondere die beiden Quasi-Zustände *noch nicht geboren* und *gestorben* umfaßt. Um Lebensverläufe vollständig zu erfassen, ist ein vollständiger Zustandsraum erforderlich.

## 1.4 Biographieschema

Unter einem *Biographieschema* verstehen wir die Festlegung einer Menge möglicher (ggf. unvollständiger) Lebensverläufe in einem Zustandsraum. Ein Biographieschema kann graphisch durch ein *Zustandsdiagramm* veranschaulicht werden. Es besteht dann aus einem gerichteten Graphen, in dem die möglichen Zustände durch Knoten, die möglichen Übergänge durch gerichtete Kanten repräsentiert werden.

AUFGABE 1.1 Konzipieren Sie einen vollständigen Zustandsraum für die Erfassung von Erwerbsverläufen, der die folgenden Zustände unterscheidet: (1) erwerbstätig, (2) arbeitslos, (3) weder erwerbstätig noch arbeitslos.

AUFGABE 1.2 Konzipieren Sie mit dem Zustandsraum aus Aufgabe 1.1 ein Biographieschema.

## 1.5 Mehrdimensionale Zustandsräume

Die Konzeption eines Zustandsraum muß durch den Modellkonstrukteur vorgegeben werden. Dies hängt davon ab, welche Aspekte realer Lebensverläufe erfaßt werden sollen, z.B. Erwerbsverläufe oder Ausbildungsverläufe oder Beziehungsverläufe. Man kann mehrere solcher Aspekte durch einen mehrdimensionalen Zustandsraum repräsentieren. Als symbolische Form eines  $m$ -dimensionalen Zustandsraum hat man dann

$$\tilde{Y} = \tilde{Y}_1 \times \cdots \times \tilde{Y}_m$$

Andererseits ist es möglich, stattdessen einen einfachen (eindimensionalen) Zustandsraum zu verwenden, bei dem jede mögliche Kombination von Zuständen in  $\tilde{Y}_1, \dots, \tilde{Y}_m$  als ein gesonderter Zustand im kombinierten Zustandsraum  $\tilde{Y}$  repräsentiert wird.

AUFGABE 1.3 Konzipieren Sie einen vollständigen Zustandsraum für die Zustände: (1) unverheiratet, (2) verheiratet. Bilden Sie dann aus diesem und dem in Aufgabe 1.1 konzipierten Zustandsraum einen zweidimensionalen Zustandsraum.

AUFGABE 1.4 Konzipieren Sie für den zweidimensionalen Zustandsraum aus Aufgabe 1.3 ein Biographieschema.

## 1.6 Die Zeitachse

Die Grundvorstellung besteht darin, Lebensverläufe als ein zeitlich geordnetes „Durchwandern“ von Zustandsräumen aufzufassen. Es ist also erforderlich, sich explizit auf eine Zeitachse zu beziehen. Hierfür gibt es zwei Möglichkeiten.

- Wir können uns eine Zeitachse als eine Folge von Zeitstellen vorstellen, z.B. Stunden, Tage, Wochen oder Monate. Man spricht dann von einer *diskreten Zeitachse*, und zur numerischen Repräsentation können die natürlichen Zahlen verwendet werden.
- Wir können uns eine Zeitachse als einen kontinuierlichen Zeitfluß vorstellen, d.h. von der Annahme ausgehen, daß Zeitstellen beliebig teilbar sind. Man spricht dann von einer *kontinuierlichen* oder *stetigen Zeitachse* und verwendet zur numerischen Repräsentation die reellen Zahlen.

Wir werden zunächst von einer diskreten Zeitachse ausgehen. Dies hat den Vorteil, daß von einer *Folge* von Zeitstellen gesprochen werden kann. Statistische Modelle verwenden jedoch häufig eine kontinuierliche Zeitachse, so daß wir uns später auch dieser Vorstellung bedienen werden.

## 1.7 Ereignisse als Zustandswechsel

Wir haben bisher Lebensverläufe als Folgen von Zuständen betrachtet, wobei die Aufenthaltsdauer in jedem der möglichen Zustände von unterschiedlicher Dauer sein kann. Stattdessen kann man das Augenmerk auch auf die Zustandswechsel richten, also auf die Übergänge von einem gegebenen in einen neuen Zustand. Diese Zustandswechsel werden auch *Ereignisse* genannt.<sup>1</sup>

AUFGABE 1.5 Geben Sie eine Liste aller Ereignisse an, die in dem Biographieschema, das in Aufgabe 1.2 konzipiert wurde, möglich sind.

## 1.8 Kalenderzeit und Prozeßzeit

Wenn man sich auf reale Individuen und deren Lebensverläufe beziehen will, muß man zunächst immer von einer Kalenderzeitachse ausgehen. Man spricht gelegentlich auch von einer historischen Zeitachse. Für die Modellbildung verwendet man stattdessen meistens eine *Prozeßzeitachse*. Es handelt sich um eine Zeitachse, bei der der Nullpunkt durch das Eintreten eines Ereignisses definiert wird. Zum Beispiel könnte man eine Prozeßzeitachse konzipieren, die mit der Geburt beginnt oder mit der Aufnahme eines Studiums oder dem Beginn einer Eheschließung.

<sup>1</sup> Es sei angemerkt, daß das Wort ‘Ereignis’ dadurch eine spezifische Bedeutung bekommt. Wer sich für eine gründlichere Diskussion interessiert, sei auf Galton (1994) verwiesen.

**Box 1.1** Datensatz 1

ID	Geburt	Beginn des Studiums	Ende des Studiums
1	1970	1990	1995
2	1975	1994	1999
3	1973	1991	1996
4	1970	1989	1995
5	1975	1993	1999
6	1973	1993	1996
7	1970	1988	1995
8	1975	1995	1999
9	1973	1992	1997

AUFGABE 1.6 Betrachten Sie die Daten in Box 1.1. Konzipieren Sie dazu einen Zustandsraum und ein Biographieschema. Stellen Sie die Daten auf einer Prozeßzeitachse dar, deren Zeiteinheiten Jahre sind und die mit dem Beginn des Studiums beginnt.

## 1.9 Verlaufsdiagramme

Ein *Verlaufsdiagramm* ist ein Diagramm, bei dem die horizontale Achse die Zeitachse und die vertikale Achse den Zustandsraum repräsentiert. Dabei kann die Zeitachse entweder eine Kalenderzeitachse oder eine Prozeßzeitachse sein. Solche Diagramme sind oft nützlich, um exemplarisch einzelne oder auch mehrere Verläufe darzustellen.

AUFGABE 1.7 Stellen Sie die ersten drei Verläufe aus dem Datensatz in Box 1.1 zunächst in einem Verlaufsdiagramm dar, bei dem die Zeitachse eine Kalenderzeitachse ist, dann in einem Verlaufsdiagramm, bei dem die Zeitachse die Prozeßzeitachse ist, die mit dem Beginn des Studiums beginnt.

## 1.10 Kohorten

In der empirischen Sozialforschung wird oft der Begriff *Kohorte* verwendet, um eine Menge von Individuen zu bezeichnen, die ein Ereignis eines bestimmten Typs in der gleichen Kalenderzeitstelle erfahren haben. Zum Beispiel bilden alle Individuen, die im Jahr 1970 geboren worden sind, eine Geburtskohorte. Dabei muß natürlich angegeben werden, auf welche Grundgesamtheit von Individuen man sich beziehen möchte. Und außerdem muß die Dauer der Zeitstelle fixiert werden, die zur Definition von Kohorten dienen soll.

AUFGABE 1.8 Betrachten sie die Daten in Box 1.1. Wieviel Geburtskohorten gibt es? Erstellen Sie eine Tabelle, in der die Individuen den Geburtskohorten zugeordnet werden. Machen Sie dann das gleiche für die Kohorten von Studienanfängern.

## 2 Statistische Beschreibungen

In diesem Kapitel beginnen wir mit einer Diskussion der Frage, wie Lebensverläufe beschrieben werden können. Zwei komplementäre Betrachtungsweisen können eingenommen werden. Man kann versuchen, Lebensverläufe spezifischer Individuen ins Auge zu fassen und in ihrer jeweils einmaligen Entwicklung zu beschreiben. Andererseits kann man eine *vergleichende Betrachtungsweise* einnehmen. Dies setzt voraus, daß man sich auf eine Mehrzahl vergleichbarer Lebensverläufe beziehen kann. Vergleichbarkeit ist allerdings kein Merkmal, das Lebensverläufen „an und für sich“ zukommt oder nicht zukommt, sondern Gesichtspunkte für einen Vergleich kommen stets durch den Sozialforscher zustande. Er ist es, der Lebensverläufe vergleichen möchte und dafür die ihm wichtig erscheinenden Gesichtspunkte definiert.

Für den hier zu behandelnden statistischen Ansatz kommen die Gesichtspunkte für einen Vergleich von Lebensverläufen durch die Definition eines Biographieschemas zustande. Wir nehmen im folgenden an, daß ein Biographieschema definiert worden ist und daß man sich auf eine vorgegebene Menge von Individuen beziehen kann, deren Lebensverläufe (meistens nur ausschnittshaft) durch das vorgegebene Biographieschema verglichen werden können. Wir bezeichnen diese Menge von Individuen mit dem Symbol  $\Omega$ .

Diese Voraussetzungen erlauben es, Lebensverläufe mit statistischen Begriffen zu beschreiben. Was damit gemeint ist, wird sogleich deutlicher werden, wenn wir die beiden Grundbegriffe, *statistische Variable* und *statistische Verteilung*, eingeführt haben.

## 2.1 Statistische Variablen

Eine statistische Variable ist eine Abbildung (auch Funktion genannt), die jedem Individuum aus einer vorgegebenen Menge einen bestimmten Wert in einem Merkmalsraum zuordnet. Zur symbolischen Repräsentation verwenden wir die Schreibweise

$$X : \Omega \longrightarrow \tilde{X}$$

Hier ist  $X$  eine statistische Variable, die jedem Individuum  $\omega \in \Omega$  einen Merkmalswert  $X(\omega)$  aus dem Merkmalsraum  $\tilde{X}$  zuordnet. Wir setzen voraus, daß es für den Merkmalsraum eine numerische Repräsentation gibt. In dieser Einführung betrachten wir zwei Arten numerischer Repräsentationen. Wenn  $\tilde{X}$  durch eine Teilmenge der natürlichen Zahlen repräsentiert werden kann, nennen wir  $X$  eine *diskrete Variable*. Wenn  $\tilde{X}$  durch einen zusammenhängenden Teilbereich der reellen Zahlen repräsentiert werden kann, nennen wir  $X$  eine *kontinuierliche Variable*. Variablen können außerdem danach

unterschieden werden, ob es sich um einen qualitativen, ordinalen oder quantitativen Merkmalsraum handelt. Eine diskrete numerische Repräsentation kann für alle drei Arten von Variablen verwendet werden, eine kontinuierliche numerische Repräsentation ist im allgemeinen nur bei quantitativen Variablen sinnvoll.

## 2.2 Zustandsvariablen

Der Begriff der statistischen Variablen kann nun verwendet werden, um Lebensverläufe zu repräsentieren. Vorausgesetzt wird ein Biographieschema, also insbesondere ein Zustandsraum  $\tilde{Y}$  und eine Zeitachse  $\tilde{T}$ , die zunächst als eine diskrete Prozeßzeitachse angenommen wird, also

$$\tilde{T} = \{0, 1, 2, 3, \dots\}$$

Weiterhin wird eine endliche Menge von Individuen,  $\Omega$ , vorausgesetzt. Dann können die Zustände, in denen sich die Individuen befinden, durch *statistische Zustandsvariablen* erfaßt werden. Für jede Zeitstelle  $t \in \tilde{T}$  gibt es eine Zustandsvariable

$$Y_t : \Omega \longrightarrow \tilde{Y}$$

$Y_t(\omega)$  ist der Zustand, in dem sich das Individuum  $\omega \in \Omega$  in der Zeitstelle  $t$  befindet. Der Lebensverlauf jedes Individuums ist dann durch eine Folge von Zuständen:

$$(Y_0(\omega), Y_1(\omega), Y_2(\omega), \dots)$$

gegeben. Da wir angenommen haben, daß Zustandsräume stets nur eine endliche Anzahl unterschiedlicher Zustände enthalten, handelt es sich bei Zustandsvariablen stets um diskrete Variablen.

## 2.3 Partielle Lebensverläufe

Die Idee, Lebensverläufe durch Folgen von Zuständen zu repräsentieren, bereitet dann keine Schwierigkeiten, wenn es sich um vollständige Lebensverläufe handelt. Jeder Lebensverlauf mündet dann in einem Endzustand, in dem Quasi-Zustand *gestorben*. In der empirischen Sozialforschung werden jedoch meistens nur partielle Lebensverläufe untersucht. Man muß dann festlegen, welchen Teil von Lebensverläufen man betrachten möchte. Dafür gibt es zwei Möglichkeiten. In beiden Fällen beginnt man mit einem Anfangsereignis, dessen Eintritt den Beginn des partiellen Lebensverlaufs markiert; zum Beispiel: Geburt eines Individuums, Beginn eines Studiums, Eintritt in das Erwerbsleben. Dies erlaubt es, eine entsprechende Prozeßzeitachse zu definieren. Um die Entwicklung partieller Lebensverläufe auf dieser Prozeßzeitachse zu erfassen und zu vergleichen, gibt es dann zwei Möglichkeiten.

- a) Man kann einen festen Zeitraum fixieren; zum Beispiel die ersten 20 Jahre seit der Geburt, oder 6 Jahre seit dem Beginn eines Studiums. Das heißt, man fixiert auf der vorgegebenen Zeitachse eine maximale Zeitstelle  $t^*$  und erhält dann für alle Individuen aus  $\Omega$  partielle Lebensverläufe gleicher Länge, nämlich

$$(Y_0(\omega), Y_1(\omega), Y_2(\omega), \dots, Y_{t^*}(\omega))$$

Hierbei muß natürlich ein geeigneter Zustandsraum vorausgesetzt werden, der es erlaubt, alle Lebensverläufe für die vorgegebene Zeitspanne zu definieren.

- b) Eine andere Möglichkeit besteht darin, daß man die partiellen Lebensverläufe enden läßt, wenn eines aus einer vorgegebenen Menge möglicher Ereignisse eintritt. Statistiker sprechen dann manchmal von einem „absorbierenden Endzustand“, der durch das Eintreten eines solchen Ereignisses erreicht wird. Analog kann man von „absorbierenden Endergebnissen“ sprechen, die einen partiellen Lebensverlauf beenden. Um eine Menge absorbierender Endzustände zu fixieren, verwenden wir das Symbol  $\tilde{Y}^*$ . Es muß gelten, daß  $\tilde{Y}^*$  eine Teilmenge des Zustandsraums ist, also  $\tilde{Y}^* \subset \tilde{Y}$ . Jeder individuelle Lebensverlauf wird dann so lange erfaßt, bis zum ersten Mal ein Zustand in  $\tilde{Y}^*$  erreicht wird.

In der empirischen Sozialforschung wird hauptsächlich die zweite Herangehensweise verwendet. Sie hat zur Folge, daß die individuellen (partiellen) Lebensverläufe im allgemeinen eine unterschiedliche zeitliche Ausdehnung bekommen. Einige Individuen erreichen einen absorbierenden Endzustand schon nach kurzer Zeit, andere brauchen dafür länger.

