

### Statistische Modellierung von Verlaufsdaten in der Analyse sozialer Probleme: T. 2: Datenauswertung

Ludwig-Mayerhofer, Wolfgang

Veröffentlichungsversion / Published Version

Zeitschriftenartikel / journal article

Zur Verfügung gestellt in Kooperation mit / provided in cooperation with:

Centaurus-Verlag

#### Empfohlene Zitierung / Suggested Citation:

Ludwig-Mayerhofer, W. (1994). Statistische Modellierung von Verlaufsdaten in der Analyse sozialer Probleme: T. 2: Datenauswertung. *Soziale Probleme*, 5(1/2), 229-263. <https://nbn-resolving.org/urn:nbn:de:0168-ssoar-247330>

#### Nutzungsbedingungen:

Dieser Text wird unter einer Deposit-Lizenz (Keine Weiterverbreitung - keine Bearbeitung) zur Verfügung gestellt. Gewährt wird ein nicht exklusives, nicht übertragbares, persönliches und beschränktes Recht auf Nutzung dieses Dokuments. Dieses Dokument ist ausschließlich für den persönlichen, nicht-kommerziellen Gebrauch bestimmt. Auf sämtlichen Kopien dieses Dokuments müssen alle Urheberrechtshinweise und sonstigen Hinweise auf gesetzlichen Schutz beibehalten werden. Sie dürfen dieses Dokument nicht in irgendeiner Weise abändern, noch dürfen Sie dieses Dokument für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, aufführen, vertreiben oder anderweitig nutzen.

Mit der Verwendung dieses Dokuments erkennen Sie die Nutzungsbedingungen an.

#### Terms of use:

This document is made available under Deposit Licence (No Redistribution - no modifications). We grant a non-exclusive, non-transferable, individual and limited right to using this document. This document is solely intended for your personal, non-commercial use. All of the copies of this documents must retain all copyright information and other information regarding legal protection. You are not allowed to alter this document in any way, to copy it for public or commercial purposes, to exhibit the document in public, to perform, distribute or otherwise use the document in public.

By using this particular document, you accept the above-stated conditions of use.

# Statistische Modellierung von Verlaufsdaten in der Analyse sozialer Probleme, Teil II: Datenauswertung

von Wolfgang Ludwig-Mayerhofer

## **Abstract**

*This paper resumes the discussion from an introductory article on event history analysis («survival analysis», «analysis of failure times») which appeared in the first part of this volume. It describes various approaches to the analysis of hazard functions and the assessment of the influence of covariates. After introductory remarks on non-parametric and semi-parametric estimation, an extensive discussion covers parametric models which may take into account variation of the hazard function over time. In addition, models for discrete time hazard functions, time-dependent covariates, competing risks, and repeated events are treated. All models are illustrated by an example from the German Socio-Economic Panel (SOEP).*

## **Zusammenfassung**

*Die Arbeit baut auf dem einführenden Artikel zur Verlaufsdatenanalyse auf, der im vorangegangenen Heft dieser Zeitschrift erschienen ist. Sie stellt verschiedene Möglichkeiten vor, Hazardfunktionen und die Einflüsse von Kovariaten auf diese zu analysieren. Nach non-parametrischen und semi-parametrischen Analyseverfahren werden ausführlich Modelle der parametrischen Analyse diskutiert, welche die Veränderlichkeit der Hazardfunktion in der Zeit berücksichtigen können. Ferner werden Modelle für diskrete Verweildauern, zeitveränderliche Kovariaten, mehrere Zielzustände und wiederholte Ereignisse erörtert. Alle Modelle werden anhand eines Beispiels aus dem Sozio-ökonomischen Panel (SOEP) erläutert.*

## **1. Übersicht zu statistischen Verfahren der Verlaufsdatenanalyse**

In diesem Teil der Arbeit wird dargestellt, welche verschiedenen Möglichkeiten zur statistischen Analyse von Verlaufsdaten, insbesondere zur Schätzung von Hazardfunktionen oder -raten in Abhängigkeit von Kovariaten, zur Verfügung stehen. Wie schon in Teil I verweise ich häufig auf weiterführende Stellen in den drei deutschsprachigen Lehrbüchern mit den Kürzeln *A* für Andreß (1992a), *BHM* für Blossfeld/Hamerle/Mayer (1986) und *DM* für Diekmann/Mitter (1984).

Die statistische Analyse von Verlaufsdaten steht grundsätzlich vor den gleichen Problemen wie jede andere Datenauswertung. Da in aller Regel keine Vollerhebung, sondern eine Stichprobe vorliegt, möchte man wissen, in welchem Bereich die »wahren«, also die in der untersuchten Population gültigen Parameter mit hinlänglicher Wahrscheinlichkeit liegen werden. Ferner wird man in sozialwissenschaftlichen Untersuchungen zumeist die Frage beantworten wollen, ob Einflüsse von erklärenden Variablen auf den untersuchten Prozeß, wie sie sich in den explo-

rativen Auswertungen in Teil I am Beispiel der Arbeitslosigkeitsdauern aus dem Sozio-ökonomischen Panel hinsichtlich des Alters gezeigt haben, inferenzstatistisch absicherbar sind. Auf dieses Beispiel greife ich auch in diesem Teil zurück.

Im folgenden soll nur kurz auf einfache, nicht-parametrische Verfahren der Analyse eingegangen werden; etwas ausführlicher sollen semi-parametrische und noch umfassender parametrische Verfahren dargestellt werden. In den *nicht-parametrischen Verfahren* geht es – wie bei der in Teil I ausführlich dargestellten Life-Table-Schätzung – um die möglichst »datennahe« Schätzung der relevanten Funktionen und um einfache Prüfungen von Unterschieden zwischen verschiedenen Gruppen. Die Einfachheit ist gleichzeitig Vorzug, aber auch Nachteil dieser Verfahren. Da sie mit sehr wenigen Annahmen verbunden sind, wird vermieden, dem Datenmaterial ein Modell zu oktroyieren, das ihm möglicherweise nicht angemessen ist. Andererseits sind sie für eine multivariate Analyse nicht geeignet.

*Semi-parametrische und parametrische Verfahren* versuchen, den Verlauf des untersuchten Prozesses durch einen oder mehrere Parameter wiederzugeben. Hier werden die Survivorfunktion oder andere Funktionen nicht mehr explizit für jeden einzelnen Zeitpunkt (bzw. jedes Zeitintervall) angegeben, sondern durch einige wenige Kennzahlen, eben Parameter, charakterisiert. Gleichzeitig ist auf diese Weise – anders als mit den nicht-parametrischen Verfahren – eine simultane Berücksichtigung einer Vielzahl möglicher Einflüsse, also eine multivariate Analyse des untersuchten Prozesses möglich. Im allgemeinen wird der Einfluß der relevanten erklärenden Variablen auf die Hazardfunktion untersucht, da diese, wie wir gesehen haben, als den übrigen Funktionen zugrundeliegend aufgefaßt werden kann und im übrigen besser zwischen verschiedenen Verläufen unterscheidet, weil aus sehr unterschiedlichen Hazardraten zumindest oberflächlich recht ähnliche Survivorfunktionen resultieren können.