## 2.4 Statistische Verteilungen

Grundlegend für statistische Beschreibungen ist der Begriff der statistischen Verteilung. Vorausgesetzt wird, daß man sich auf eine statistische Variable beziehen kann, also auf ein Kollektiv  $\Omega$  und eine Abbildung  $X$ , die jedem Mitglied des Kollektivs einen Wert in einem Merkmalsraum,  $\tilde{X}$ , zuordnet. Die Idee ist, daß man sich bei einer statistischen Beschreibung nicht für die jeweils individuellen Merkmalswerte der Mitglieder des Kollektivs interessiert, sondern nur dafür, wie sich die Mitglieder auf die möglichen Merkmalswerte verteilen. Diese Betrachtungsweise kommt gut in folgenden Worten der „Declaration on Professional Ethics“ zum Ausdruck, die vom *International Statistical Institute* erstellt worden ist:

“Statistical data are unconcerned with individual identities. They are collected to answer questions such as ‘how many?’ or ‘what proportions?’, not ‘who?’. The

identities and records of cooperating (or non-cooperating) subjects should therefore be kept confidential, whether or not confidentiality has been explicitly pledged.”<sup>1</sup>

Eine statistische Verteilung wird deshalb als eine Funktion

$$P : \mathcal{A}(\tilde{X}) \longrightarrow [0, 1]$$

definiert.  $\mathcal{A}(\tilde{X})$  ist eine Menge von Teilmengen des Merkmalsraums  $\tilde{X}$ . Dabei wird üblicherweise vorausgesetzt, daß es sich um eine Mengenalgebra handelt, die bezüglich der mengentheoretischen Basisoperationen (Vereinigung, Durchschnitt und Komplement) abgeschlossen ist. Die Elemente von  $\mathcal{A}(\tilde{X})$  werden wir *Merkmalsmengen* nennen. Die Funktion P kann dann folgendermaßen spezifiziert werden: Sie soll für jede Merkmalsmenge  $\tilde{x} \in \mathcal{A}(\tilde{X})$  den Anteil der Mitglieder von  $\Omega$  angeben, deren Merkmalswerte in dieser Merkmalsmenge liegen. Also in einer expliziten Definition:

$$P(\tilde{x}) := |\{\omega \in \Omega \mid X(\omega) \in \tilde{x}\}| / |\Omega|$$

Es ist erkennbar, wie durch diese Definition eine Bezugnahme auf individuelle Mitglieder von  $\Omega$  verschwindet und es nur noch darauf ankommt, wieviele Mitglieder an den jeweiligen Merkmalsmengen teilhaben.

Um uns flexibler auf Merkmalsmengen beziehen zu können, werden wir auch noch einige abkürzende Schreibweisen verwenden; insbesondere die folgenden:

$$P(X \in \tilde{x}) := P(\tilde{x})$$

$$P(X = x) := P(\{x\})$$

Bei quantitativen Variablen wird auch noch die Schreibweise

$$P(X \leq x) := P(\{\omega \in \Omega \mid X(\omega) \leq x\})$$

verwendet und im allgemeinen als (*kumulative*) *Verteilungsfunktion* von X bezeichnet. Die meistens verwendete Symbolik is

$$F(x) := P(X \leq x)$$

**AUFGABE 2.1** Es sei  $\Omega$  ein Kollektiv mit 10 Mitgliedern und es gebe die folgenden Merkmalswerte einer Variablen X:

$$3, 2, 3, 1, 4, 3, 1, 3, 4, 2$$

(a) Geben sie den Merkmalsraum an. (b) Definieren Sie eine Algebra von Merkmalsmengen durch die Potenzmenge des Merkmalsraums. (c) Berechnen

### Box 2.1 Datensatz 2

ID	t = 1	2	3	4	5	6
1	0	0	1	1	0	0
2	1	0	0	0	1	1
3	1	1	0	0	0	0
4	0	0	0	1	1	1
5	0	1	1	1	0	0
6	1	1	0	0	1	1

Sie die statistische Verteilung der Variablen X und geben Sie das Resultat für alle möglichen Merkmalsmengen in einer Tabelle an. (d) Nehmen Sie an, daß es sich um eine quantitative Variable handelt. Berechnen Sie dann die Verteilungsfunktion der Variablen und geben Sie das Resultat in einer Tabelle an.

**AUFGABE 2.2** Zeigen Sie, daß die Verteilungsfunktion P additiv ist, d.h. daß folgendes gilt: Wenn  $\tilde{x}_1$  und  $\tilde{x}_2$  zwei disjunkte Merkmalsmengen sind, dann gilt

$$P(\tilde{x}_1 \cup \tilde{x}_2) = P(\tilde{x}_1) + P(\tilde{x}_2)$$

## 2.5 Zustandsverteilungen

Zu überlegen ist, wie statistische Beschreibungen von Lebensverläufen entwickelt werden können. Das kann man zunächst auf ganz einfache Weise dadurch machen, daß man sich auf die Zustandsvariablen  $Y_t$  bezieht, die in Abschnitt 2.2 zur Repräsentation von Lebensverläufen eingeführt worden sind. D.h. man kann für jede Zeitstelle  $t \in \tilde{T}$  die statistische Verteilung der Zustandsvariablen  $Y_t$  berechnen. Bezieht man sich nur auf eine einzige Zeitstelle, spricht man von einer *Querschnittsverteilung*. Eine Querschnittsverteilung ergibt natürlich noch kein Bild der Entwicklung von Lebensverläufen. Eine Möglichkeit, hier weiterzukommen, besteht darin, die Querschnittsverteilungen für alle Zeitstellen der Zeitachse zu berechnen. Wir sprechen dann von *diachronen Zustandsverteilungen*. Man kann das Ergebnis in einer Tabelle oder in einem Schaubild darstellen.

**AUFGABE 2.3** Betrachten Sie die Daten in Box 2.1. Es handelt sich um Erwerbsverläufe bei 6 Individuen. Es gibt zwei Zustände: 1 = erwerbstätig, 0 = nicht erwerbstätig. Berechnen sie die diachrone Zustandsverteilung und stellen Sie diese Verteilung (a) in einer Tabelle und (b) in einem Schaubild dar.

<sup>1</sup> International Statistical Institute 1986, S. 238.

*Problematik.* Diachrone Zustandsverteilungen liefern sinnvolle statistische Beschreibungen, wenn es sich um nicht wiederholbare Zustände handelt. Zum Beispiel: 0 = noch nie verheiratet gewesen, 1 = verheiratet oder mindestens einmal verheiratet gewesen. Wenn es sich jedoch um wiederholbare Zustände handelt, wie z.B. bei Erwerbsverläufen, können diachrone Zustandsverteilungen irreführend werden, weil sie keine Rückschlüsse auf die individuellen Verläufe gestatten.

AUFGABE 2.4 Konstruieren Sie ein Beispiel, um diese Problematik sichtbar zu machen. Es soll zwei Zustände geben: 1 = arbeitslos, 0 = nicht arbeitslos. Konstruieren Sie dann zwei Varianten für 6 individuelle Verläufe, so daß der Anteil der arbeitslosen Personen in jeder Zeitstelle  $1/3$  beträgt. Bei der ersten Variante sollen 2 Personen immer, 4 Personen nie arbeitslos sein. Bei der zweiten Variante sollen alle Personen gleichmäßig von Arbeitslosigkeit betroffen sein.

### 3 Verweildauerverteilungen

In diesem Kapitel werden einige Begriffe diskutiert, die dazu dienen können, die Verweildauern in den durch ein Biographieschema vorgegebenen Zuständen statistisch darzustellen. Soweit wir uns dabei auf Daten beziehen, wird angenommen, daß vollständige Beobachtungen verfügbar sind. Die Problematik unvollständiger (zensierter) Beobachtungen wird im nächsten Teil behandelt.

#### 3.1 Episoden

Gegeben ein Biographieschema, stellen wir uns einen Lebensverlauf als ein sequentielles Durchwandern des zugehörigen Zustandsraums vor. Ein Individuum beginnt in einem gewissen Zustand und hält sich eine mehr oder weniger lange Zeit in diesem Zustand auf, dann wechselt es in einen neuen Zustand und hält sich in diesem neuen Zustand mehr oder weniger lange auf, usw. Wir können uns einen Lebensverlauf also auch als eine Folge von *Episoden* vorstellen, d.h. Aufenthaltsdauern in einem gegebenen Zustand bis ein Wechsel in einen neuen Zustand erfolgt. Eine einzelne Episode läßt sich durch vier Angaben charakterisieren:

- durch einen *Anfangszustand*, mit dessen Auftreten die Episode beginnt;
- durch einen *Endzustand*, oder *Folgezustand*, mit dessen Auftreten die Episode beendet wird;
- durch eine *Anfangszeitstelle*, die angibt, wann der Anfangszustand zum erstenmal eingenommen wird; und
- durch eine *Endzeitstelle*, die angibt, wann der Endzustand zum erstenmal eingenommen wird.

Der Begriff der Episode (verwendet wird auch gelegentlich das englische Wort *Spell*) erlaubt es, ein allgemeines Schema für die Representation von Lebensverlaufsdaten zu definieren. Box 3.1 illustriert dies Schema anhand von vier Verläufen. Der Zustandsraum umfaßt vier Zustände; 1 ist der Anfangszustand, 4 ist der (absorbierende) Endzustand. Jede Zeile in dem Schema bezieht sich auf eine Episode, und für jedes Individuum gibt es also so viele Zeilen, wie ihr Lebensverlauf Episoden aufweist. Die die Spalten benennenden Abkürzungen sind folgendermaßen zu verstehen:

- ID ist die Identifikationsnummer der Individuen,
- SN ist die laufende Nummer der Episode,

**Box 3.1** Schema für Episodendaten (Datensatz 3)

ID	SN	ORG	DES	TS	TF
1	1	1	2	0	10
1	2	2	3	10	15
1	3	3	4	15	20
2	1	1	4	0	15
3	1	1	3	0	16
3	2	3	4	19	18
4	1	1	2	0	6
4	2	2	3	6	11
4	3	3	2	11	17
4	4	2	4	17	23

- ORG ist der Anfangszustand der Episode,
- DES ist der Endzustand der Episode,
- TS ist die Anfangszeitstelle der Episode,
- TF ist die Endzeitstelle der Episode.

Wir werden ein solches Schema ein *Episodendatenschema* nennen.

**AUFGABE 3.1** Konstruieren Sie für den Datensatz 1 (Box 1.1) zunächst ein Biographieschema und stellen Sie die Daten dann in einem Episodendatenschema dar.

**AUFGABE 3.2** Konstruieren Sie für den Datensatz 2 (Box 2.1) zunächst ein vollständiges Biographieschema und stellen Sie die Daten dann in einem Episodendatenschema dar.

### 3.2 Statistischer Begriffsrahmen

Wir setzen jetzt die in Abschnitt 2 begonnene Diskussion fort, wie Lebensverläufe statistisch beschrieben werden können. Die Idee, die wir im weiteren verfolgen, besteht darin, sich zunächst auf einzelne Episoden zu konzentrieren, genauer gesagt, auf die Gesamtheit der Episoden, die in einem bestimmten, der Beschreibung vorausgesetzten Anfangszustand beginnen. Wir setzen außerdem voraus, daß wir diese Episoden auf einer Prozeßzeitachse beschreiben wollen, die mit dem Eintritt des Anfangszustands beginnt. Die Gesamtheit der Episoden, auf die wir uns beziehen wollen, kann dann durch eine zweidimensionale statistische Variable

$$(T, D)$$

repräsentiert werden.  $T$  erfaßt die Zeitdauer der Episode, d.h. die Verweildauer im Ausgangszustand, und  $D$  erfaßt den Folgezustand, dessen Eintreten die Episode abschließt.

Es ist klar, daß sich die Darstellung vereinfacht, wenn eine Episode in nur einem möglichen Folgezustand enden kann. Dann kann  $D$  nur einen möglichen Wert annehmen und braucht nicht explizit erfaßt zu werden. Oder anders gesagt, eine Episode wird dann vollständig durch ihre Dauer, den Wert von  $T$ , charakterisiert.

### 3.3 Ein möglicher Folgezustand

Wenn es nur einen möglichen Folgezustand gibt, genügt es, die Verweildauervariable  $T$  zu betrachten. Eine statistische Beschreibung zielt dann darauf, die statistische Verteilung dieser Verweildauervariablen zu ermitteln und darzustellen. Die begrifflichen Hilfsmittel hängen davon ab, ob man sich die Zeitachse als diskret oder stetig vorstellen will. In beiden Fällen können wir die Verteilung durch eine (kumulative) Verteilungsfunktion

$$F(t) = P(T \leq t)$$

charakterisieren. Ebenfalls unabhängig von der Art der Zeitachse kann man einen weiteren in der Verweildaueranalyse oft verwendeten Begriff definieren, die *Survivorfunktion*. Sie ergibt sich unmittelbar aus der Verteilungsfunktion durch die Definition

$$G(t) = 1 - F(t)$$

Eine Unterscheidung wird allerdings erforderlich, wenn wir uns auf eine zeitstellenbezogene Ereignisdichte beziehen wollen. Im diskreten Fall kann man dann eine diskrete Dichtefunktion

$$f(t) = P(T = t)$$

verwenden. Im stetigen Fall wird der Ausdruck  $P(T = t)$  problematisch, und es ist zweckmäßig, zunächst von Zeitintervallen auszugehen, also Ausdrücken der Art

$$P(t \leq T < t + \Delta)$$

wobei  $\Delta$  die Dauer des Zeitintervalls angibt, das an der Stelle  $t$  beginnt. Es ist klar, daß der Wert eines solchen Ausdrucks von  $\Delta$  abhängt, und man definiert deshalb die Ereignisdichte *pro Zeiteinheit* durch

$$f(t) = \lim_{\Delta \rightarrow 0} \frac{P(t \leq T < t + \Delta)}{\Delta}$$

Schließlich ist die Unterscheidung auch noch für den Begriff der *Übergangsrate* relevant, der in vielen Ansätzen der Verweildaueranalyse eine zentrale Rolle spielt. Die Idee ist, eine zeitstellenbezogene Ereignisdichte *unter der Bedingung* zu betrachten, daß das Ereignis noch nicht eingetreten ist. Im diskreten Fall lautet die Definition

$$r(t) = P(T = t | T \geq t)$$

Im stetigen Fall verwendet man die Definition

$$r(t) = \lim_{\Delta \rightarrow 0} \frac{P(t \leq T < t + \Delta | T \geq t)}{\Delta}$$

AUFGABE 3.3 Zeigen Sie zunächst für den diskreten, dann für den stetigen Fall, daß die Begriffe ‘Verteilungsfunktion’, ‘Survivorfunktion’, ‘Dichtefunktion’ und ‘Übergangsrate’ äquivalent sind, d.h. daß sie wechselseitig auseinander abgeleitet werden können. Zeigen Sie insbesondere, daß folgende Zusammenhänge gelten. Im diskreten Fall:

$$r(t) = f(t)/G(t-1)$$

und

$$G(t) = \prod_{\tau=1}^t (1 - r(\tau))$$

Und im stetigen Fall:

$$r(t) = f(t)/G(t)$$

und

$$G(t) = \exp \left\{ - \int_0^t r(\tau) d\tau \right\}$$

AUFGABE 3.4 Berechnen Sie mit dem Datensatz 1 (Box 1.1) die diskrete Übergangsrate für die Beendigung des Studiums.

AUFGABE 3.5 Betrachten Sie im Datensatz 2 (Box 2.1) zwei Gruppen von Episoden: Episoden, die im Zustand 0 beginnen, und Episoden, die im Zustand 1 beginnen. Verwenden Sie nur die nicht-zensierten Episoden, d.h. diejenigen Episoden, für die aus dem Datenbestand erkennbar ist, daß sie durch den Übergang in einen neuen Zustand beendet werden. Berechnen Sie dann die Übergangsraten für den Übergang in den Zustand 1 und für den Übergang in den Zustand 0.

### 3.4 Mehrere mögliche Folgezustände

Wenn eine Episode in zwei oder mehr möglichen Folgezuständen enden kann, genügt es nicht, nur die Verweildauervariable  $T$  zu betrachten, sondern man muß sich direkt auf die zweidimensionale Variable  $(T, D)$  beziehen. Die Aufgabe besteht dann darin, eine zweidimensionale Verteilung zu ermitteln und darzustellen. Um einen Zugang zu dieser Aufgabe zu finden, ist es zweckmäßig, mit der Idee einer *zielzustandsspezifischen Übergangsrate* zu beginnen. Im diskreten Fall lautet die Definition

$$r_d(t) = P(T = t, D = d | T \geq t)$$

wobei  $d$  einen der möglichen Folgezustände bezeichnet. Im stetigen Fall lautet die Definition

$$r_d(t) = \lim_{\Delta \rightarrow 0} \frac{P(t \leq T < t + \Delta, D = d | T \geq t)}{\Delta}$$

Die Menge der möglichen Folgezustände werden wir im folgenden stets mit dem Symbol  $\tilde{D}$  bezeichnen und dabei als Konvention annehmen, daß

$$\tilde{D} = \{1, \dots, m\}$$

ist, wenn es  $m$  mögliche Folgezustände gibt.

AUFGABE 3.6 Betrachten Sie in Box 3.1 alle Episoden, die im Zustand 1 beginnen. Bestimmen Sie die Menge  $\tilde{D}$  der möglichen Folgezustände und berechnen Sie für jeden Zustand  $d \in \tilde{D}$  die Übergangsrate  $r_d(t)$ .

AUFGABE 3.7 Wenn Episoden in mehreren möglichen Folgezuständen enden können, kann man auch von den Unterscheidungen abstrahieren und stattdessen nur einen möglichen Folgezustand betrachten, nämlich das Verlassen des Anfangszustands. Man kann dann die Episoden so betrachten, als ob es nur einen möglichen Folgezustand gibt und die in Abschnitt 3.3 eingeführten Begriffsbildungen verwenden. Zeigen Sie, daß folgender Zusammenhang gilt:

$$r(t) = \sum_{d \in \tilde{D}} r_d(t)$$

wobei  $\tilde{D}$  die Menge der möglichen Folgezustände bezeichnet. Verifizieren Sie diesen Zusammenhang an den Rechenergebnissen der Aufgabe 3.6.

AUFGABE 3.8 Bei Episoden mit mehreren möglichen Folgezuständen kann man folgendermaßen sog. *Sub-Survivorfunktionen* definieren:

$$G_d(t) = \exp \left\{ - \int_0^t r_d(\tau) d\tau \right\}$$

(a) Überlegen Sie sich, ob bzw. wie man diese Sub-Survivorfunktionen inhaltlich interpretieren kann. (b) Zeigen Sie, daß folgender Zusammenhang zum Begriff der Survivorfunktion gilt:

$$G(t) = \prod_{d \in \tilde{D}} G_d(t)$$

## 4 Zensierte Beobachtungen

Bisher haben wir angenommen, daß für die Verweildauervariable  $T$ , bzw.  $(T, D)$  bei mehreren möglichen Folgezuständen, vollständige Beobachtungen verfügbar sind, daß also die Episoden für alle Individuen abgeschlossen sind und wir die Verweildauern und Folgezustände kennen. Das ist bei den in der Praxis ermittelbaren Daten oft nicht der Fall. In diesem Kapitel behandeln wir einen wichtigen Spezialfall unvollständiger Daten, sog. rechts zensierte Beobachtungen.

### 4.1 Rechts zensierte Beobachtungen

Man sagt, daß die Beobachtung einer Episode bei einem Individuum *rechts zensiert* ist, wenn man zwar weiß, wie lange sich das Individuum schon im Anfangszustand aufhält, aber nicht weiß, wie lange es noch in diesem Zustand bleiben wird und welcher der möglichen Folgezustände dann eintreten wird. Die Situation ist dann folgende: Wir unterstellen eine statistische Variable  $(T, D)$  mit einer Menge möglicher Folgezustände  $\tilde{D}$ . Unsere Beobachtungen für  $i = 1, \dots, n$  Individuen liefern uns jedoch nicht unmittelbar Werte von  $(T, D)$ , sondern Werte einer Variablen  $(T^*, D^*)$ .  $D^*$  kann Werte in einer Menge

$$\tilde{D}^* = \tilde{D} \cup \{0\}$$

annehmen, wobei 0 der Anfangszustand der Episode ist und infolgedessen kein Element von  $\tilde{D}$  sein kann.<sup>1</sup> Die Beobachtungen sind in Form von Werten

$$(t_i^*, d_i^*) \quad \text{für } i = 1, \dots, n$$

gegeben, und der Zusammenhang mit den unterstellten Werten  $(t_i, d_i)$ , also den Werten der als theoretischer Rahmen angenommenen Variablen  $(T, D)$ , wird folgendermaßen hergestellt:

- Wenn  $d_i^* \in \tilde{D}$ , liegt eine nicht zensierte Beobachtung vor, und es gilt:  $t_i = t_i^*$  und  $d_i = d_i^*$ .
- Wenn  $d_i^* = 0$ , liegt eine zensierte Beobachtung vor; über den Folgezustand ist also nichts bekannt, es gilt jedoch  $t_i > t_i^*$ .

Diese Form der Repräsentation zensierter Beobachtungen erlaubt es, sie auf einfache Weise in einem Episodendatenschema (vgl. Abschnitt 3.1) kenntlich zu machen. Sie werden dadurch kenntlich gemacht, daß man für den Endzustand der Episode ihren Anfangszustand einsetzt, und für

<sup>1</sup> Entsprechend unserer Konvention, für  $\tilde{D}$  positive natürliche Zahlen zu verwenden, ist also  $\tilde{D}^* = \{0, 1, \dots, m\}$ , wenn es  $m$  mögliche Folgezustände gibt.

die Endzeitstelle diejenige Zeitstelle, bis zu der man weiß, daß sich das Individuum im Anfangszustand der Episode aufgehalten hat.

AUFGABE 4.1 Stellen Sie die Daten des Datensatzes 2 (Box 2.1) in einem Episodendatenschema dar, wobei rechts zensierte Episoden durch die eben genannte Konvention kenntlich gemacht werden.

## 4.2 Berechnung von Survivorfunktionen

Wir behandeln zunächst eine Situation, in der es nur einen möglichen Folgezustand gibt. Wir können also  $\tilde{D}^* = \{0, 1\}$  annehmen, wobei 0 zensierte, 1 unzensierte Beobachtungen kennzeichnet. Wie lassen sich dann Survivorfunktionen berechnen, wenn einige Beobachtungen rechts zensiert sind? Eine genaue Berechnung ist offenbar nicht möglich, denn bei den zensierten Beobachtungen kennt man nur  $t_i^*$ , nicht jedoch  $t_i$ . Wir können jedoch untere und obere Grenzen für die unbekannte Survivorfunktion  $G(t)$  berechnen.

- Eine untere Grenze, wir bezeichnen sie mit  $G^-(t)$ , erhält man, wenn man für die zensierten Beobachtungen annimmt, daß der Anfangszustand unmittelbar nach dem Zensierungszeitpunkt verlassen wird, also  $t_i = t_i^*$  oder, bei einer diskreten Zeitachse,  $t_i = t_i^* + 1$ .
- Eine obere Grenze, durch  $G^+(t)$  bezeichnet, erhält man, wenn man für die zensierten Beobachtungen annimmt, daß der Anfangszustand erst nach einer „beliebig langen“ Verweildauer verlassen wird. Es genügt jedoch, die Verweildauern der zensierten Episoden so anzusetzen, daß sie länger sind als die längste unzensierte Verweildauer.

Die unbekannte Survivorfunktion  $G(t)$  liegt sicherlich zwischen diesen Grenzen, d.h.

$$G^-(t) \leq G(t) \leq G^+(t)$$

Die Breite der Intervalle (abhängig von  $t$ ) hängt natürlich von dem Anteil zensierter Beobachtungen ab und davon, wie sie sich auf der Zeitachse verteilen. Je nachdem liefern die Daten mehr oder weniger viel Information über die Survivorfunktion  $G(t)$ .

AUFGABE 4.2 Berechnen Sie für die Daten in Box 4.1 untere und obere Grenzen der Survivorfunktion. Stellen Sie dann das Ergebnis in einem Schaubild dar.

## 4.3 Das Kaplan-Meier-Verfahren

Wenn man etwas nicht genau kennt, wie in diesem Fall die Survivorfunktion  $G(t)$ , neigen Statistiker dazu, sich Verfahren auszudenken, wie man das, was man nicht kennt, trotzdem möglichst sinnvoll *schätzen* kann. Ein für

### Box 4.1 Datensatz 4

ID	DUR	CEN
1	17	1
2	5	0
3	22	1
4	13	1
5	2	0
6	9	1
7	12	0
8	15	1

diesen Zweck ausgedachtes Verfahren stammt von E. L. Kaplan und P. Meier (1958). Um das Verfahren darzustellen, wird zunächst eine diskrete Zeitachse angenommen. Dann gibt es, wie in Abschnitt 3.3 gezeigt worden ist, folgenden Zusammenhang zwischen der Survivorfunktion und der Übergangsrate:

$$G(t) = \prod_{\tau=1}^t (1 - r(\tau))$$

Die Idee ist nun, zunächst die Übergangsraten  $r(t)$  zu schätzen und dann daraus die Survivorfunktion  $G(t)$  zu berechnen. Wenn es keine zensierten Beobachtungen gibt, ist unmittelbar einsichtig, wie man die Übergangsraten berechnen kann, nämlich durch

$$r(t) = \frac{E(t)}{R(t)}$$

Dabei ist  $E(t)$  die Anzahl der Individuen, die in der Zeitstelle  $t$  den Ausgangszustand der Episode verlassen; und  $R(t)$  ist die Anzahl der Individuen, bei denen es in der Zeitstelle  $t$  noch möglich ist, daß sie den Ausgangszustand verlassen, also die Anzahl derjenigen Individuen, die den Ausgangszustand nicht schon vorher verlassen haben.

Wenn es zensierte Beobachtungen gibt, kennen wir zwar weder  $E(t)$  noch  $R(t)$ , jedoch zwei vergleichbare Größen. Nämlich  $E^*(t)$ , die Anzahl der Individuen, deren Verlassen des Ausgangszustands in der Zeitstelle  $t$  wir beobachten können; und  $R^*(t)$ , die Anzahl der Individuen, bei denen ein Verlassen des Ausgangszustands in  $t$  noch beobachtet werden könnte, weil sie nicht schon vorher den Ausgangszustand verlassen und/oder rechts zensiert sind. Mithilfe dieser beobachteten Größen kann dann eine beobachtete Übergangsrate

$$r^*(t) = \frac{E^*(t)}{R^*(t)}$$

und daraus schließlich durch Anwendung der Formel (die jetzt eine Definition ist)

$$G^*(t) = \prod_{\tau=1}^t (1 - r^*(\tau))$$

eine Survivorfunktion  $G^*(t)$  berechnet werden. Offenbar ist  $G^*(t)$  eine sinnvolle Schätzung für  $G(t)$ , wenn man voraussetzen kann, daß  $r(t)$  sinnvoll durch  $r^*(t)$  geschätzt werden kann.

Das gleiche Verfahren kann natürlich angewendet werden, wenn man annimmt, daß die beobachteten (zensierten und nicht zensierten) Verweildauern als exakte Zeitangaben auf einer kontinuierlichen Zeitachse interpretiert werden können. Man erhält dann eine Treppenfunktion, die genau in denjenigen Zeitpunkten Sprungstellen aufweist, in denen mindestens ein Ereignis stattfindet.

**AUFGABE 4.3** Berechnen Sie für die Daten in Box 4.1 die Survivorfunktion  $G^*(t)$  mit dem Kaplan-Meier-Verfahren. Stellen Sie dann das Ergebnis in einem Schaubild dar, das außerdem die unteren und oberen Schranken,  $G^-(t)$  und  $G^+(t)$ , zeigt. Beachten Sie, daß  $r^*(t)$  nur für diejenigen Zeitstellen berechnet zu werden braucht, in denen mindestens ein Ereignis stattfindet, also  $E^*(t) \neq 0$  ist.

#### 4.4 Mehrere Folgezustände

Das Kaplan-Meier-Verfahren läßt sich auch dann verwenden, wenn die Episoden in zwei oder mehr möglichen Folgezuständen enden können. Es werden dann Sub-Survivorfunktionen geschätzt, also

$$G_d^*(t) = \prod_{\tau=1}^t (1 - r_d^*(\tau))$$

wobei  $d \in \tilde{D}$ . Die zielzustandsspezifischen Übergangsraten können durch

$$r_d^*(t) = \frac{E_d^*(t)}{R^*(t)}$$

geschätzt werden, wobei jetzt  $E_d^*(t)$  die Anzahl der Individuen ist, bei denen in der Zeitstelle  $t$  ein Übergang in den Folgezustand  $d$  festgestellt werden kann. Man beachte, daß in diesem Fall der multiplikative Zusammenhang

$$G^*(t) \approx \prod_{d \in \tilde{D}} G_d^*(t)$$

nur näherungsweise gilt.

#### 4.5 Selbst-Konsistenz

Das Kaplan-Meier-Verfahren kann auch mit der Idee einer Selbst-Konsistenz begründet werden, die wir kurz diskutieren wollen. Die Idee ist nicht auf rechts zensierte Daten beschränkt, sondern allgemeiner, und wir besprechen sie deshalb zunächst für eine beliebige diskrete Variable

$$X : \Omega \longrightarrow \tilde{X}$$

Wenn uns für alle Mitglieder von  $\Omega$  genaue Beobachtungen vorliegen, kann natürlich ohne weiteres die Verteilungsfunktion

$$P(X = x) \quad \text{für alle } x \in \tilde{X}$$

berechnet werden (vgl. Abschnitt 2.4). Jetzt nehmen wir jedoch an, daß wir die genauen Werte nicht kennen, sondern für jedes  $\omega \in \Omega$  nur eine Teilmenge von  $\tilde{X}$ , in der der Variablenwert  $X(\omega)$  liegt. Um den Gedankengang einfacher darstellen zu können, stellen wir uns vor, daß es für die Mitglieder von  $\Omega$  Nummern,  $i = 1, \dots, n$ , gibt. Die beobachteten Werte der Variablen seien durch Merkmalsmengen

$$\tilde{x}_i \subseteq \tilde{X}$$

gegeben. Dann können wir zwar die Verteilungsfunktion  $P$  nicht genau berechnen; wir können jedoch zunächst untere und obere Grenzen ermitteln. In einem ersten Schritt definieren wir:

$$\begin{aligned} \text{pmin}(\tilde{x}_i, \tilde{x}) &:= \begin{cases} 1 & \text{wenn } \tilde{x}_i \subseteq \tilde{x} \\ 0 & \text{andernfalls} \end{cases} \\ \text{pmax}(\tilde{x}_i, \tilde{x}) &:= \begin{cases} 0 & \text{wenn } \tilde{x}_i \cap \tilde{x} = \emptyset \\ 1 & \text{andernfalls} \end{cases} \end{aligned}$$

wobei  $\tilde{x}$  eine beliebige Teilmenge von  $\tilde{X}$  sein kann. Dann ergeben sich untere und obere Grenzen für  $P$  durch die Definitionen:

$$\begin{aligned} P^+(\tilde{x}) &:= \frac{1}{n} \sum_{i=1}^n \text{pmin}(\tilde{x}_i, \tilde{x}) \\ P^-(\tilde{x}) &:= \frac{1}{n} \sum_{i=1}^n \text{pmax}(\tilde{x}_i, \tilde{x}) \end{aligned}$$

Wie man sich leicht überlegen kann, gilt

$$P^+(\tilde{x}) \leq P(\tilde{x}) \leq P^-(\tilde{x})$$

Die Frage ist nun, wie man sinnvoll eine „mittlere“ Verteilungsfunktion definieren kann, die zwischen den beiden Grenzen liegt; denn die Verteilung  $P$  kennt man nicht, und man kann sie (ohne weitere Annahmen) auch nicht aus den Daten schätzen. Eine Überlegung wäre die folgende. Man nimmt an, daß das Individuum  $i$  an der Merkmalsmenge  $\tilde{x}$  in dem Maße partizipiert, wie sich  $\tilde{x}_i$  und  $\tilde{x}$  überschneiden. Diese Idee führt zu folgender Definition einer Verteilungsfunktion:

$$\bar{P}(\tilde{x}) := \frac{1}{n} \sum_{i=1}^n \frac{|\tilde{x}_i \cap \tilde{x}|}{|\tilde{x}_i|}$$

Wir nennen sie *mittlere Verteilungsfunktion*. Dies ist jedoch ein Spezialfall einer allgemeineren Idee. Wenn wir uns nämlich auf eine beliebige Verteilungsfunktion  $\tilde{P}$  beziehen, erscheint es sinnvoll zu sagen, daß das Individuum  $i$  an  $\tilde{x}$  entsprechend der Größe von

$$\frac{\tilde{P}(\tilde{x}_i \cap \tilde{x})}{\tilde{P}(\tilde{x}_i)}$$

partizipiert. Die mittlere Verteilungsfunktion resultiert dann als ein Spezialfall, wenn man nämlich für  $\tilde{P}$  eine Gleichverteilung annimmt.