Die grundsätzliche Logik der Verfahren ist ganz ähnlich wie in den üblichen multivariaten Analysetechniken, etwa der linearen oder logistischen Regression. Untersucht wird, ob bestimmte Merkmale der Untersuchungseinheiten *ceteris paribus*, also unter statistischer Kontrolle der übrigen Merkmale, einen Einfluß auf den untersuchten Prozeß, hier: auf die Hazardfunktion, haben. Es werden Koeffizienten geschätzt, die angeben, ob bzw. wie die erklärende(n) Variable(n) - im Rahmen der Verlaufsdatenanalyse spricht man üblicherweise von Kovariaten - eine höhere oder niedrigere Hazardfunktion, also einen schnelleren oder langsameren Übergang in den Zielzustand, bewirken. Darüber hinaus lassen sich mit Hilfe der *parametrischen* Verfahren auch Koeffizienten schätzen, die die spezifische Verlaufsform der Hazardrate charakterisieren. Als *semi-parametrisch* wird dagegen ein von Cox (1972) entwickeltes Schätzverfahren bezeichnet, in welchem nur der Einfluß der Kovariaten geprüft wird, die dem Prozeß zugrundeliegende »Basisrate« jedoch unspezifiziert bleibt (weshalb das Verfahren auch nicht vollständig parametrisch, sondern eben nur semi-parametrisch ist).<sup>1</sup> Der Anspruch des letztgenannten Modells ist also bescheidener als der der parametrischen Verfahren. Je nach Forschungsinteresse kann hierin ein Vor- oder ein Nachteil gesehen werden. Geht man

von der Voraussetzung aus, daß Hypothesen über den Verlauf der Basisrate nur dann getestet werden sollen, wenn man dafür begründete theoretische Annahmen hat, so bietet das Cox-Modell eine ausgezeichnete Möglichkeit der multivariaten Analyse, wenn keine solchen Annahmen vorliegen. Umgekehrt läßt sich die Position vertreten, daß man mit dem Cox-Modell gerade die Möglichkeit verschenkt, die Zeitabhängigkeit des Prozesses zu analysieren, daß damit also ein unnötiger Verzicht auf relevante Informationen verbunden ist (Brüderl/Diekmann 1995).

## 2. Nicht-parametrische Verfahren

In Zusammenhang mit den nicht-parametrischen Verfahren der Verlaufsdatenanalyse sind drei Aspekte wichtig: 1. Die Schätzung der einschlägigen Funktionen, 2. die Berechnung von Konfidenzintervallen und 3. die Möglichkeit einfacher Gruppenvergleiche.

1. In Teil I, Abschnitt 3.2, wurde bereits ausführlich *eine* Möglichkeit einer einfachen Schätzung der grundlegenden Funktionen, also der Survivorfunktion  $S(t)$ , der Dichtefunktion  $f(t)$  und der Hazardfunktion  $r(t)$ , erläutert: die Life-Table- oder Sterbetafel-Methode. Ein zweites Verfahren, der sog. *Kaplan-Meier-* oder *Product-Limit-Schätzer* für  $S(t)$ , basiert auf ganz ähnlichen Überlegungen, allerdings wird hier davon ausgegangen, daß die Verweildauern exakt, also nicht gruppiert gemessen wurden. Wegen der Annahme exakter Messungen wird  $S(t)$  für jeden Zeitpunkt berechnet, zu dem ein oder mehrere Ereignisse eingetreten sind. Angesichts dieser »punktuellen« Betrachtungsweise lassen sich  $f(t)$  und  $r(t)$  nicht unmittelbar schätzen, jedoch »Hazardkomponenten« für die Sprungstellen von  $S(t)$ , aus denen sich auch eine kumulierte Hazardrate schätzen läßt.<sup>2</sup> Der wesentliche Unterschied zum Life-Table-Schätzer ist darin zu sehen, daß wegen der Annahme exakter Verweildauern eine Korrektur der »Risikomenge« bei Ereignissen nicht erfolgt. Es wird vielmehr angenommen, daß die Zensierungen und dementsprechend die Verringerung der Risikomenge jeweils zwischen Ereignissen fallen. Werden doch Zensierungen und Ereignisse zum gleichen Zeitpunkt beobachtet, werden die Zensierungen so behandelt, als wären sie nach den Ereignissen aufgetreten. Der Kaplan-Meier-Schätzer führt daher im vorliegenden Beispiel im Detail zu etwas anderen Ergebnissen als der Life-Table-Schätzer, die wesentlichen Schlußfolgerungen hinsichtlich des Arbeitslosigkeitsverlaufs und der Einflüsse des Alters sind aber in unserem Beispiel bei beiden Verfahren - wie auch sonst in aller Regel - identisch. Ausführlichere Darstellungen finden sich in den Lehrbüchern (A: 147 ff.; BHM: 44 ff., 124 ff.; DM: 76 ff.).

2. Im Rahmen der beiden genannten Schätzverfahren lassen sich jeweils Standardfehler und hieraus Vertrauensintervalle für  $S(t)$  bzw. - nur im Life-Table-Schätzer - für  $f(t)$  und  $r(t)$  berechnen. Hierzu sei wiederum auf die Lehrbuchliteratur verwiesen (A: 156 ff.; BHM: 45; DM: 66, 78).

3. Wichtiger als die Berechnung von Konfidenzintervallen für einzelne Funktionen sind in sozialwissenschaftlichen Anwendungen die Vergleiche zwischen Grup-

pen, wie etwa in unserem Beispiel zwischen den Altersgruppen.<sup>3</sup> Hierfür sind verschiedene nicht-parametrische Teststatistiken entwickelt worden. Die wichtigsten unter diesen sind zwei »klassische« Testverfahren, der *Log-Rank-Test*, auch als Mantel-Cox-Test oder Verallgemeinerter Savage-Test bezeichnet, sowie eine Teststatistik nach Gehan und Breslow, auch als *Verallgemeinerter Wilcoxon-Test* bezeichnet (vgl. A: 159; BHM: 48, Anwendung S. 128 ff.; DM: 86 ff.).<sup>4</sup> Der erstgenannte Test kann tendenziell eher Unterschiede am rechten Ende der Survivorfunktion entdecken, der zweite reagiert eher auf Unterschiede zu Beginn der Survivorfunktion, so daß in der Praxis unterschiedliche Entscheidungen über die Signifikanz von Unterschieden zustandekommen können. Zwei neuere Statistiken, die ebenfalls in einigen Programmen implementiert sind, wurden von Tarone/Ware (1977) und Prentice (1978; siehe auch Prentice/Marek 1979) entwickelt und liefern im allgemeinen Werte zwischen denjenigen des Log-Rank-Tests und der Gehan/Breslow-Statistik.