Ein zweiter Ansatz ergibt sich durch die Idee der Selbst-Konsistenz. Die Idee ist, nach einer Verteilung  $\tilde{P}$  zu suchen, so daß die folgende Funktionalgleichung erfüllt ist:

$$\tilde{P}(\tilde{x}) := \frac{1}{n} \sum_{i=1}^n \frac{\tilde{P}(\tilde{x}_i \cap \tilde{x})}{\tilde{P}(\tilde{x}_i)} \quad (4.5.1)$$

Eine Lösung kann meistens mit einem iterativen Verfahren gefunden werden: Man beginnt mit einer beliebigen, z.B. der mittleren Verteilungsfunktion; dann wendet man die Formel an, um daraus eine neue Verteilungsfunktion zu berechnen; dann wiederholt man die Berechnung mit der neuen Verteilungsfunktion, usw., bis sich keine wesentlichen Änderungen mehr ergeben. Wenn man eine Lösung gefunden hat, wird sie als *selbst-konsistente Verteilung* bezeichnet.

AUFGABE 4.4 Berechnen Sie mit den Daten

$$\tilde{x}_1 = \{1\}, \tilde{x}_2 = \{2\}, \tilde{x}_3 = \{3\}, \tilde{x}_4 = \{1, 2\}, \tilde{x}_5 = \{2, 3\}$$

die Verteilungen  $P^+$ ,  $P^-$ ,  $\bar{P}$  und  $\tilde{P}$ .

AUFGABE 4.5 Man kann zeigen, daß das Kaplan-Meier-Verfahren eine selbst-konsistente Verteilung erzeugt. Überlegen Sie sich, wie mit dem beschriebenen Verfahren eine selbst-konsistente Verteilungs- oder Survivorfunktion für den Datensatz 4.1 berechnet werden kann.

## 5 Regressionsmodelle für Zustände

In den weiteren Kapiteln beschäftigen wir uns mit statistischen Modellen für die Analyse von Abhängigkeiten zwischen Variablen. Um den Ausgangspunkt zu fixieren, erinnern wir uns an die beiden Formen, die wir zur Repräsentation von Verlaufsdaten eingeführt haben. Einerseits haben wir Folgen von Zustandsvariablen betrachtet:  $Y_t$ , wobei der Index  $t$  sich auf einer diskreten Zeitachse bewegen kann. Andererseits haben wir uns auf einzelne Episoden bezogen und diese durch eine zweidimensionale Variable  $(T, D)$  repräsentiert. Beide Varianten können als Ausgangspunkt für Modellkonstruktionen dienen. In diesem Kapitel gehen wir von der ersten Variante aus.

### 5.1 Der Modellansatz

Wir beziehen uns auf eine Prozeßzeitachse  $t = 0, 1, 2, \dots$  und nehmen an, daß eine Folge von Zustandsvariablen

$$Y_t : \Omega \longrightarrow \tilde{Y}$$

gegeben ist.  $\tilde{Y}$  ist der Zustandsraum, der zwei oder mehr unterschiedliche Zustände enthalten kann. Wenn sich der Prozeß bis zu einem Zeitpunkt  $t$  entwickelt hat, wird er statistisch durch die Verteilung

$$P(Y_t = y_t, Y_{t-1} = y_{t-1}, \dots, Y_0 = y_0)$$

erfaßt. Dabei sind  $y_0, \dots, y_t$  mögliche Zustände im Zustandsraum  $\tilde{Y}$ .

Ein Modell soll es erlauben, über Abhängigkeiten zwischen den Zustandsvariablen nachzudenken zu können. Um unseren Modellansatz einfach zu schreiben, verwenden wir folgende Abkürzung:

$$\bar{Y}_t := (Y_t, Y_{t-1}, \dots, Y_0)$$

und nennen dies eine *Prozeßvariable*. Auf mögliche Werte wird durch entsprechende Kleinbuchstaben, also

$$\bar{y}_t := (y_t, y_{t-1}, \dots, y_0)$$

verwiesen; es handelt sich um mögliche Zustandsfolgen. Der Ausgangspunkt für die Modellbildung kann dann durch

$$P(\bar{Y}_t = \bar{y}_t)$$

fixiert werden. Um Abhängigkeiten zwischen den Zustandsvariablen zu erfassen, werden bedingte Verteilungen verwendet. Dabei nehmen wir an,

daß die Ausgangsverteilung, also die Verteilung von  $Y_0$ , vorgegeben ist und sich der Prozeß dann sequentiell entwickelt; in symbolischer Notation:

$$\begin{aligned} & Y_0 \\ & Y_1 | Y_0 \\ & Y_2 | Y_0, Y_1 \\ & \vdots \\ & Y_t | Y_0, Y_1, \dots, Y_{t-1} \end{aligned}$$

Durch sukzessive Anwendung der Regel zur Bildung von bedingten Verteilungen erhält man dann:

$$P(\bar{Y}_t = \bar{y}_t) = \prod_{\tau=1}^t P(Y_\tau = y_\tau | \bar{Y}_{\tau-1} = \bar{y}_{\tau-1}) P(Y_0 = y_0) \quad (5.1.1)$$

## 5.2 Spekulation und Empirie

Der allgemeine Modellansatz (5.1.1) kann in zwei unterschiedlichen Weisen als Ausgangspunkt für weitere Überlegungen dienen. Er kann einerseits als Ausgangspunkt für spekulative Überlegungen zur Prozeßentwicklung, andererseits als ein formaler Rahmen für die Repräsentation von Daten über die Prozeßentwicklung verwendet werden. Um mit dem Modellansatz etwas vertrauter zu werden, beginnen wir mit spekulativen Überlegungen.

Da wir annehmen, daß die Verteilung von  $Y_0$  vorgegeben ist, konzentriert sich die Spekulation auf Annahmen über

$$P(Y_\tau = y_\tau | \bar{Y}_{\tau-1} = \bar{y}_{\tau-1})$$

Eine einfache Annahme könnte zum Beispiel darin bestehen, daß es nur ein einstufiges Gedächtnis gibt:

$$P(Y_\tau = y_\tau | \bar{Y}_{\tau-1} = \bar{y}_{\tau-1}) = P(Y_\tau = y_\tau | Y_{\tau-1} = y_{\tau-1})$$

Hier wird also angenommen, daß der Zustand, der in einer Zeitstelle angenommen wird, nur davon abhängt, welcher Zustand in der vorangegangenen Zeitstelle angenommen worden ist. Etwas komplizierter wäre ein zweistufiges Gedächtnis, das man durch die Annahme

$$\begin{aligned} & P(Y_\tau = y_\tau | \bar{Y}_{\tau-1} = \bar{y}_{\tau-1}) = \\ & P(Y_\tau = y_\tau | Y_{\tau-1} = y_{\tau-1}, Y_{\tau-2} = y_{\tau-2}) \end{aligned}$$

ausdrücken kann.

AUFGABE 5.1 Betrachten Sie einen Prozeß mit einem einstufigen Gedächtnis. Der Zustandsraum sei  $\tilde{Y} = \{0, 1\}$ , zum Prozeßbeginn befinden sich alle Individuen im Zustand 0, und die Übergangswahrscheinlichkeiten werden durch

$$\begin{aligned} P(Y_\tau = 1 | Y_{\tau-1} = 0) &= 1/2 \\ P(Y_\tau = 1 | Y_{\tau-1} = 1) &= 1/3 \end{aligned}$$

angenommen. Konstruieren Sie mithilfe eines Würfels 10 Realisationen dieses Prozesses, für  $t = 0, \dots, 8$ , und stellen Sie in einer Tabelle dar, wie sich die Zustandsverteilung im Zeitablauf entwickelt.

AUFGABE 5.2 Nehmen Sie an, daß es 5 unterschiedliche Zustände gibt. Wieviele Übergangswahrscheinlichkeiten sind erforderlich, um einen Prozeß mit einem zweistufigen Gedächtnis vollständig darzustellen?

## 5.3 Modelle für zwei Zustände

Wir betrachten Prozesse, bei denen es nur zwei unterschiedliche Zustände gibt, also  $\tilde{Y} = \{0, 1\}$ . Außerdem nehmen wir an, daß es nur ein einstufiges Gedächtnis gibt. Wenn man keine weiteren Einschränkungen vornimmt, müssen für jede Zeitstelle zwei Parameter ermittelt werden:

$$\begin{aligned} P(Y_t = 1 | Y_{t-1} = 0) &= \theta_{10,t} \\ P(Y_t = 1 | Y_{t-1} = 1) &= \theta_{11,t} \end{aligned}$$

Die Größen  $\theta_{10,t}$  und  $\theta_{11,t}$  nennen wir *Parameter des Prozesses*. Im allgemeinen muß angenommen werden, daß sich diese Prozeßparameter während der Entwicklung des Prozesses verändern können. Eine radikal vereinfachende, aber auch problematische Annahme ist, daß die Prozeßparameter zeitkonstant sind; man spricht dann von einem *stationären Prozeß*. Der Modellansatz ist dann

$$\begin{aligned} P(Y_t = 1 | Y_{t-1} = 0) &= \theta_{10} \\ P(Y_t = 1 | Y_{t-1} = 1) &= \theta_{11} \end{aligned}$$

AUFGABE 5.3 Verwenden Sie die Daten aus Box 2.1. Berechnen Sie zunächst die zeitstellenspezifischen Prozeßparameter

$$\theta_{ij,t} \quad \text{für } i, j \in \{0, 1\}, t = 2, \dots, 6$$

Nehmen Sie dann an, daß die Daten aus einem stationären Prozeß stammen und berechnen Sie die zeitkonstanten Prozeßparameter

$$\theta_{ij} \quad \text{für } i, j \in \{0, 1\}$$

## 5.4 Modelle mit Kovariablen

Bisher haben wir nur eine Prozeßvariable,  $Y_t$ , betrachtet, und die Modellbildung bezog sich darauf, herauszufinden, wie der jeweils gegenwärtige Zustand von vorangegangenen Zuständen abhängt. Bei praktischen Anwendungen ist man oft daran interessiert, noch weitere Variablen, sog. *Kovariablen*, zu berücksichtigen. Zwei Arten von Kovariablen können dabei unterschieden werden:

- *Zeitunabhängige Kovariablen*, deren Werte zu Beginn des Prozesses feststehen und sich während des Prozesses nicht ändern können; und
- *Zeitabhängige Kovariablen*, deren Werte sich während des Prozesses verändern können.

Offenbar können zeitunabhängige Kovariablen als ein Spezialfall zeitabhängiger Kovariablen betrachtet werden. Wir betrachten deshalb im folgenden nur zeitabhängige Kovariablen. Einen sinnvollen Begriffsrahmen liefert dann die Vorstellung *paralleler Prozesse*. Den primär interessierenden Prozeß repräsentieren wir wie bisher durch die Zustandsvariablen  $Y_t$ , den parallelen Kovariablenprozeß durch eine Folge von Kovariablen  $X_t$ ; in expliziter Schreibweise:

$$(X_t, Y_t) : \Omega \longrightarrow \tilde{X} \times \tilde{Y}$$

Wie bisher nehmen wir an, daß  $Y_t$  eine diskrete eindimensionale Zustandsvariable ist. Bei  $X_t$  kann es sich um eine mehrdimensionale Variable handeln, z.B. um eine  $m$ -dimensionale Variable

$$X_t = (X_{t1}, \dots, X_{tm})$$

Zur Vereinfachung werden wir jedoch annehmen, daß auch  $X_t$  eine diskrete Variable ist.

Die Idee ist nun, daß der Zustand in einer Zeitstelle  $t$  nicht nur von Zuständen abhängen kann, die in vorangehenden Zeitstellen eingenommen worden sind, sondern auch von den bisher realisierten Werten der Kovariablen. Für die Modellbildung nehmen wir an, daß die Werte der Kovariablen ihrerseits nicht von den bisher realisierten Werten der Prozeßvariablen  $Y_t$  abhängig sind.<sup>1</sup> Der allgemeine Modellansatz (5.1.1) kann dann folgendermaßen erweitert werden:

$$P(\bar{Y}_t = \bar{y}_t) = \quad (5.4.1)$$

<sup>1</sup> Die Kovariablen werden dann *exogen* genannt. Wenn diese Annahme nicht erfüllt ist, spricht man gelegentlich von *interdependenten* Prozessen. Damit werden wir uns in dieser Einführung jedoch nicht näher beschäftigen.

### Box 5.1 Datensatz 5

ID	t =	0	1	2	3	4	5	6
1	Y	0	0	0	1	1	0	0
	X1	0	0	0	0	0	0	0
	X2	20	21	22	23	24	25	26
2	Y	1	1	0	0	0	1	1
	X1	0	0	0	0	0	0	0
	X2	22	23	24	25	26	27	28
3	Y	1	1	1	0	0	0	0
	X1	0	0	0	0	0	0	0
	X2	21	22	23	24	25	26	27
4	Y	0	0	0	0	1	1	1
	X1	1	1	1	1	1	1	1
	X2	20	21	22	23	24	25	26
5	Y	0	0	1	1	1	0	0
	X1	1	1	1	1	1	1	1
	X2	22	23	24	25	26	27	28
6	Y	1	1	1	0	0	1	1
	X1	1	1	1	1	1	1	1
	X2	21	22	23	24	25	26	27

$$\prod_{\tau=1}^t P(Y_\tau = y_\tau | \bar{Y}_{\tau-1} = \bar{y}_{\tau-1}, \bar{X}_{\tau-1} = \bar{x}_{\tau-1}) \\ P(Y_0 = y_0, X_0 = x_0)$$

Wiederum kann dieser Modellansatz auf vielfältige Weisen vereinfacht werden. Denkt man an die Idee eines Prozesses mit einem einstufigen Gedächtnis, kann man das zum Beispiel auch für die Kovariablen annehmen und erhält dann den Modellansatz

$$P(Y_\tau = y_\tau | \bar{Y}_{\tau-1} = \bar{y}_{\tau-1}, \bar{X}_{\tau-1} = \bar{x}_{\tau-1}) = \\ P(Y_\tau = y_\tau | Y_{\tau-1} = y_{\tau-1}, X_{\tau-1} = x_{\tau-1})$$

Als Beispiel betrachten wir den Datensatz 5 in Box 5.1. Er enthält Angaben über die Entwicklung von Zuständen bei 6 Personen. Es gibt zwei mögliche Zustände,  $\tilde{Y} = \{0, 1\}$ , und zwei Kovariablen. Die Kovariable  $X_1$  ist zeitunabhängig, z.B. das Geschlecht der Personen (0 = Männer, 1 = Frauen); die Kovariable  $X_2$  ist zeitabhängig, z.B. das Alter der Personen.

## 5.5 Binäre Logitmodelle

Zur Modellierung von Prozessen mit Kovariablen werden in der Praxis oft Logitmodelle verwendet, die rechentechnisch verhältnismäßig einfach handhabbar sind. Wir besprechen hier diese Modelle für Prozesse, bei denen

es nur zwei mögliche Zustände gibt, also  $\tilde{Y} = \{0, 1\}$ . Verzichtet man zunächst auf die Annahme, daß es sich um einen stationären Prozeß handelt, ist der Modellansatz:

$$P(Y_t = 1 | Y_{t-1} = 0, X_{t-1} = x_{t-1}) = \frac{\exp(\alpha_{10,t} + x_{t-1}\beta_{10,t})}{1 + \exp(\alpha_{10,t} + x_{t-1}\beta_{10,t})}$$

$$P(Y_t = 1 | Y_{t-1} = 1, X_{t-1} = x_{t-1}) = \frac{\exp(\alpha_{11,t} + x_{t-1}\beta_{11,t})}{1 + \exp(\alpha_{11,t} + x_{t-1}\beta_{11,t})}$$

Nimmt man an, daß es sich um einen stationären Prozeß handelt, entfällt auf den rechten Seiten der Zeitindex bei den Modellparametern. Das in Abschnitt 5.3 behandelte Modell ohne Kovariablen ist offenbar ein Spezialfall dieses Logitmodells.

Leider können die Parameter eines Logitmodells im allgemeinen nicht mit einfachen Rechenverfahren aus Daten berechnet werden. Wir werden uns im nächsten Abschnitt mit einem Schätzverfahren beschäftigen. Zuvor behandeln wir einige Aufgaben, um mit dem Modellansatz etwas vertrauter zu werden.

AUFGABE 5.4 Zeichnen Sie den Verlauf der Logitfunktion

$$z = \frac{\exp(x)}{1 + \exp(x)}$$

im Bereich  $-3 \leq x \leq 3$ .

AUFGABE 5.5 Unter Verwendung der Notation aus Aufgabe 5.4, entwickeln Sie die Umkehrfunktion, die zeigt, wie  $x$  von  $z$  abhängt. Berechnen Sie die  $x$ -Werte, die den  $z$ -Werten 0.5, 0.6 und 0.7 entsprechen.

AUFGABE 5.6 Formulieren Sie ein Logitmodell für die Daten aus Box 5.1 unter der Annahme, daß es sich um Daten aus einem stationären Prozeß handelt. Überlegen Sie sich insbesondere, welche Modellparameter berechnet werden müßten.

## 5.6 Maximum-Likelihood-Schätzung

In diesem Abschnitt besprechen wir, wie die Parameter eines binären Logitmodells mit der ML-Methode berechnet werden können. Wenn angenommen wird, daß es sich um einen stationären Prozeß handelt, müssen vier Parameter berechnet werden:

$$\alpha_{10}, \beta_{10}, \alpha_{11}, \beta_{11}$$

Die Daten seien in der Form

$$(x_{it}, y_{it}) \quad \text{für } i = 1, \dots, N, t = 0, \dots, T$$

gegeben. Dann folgt aus dem Modellansatz:

$$P(Y_t = y_{i,t} | Y_{t-1} = y_{i,t-1}, X_{t-1} = x_{i,t-1}) = \begin{cases} \frac{\exp(\alpha_{10} + x_{i,t-1}\beta_{10})}{1 + \exp(\alpha_{10} + x_{i,t-1}\beta_{10})} & \text{wenn } y_{i,t} = 1, y_{i,t-1} = 0 \\ \frac{1}{1 + \exp(\alpha_{10} + x_{i,t-1}\beta_{10})} & \text{wenn } y_{i,t} = 0, y_{i,t-1} = 0 \\ \frac{\exp(\alpha_{11} + x_{i,t-1}\beta_{11})}{1 + \exp(\alpha_{11} + x_{i,t-1}\beta_{11})} & \text{wenn } y_{i,t} = 1, y_{i,t-1} = 1 \\ \frac{1}{1 + \exp(\alpha_{11} + x_{i,t-1}\beta_{11})} & \text{wenn } y_{i,t} = 0, y_{i,t-1} = 1 \end{cases}$$

Das kann man in einer Formel folgendermaßen ausdrücken:

$$P(Y_t = y_{i,t} | Y_{t-1} = y_{i,t-1}, X_{t-1} = x_{i,t-1}) = \left( \frac{\exp(\alpha_{10} + x_{i,t-1}\beta_{10})}{1 + \exp(\alpha_{10} + x_{i,t-1}\beta_{10})} \right)^{1-y_{i,t-1}} \left( \frac{\exp(\alpha_{11} + x_{i,t-1}\beta_{11})}{1 + \exp(\alpha_{11} + x_{i,t-1}\beta_{11})} \right)^{y_{i,t-1}}$$

Betrachtet man diesen Ausdruck als eine Funktion der Modellparameter, wird er als *Likelihood* der Beobachtung  $i$  in der Zeitstelle  $t$  bezeichnet. Die ML-Methode besteht nun darin, diese Likelihood für alle Beobachtungen und Zeitstellen zusammenzufassen und dann als eine Funktion der Modellparameter zu maximieren. Die *Likelihoodfunktion* sieht also folgendermaßen aus:

$$\mathcal{L}(\alpha_{10}, \beta_{10}, \alpha_{11}, \beta_{11}) = \prod_{i=1}^N \prod_{t=1}^T \left( \frac{\exp(\alpha_{10} + x_{i,t-1}\beta_{10})}{1 + \exp(\alpha_{10} + x_{i,t-1}\beta_{10})} \right)^{1-y_{i,t-1}} \left( \frac{\exp(\alpha_{11} + x_{i,t-1}\beta_{11})}{1 + \exp(\alpha_{11} + x_{i,t-1}\beta_{11})} \right)^{y_{i,t-1}}$$

Aus der Maximierung dieser Likelihoodfunktion ergeben sich die ML-Schätzwerte der Modellparameter, die mit

$$\hat{\alpha}_{10}, \hat{\beta}_{10}, \hat{\alpha}_{11}, \hat{\beta}_{11}$$

bezeichnet werden.

AUFGABE 5.7 Für die praktische Berechnung verwendet man meistens die *Log-Likelihoodfunktion*, also

$$\ell(\alpha_{10}, \beta_{10}, \alpha_{11}, \beta_{11}) = \log(\mathcal{L}(\alpha_{10}, \beta_{10}, \alpha_{11}, \beta_{11}))$$

Berechnen Sie die Log-Likelihoodfunktion für die oben angegebene Likelihoodfunktion.

**Box 5.2** Datensatz 5a

ID	t	Y(t)	Y(t-1)	X1(t-1)	X2(t-1)
1	1	0	0	0	20
1	2	0	0	0	21
1	3	1	0	0	22
1	4	1	1	0	23
1	5	0	1	0	24
1	6	0	0	0	25
2	1	1	1	0	22
2	2	0	1	0	23
2	3	0	0	0	24
2	4	0	0	0	25
2	5	1	0	0	26
2	6	1	1	0	27
3	1	1	1	1	21
3	2	1	1	1	22
3	3	1	0	1	23
3	4	0	0	1	24
3	5	0	0	1	25
3	6	0	0	1	26
4	1	0	0	1	20
4	2	0	0	1	21
4	3	0	0	1	22
4	4	1	0	1	23
4	5	1	1	1	24
4	6	1	1	1	25
5	1	0	0	1	22
5	2	1	0	1	23
5	3	1	1	1	24
5	4	1	1	1	25
5	5	0	1	1	26
5	6	0	0	1	27
6	1	1	1	1	21
6	2	1	1	1	22
6	3	0	1	1	23
6	4	0	0	1	24
6	5	1	0	1	25
6	6	1	1	1	26

AUFGABE 5.8 Berechnen Sie den Wert der in Aufgabe 5.7 gefundenen Log-Likelihoodfunktion für den Datensatz in Box 5.1 und für die Parameterwerte

$$\alpha_{10} = \beta_{10} = \alpha_{11} = \beta_{11} = 0$$

AUFGABE 5.9 Die ML-Schätzwerte der Modellparameter ergeben sich aus der Maximierung der Log-Likelihoodfunktion. Eine Lösung kann mithilfe eines Statistik-Programms berechnet werden, das in der Lage ist, Logitmodelle zu

schätzen.<sup>2</sup> Zunächst erkennt man, daß sich die Likelihoodfunktion in zwei Faktoren separieren läßt, die unabhängig voneinander maximiert werden können. Überlegen Sie sich, wie dementsprechend der Datensatz aus Box 5.1 reorganisiert und aufgeteilt werden kann (einen Hinweis gibt Box 5.2). Überlegen Sie sich dann, wie die Modellparameter mit einem Programm berechnet werden können, das einfache Logitmodelle für eine binäre abhängige Variable schätzen kann.

AUFGABE 5.10 Verwendet man für die Variable  $X$  das Alter ( $X_2$  in Box 5.1), erhält man folgende ML-Schätzwerte der Modellparameter:

$$\hat{\alpha}_{10} = -3.14, \hat{\beta}_{10} = 0.10, \hat{\alpha}_{11} = 4.91, \hat{\beta}_{11} = -0.16$$

Das für diese Berechnungen verwendete Computerprogramm (TDA) gibt als Wert der maximierten Log-Likelihoodfunktion die Werte -12.14 (für die erste Modellhälfte,  $Y_{t-1} = 0$ ) und -8.88 (für die zweite Modellhälfte,  $Y_{t-1} = 1$ ) an. Berechnen Sie den Wert der Log-Likelihoodfunktion für das Gesamtmodell.

AUFGABE 5.11 Unter Verwendung der in Aufgabe 5.10 angegebenen Schätzwerte für die Modellparameter: Erstellen Sie eine Tabelle, die für die Alterswerte  $X = 20, \dots, 26$  und die Werte  $Y_{t-1} = 0, 1$  die geschätzten Wahrscheinlichkeiten

$$P(Y_t = 0 | Y_{t-1} = \dots, X = \dots) \quad \text{und}$$

$$P(Y_t = 1 | Y_{t-1} = \dots, X = \dots)$$

angibt.

<sup>2</sup> Wie man das praktisch machen kann, wird im Rahmen der praktischen Übungen im Anhang besprochen.

## 6 Modelle für Verweildauern

Entsprechend der zu Beginn des vorangegangenen Kapitels getroffenen Unterscheidung beziehen wir uns für die weitere Diskussion statistischer Modelle auf einzelne Episoden. Ausgangspunkt ist also die Repräsentation einer einzelnen Episode durch eine zweidimensionale statistische Variable

$$(T, D) \quad \text{wobei } T \in \tilde{T}, D \in \tilde{D}$$

$\tilde{D}$  ist die Menge der möglichen Folgezustände der Episode,  $\tilde{T}$  ist die Zeitachse. Für die Zeitachse kann man entweder eine diskrete oder eine kontinuierliche numerische Repräsentation wählen. Wir verwenden im folgenden eine kontinuierliche Repräsentation, da die meisten in der Literatur diskutierten und praktisch verwendeten Modelle von dieser Konvention ausgehen. Wir identifizieren  $\tilde{T}$  also mit den nichtnegativen reellen Zahlen.

Statistische Modelle können sowohl spekulativen als auch repräsentativen Zwecken dienen. Im ersten Fall geht es darum, Vorstellungen darüber zu bilden, wie Episoden ablaufen *könnten*. Damit beschäftigen wir uns in diesem Kapitel. Wie schon in früheren Kapiteln können wir zwei Situationen unterscheiden: Erstens eine Betrachtungsweise, bei der Episoden nur in einem möglichen Folgezustand enden können; es genügt dann, die Episode durch eine Verweildauervariable ( $T$ ) zu erfassen. Zweitens eine Betrachtungsweise, bei der es zwei oder mehr mögliche Folgezustände gibt; dann muß von der zweidimensionalen Variablen ( $T, D$ ) ausgegangen werden.

### 6.1 Zeitkonstante Raten

Beginnen wir mit einer Episode, für die es nur einen möglichen Folgezustand gibt. Es genügt dann, die Verweildauervariable  $T$  zu betrachten, und ein statistisches Modell besteht darin, Annahmen über die Verteilung dieser Variablen zu machen. Wie wir uns schon überlegt haben, gibt es zur Charakterisierung dieser Verteilung vier äquivalente Konzepte: die Verteilungsfunktion  $F(t)$ , die Survivorfunktion  $G(t)$ , die Dichtefunktion  $f(t)$  und die Rate  $r(t)$ . Um Annahmen über die Verteilung von  $T$  zu formulieren, kann man also wahlweise einen dieser vier Begriffe verwenden und die jeweils übrigen daraus berechnen.

Um inhaltliche Vorstellungen über den Ablauf einer Episode zu bilden, ist es oft hilfreich, vom Begriff der Rate auszugehen. Die allereinfachste Vorstellung besteht darin, von einer zeitkonstanten Rate auszugehen. Diese Annahme kann so formuliert werden:

$$r(t) = \theta$$

wobei  $\theta$  ein Modellparameter ist, der in einem vorgegebenen Parameterraum variieren kann. Wir bezeichnen den Parameterraum mit  $\Theta$ ; und da Raten keine negativen Werte annehmen können, können wir  $\Theta$  mit den nichtnegativen reellen Zahlen identifizieren.

Ein Verteilungsmodell, das auf der Annahme einer zeitkonstanten Rate beruht, wird *Exponentialmodell* genannt. Man spricht von einer *Exponentialverteilung* mit dem Parameter  $\theta$ .

Die Survivorfunktion für eine Exponentialverteilung ergibt sich aus unserer Basisformel für den Zusammenhang zwischen Rate und Survivorfunktion in folgender Weise:

$$G(t) = \exp \left\{ - \int_0^t r(\tau) d\tau \right\} = \exp(-\theta t)$$

AUFGABE 6.1 Berechnen Sie die Verteilungs- und Dichtefunktionen für eine Exponentialverteilung mit dem Parameter  $\theta$ .

AUFGABE 6.2 Zeichnen Sie den Verlauf der Survivor- und Dichtefunktionen für eine Standard-Exponentialverteilung, d.h. für eine Exponentialverteilung mit dem Parameter  $\theta = 1$ .

## 6.2 Weibull-Verteilung

Interessanter als die Exponentialverteilung sind Verteilungen, bei denen sich die Rate im Zeitablauf verändern kann. Eine in Anwendungen oft verwendete Verteilung ist die *Weibull-Verteilung*, bei der die Rate monoton steigen oder fallen kann. Die Weibull-Verteilung hat zwei Parameter, die Survivorfunktion sieht folgendermaßen aus:

$$G(t) = \exp\{-(\alpha t)^\beta\}$$

$\alpha$  und  $\beta$  sind die Verteilungsparameter, und es wird vorausgesetzt, daß beide Parameter nur positive Werte annehmen können.

Aus der Survivorfunktion erhält man durch Differenzieren die Dichtefunktion

$$f(t) = \beta \alpha^\beta t^{\beta-1} \exp\{-(\alpha t)^\beta\}$$

und daraus die Rate

$$r(t) = \beta \alpha^\beta t^{\beta-1}$$

AUFGABE 6.3 Zeigen Sie schrittweise, wie sich die Dichtefunktion und die Rate aus der Survivorfunktion der Weibull-Verteilung berechnen lassen.

AUFGABE 6.4 Für welche Parameterwerte der Weibull-Verteilung erhält man die Exponentialverteilung als einen Spezialfall?

AUFGABE 6.5 Zeichnen Sie den Verlauf der Rate der Weibull-Verteilung auf einer Zeitachse von 0 bis 3, und zwar für die Parameterwerte  $\alpha = 1$  und  $\beta = 0.5, 1.0, 1.5$ .

## 6.3 Loglogistische Verteilung

Die Exponentialverteilung kann nur zeitkonstante, die Weibull-Verteilung nur monotone Ratenverläufe ausdrücken. Für nichtmonotone Ratenverläufe eignet sich manchmal die sog. *loglogistische Verteilung*. Sie hat, wie die Weibull-Verteilung, zwei Parameter,  $\alpha$  und  $\beta$ , die nur positive Werte annehmen können. Die Survivorfunktion sieht folgendermaßen aus:

$$G(t) = \frac{1}{1 + (\alpha t)^\beta}$$

Daraus erhält man durch Differenzieren die Dichtefunktion

$$f(t) = \frac{\beta \alpha^\beta t^{\beta-1}}{(1 + (\alpha t)^\beta)^2}$$

und schließlich die Rate

$$r(t) = \frac{\beta \alpha^\beta t^{\beta-1}}{1 + (\alpha t)^\beta}$$

AUFGABE 6.6 Zeigen Sie schrittweise, wie sich die Dichtefunktion und die Rate aus der Survivorfunktion der loglogistischen Verteilung berechnen lassen.

AUFGABE 6.7 Zeigen Sie, daß die Exponentialverteilung *kein* Spezialfall der loglogistischen Verteilung ist.

AUFGABE 6.8 Zeichnen Sie den Verlauf der Rate der loglogistischen Verteilung auf einer Zeitachse von 0 bis 3, und zwar für die Parameterwerte  $\alpha = 1$  und  $\beta = 1, \beta = 2$ .

AUFGABE 6.9 Zeigen Sie, daß die Rate der loglogistischen Verteilung (wenn sie konkav ist) ihr Maximum in der Zeitstelle

$$t_{\max} = \frac{1}{\alpha} (\beta - 1)^{\frac{1}{\beta}}$$

annimmt. Für welche Parameterwerte ist der Verlauf der Rate konkav? Welchen Wert nimmt die Rate in ihrem Maximum an?