Abschließend ist festzuhalten, daß die Anwendungsmöglichkeiten nicht-parametrischer Verfahren für viele sozialwissenschaftliche Fragestellungen sicher begrenzt sind. Die hier angesprochenen Teststatistiken sind eher im Rahmen von kontrollierten Experimenten von Bedeutung, wo durch die randomisierte Zuteilung zu Kontroll- und Experimentalgruppen weitere Einflüsse ausgeschaltet werden können. Bei Untersuchungen, die nicht am Experimental-Paradigma ausgerichtet sind, wird dagegen häufig eine Vielzahl von potentiellen Einflüssen erfaßt, die simultan nur mit multivariaten Verfahren analysiert werden können. Es ist in einer solchen Situation sicherlich auch nicht sinnvoll, zunächst eine Vielzahl nicht-parametrischer Signifikanztests zu berechnen und im Anschluß nur die in diesem Schritt signifikanten Einflüsse mit multivariaten Modellen zu prüfen, da wegen Suppressor-Effekten möglicherweise bedeutsame Einflüsse so nicht erkannt werden.

Dagegen sind einfache Berechnungen bzw. graphische Darstellungen von  $S(t)$ ,  $r(t)$  oder  $\log S(t)$ , wie sie in Teil I dargestellt wurden, wichtig für eine *explorative Datenanalyse*. Sie erlauben eine optische Inspektion der Daten, durch die eventuell vorhandene Datenfehler erkannt werden können. Ferner führen sie zu ersten Aufschlüssen über die Form der Survivor- bzw. Hazardfunktion, also über eine mögliche Zeitabhängigkeit des untersuchten Prozesses. Metrische Variablen können in mehrere Gruppen zerlegt werden, so daß - vorbehaltlich einer genaueren Prüfung mit multivariaten Verfahren - u.U. nicht-lineare Einflüsse erkannt werden können, wie sie in unserem Beispiel wohl hinsichtlich des Alters vorliegen.

Ob man ein solches exploratives Vorgehen für sinnvoll hält, hängt davon ab, welche Auffassung man von der Datenanalyse hat. Wer diese ausschließlich für den Test von Hypothesen als zulässig erachtet, muß ein exploratives Vorgehen als fragwürdig beurteilen. Nun besteht gewiß die Gefahr, daß man nach einer ausführlichen Datenexploration nur mehr die Zusammenhänge »bestätigt« findet, die man beim Screenen der Daten erst entdeckt hat. Andererseits sollte man gegenüber einer Praxis empirischer Sozialforschung, die abstrakte Modelle auf Daten anwendet,

ohne sich zu vergewissern, ob jene diesen in irgendeiner Weise angemessen sind, ebenfalls skeptisch sein. Daher sind explorative Analysen auch nicht nur eine Vorstufe der Datenauswertung, sondern können sehr wichtig sein, um etwa auffälligen oder unerwarteten Ergebnissen (und Nicht-Ergebnissen) auf die Spur zu kommen. »Hypothesentestende« und »explorative« Datenanalyse, das hat soeben Schnell (1994, vor allem Kap. 11) einmal mehr verdeutlicht, lassen sich ohnehin kaum sinnvoll voneinander abgrenzen; auch wer »explorativ« seine Daten durchforstet, hat dabei im allgemeinen bestimmte Ideen im Kopf, verfolgt also Hypothesen. Sogar der in Lehrbüchern als warnendes Beispiel dargestellte *Idealtypus* des Forschers, der »Alles mit Allem« korreliert, wird als *Realtypus* nur selten völlig willkürlich handeln, da ja schon die Auswahl der erhobenen Variablen nicht zufällig, sondern nach wenigstens impliziten Hypothesen erfolgte.

### 3. Ein semi-parametrischer Ansatz: Das »Partial-Likelihood«-Verfahren

In dem semi-parametrischen Modell von Cox wird folgende Gleichung zur Charakterisierung der Hazardfunktion geschätzt:<sup>5</sup>

$$r(t; \mathbf{X}) = r_0(t) \exp(\mathbf{X}\boldsymbol{\beta}) \quad (1)$$

$r(t; \mathbf{X})$  ist die Hazardrate für ein Individuum mit gegebenem Kovariaten-Vektor  $\mathbf{X}$ . Diese Hazardrate ergibt sich also aus einer Basisrate  $r_0(t)$  und den mit einem (Spalten-)Vektor  $\boldsymbol{\beta}$  gewichteten individuellen Ausprägungen der Kovariaten. Die Größe  $\mathbf{X}\boldsymbol{\beta}$  wird mit der Basisrate exponentiell verknüpft. Wie oben erwähnt, wird die Basisrate nicht explizit geschätzt. Die exponentielle Verknüpfung mit  $\mathbf{X}\boldsymbol{\beta}$  ist rein technisch motiviert; sie stellt sicher, daß keine negativen Werte für die Hazardfunktion geschätzt werden können. Die Parameter des Modells können aber leicht interpretiert werden, indem ihr Antilogarithmus  $\alpha = \exp(\boldsymbol{\beta})$  gebildet wird. Diese  $\alpha$ -Parameter geben an, um welchen Faktor das Risiko, in den Zielzustand zu wechseln, erhöht ( $\alpha > 1$ ) oder verringert ( $\alpha < 1$ ) wird (BHM: 147; DM: 98 f. und 128 f.). Das Modell geht übrigens von der Annahme kontinuierlich gemessener Verweildauern aus, es gibt aber auch Vorschläge zu seiner Modifikation für diskrete bzw. gruppierte Dauern (siehe Abschnitt 5).

Aus Gleichung (1) und der soeben erläuterten Interpretation der Parameter ergibt sich, daß das Cox-Modell davon ausgeht, daß die Hazardfunktionen für verschiedene Werte von Kovariaten jeweils *proportional* zueinander sind. Angenommen, wir würden für eine Kovariate mit drei Ausprägungen 0, 1, und 2 einen Koeffi

zienten  $\beta$  von 0,405 schätzen. Das Modell nimmt an, daß die Hazardfunktion für die Individuen mit der Ausprägung 1 in der betreffenden Variablen *stets* das  $\alpha = \exp(0,405) = 1,5$ -fache der Hazardfunktion der Individuen mit der Ausprägung 0 beträgt, für Individuen mit der Ausprägung 2 *stets* das  $\alpha = \exp(2 \times 0,405) = 2,25$ -

fache. Man spricht daher häufig im Zusammenhang mit dem Cox-Modell vom Proportional-Hazards-Modell. Tatsächlich gehen aber auch einige der später geschilderten parametrischen Modelle von proportionalen Hazardfunktionen aus. Die Annahme der proportionalen Hazardfunktionen kann explorativ mit einfachen Verfahren (siehe z.B. Teachman 1983), aber auch auf komplexere Weise im Rahmen der Modell-Schätzung geprüft werden (A: 254 ff.; BHM: 143 ff.). Allerdings dürfte das Cox-Modell relativ robust gegenüber nicht allzu starken Abweichungen von der Proportional-Hazards-Annahme sein. Liegen stärkere Abweichungen vor, ist eine Analyse im Rahmen des Cox-Modells unter Umständen trotzdem möglich, indem eine sog. stratifizierte Schätzung vorgenommen wird. Konkret: Ist für eine Variable die Annahme der proportionalen Hazards verletzt, so kann die gesamte Stichprobe in Gruppen oder »Schichten« (engl. strata) zerlegt werden, die den einzelnen Ausprägungen dieser Variablen entsprechen. Ist nun innerhalb dieser Gruppen die Annahme proportionaler Hazards gültig, kann ein Modell geschätzt werden, in dem für die verschiedenen Gruppen (Schichten) eine unterschiedliche Basisrate angenommen wird. Ausführliche Beispiele finden sich bei A: 254 ff. und BHM: 139 ff.

Die Modellschätzung erfolgt über ein Verfahren, welches als Partial-Likelihood-(PL-)Schätzung bezeichnet wird (grob gesagt deshalb, weil wegen der nicht explizit geschätzten Basisrate nicht die gesamte Information der Daten für die Schätzung benutzt wird). Trotzdem ist die grundsätzliche »Logik« des Schätzverfahrens nicht anders als in der Maximum-Likelihood-(ML-)Schätzung, wie sie den nachfolgend geschilderten parametrischen Modellen zugrundeliegt. Daher sei hier ganz kurz auf die wesentlichen inferenzstatistischen Gesichtspunkte eingegangen, also auf die Frage der statistischen Signifikanz der geschätzten Parameter.

Die Modellschätzung mit PL oder ML beruht auf iterativen Verfahren, in welchen der Wert einer Likelihood-Funktion maximiert wird, die angibt, wie wahrscheinlich man die gegebenen Daten in der Stichprobe erhalten würde, wenn die geschätzten Parameter den »wahren« Parametern entsprächen. Die Parameter mit der größten Wahrscheinlichkeit werden als die besten Schätzungen der wahren Parameter aufgefaßt. Die mathematische Logik des Verfahrens braucht uns hier nicht zu interessieren,<sup>6</sup> wohl aber die Frage, wie man aus den Resultaten Aufschluß über die inferenzstatistische Absicherung des gesamten Modells bzw. einzelner Effekte erhält.

Zunächst werden in allen Verfahren Standardfehler (engl. standard error, meist abgekürzt als S.E.) für die einzelnen Koeffizienten berechnet, welche bei hinreichend großen Stichproben als normalverteilt gelten können.<sup>7</sup> Daher kann - grob gesagt - ein Koeffizient als auf dem 5-Prozent-Niveau signifikant von Null verschieden angesehen werden, wenn das Verhältnis des Koeffizienten zu seinem Standardfehler (in den Programmen häufig auch als T-Statistik bezeichnet) mindestens 1,96 beträgt. Die entsprechenden Signifikanzniveaus werden aber von den meisten Programmen explizit berechnet.<sup>8</sup>

Noch zuverlässiger (und gleichzeitig universeller einsetzbar) sind jedoch Tests, die auf der Likelihood-Funktion aufbauen bzw. auf deren Logarithmus, im folgenden als Log-Likelihood bezeichnet. Grundsätzlich geht es immer darum, ein gegebenes Modell - nennen wir es  $M_1$  - bzw. dessen Log-Likelihood  $LL_1$  mit einem anderen - Modell  $M_0$  mit Log-Likelihood  $LL_0$  - zu vergleichen, aus welchem im Vergleich zu Modell  $M_1$  einer oder mehrere Parameter weggelassen wurden. Die Größe

$$2(LL_1 - LL_0) \quad (2)$$

ist  $\chi^2$ -verteilt mit  $s$  Freiheitsgraden, wobei  $s$  der Zahl der weggelassenen Parameter entspricht. Dieser Likelihood-Verhältnis- oder Likelihood-Ratio-Test<sup>9</sup> (abgekürzt: LR-Test) kann also eingesetzt werden, um

- zu testen, ob ein einzelner Parameter signifikant von Null verschieden ist, d.h. Modell  $M_0$  wäre in diesem Fall ein Modell, aus dem im Vergleich zu  $M_1$  eine einzelne Variable weggelassen wurde;
- zu testen, ob die Gesamtheit aller Variablen zusammengenommen eine signifikante Erklärungskraft besitzt, d.h. Modell  $M_0$  wäre hier ein sog. Null-Modell, welches nur eine Modellkonstante enthält,<sup>10</sup> oder
- um den Einfluß einer beliebigen Zahl von Variablen zu testen, was z.B. geboten sein kann, wenn eine Variable in mehrere Dummy-Variablen zerlegt wurde und der Einfluß all dieser Dummy-Variablen zusammen geprüft werden soll (A: 205; BHM: 89; DM: 106 u. 141).

Weitere häufig verwendete Teststatistiken, die asymptotisch - d.h. mit zunehmendem Stichprobenumfang - mit dem LR-Test identisch sind, sind die Wald- und die Score-Teststatistik (vgl. A: 203 ff.; BHM: 89; DM: 106). Besonders hinzuweisen ist auf die Möglichkeit im Programm TDA (Rohwer 1994; vgl. Anhang), den Modell-Parametern beliebige Restriktionen aufzuerlegen. So kann z.B. leicht getestet werden, ob zwei oder mehrere Parameter sich signifikant voneinander unterscheiden.<sup>11</sup>

Nun aber zu den Ergebnissen des Cox-Modells für unseren konkreten Fall.<sup>12</sup> In *Darstellung 1* werden, wie im folgenden für weitere Beispiele, die Log-Likelihood für das geprüfte Modell mit den zwei (Dummy-)Variablen für die Altersgruppen »31 bis 50« und »über 50« sowie die  $\beta$ -Koeffizienten für diese beiden Variablen, deren Standardfehler, die sich hieraus ergebende T-Statistik und das (zweiseitige) Signifikanzniveau angegeben.<sup>13</sup> Zunächst zur inhaltlichen Interpretation: Da im Cox-Modell die Basisrate nicht geschätzt wird, wird keine Regressionskonstante berechnet; das Modell sagt also nichts über die Hazardrate der jüngsten Altersgruppe, sondern nur darüber, wie die Raten der beiden aufgeführten Altersgruppen sich von der der jüngsten Gruppe unterscheiden. Der Koeffizient von -0,0494 für die 31- bis 50jährigen ergibt in den Exponenten erhoben einen Wert von etwa 0,95, d.h., die Rate dieser Gruppe beträgt im Schnitt das 0,95fache der Rate der Vergleichsgruppe, liegt also nur geringfügig unter dieser. Dagegen beträgt die Rate der



ältesten Gruppe das  $\exp(-1,8977) = 0,15$ fache der jüngsten Gruppe, also nicht einmal ein Sechstel. Wie auch nach den explorativen Ergebnissen zu erwarten, ist der Unterschied zwischen der jüngsten Gruppe als Bezugsgruppe und den 31- bis 50jährigen nicht signifikant, dagegen derjenige zwischen der höchsten und der jüngsten Altersgruppe höchst signifikant.