## 6.4 Lognormal-Verteilung

Bei Regressionsmodellen wird oft eine Normalverteilung unterstellt. Da unsere Verweildauervariable  $T$  nur positive Werte annehmen kann, liegt es nahe, ihren Logarithmus zu betrachten und dafür eine Normalverteilung zu unterstellen. Allgemein sagt man: eine Variable ist *lognormal* verteilt, wenn ihr Logarithmus normalverteilt ist.<sup>1</sup>

Um diesen Gedanken zu verfolgen, betrachten wir zunächst ein allgemeineres Problem. Es sei  $X$  eine kontinuierliche Variable mit der Verteilungsfunktion  $F_X(x)$  und der Dichtefunktion  $f_X(x)$ . Außerdem sei  $g$  eine beliebige monoton steigende Funktion. Wir können dann eine neue Variable

$$Y = g(X)$$

bilden. Die Frage ist, wie man die Verteilungs- und Dichtefunktionen von  $Y$ , also  $F_Y(y)$  und  $f_Y(y)$ , aus den entsprechenden Funktionen für  $X$  berechnen kann; oder umgekehrt, wie man die Verteilung von  $X$  aus der Verteilung von  $Y$  finden kann. Da die Transformationsfunktion  $g$  monoton ist, kann man natürlich auch

$$X = g^{-1}(Y)$$

betrachten, insofern ist das Problem symmetrisch. Zur Beantwortung der Fragen gibt es zunächst folgende Beziehung zwischen den Verteilungsfunktionen:

$$F_X(x) = P(X \leq x) = P(Y \leq g(x)) = F_Y(g(x))$$

Beziehungen zwischen den Dichtefunktionen ergeben sich aus Ableitungen der Verteilungsfunktionen. Man findet

$$f_X(x) = \left. \frac{dF_X(u)}{du} \right|_{u=x} = \left. \frac{dF_Y(g(u))}{du} \right|_{u=x} = f_Y(g(x))g'(x)$$

wobei  $g'(x)$  die Ableitung der Transformationsfunktion  $g$  bezeichnet.

Wenden wir jetzt diese Überlegung auf eine Situation an, in der die Variable  $Y$  normalverteilt und der Zusammenhang zwischen  $T$ , unserer Verweildauervariablen, und  $Y$  durch die Transformation

$$Y = \log(T)$$

definiert ist. Für die Dichtefunktion von  $Y$  haben wir also

$$f_Y(y) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{1}{2} \left( \frac{y - \mu}{\sigma} \right)^2 \right\}$$

<sup>1</sup> Wir verwenden hier stets den natürlichen Logarithmus, also die Umkehrfunktion zur Exponentialfunktion.

Wenden wir jetzt die oben abgeleitete Transformationsformel an, erhalten wir für die Dichtefunktion von  $T$ , die wir jetzt wieder mit  $f(t)$  bezeichnen, den Ausdruck

$$f(t) = \frac{1}{\sqrt{2\pi}\sigma t} \exp \left\{ -\frac{1}{2} \left( \frac{\log(t) - \mu}{\sigma} \right)^2 \right\}$$

Dies ist die Dichtefunktion der Lognormal-Verteilung mit den Parametern  $\mu$  und  $\sigma$ . Dabei kann  $\mu$  beliebige,  $\sigma$  nur positive Werte annehmen.

Ein Nachteil dieser Verteilung ist, daß es für die Verteilungsfunktion und infolgedessen für die Rate keinen geschlossenen, einfach berechenbaren Ausdruck gibt. Die Lognormal-Verteilung ist trotzdem wichtig, weil sie es erlaubt, einen Zusammenhang zwischen Ratenmodellen und üblichen Regressionsmodellen herzustellen.

**AUFGABE 6.10** Zeichnen Sie den Verlauf der Dichtefunktion der Lognormal-Verteilung auf einer Zeitachse von 0 bis 3, und zwar für die Parameterwerte  $\mu = 0$  und  $\sigma = 1$ .

## 6.5 Mehrere Zielzustände

Bisher haben wir Episoden betrachtet, bei denen es nur einen möglichen Folgezustand gibt. Wenn es zwei oder mehr mögliche Folgezustände gibt, muß die zweidimensionale Variable  $(T, D)$  betrachtet werden. Wie wir uns schon überlegt haben, kann ihre Verteilung durch zielzustandsspezifische Raten charakterisiert werden:

$$r_d(t) \quad \text{für } d \in \tilde{D}$$

Um zu Verteilungsmodellen zu gelangen, kann man also die in den vorangegangenen Abschnitten diskutierten Annahmen über Raten einfach auf die zielzustandsspezifischen Raten übertragen. Dabei hat man die Möglichkeit, für alle Raten das gleiche Verteilungsmodell, aber mit jeweils unterschiedlichen Verteilungsparametern, zu verwenden; man kann aber auch für jede zielzustandsspezifische Rate ein eigenes Verteilungsmodell annehmen.

Im einfachsten Fall könnte man zum Beispiel für alle zielzustandsspezifischen Raten zeitkonstante Werte annehmen, also

$$r_d(t) = \theta_t$$

## 6.6 Mischungen

Mischungsmodelle ergeben sich, wenn man annimmt, daß die Gesamtheit  $\Omega$  aus zwei oder mehr Teilgesamtheiten besteht, bei denen es jeweils unterschiedliche

Ratenverläufe gibt. Nehmen wir an, daß es  $m$  Teilgesamtheiten gibt. Zum Zeitpunkt  $t = 0$  sei der Anteil der  $j$ .ten Teilgesamtheit durch  $\pi_j$  gegeben, also

$$\sum_{j=1}^m \pi_j = 1$$

$r_j(t)$  sei die Rate,  $f_j(t)$  die Dichtefunktion und  $G_j(t)$  die Survivorfunktion für die  $j$ .te Teilgesamtheit. Der Anteil der  $j$ .ten Teilgesamtheit entwickelt sich dann im Zeitablauf entsprechend

$$\pi_j G_j(t) / \sum_{k=1}^m \pi_k G_k(t)$$

wobei

$$G(t) = \sum_{j=1}^m \pi_j G_j(t)$$

die Survivorfunktion in der Gesamtheit ist. Daraus findet man die Dichtefunktion

$$f(t) = -\frac{dG(t)}{dt} = \sum_{j=1}^m \pi_j \left( -\frac{dG_j(t)}{dt} \right) = \sum_{j=1}^m \pi_j f_j(t)$$

Schließlich findet man für die durchschnittliche Rate in der Gesamtheit den Ausdruck

$$r(t) = \frac{\sum_{j=1}^m \pi_j f_j(t)}{\sum_{j=1}^m \pi_j G_j(t)}$$

**AUFGABE 6.11** Nehmen Sie an, daß es zwei Teilgesamtheiten gibt und  $\pi_1 = \pi_2 = 0.5$  ist. Nehmen Sie außerdem an, daß es in der ersten Teilgesamtheit eine konstante Rate  $r_1(t) = 1$ , in der zweiten Teilgesamtheit eine konstante Rate  $r_2(t) = 2$  gibt. Berechnen Sie die Entwicklung der durchschnittlichen Rate im Zeitablauf und zeigen Sie, daß diese Rate immer kleiner wird.

## 7 Ratenmodelle mit Kovariablen

In diesem Kapitel werden Ratenmodelle mit zeitunabhängigen Kovariablen behandelt. Wie im vorangegangenen Kapitel beziehen wir uns auf Episoden, die wir durch eine zweidimensionale Variable  $(T, D)$  repräsentieren. Allerdings müssen wir jetzt berücksichtigen, daß es noch Kovariablen gibt. Als Ausgangspunkt haben wir also eine Variable  $(T, D, X)$ , wobei  $X$  die Kovariablen bezeichnet. Im allgemeinen kann es sich bei  $X$  um eine ein- oder mehrdimensionale Variable handeln. Zunächst behandeln wir Modelle für Episoden, bei denen es nur einen möglichen Folgezustand gibt; Modelle für Episoden mit mehreren Folgezuständen werden in Abschnitt 7.4 besprochen.

### 7.1 Das Exponentialmodell

Beim Exponentialmodell wird für die Verweildauerverteilung eine Exponentialverteilung unterstellt, also eine zeitkonstante Rate

$$r(t) = \theta$$

Die Idee ist, daß diese Rate von den Werten von Kovariablen abhängig sein kann. Da die Rate nur positive Werte annehmen kann, sollte eine Link-Funktion verwendet werden, die dies garantiert. Beim Standardmodellansatz wird als Link-Funktion eine Exponentialfunktion verwendet. Wenn  $X$  eine eindimensionale Kovariable ist, sieht das Modell dann so aus:

$$r(t | X = x) = \exp(\beta_0 + x\beta_1)$$

Wenn es  $m$  Kovariablen  $(X_1, \dots, X_m)$  gibt, kann man folgende allgemeine Formulierung verwenden:

$$r(t | X_1 = x_1, \dots, X_m = x_m) = \exp(\beta_0 + x_1\beta_1 + \dots + x_m\beta_m)$$

**AUFGABE 7.1** Beweisen Sie folgende Formel für den Mittelwert einer Exponentialverteilung mit dem Parameter  $\theta$ :<sup>1</sup>

$$E(T) = \int_0^{\infty} t f(t) dt = \int_0^{\infty} \theta t \exp(-\theta t) dt = \frac{1}{\theta}$$

**AUFGABE 7.2** Finden Sie eine Formel für den Median einer Exponentialverteilung mit dem Parameter  $\theta$ .

<sup>1</sup> Verwenden Sie folgende Regel für partielle Integration:

$$\int F(t)g(t) dt = F(t)G(t) - \int f(t)G(t) dt$$

wobei  $f(t) = dF(t)/dt$  und  $g(t) = dG(t)/dt$ . Setzen Sie  $F(t) = \theta t$ ,  $g(t) = \exp(-\theta t)$ .

**Box 7.1** Datensatz 6

ID	T	D	X1	X2
1	17	1	8	1
2	5	0	5	0
3	22	1	9	1
4	13	1	7	0
5	2	0	5	0
6	9	1	6	1
7	12	0	5	1
8	15	1	7	1

## 7.2 Parameterschätzungen

Wir nehmen an, daß uns  $n$  Beobachtungen

$$(t_i, d_i, x_{i1}, \dots, x_{im}) \quad \text{für } i = 1, \dots, n$$

gegeben sind.  $t_i$  ist der Wert der Verweildauer  $T$ ,  $d_i$  gibt an, ob es sich um eine zensierte ( $d_i = 0$ ) oder unzensierte ( $d_i = 1$ ) Beobachtung handelt, und  $x_{i1}, \dots, x_{im}$  sind die Werte der Kovariablen. Zur Modellschätzung wird die ML-Methode verwendet. Für die Bildung der Likelihood-Funktion wird bei nicht-zensierten Beobachtungen die Dichtefunktion

$$f(t_i | X_1 = x_{i1}, \dots, X_m = x_{im})$$

und bei zensierten Beobachtungen die Survivorfunktion

$$G(t_i | X_1 = x_{i1}, \dots, X_m = x_{im})$$

verwendet; sie sieht also folgendermaßen aus:

$$\mathcal{L}(\beta_0, \dots, \beta_m) = \prod_{i=1}^n f(t_i | X_1 = x_{i1}, \dots, X_m = x_{im})^{d_i} G(t_i | X_1 = x_{i1}, \dots, X_m = x_{im})^{1-d_i}$$

Wegen der Beziehung  $r(t) = f(t)/G(t)$  kann man diese Likelihood-Funktion auch folgendermaßen schreiben:

$$\mathcal{L}(\beta_0, \dots, \beta_m) = \prod_{i=1}^n r(t_i | X_1 = x_{i1}, \dots, X_m = x_{im})^{d_i} G(t_i | X_1 = x_{i1}, \dots, X_m = x_{im})$$

Wenn wir uns jetzt auf das in Abschnitt 7.1 besprochene Exponentialmodell beziehen, finden wir zunächst:

$$r(t_i | X_1 = x_{i1}, \dots, X_m = x_{im}) = \exp(\beta_0 + x_{i1}\beta_1 + \dots + x_{im}\beta_m)$$

$$G(t_i | X_1 = x_{i1}, \dots, X_m = x_{im}) = \exp\{-\exp(\beta_0 + x_{i1}\beta_1 + \dots + x_{im}\beta_m) t_i\}$$

Daraus findet man die Likelihood-Funktion

$$\mathcal{L}(\beta_0, \dots, \beta_m) = \prod_{i=1}^n \exp(\beta_0 + x_{i1}\beta_1 + \dots + x_{im}\beta_m)^{d_i} \exp\{-\exp(\beta_0 + x_{i1}\beta_1 + \dots + x_{im}\beta_m) t_i\}$$

Einen einfacheren Ausdruck erhält man, wenn man zur Log-Likelihood-Funktion übergeht:

$$\ell(\beta_0, \dots, \beta_m) = \sum_{i=1}^n d_i(\beta_0 + x_{i1}\beta_1 + \dots + x_{im}\beta_m) - \exp(\beta_0 + x_{i1}\beta_1 + \dots + x_{im}\beta_m) t_i$$

Im allgemeinen, wenn das Modell Kovariablen enthält, kann man das Maximum dieser Funktion nicht auf analytischem Wege finden, sondern benötigt ein iteratives Verfahren.<sup>2</sup> Einfach geht es nur, wenn das Modell keine Kovariablen enthält. Die Log-Likelihood-Funktion sieht dann folgendermaßen aus:

$$\ell(\beta_0) = \sum_{i=1}^n d_i\beta_0 - t_i \exp(\beta_0)$$

Daraus gewinnt man den Gradienten, also die erste Ableitung

$$\frac{\partial \ell(\beta_0)}{\partial \beta_0} = \sum_{i=1}^n d_i - t_i \exp(\beta_0)$$

Den ML-Schätzwert für  $\beta_0$  findet man schließlich an der Nullstelle des Gradienten, also durch

$$\hat{\beta}_0 = \log\left(\frac{\sum_{i=1}^n d_i}{\sum_{i=1}^n t_i}\right)$$

**AUFGABE 7.3** Bilden Sie die zweite Ableitung der Log-Likelihoodfunktion für das einfache Exponentialmodell ohne Kovariablen und zeigen Sie anhand dieser Ableitung, daß diese Funktion genau ein globales Maximum hat.

**AUFGABE 7.4** Nehmen Sie für die Daten in Box 7.1 ein Exponentialmodell ohne Kovariablen an. Schätzen Sie dafür die zeitkonstante Rate und gewinnen Sie daraus einen Schätzwert für die mittlere Verweildauer.

<sup>2</sup> Das wird in einigen Übungen im Anhang besprochen.

### 7.3 Ein allgemeiner Modellansatz

Wie im vorangegangenen Kapitel besprochen worden ist, kann man zahlreiche unterschiedliche Vorstellungen über den Ratenverlauf einer Episode bilden. Bei parametrischen Modellansätzen geht man von einer Verweildauerverteilung aus, die im allgemeinen von einem ein- oder mehrdimensionalen Parameter  $\theta$  abhängt. Rate, Dichtefunktion und Survivorfunktion können also allgemein folgendermaßen geschrieben werden:

$$r(t|\theta), f(t|\theta), G(t|\theta)$$

Kovariablen können in einen solchen Modellansatz eingebaut werden, indem man mittels einer Link-Funktion die Verteilungsparameter von den Kovariablen abhängig macht. Wenn also die Verteilung zum Beispiel zwei Parameter hat, also  $\theta = (\alpha, \beta)$ , kann man sie von den Kovariablen auf folgende Weise abhängig machen:

$$\begin{aligned}\alpha &= g_\alpha(\alpha_0 + x_1\alpha_1 + \dots + x_m\alpha_m) \\ \beta &= g_\beta(\beta_0 + x_1\beta_1 + \dots + x_m\beta_m)\end{aligned}$$

Hierbei ist  $g_\alpha$  die Link-Funktion für den Verteilungsparameter  $\alpha$ ,  $g_\beta$  die Link-Funktion für den Verteilungsparameter  $\beta$ . Zur Unterscheidung von den Verteilungsparametern werden  $\alpha_0, \dots, \alpha_m$  und  $\beta_0, \dots, \beta_m$  als *Modellparameter* bezeichnet. Natürlich ist es möglich, die Verteilungsparameter von unterschiedlichen Kovariablen abhängig zu machen oder auch vollständig auf Kovariablen zu verzichten.