**Darstellung 1:** *Ergebnisse eines Cox-Modells mit zwei Kovariaten  
(Erweiterte Beispieldaten aus Teil I, Darstellung 3)*

Variable	Coeff	Error	T-Stat	Signif
Alter 31 bis 50	-0.0494	0.1033	-0.4778	0.3672
Alter über 50	-1.8977	0.2404	-7.8928	1.0000

Log-likelihood: -2361.52

Die Log-Likelihood für das Null-Modell ohne Kovariaten beträgt -2417,26, so daß der LR-Test nach Formel (2) für das Gesamtmodell einen Wert von  $2(-2361,52 - (-2417,26)) = 111,49$  hat, der mit 2 Freiheitsgraden hoch signifikant ist.<sup>14</sup> Im Rahmen des Cox-Modells wird ferner häufig eine »Globale Chi-Quadrat-Statistik« berechnet, die einem Score-Test entspricht (vgl. dazu BHM: 89 und zur Anwendung S. 145 ff.). Deren Wert beträgt 82,70 und zeigt ebenfalls, daß das Gesamtmodell sich signifikant von einem Modell ohne Kovariaten unterscheidet. Will man z.B. den LR-Test auf die Variable »Alter 31 bis 50« anwenden und prüfen, ob deren Effekt signifikant von Null verschieden ist, so ergibt sich für das Modell, aus dem der Effekt dieser Variablen weggelassen (und damit auf Null gesetzt wurde), eine Log-Likelihood von 2361,63. Der LR-Test, ob dieses Modell signifikant von dem Ausgangsmodell verschieden ist, ergibt  $2(-2361,52 - (-2361,63)) = 0,22$ , d.h., das Modell, welches keinen Unterschied zwischen den beiden jüngeren Altersgruppen postuliert, ist nicht signifikant schlechter und könnte somit als einfacheres Modell dem Ausgangsmodell vorgezogen werden.<sup>15</sup>

#### 4. Parametrische Modelle (kontinuierliche Zeit)

Die parametrischen Modelle für kontinuierliche Verweildauern, die im folgenden vorgestellt werden, haben den Vorzug, nicht nur den Einfluß von Kovariaten - relativ zu einer unspezifizierten »Basisrate« - zu prüfen, sondern auch die zugrundeliegende Hazardfunktion selbst zu modellieren. Im folgenden werden exemplarisch die wichtigsten Modelle dargestellt, wobei weiterhin die schon verwendeten einfachen Beispieldaten herangezogen werden. Im Vordergrund stehen die Möglichkeiten, unterschiedliche Verläufe der Hazardrate zu modellieren.

Zunächst ist hier eine Warnung auszusprechen. Definitive Aussagen über den zeitlichen Verlauf sind immer mit Fragezeichen zu versehen, da sie auch auf nicht bzw. nicht ausreichend berücksichtigte Einflüsse (sog. »unbeobachtete Heterogenität«) zurückgehen können. Konkret: Wenn sich in einem Datensatz zwei (oder

mehr) Gruppen befinden, von denen die eine eine hohe und die andere eine niedrige Hazardrate aufweist, so ereignen sich bei den Personen aus der ersten Gruppe die Übergänge in den Zielzustand natürlich wesentlich schneller, so daß am Schluß vor allem die Personen aus der anderen Gruppe mit sehr langen Verweildauern verbleiben. Würde man nicht zwischen diesen Gruppen unterscheiden (also z.B. in der Arbeitslosenstichprobe nicht nach Arbeitslosen bis 50 und solchen über 50 Jahren differenzieren), hätte man den Eindruck, daß die Hazardrate insgesamt immer mehr abnehmen würde, auch wenn sie in Wirklichkeit in beiden Gruppen möglicherweise konstant ist. Auf diesen - schon lange bekannten - Sachverhalt hat kürzlich wieder Klein (1992) hingewiesen (vgl. außerdem Andreß 1988; Arminger 1984). Auf der anderen Seite sollte man sich hierdurch nicht entmutigen lassen, nach Zeitabhängigkeiten zu suchen - denn schließlich gilt die Feststellung, daß unbeobachtete Einflüsse die Resultate statistischer Modellbildung in Frage stellen, ganz grundsätzlich (Lieberson 1985). Daß man keine definitive Aussage machen kann, welches Modell den Daten angemessen ist, sollte sich eigentlich von selbst verstehen; da aber, wie wir sehen werden, die verschiedenen Modelle den Daten eine bestimmte Struktur »aufzwingen«, sollte man in jedem Fall versuchen, auch alternative Modelle zu testen. Gegebenenfalls muß es weiteren Untersuchungen überlassen bleiben, zwischen möglicherweise divergierenden Schlußfolgerungen zu entscheiden.

Im übrigen wird häufig angenommen, daß abgesehen von der Frage der Zeitabhängigkeit die Wahl des konkreten Modells für die Analyse von geringer Bedeutung ist, da es hinsichtlich der zumeist wichtigeren Informationen über die Einflüsse der Kovariaten keine großen Divergenzen zwischen verschiedenen Modellen gibt. Das zeigen empirische Beispiele (Diekmann/Klein 1991; Ziegler/Brüderl/Diekmann 1988) ebenso wie theoretische Überlegungen (Galler/Pötter 1992). Allerdings werden wir gerade in unserem Beispiel sehen, daß das nicht in allen Fällen zutrifft und ein - allerdings wahrscheinlich unangemessenes - Modell auch zu inhaltlich anderen Schlußfolgerungen führen kann!