Es sollte evident sein, wie sich dieser Ansatz für Verteilungen, die nur einen oder mehr als zwei Verteilungsparameter enthalten, modifizieren läßt. Zum Beispiel ist das in Abschnitt 7.1 behandelte Exponentialmodell ein einfacher Spezialfall dieses allgemeinen Modellansatzes.

Die Likelihood-Funktion kann dann in jedem Fall so formuliert werden, wie wir es für das Exponentialmodell gezeigt haben.

**AUFGABE 7.5** Entwickeln Sie die Log-Likelihood-Funktion für ein Weibull-Modell ohne Kovariablen.

**AUFGABE 7.6** Entwickeln Sie die Log-Likelihood-Funktion für ein log-logistisches Modell ohne Kovariablen.

### 7.4 Mehrere Folgezustände

Der im vorangegangenen Abschnitt behandelte allgemeine Modellansatz läßt sich leicht für Episoden verallgemeinern, bei denen es mehrere mögliche Folgezustände gibt. Ausgangspunkt ist eine Verteilungsannahme über die

zielzustandsspezifischen Raten, die wir allgemein folgendermaßen schreiben können:

$$r_d(t|\theta_d) \quad \text{für } d \in \tilde{D}$$

Daraus gewinnt man zunächst die Gesamtrate

$$r(t|\theta) = \sum_{d \in \tilde{D}} r_d(t|\theta_d)$$

wobei hier  $\theta$  für die Gesamtheit der zustandsspezifischen Verteilungsparameter  $\theta_d$  steht. Die Gesamtrate liefert die Survivorfunktion

$$G(t|\theta) = \exp \left\{ - \int_0^t r(\tau|\theta) d\tau \right\}$$

oder

$$G(t|\theta) = \prod_{d \in \tilde{D}} G_d(t|\theta_d)$$

wobei

$$G_d(t|\theta_d) = \exp \left\{ - \int_0^t r_d(\tau|\theta_d) d\tau \right\}$$

ist. Um schließlich die Likelihood-Funktion zu finden, nehmen wir an, daß die Daten so gegeben sind, wie in Abschnitt 7.2 angegeben worden ist, nur daß jetzt  $d_i$  mehr als zwei mögliche Werte annehmen kann: Wenn  $d_i = 0$  ist, handelt es sich um eine zensierte Beobachtung, wenn  $d_i > 0$  ist, handelt es sich um einen Übergang in den Folgezustand  $d_i$ . Also kann man die Likelihood-Funktion folgendermaßen schreiben:

$$\mathcal{L}(\theta) = \prod_{i=1}^n G(t_i|\theta) \prod_{d \in \tilde{D}} r(t_i|\theta_d)^{I(d=d_i)}$$

wobei die Indikatorfunktion  $I(d = d_i)$  den Wert 1 annimmt, wenn  $d = d_i$  ist, andernfalls den Wert 0.

**AUFGABE 7.7** Entwickeln Sie zunächst eine Likelihood-Funktion und dann eine Log-Likelihood-Funktion für ein Exponentialmodell, bei dem es drei verschiedene Folgezustände gibt.

## 7.5 Pseudo-Residuen

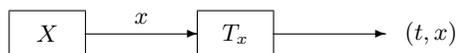
Im Unterschied zu gewöhnlichen Regressionsmodellen mit einer quantitativen abhängigen Variablen gibt es bei der ML-Schätzung von Ratenmodellen keine einfache Methode, um zu beurteilen, wie gut ein Modell zu den Daten paßt. Ein gewisses Hilfsmittel für die Modelldiagnostik liefern jedoch sog. Pseudo-Residuen. Um die Idee zu verdeutlichen, beziehen wir uns auf eine Episode mit einem möglichen Folgezustand. Daten seien in der Form

$$(t_i, d_i, x_i) \quad \text{für } i = 1, \dots, n$$

gegeben. Außerdem wird angenommen, daß bereits ein Modell geschätzt worden ist und die Schätzergebnisse durch

$$r(t | x; \hat{\theta}), f(t | x; \hat{\theta}), G(t | x; \hat{\theta})$$

gegeben sind, wobei  $\hat{\theta}$  die geschätzten Modellparameter bezeichnet. Dieses Modell kann nun als Beschreibung eines zweistufigen Zufallsgenerators betrachtet werden, der folgendermaßen aussieht:



Der erste Zufallsgenerator liefert einen Wert  $x$  für den Kovariablenvektor entsprechend der in den Daten gegebenen Verteilung von  $X$ . Der zweite Zufallsgenerator liefert eine Verweildauer  $t$ , wobei vom Modell  $f(t | x; \hat{\theta})$  ausgegangen wird, also konditional auf den im ersten Schritt realisierten Wert des Kovariablenvektors. Auf diese Weise kann eine beliebige Menge von Pseudo-Beobachtungen  $(t_j^*, x_j^*)$  erzeugt werden ( $j = 1, 2, 3, \dots$ ), so daß die Verteilung der  $x_j^*$  mit der durch die Daten gegebenen Verteilung von  $X$  und die konditionale Verteilung von  $t_j^*$  mit dem Modell übereinstimmt. Für jeden möglichen Wert  $x$  aus dem Merkmalsraum der Kovariablen  $X$  kann man die auf diese Weise erzeugten Werte als Realisierungen einer Zufallsvariablen  $T_x$  ansehen, deren Verteilung durch das Modell, d.h. durch  $f(t | x; \hat{\theta})$  definiert ist.

Im nächsten Schritt wird eine Transformation der Zufallsvariablen  $T_x$  betrachtet, so daß die Abhängigkeit vom jeweils realisierten Wert des Kovariablenvektors verschwindet. Als Transformation wird

$$T_x \longrightarrow J(T_x), \quad \text{definiert durch } t \longrightarrow J(t) = \int_0^t r(\tau | x; \hat{\theta}) d\tau$$

verwendet, denn die Verteilung von  $J(T_x)$  ist dann eine Standard-Exponentialverteilung, d.h. eine Exponentialverteilung mit der konstanten Rate 1.

Man sieht das folgendermaßen, wobei ausgenutzt wird, daß es sich um eine monotone Transformation handelt:

$$\begin{aligned} P(J(T_x) > t) &= P(T_x > J^{-1}(t)) \\ &\equiv G(J^{-1}(t) | x; \hat{\theta}) \\ &= \exp \left\{ - \int_0^{J^{-1}(t)} r(\tau | x; \hat{\theta}) d\tau \right\} \\ &= \exp \left\{ -J(J^{-1}(t)) \right\} \\ &= \exp(-t) \end{aligned}$$

Die Survivorfunktion für die transformierte Zufallsvariable  $J(T_x)$  ist also die Survivorfunktion einer Standard-Exponentialverteilung und mithin unabhängig von  $x$ .<sup>3</sup>

Diese Tatsache kann ausgenutzt werden, um zu prüfen, ob es plausibel erscheint, daß die als Stichprobe vorliegenden Daten aus dem durch das Modell beschriebenen Zufallsgenerator stammen könnten. Wenn dies der Fall ist, müßte die Anwendung der oben beschriebenen Transformation auf die vorliegenden Daten zu einer Menge transformierter Größen

$$e_i = J(t_i)$$

führen, die näherungsweise einer Standard-Exponentialverteilung folgen. Die Größen  $e_i$  werden als *Pseudo-Residuen*, gelegentlich auch als *verallgemeinerte Residuen*, bezeichnet. Um zu prüfen, ob sie näherungsweise einer Standard-Exponentialverteilung folgen, kann ihre Survivorfunktion berechnet werden. Dafür sollte das Kaplan-Meier-Verfahren verwendet werden, um berücksichtigen zu können, daß einige Beobachtungen und mithin die ihnen korrespondierenden Residuen rechts zensiert sein können. Wenn  $G_r(t)$  die auf diese Weise berechnete Survivorfunktion bezeichnet, kann man schließlich die Abbildung

$$t \longrightarrow -\log \{G_r(t)\}$$

betrachten. Wenn die Residuen näherungsweise einer Standard-Exponentialverteilung folgen, sollte diese Abbildung näherungsweise einer 45°-Linie entsprechen.

**AUFGABE 7.8** Entwickeln Sie eine Formel zur Berechnung von Pseudo-Residuen für ein Exponentialmodell mit Kovariablen.

<sup>3</sup> Es sei angemerkt, daß die Verteilung von  $T_x$  durch  $G(\cdot | x, \theta)$ , nicht durch  $G(\cdot | x, \hat{\theta})$ , definiert ist. Das Verfahren beruht jedoch darauf, daß  $\theta$  durch die mit den Daten geschätzten Parameterwerte ersetzt wird.

AUFGABE 7.9 Berechnen Sie die Pseudo-Residuen für das in Aufgabe 7.4 geschätzte Exponentialmodell ohne Kovariablen.

## 8 Zeitveränderliche Kovariablen

Im vorangegangenen Kapitel wurde diskutiert, wie Ratenmodelle mit zeitunabhängigen Kovariablen formuliert und geschätzt werden können. Für viele Anwendungen ist es jedoch erforderlich, auch Kovariablen zu berücksichtigen, die ihre Werte während einer laufenden Episode verändern können. Man möchte dann herausfinden, wie die Rate für den Übergang in einen neuen Folgezustand von den jeweils aktuellen, sich während des Episodenverlaufs verändernden Kovariablenwerten bedingt wird. In diesem Kapitel wird besprochen, wie dafür geeignete Ratenmodelle formuliert werden können.

### 8.1 Konditionale Survivorfunktionen

Als Ausgangspunkt erinnern wir an unsere Basisformel für den Zusammenhang zwischen Rate und Survivorfunktion, nämlich

$$G(t) = \exp \left\{ - \int_0^t r(\tau) d\tau \right\}$$

Eine *konditionale Survivorfunktion* wird dementsprechend durch

$$G(t | s) = \exp \left\{ - \int_s^t r(\tau) d\tau \right\}$$

definiert. Offenbar gilt:

$$G(t) = G(t | s) G(s)$$

Dies läßt sich beliebig wiederholen. Wir können zum Beispiel die Zeitspanne von 0 bis  $t$  in  $k$  beliebige Subintervalle aufteilen:

$$0 = t_0 < t_1 < \dots < t_{k-1} < t_k = t$$

Dann finden wir

$$G(t) = \prod_{j=1}^k G(t_j | t_{j-1})$$

### 8.2 Reformulierte Likelihoodfunktion

Erinnern wir uns jetzt an die Likelihoodfunktionen zur Schätzung von Ratenmodellen. Bei einer Episode mit einem möglichen Folgezustand ist die allgemeine Formulierung:

$$\mathcal{L}(\theta) = \prod_{i=1}^n r(t_i | \theta)^{d_i} G(t_i | \theta)$$

Für jedes Individuum  $i$  kann man jetzt die Verweildauer  $t_i$  in beliebig viele Subintervalle aufteilen, etwa

$$0 = t_{i0} < t_{i1} < \dots < t_{i,k_i-1} < t_{ik_i} = t_i$$

Dann kann man die Likelihoodfunktion folgendermaßen reformulieren:

$$\mathcal{L}(\theta) = \prod_{i=1}^n r(t_i | \theta)^{d_i} \prod_{j=1}^{k_i} G(t_{ij} | t_{i,j-1}, \theta)$$

Wichtig ist, daß sich durch diese Reformulierung tatsächlich nur die äußere Form der Likelihoodfunktion verändert, nicht jedoch der funktionale Zusammenhang zwischen Modellparametern und Werten der Likelihoodfunktion. Zur Modellschätzung kann man deshalb ebensogut die reformulierte Likelihoodfunktion verwenden.

**AUFGABE 8.1** Entwickeln Sie eine reformulierte Likelihoodfunktion zur Schätzung eines einfachen Exponentialmodells ohne Kovariablen. Zeigen Sie, daß man mit der reformulierten Likelihoodfunktion die gleichen Schätzwerte erhält wie mit der ursprünglichen Likelihoodfunktion.

### 8.3 Zeitveränderliche Indikatorvariablen

Jetzt nehmen wir an, daß die Daten zeitveränderliche Kovariablen enthalten und in folgender Form gegeben sind:

$$(t_i, d_i, x_i, z_{i1}(t), \dots, z_{im}(t)) \quad \text{für } i = 1, \dots, n$$

$x_i$  ist wie bisher ein Vektor mit zeitunabhängigen Kovariablen, deren Werte also spätestens zu Beginn der Episode feststehen. Die Variablen  $z_{ij}(t)$  sind dagegen zeitveränderlich. Und zwar nehmen wir an, daß es sich um zeitveränderliche 0-1-Variablen handelt, das heißt: bis zu einem gewissen Zeitpunkt haben sie den Wert 0, dann den Wert 1. Es sei  $t_{ij}$  der Zeitpunkt, bei dem die Variable  $z_{ij}(t)$  ihren Wert von 0 auf 1 wechselt. Wir können dann die Gesamtheit der Zeitpunkte innerhalb der Episoden  $[0, t_i]$  betrachten, zu denen mindestens eine der zeitveränderlichen Kovariablen ihren Wert verändert.<sup>1</sup> Wir stellen uns vor, daß diese Zeitpunkte folgendermaßen der Größe nach geordnet sind:

$$0 = \tau_{i0} < \tau_{i1} < \dots < \tau_{i,k_i-1} < \tau_{ik_i} = t_i$$

<sup>1</sup> Es ist klar, daß zeitveränderliche Kovariablen, die ihre Werte während des Episodenverlaufs nicht verändern, wie zeitunabhängige Kovariablen behandelt werden können.

#### Box 8.1 Datensatz 7: Episodensplitting

ID	ORG	DES	TS	TF	ID	SPN	ORG	DES	TS	TF
1	0	1	0	10	1	1	0	0	0	3
2	0	0	0	12	1	2	0	0	3	6
					1	3	0	1	6	10
					2	1	0	0	0	5
					2	2	0	0	5	9
					2	3	0	0	9	11
					2	4	0	0	11	12

#### Box 8.2 Datensatz 8: Episodensplitting

ID	ORG	DES	TS	TF	X	ID	SPN	ORG	DES	TS	TF	D
1	0	1	0	10	7	1	1	0	0	0	7	0
2	0	1	0	8	-1	1	2	0	1	7	10	1
3	0	1	0	5	8	2	1	0	1	0	0	1
4	0	0	0	12	9	3	1	0	1	0	5	0
						4	1	0	0	0	9	0
						4	2	0	0	9	12	1

Dann kann man die Likelihood unter Berücksichtigung der zeitunabhängigen und der zeitveränderlichen Kovariablen folgendermaßen schreiben:

$$\mathcal{L}(\theta) = \prod_{i=1}^n r(t_i | z_i, x_{i1}(t_i), \dots, x_{im}(t_i), \theta)^{d_i} \prod_{j=1}^{k_i} G(\tau_{ij} | \tau_{i,j-1}, z_i, x_{i1}(\tau_{i,j-1}), \dots, x_{im}(\tau_{i,j-1}), \theta)$$

Offenbar berücksichtigt man dadurch für jedes Teilstück des Episodenverlaufs die dafür gegebenen aktuellen Werte der Kovariablen.

### 8.4 Episodensplitting

Die praktische Umsetzung der eben beschriebenen Idee zur Berücksichtigung zeitveränderlicher Kovariablen geschieht mit der Methode des Episodensplitting. Wir erklären die Methode zunächst an einem Datensatz ohne Berücksichtigung von Kovariablen. Box 8.1 zeigt den Datensatz 7, zunächst in ungesplitteter Form auf der linken Seite. Es gibt zwei Episoden. Die Episode für das erste Individuum endet mit einem Ereignis, die für das zweite Individuum ist rechts zensiert. Die rechte Seite zeigt, wie die Episoden

gesplittet worden sind; in diesem Beispiel ganz willkürlich, die erste Episode in drei, die zweite in vier Splits.