Zur Überprüfung, ob unbeobachtete Heterogenität vorliegt, wurde ein Verfahren entwickelt, welches unter der Annahme einer bestimmten Verteilung (der Gamma-Verteilung) für die Varianz der Fehler des Modells - also der nicht erklärten Anteile - prüfen kann, wie groß die unbeobachtete Heterogenität ist (A: 266 ff.; BHM: 251 ff.). Allerdings wird man realistischerweise davon ausgehen müssen, daß mit

diesem Verfahren fast immer ein relevanter Anteil an unbeobachteter Heterogenität entdeckt wird (vgl. auch Galler/Pötter 1987; Petersen 1993).<sup>16</sup>

Wir beginnen nunmehr mit dem einfachsten parametrischen Modell, dem sog. *Exponentialmodell*. Hier besteht zwischen den Kovariaten  $\mathbf{X}$  und der Rate  $r(t; \mathbf{X})$  folgende einfache Beziehung:

$$r(t; \mathbf{X}) = \exp(\beta_0 + \mathbf{X}\beta) \quad (4)$$

**Darstellung 2:** *Ergebnisse verschiedener parametrischer Exponentialmodelle (Erweiterte Beispieldaten aus Teil I, Darstellung 3)*

Variable	Coeff	Error	T-Stat	Signif
<i>Einfaches Exponentialmodell (Log-likelihood: -1303.72)</i>				
Konstante	-1.9351	0.0639	-30.2885	1.0000
Alter 31 bis 50	-0.2175	0.1028	-2.1153	0.9656
Alter über 50	-2.2821	0.2381	-9.5829	1.0000
<i>Piecewise Constant (Log-likelihood: -1244.53)</i>				
Konstante für Monat 1-2	-1.8447	0.0914	-20.1773	1.0000
Konstante für Monat 3	-1.2781	0.1143	-11.1787	1.0000
Konstante für Monat 4-6	-1.7875	0.1087	-16.4466	1.0000
Konstante für Monat 7-12	-2.1761	0.1372	-15.8641	1.0000
Konstante für Monat 13-24	-3.0400	0.2305	-13.1912	1.0000
Konstante für Monat 25 u. höher	-3.6244	0.3395	-10.6756	1.0000
Alter 31 bis 50	-0.0812	0.1034	-0.7853	0.5677
Alter über 50	-1.9874	0.2402	-8.2755	1.0000
<i>Polynom 1. Grades (Log-likelihood: -1254.95)</i>				
Konstante	-1.5832	0.0725	-21.8417	1.0000
Alter 31 bis 50	-0.0686	0.1033	-0.6643	0.4935
Alter über 50	-1.9894	0.2399	-8.2925	1.0000
$\beta_{t1}$	-0.0671	0.0087	-7.7399	1.0000
<i>Polynom 4. Grades (Log-likelihood: -1247.23)</i>				
Konstante	-1.8783106	0.0912149	-20.5921450	1.0000
Alter 31 bis 50	-0.0712948	0.1032454	-0.6905368	0.5101
Alter über 50	-1.9932917	0.2402498	-8.2967483	1.0000
$\beta_{t1}$	0.1560741	0.0305186	5.1140619	1.0000
$\beta_{t2}$	-0.0303144	0.0043719	-6.9338968	1.0000
$\beta_{t3}$	0.0012134	0.0002318	5.2344372	1.0000
$\beta_{t4}$	-0.0000145	0.0000036	-3.9948453	0.9999

Log-likelihood des *Exponentialmodelles nur mit Konstante*: -1393.38

Wie im Cox-Modell wird die Rate hier (wie auch in den folgenden Modellen) exponentiell mit dem durch  $\mathbf{B}$  gewichteten Kovariatenvektor verknüpft, um eine negative Schätzung der Rate zu vermeiden.<sup>17</sup> Im Gegensatz zu jenem Modell enthält der Koeffizientenvektor jetzt aber (wie ebenfalls in allen folgenden Modellen) eine Modellkonstante - hier als  $\beta_0$  bezeichnet -, die die »Basisrate« angibt, also die Rate für diejenigen Fälle, die in allen Kovariaten den Wert 0 aufweisen (in unserem einfachen Beispiel die Altersgruppe bis 30 Jahre). Damit lassen sich nicht nur Aussagen über die *relativen Chancen* der verschiedenen Gruppen treffen, die Arbeitslosigkeit zu verlassen, sondern diese Chancen lassen sich als Hazardrate für die verschiedenen Gruppen explizit »bezziffern«. Wie aus der Modellformulierung ersichtlich wird, wird  $r(t; \mathbf{X})$  als zeitlich konstant aufgefaßt.

Eine Schätzung dieses einfachen Exponentialmodells kommt zu dem etwas überraschenden Ergebnis, daß der Unterschied in den Hazardfunktionen (und damit in den Arbeitslosigkeitsdauern) zwischen der jüngsten und der mittleren Altersgruppe signifikant ist (vgl. *Darstellung 2*)!<sup>18</sup> Im einzelnen lassen sich die Ergebnisse so interpretieren: Für die jüngste Altersgruppe wird eine Hazardrate von  $\exp(-1,9351) = 0,144$  geschätzt, für die mittlere Gruppe eine Rate von  $\exp(-1,9351 + (-0,2175)) = 0,116$  und für die älteste Gruppe eine Rate von  $\exp(-1,9351 + (-0,2812)) = 0,015$ .

Wie schneidet dieses Modell im Vergleich mit einem Null-Modell, also einem Exponentialmodell mit einer Konstanten, aber ohne Kovariaten ab? (Da das einfache Exponentialmodell ohne Kovariaten den allereinfachsten Fall eines parametrischen Modells darstellt, wird es im allgemeinen als Null-Modell für alle parametrischen Modellklassen eingesetzt.) Die Log-Likelihood des letzteren Modells beträgt -1393,39, das Modell mit Kovariaten bringt also durchaus eine entscheidende Verbesserung mit  $2(-1303,72 - (-1393,38)) = 179,32$ .

Im vorliegenden Fall müssen wir allerdings aufgrund der explorativen Ergebnisse damit rechnen, daß das einfache Exponentialmodell, auch wenn es gegenüber dem Null-Modell eine erhöhte Erklärungskraft besitzt, den Arbeitslosigkeitsverlauf nicht adäquat modelliert, da die Hazardrate für die einzelnen Gruppen als konstant angenommen wird. Nach den explorativen Analysen in Teil I dürfte die Annahme der Konstanz jedoch kaum gerechtfertigt sein. Man beachte auch, daß nach dem geschätzten Modell die Rate der Altersgruppe 31 bis 50 Jahre nur etwa das 0,8fache der Basisrate (also der Rate für die Vergleichsgruppe bis 30 Jahre) beträgt; auch dies scheint nach den explorativen Darstellungen des Life-Table-Schätzers kaum gerechtfertigt.

Eine variable Hazardrate kann man durch verschiedene andere Verteilungsannahmen modellieren, auf die wir noch zurückkommen. Allerdings läßt sich auch das einfache Exponentialmodell durch verschiedene Zusatzannahmen zu folgenden zwei Modellen erweitern:<sup>19</sup>

*Piecewise Constant Exponentialmodell:*

$$r(t; \mathbf{X}) = \exp(\beta_{0t} + \mathbf{X}\boldsymbol{\beta}) \quad (5)$$

*Exponentialmodell mit Polynom-Term für die Zeit:*

$$r(t; \mathbf{X}) = \exp(\beta_0 + \mathbf{X}\boldsymbol{\beta} + \beta_{t1}t + \beta_{t2}t^2 \dots + \beta_{tm}t^m) \quad (6)$$

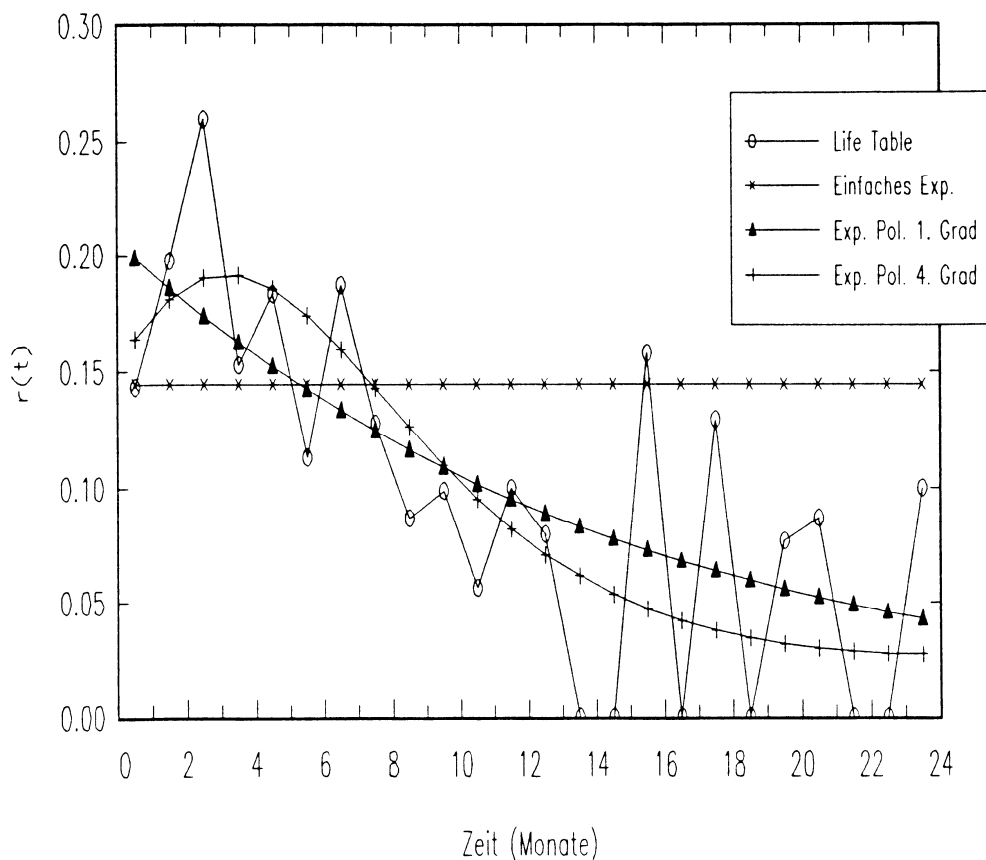
Im *Piecewise Constant Exponentialmodell* (oder periodenspezifischen Exponentialmodell) gibt es nicht nur eine Modellkonstante  $\beta_0$ , sondern mehrere *periodenspezifische* Konstanten  $\beta_{0t}$  (vgl. A: 231 f. und 235 ff.). Die Perioden beziehen sich auf die Prozeßzeit und können vom Benutzer selbst beliebig spezifiziert werden. Das legt natürlich die Gefahr des »curve fitting« nahe, erlaubt andererseits eine recht flexible Modellstruktur. Für die Beispieldaten wurde ein Modell getestet mit der Annahme von sechs unterschiedlichen Basisraten in den Monaten 1 und 2, 3, 4

bis 6, 7 bis 12, 13 bis 24 sowie über 24 Monate. Offensichtlich entspricht dieses Modell dem beobachtbaren Verlauf der Hazardrate deutlich besser, was sich in einer erheblichen Erhöhung der Log-Likelihood niederschlägt (vgl. Darstellung 2).

Im *Exponentialmodell mit Polynom-Term für die Zeit* werden dagegen zusätzliche  $\beta$ -Koeffizienten geschätzt - hier als  $\beta_{tn}$  bezeichnet -, die mit der Zeit in der Form von Polynomen beliebigen Grades verknüpft werden können. Wieder kann bzw. muß vom Benutzer spezifiziert oder ausprobiert werden, welcher Grad des Polynoms zu sinnvollen Ergebnissen führt. Bei den Beispieldaten ergibt sich in einem Modell mit dem Zeitglied 1. Grades eine fallende Rate (vgl. Darstellung 2). Modelle mit einem Glied 2. oder 3. Grades, die grundsätzlich in der Lage sein müßten, einen zuerst steigenden und dann fallenden Verlauf zu modellieren - wie er nach den explorativen Analysen gegeben ist -, haben im vorliegenden Fall weder einen signifikanten Erklärungszuwachs erbracht noch auch den entsprechenden Verlauf modellieren können. Erst ein Polynom 4. Grades führt zur Modellierung einer zunächst steigenden ( $\beta_{t1}$  ist positiv!) und dann fallenden Rate; der Erklärungszuwachs dieses Modells gegenüber dem Polynom 1. Grades in Höhe von  $2(-1247,23 - (-1254,95)) = 15,44$  ist auch bei drei Freiheitsgraden (für die drei zusätzlichen Parameter) nach der  $\chi^2$ -Verteilung signifikant. Trotzdem sollte man nicht annehmen, hiermit »bewiesen« zu haben, daß das Polynom 4. Grades ein besseres Modell impliziert; eine gewisse Plausibilität hierfür wird allerdings durch den höheren Modellfit nahegelegt.

Da dieses Modell in der Lehrbuchliteratur (mit Ausnahme einer kurzen Erwähnung bei A: 233) nicht dargestellt wird, seien noch einmal exemplarisch für einige Zeitpunkte die geschätzten Hazardraten für die Referenzgruppe bis 30 Jahre berechnet. Für den Zeitpunkt 1 ergibt sich - annäherungsweise berechnet - eine Rate von  $\exp(-1,8783 + 0,1561 \times 1 - 0,0303 \times 1^2 + 0,0012 \times 1^3 - 0,0000145 \times 1^4) = 0,1735$ , für den Zeitpunkt 3 erhält man eine etwas höhere Rate von 0,1918, nach 12 Monaten beträgt sie aber nur mehr 0,0746 und nach 24 Monaten schließlich  $\exp(-1,8783 + 0,1561 \times 24 - 0,0303 \times 24^2 + 0,0012 \times 24^3 - 0,0000145 \times 24^4) = 0,0222$ . Bei den anderen beiden Altersgruppen sind die jeweiligen gruppenspezifischen Koeffizienten noch mit in den Ausdruck in der Klammer aufzunehmen.