Box 8.2 zeigt ein zweites Beispiel mit einer zeitveränderlichen Variablen ( $X$ ); diese Variable enthält zunächst den Zeitpunkt, zu dem die korrespondierende 0-1-Variable ( $D$ ) ihren Wert von 0 auf 1 verändert. Die rechte Seite der Box zeigt dann den gesplitteten Datensatz.

## A Übungen mit R

AUFGABE A.1 Installieren Sie das Programm R. Es gibt Versionen für alle gängigen Betriebssysteme (Linux, Windows, Mac) auf <http://cran.r-project.org/>. R ist ein Interpreter, d.h. es führt Befehle nach Eingabe aus und präsentiert das Ergebnis. Daher kann man es wie einen Taschenrechner benutzen: Tippen Sie

```
> 2+2
```

antwortet R mit

```
[1] 4
```

Hilfe erhält man mit

```
> ?"+"
```

```
> ?matrix
```

```
> help(log)
```

etc. Die HTML-Version kann mit

```
> options(htmlhelp=TRUE)
```

```
> help.start()
```

gestartet werden. Alternativ kann die HTML-Hilfe direkt in einem Browser geöffnet werden. Sie können R mit

```
> q()
```

verlassen. Sie werden dann gefragt, ob Sie den sogenannten „workspace“ speichern wollen. I.d.R. sollten Sie mit **n** (no) antworten.

AUFGABE A.2 Die meisten Statistikprogramme benutzen zur Berechnung numerischer Ausdrücke „double precision reals“. Berechnungen in diesem „Zahlensystem“ führen immer zu Rundungs- und Abschneidefehlern. Die maximale Genauigkeit dieser Arithmetik entspricht etwa 15 Dezimalstellen. Die Druckdarstellung der Zahlen kann in R durch `options(digits=5)` verändert werden (7 ist die Voreinstellung). Dies verändert nicht die Rechengenauigkeit!

a) Berechnen Sie  $2 + 3 * (4 + 5 * (6 + 7 * (8 + 9)))$ ,  $2^{11}$ ,  $\log_e(2)$ , und  $\sqrt{333}$ .

AUFGABE A.3 Variable in R können beliebige Namen haben, die aber mit einem Buchstaben beginnen müssen. Groß- und Kleinschreibung wird unterschieden. Die Zuweisung von Werten zu Variablen erfolgt durch `<-`:

```
a <- 3; a
```

Die wichtigste Datenstruktur in R ist ein „Vektor“, eine Liste von Elementen. Vektoren können durch den Befehl `c()` erzeugt werden:

```
a <- c(1,3,5)
```

```
b <- c(a,1,a)
```

Die Länge eines Vektors (die Anzahl seiner Elemente) wird durch `length(b)` berechnet. Fast alle Befehle in R sind „vektoriert“ und operieren auf allen

Elementen eines Vektors.

```
sin(a)+1
```

Wenn Vektoren unterschiedlicher Länge in einem Ausdruck verwandt werden, dann werden in den Berechnungen die kürzeren Vektoren durch Wiederholung ihrer Werte verlängert, bis sie die Länge des längsten Vektors erreicht haben.

```
d <- 2*a + b + 1
```

AUFGABE A.4 Um mit R zu arbeiten, wird ein Editor benötigt, in dem man Programme korrigieren und editieren kann. Wir schlagen vor, mit TINN-R zu arbeiten (<http://www.sciviews.org/Tinn-R/>). Installieren Sie TINN-R.

Legen Sie eine Datei mit dem Namen „test.R“ in einem Arbeitsverzeichnis an und schreiben Sie die folgenden Kommandos in die Datei:

```
a <- c(1,3,5); b <- c(a,1,a)
```

```
d <- 2*a+b+1
```

```
sin(d)
```

Speichern Sie die Datei und versuchen Sie, den Inhalt der Datei aus TINN-R heraus an R zu senden. Beachten Sie dabei, daß R beim Start zunächst ein eigenes Arbeitsverzeichnis wählt. Sie müssen gegebenenfalls R im Menü **Verzeichnis wechseln** Ihr Arbeitsverzeichnis angeben. Der R Befehl `getwd()` zeigt das gegenwärtig benutzte Arbeitsverzeichnis an.

AUFGABE A.5 Eine  $n \times m$  Matrix  $A$  besteht aus  $nm$  rechteckig angeordneten Zahlen mit  $n$  Zeilen und  $m$  Spalten:

$$A = \begin{pmatrix} a_{11} & a_{12} & \dots & a_{1m} \\ a_{21} & a_{22} & \dots & a_{2m} \\ \dots & \dots & \dots & \dots \\ a_{n1} & a_{n2} & \dots & a_{nm} \end{pmatrix}$$

Analog kann man Felder beliebiger Dimension definieren.

Matrizen und Felder beliebiger Dimension werden durch

```
A <- matrix(c(1,2,3,4), nrow=2, ncol=2)
```

```
B <- array(c(1,2,3,4,1,2,3,4), dim=c(2,2,2))
```

erzeugt.

Im Gegensatz zur mathematischen Sprechweise haben in R Vektoren keine Dimension, wohl aber Matrizen und Felder.

```
dim(a); dim(B)
```

Elementweise Operationen werden auch bei Matrizen und Feldern für alle Elemente gleichzeitig durchgeführt:

```
2*A + 1
```

```
C <- matrix(c(3,4,1,2),nrow=2)
```

```
A + C
```

Elemente von Vektoren, Matrizen und Feldern können durch die Angabe ihrer Indizes angesprochen werden.

```
a[2]
```

```
a[c(1,2)]
```

```
A[1,1]
```

```
A[,1]
```

```
C[C>2]
```

```
B[1,2,1]
```

Negative Indizes entfernen die entsprechenden Elemente.

```
a[-2]
```

AUFGABE A.6 Listen sind geordnete Mengen von verschiedenen Objekten:

```
d <- list(a, c(-1,3,5), "Gesundheit"); d
```

`unlist()` macht aus einer Liste einen Vektor. Elemente einer Liste können durch Angabe ihres Indexes ausgewählt werden.

```
d[1]
```

```
d[[1]]
```

Die erste Form ergibt eine Liste mit den ausgewählten Elementen, die zweite Form die ausgewählten Elemente. Da das erste Element von `d` ein Vektor ist, ist

```
d[[1]][2]
```

das zweite Element des Vektors `a`. Elemente einer Liste können Namen haben. Einzelne Elemente einer Liste können auch über ihren Namen ausgewählt werden:

```
d <- list(eins=a, zwei=c(-1,3,5), drei="Gesundheit")
```

```
names(d)
```

```
names(d) <- c("alpha", "beta", "Wort")
```

```
names(d); d$Wort
```

AUFGABE A.7 Die Funktion `runif()` liefert gleichverteilte Pseudozufallszahlen auf dem 0-1 Intervall. Die Funktion `rnorm()` berechnet normalverteilte Pseudozufallszahlen. Insbesondere liefert

```
X <- runif(100)
```

100 gleichverteilte Pseudozufallszahlen, die der Variablen mit dem Namen `X` zugewiesen werden. Man kann die Startwerte des Zufallszahlengenerators durch `set.seed(n)` festlegen, wobei `n` eine ganze Zahl ist.

a) Erzeugen Sie 100 gleichverteilte Pseudozufallszahlen auf dem Intervall  $(-1, 1)$ .

b) Erzeugen Sie 100 Pseudozufallszahlen mit Werten in  $\{1, 2\}$  und den Wahrscheinlichkeiten  $P(\{1\}) = P(\{2\}) = 0.5$ . Hinweis: Der Operator `X < Y` erzeugt den Wert `TRUE`, falls `X < Y` ist, `FALSE` sonst. `as.numeric(X)` verwandelt diese logischen Werte in die Zahlen 1 bzw. 0.

AUFGABE A.8 `sum()`, `mean()`, `var()`, `sd()` berechnen die Summe, den Mittelwert, die Varianz und die Standardabweichung ihrer Argumente. `cov()` und `cor()` berechnen Kovarianzen und Korrelationen. `summary()` berechnet die Quartile, den Mittelwert und die Extrema, wenn das Argument der Funktion ein numerischer Vektor ist. `plot()` interpretiert die Elemente seiner Argumente als Koordinaten und malt sie in einem Graphikfenster. `hist()` malt ein Histogramm. `density()` berechnet einen Kern-Dichte-Schätzer der Daten.

```
X <- runif(100)
mean(X)
sum(X)/length(X)
var(X)
Y <- X/2 + runif(100)
cov(X,Y)
plot(X,Y)
summary(X)
hist(X)
plot(density(X))
```

- Erzeugen Sie zwei Variable wie in den vorletzten Aufgaben mit je 1000 Beobachtungen und berechnen Sie deskriptive Statistiken.
- Berechnen Sie einen Kern-Dichte-Schätzer für die Pseudozufallszahlen mit der Verteilung  $F(x) = 1 - \exp(-0.5x)$  und zeichnen Sie das Ergebnis.
- Bilden Sie neue Variable durch

```
X <- rnorm(1000)
Y <- rnorm(1000) + X
Z <- rnorm(1000) + X
```

Was ist die Korrelation zwischen X und Y? Zwischen Y und Z?

AUFGABE A.9 Erzeugen Sie (etwa mit `TINN-R`) eine Datei mit den Daten aus A.1. Schreiben (oder kopieren) Sie auch die erste Zeile mit den Variablennamen in die Datei. Speichern Sie die Datei unter dem Namen `eha1.dat` in Ihrem Arbeitsverzeichnis. Die Daten können dann in R durch `dat <- read.table("eha1.dat", header=TRUE)` in eine Matrix (genauer: ein `data.frame`) mit dem Namen `dat` eingelesen werden. Die Option `header=TRUE` steht dabei dafür, die erste Zeile als die Namen der entsprechenden Spalten (Variablen) zu interpretieren (der logische Wert `TRUE` kann durch `T` abgekürzt werden).

### Box A.1 Datei `eha1.dat`

ID	DUR	CEN
1	17	1
2	5	0
3	22	1
4	13	1
5	2	0
6	9	1
7	12	0
8	15	1

- Berechnen Sie den Mittelwert der Variablen `DUR`. Benutzen Sie dazu den `mean()` Befehl. Sprechen Sie den Vektor `DUR` sowohl in der Form `dat$DUR` als auch in der Form `dat[,2]` an.

AUFGABE A.10 Die Ansprache einzelner Variabler in der Form `dat$DUR` ist umständlich, und die Variante `dat[,2]` ist bei vielen Variablen sehr fehleranfällig. Man kann die Variablennamen eines `data.frame` in einer R Sitzung durch den Befehl

```
attach(dat)
```

zugänglich machen. Dadurch wird eine Kopie der Daten erzeugt, für die die Namen der Variablen (Spalten) direkt zugänglich sind.

- Schicken Sie den Befehl `attach(dat)` an R. Sie können anschließend direkt auf die Namen der Variablen zugreifen. Versuchen Sie z.B. `mean(DUR)`.
- Verändern Sie den zweiten Wert von `DUR` zum Wert 2, etwa durch `DUR[2] <- 2`. Dies verändert den Wert in der durch `attach()` erstellten Kopie der Daten. Überprüfen Sie die Werte, indem Sie `DUR` aufrufen.
- Da `attach()` eine Kopie erstellt, können Sie die ursprünglichen Daten wiedererhalten, ohne sie erneut einzulesen. Probieren Sie: `detach(dat)`  
`dat$DUR[2]`
- Sie überschreiben und verändern die ursprünglichen Daten nur durch direkten Verweis auf die Elemente der Datenmatrix (oder auf die Elemente des entsprechenden `data.frame`). Versuchen Sie, den Wert von `dat[2,1]` in einen anderen Wert zu verändern. Was passiert, wenn Sie nicht numerische Werte, etwa Wörter oder logische Werte wie `T` oder `F` verwenden?

AUFGABE A.11 In R sind viele spezielle Funktionen und statistische Prozeduren in sogenannte `libraries` ausgelagert, sowohl um den Überblick in den Hilfsfunktionen als auch die Entwicklung eigener Programme zu erleichtern und um Probleme mit möglichen Überschneidungen mit gleichen Befehlsnamen zu minimieren. Die grundlegende `library` für Verlaufsdaten ist die `library survival`. Sie können ihre Funktionen dem Interpreter von R zur Verfügung stellen, indem Sie

```
library(survival)
```

eingeben. Ihre wichtigste Funktion ist die Bereitstellung eines Befehls, der die Angabe der Form zensierter Daten ermöglicht. Im einfachsten Fall muß nur angegeben werden, ob eine Zeitstelle als zensiert oder als Ereignis gewertet werden soll. Das erfolgt durch die Angabe der Form (der Datensatz `eha1.dat` sollte `attach()`ed sein):

```
library(survival)
```

```
s <- Surv(dat$DUR, dat$CEN)
```

Die Konvention in dieser Angabe ist, daß der Wert 1 (oder `TRUE`) für unzensierte Beobachtungen steht, der Wert 0 (oder `FALSE` aber für zensierte Beobachtungen).

- Untersuchen Sie die Struktur des Objekts `s` (z.B. mit dem Befehl `str()`, welcher die ersten 3–4 Werte der Variablen und ihre Typen darstellt).
- Der Befehl `erg <- survfit(s)` berechnet einen Kaplan–Meier Schätzer.
- `summary(erg)` liefert einige erste Ergebnisse. Können Sie erkennen, welche Maßzahlen ausgegeben werden? Wie erhalten Sie die Ereigniszeitpunkte aus dem Objekt `erg` zurück?
- Der Befehl `plot(erg)` ergibt eine Darstellung des Kaplan–Meier–Schätzers der Survivorfunktion einschließlich der 95% Konfidenzintervalle.
- Versuchen Sie, das Ergebnis (und damit die Survivorfunktion) zu plotten (etwa durch `plot(erg)`).
- Können Sie die Kurven interpretieren?

AUFGABE A.12 Berechnen Sie Ober- und Unterschranken des Kaplan–Meier–Schätzers, indem Sie entweder annehmen, zensierte Beobachtungen hätten ein Ereignis genau zur Zeit der Zensur bzw. am Zeitpunkt der spätesten Beobachtung.

AUFGABE A.13 Interpretieren Sie die Daten aus A.1 als links abgeschnitten, indem Sie nur Fälle mit einer Mindestdauer von 10 betrachten. Wie können solche *links abgeschnittene* Beobachtungen benutzt werden?

AUFGABE A.14 Wie können zensierte Beobachtungen simuliert werden? Machen Sie einige Vorschläge und diskutieren Sie das Konzept *unabhängiger* Zensuren.

## Literatur

- Andersen, P. K., Borgan, Ø, Gill, R. D., Keiding, N. 1993. *Statistical Models Based on Counting Processes*. Springer: Berlin.
- Blossfeld, H.-P., Rohwer, G. 2002<sup>2</sup>. *Techniques of Event History Modeling. New Approaches to Causal Analysis*, Mahwah, NJ: Lawrence Erlbaum.
- Cox, D. R., Oakes, D. 1984. *Analysis of Survival Data*. London: Chapman and Hall.
- Galton, A. 1984. *The Logic of Aspect*. Oxford: Clarendon.
- International Statistical Institute 1986. Declaration on Professional Ethics. *International Statistical Review* 54, 227–242.
- Kalbfleisch, J. D., Prentice, R. L. 2002<sup>2</sup>. *The Statistical Analysis of Failure Time Data*. Wiley: New York.
- Kaplan, E. L., Meier, P. 1958. Nonparametric Estimation from Incomplete Observations. *Journal of the American Statistical Association* 53, 457–481.
- Lawless, J. F. 1982. *Statistical Models and Methods for Lifetime Data*. New York: Wiley.
- Martinussen, T., Scheike, T. H. 2006. *Dynamic Regression Models for Survival Data*. Springer: Berlin.
- Tableman, M., Jong Sung Kim 2004. *Survival Analysis using S. Analysis of Time-to-Event Data*. Chapman & Hall: London.
- Therneau, T. M., Grambsch, P. M. 2000. *Modeling Survival Data. Extending the Cox Model*. Springer: Berlin.