**Darstellung 3:** Hazardfunktionen für die Altersgruppe bis 30 (Schätzungen nach Life Table, Exponentialmodell, Exponentialmodell mit Polynom 1. und Polynom 4. Grades)



Schließlich ist noch darauf hinzuweisen, daß die Modelle mit Polynom in der Praxis teilweise schwer zu handhaben sind, da im allgemeinen spätestens mit dem Polynom-Glied vierten Grades die iterative Modellschätzung numerisch sehr schwierig wird. Auch im vorliegenden Fall ist nicht ganz sicher, ob der Schätzalgorithmus tatsächlich das globale Maximum erreicht hat. Im folgenden werden einige Modelle vorgestellt, die jedenfalls teilweise leichter zu handhaben sind. Doch zuvor möchte ich einige der von den bislang erörterten Modellen geschätzten Hazardraten graphisch zeigen, und zwar exemplarisch für die jüngste Altersgruppe (*Darstellung 3*). So kann besser verdeutlicht werden, welche unterschiedlichen Implikationen die verschiedenen Modellschätzungen für die Hazardrate haben: die konstante Rate des einfachen Exponentialmodells, die fallende Rate des Exponentialmodells mit einem Polynom-Glied und die erst steigende und dann fallende des Exponentialmodells mit einem vierfachen Polynom.<sup>20</sup> Auch anhand dieser Darstellung wird man zu der Schlußfolgerung kommen, daß das letztgenannte Modell

noch am ehesten geeignet ist, den »datennahen« Life-Table-Schätzer nachzuvollziehen. Angemerkt sei abschließend auch, daß alle bislang besprochenen komplexeren Modelle den Alterseffekt für die mittlere Gruppe im Vergleich zur jüngsten Gruppe als nicht signifikant ausweisen. Das wird auch für die folgenden Modelle gelten, so daß dort nicht mehr eigens darauf hingewiesen wird.

Nunmehr möchte ich noch einige weitere sehr häufig gebrauchte Verteilungen vorstellen, die geeignet sind, nicht konstante Verläufe der Hazardrate zu modellieren (vgl. zu den Ergebnissen *Darstellung 4*).

Das *Weibull-Modell* läßt sich in den äquivalenten Formulierungen

$$r(t; \mathbf{X}) = \exp(\beta_0 + \mathbf{X}\boldsymbol{\beta}) p \left( \exp(\beta_0 + \mathbf{X}\boldsymbol{\beta}) t \right)^{p-1} \quad (7)$$

bzw.

$$r(t; \mathbf{X}) = \exp(\beta_0 + \mathbf{X}\boldsymbol{\beta})^p p t^{p-1} \quad (8)$$

darstellen. Der Parameter  $p^{21}$  ist als sog. »Shape Parameter« u.a. mit der Prozeßzeit  $t$  verknüpft, und somit ist leicht zu sehen, daß  $r(t)$  mit zunehmender Zeit monoton ansteigt, wenn  $p > 1$ , und monoton absinkt, wenn  $p < 1$ . Im Falle von  $p = 1$  erhält man wieder ein einfaches Exponentialmodell. Zu beachten ist, daß in den meisten Fällen - und so auch hier - von den Statistikprogrammen nicht  $p$  ausgegeben wird, sondern  $\ln(p)$ , und das heißt, daß  $\ln(p) < 0$  eine fallende und  $\ln(p) > 0$  eine steigende Hazardrate impliziert. Im vorliegenden Fall wird aufgrund der längerfristig fallenden Rate ein signifikant negativer Wert für  $\ln(p)$  geschätzt.

Ebenfalls eine monoton steigende oder fallende Rate modellieren kann die sog. *Gompertz-Verteilung*, die sich in unserem Kontext folgendermaßen formulieren läßt:

$$r(t; \mathbf{X}) = \exp(\beta_0 + \mathbf{X}\boldsymbol{\beta}) + \exp((\gamma_0 + \mathbf{X}\boldsymbol{\gamma}) t) \quad (9)$$

Hier wird die Zeitabhängigkeit modelliert, indem die Prozeßzeit mit einem weiteren Parameter  $\gamma_0$  und gegebenenfalls zusätzlich mit einem Parametervektor  $\boldsymbol{\gamma}$  und Kovariaten verknüpft wird. Wird nur eine Konstante  $\gamma_0$  geschätzt, wird der zeitlich variable Hazardratenverlauf für alle Untersuchungseinheiten (bzw. alle Kovariatenkonstellationen) als gleich angenommen, ansonsten können durch die verschiedenen mit dem Parametervektor  $\boldsymbol{\gamma}$  verknüpften Kovariaten auch spezifische Verläufe modelliert werden (wobei sich dann u.U. für bestimmte Konstellationen steigende und für andere Konstellationen fallende Verläufe ergeben können) (vgl. A: 228 ff.; BHM: 211 ff.). Bei unseren Beispieldaten ergibt sich (vgl. *Darstellung 4*), daß nur die Konstante  $\gamma_0$  signifikant von Null verschieden ist, nicht aber die  $\boldsymbol{\gamma}$ -Koeffizienten für die beiden Kovariaten, so daß man davon ausgehen kann, daß –

sofern man eine konstant sinkende Rate überhaupt als sinnvolle Modellierung erachtet - der Ratenverlauf für alle Subgruppen gleich ist.

**Darstellung 4:** *Ergebnisse weiterer parametrischer Modelle (erweiterte Beispieldaten aus Teil I, Darstellung 3)*

Variable	Coeff	Error	T-Stat	Signif
<i>Weibull (Log-likelihood: -1291.22)</i>				
Konstante	-1.9107	0.0758	-25.2143	1.0000
Alter 31 bis 50	-0.1755	0.1223	-1.4354	0.8488
Alter über 50	-2.5326	0.2878	-8.8005	1.0000
ln(p)	-0.1702	0.0359	-4.7380	1.0000
<i>Gompertz, Kovariaten in <math>\beta</math>- und g-term (Log-likelihood: -1254.74)</i>				
$\beta$ -Konstante	-1.6039	0.0813	-19.7368	1.0000
$\beta$ -Alter 31 bis 50	-0.0317	0.1313	-0.2417	0.1910
$\beta$ -Alter über 50	-1.8592	0.3417	-5.4406	1.0000
g-Konstante	-0.0619	0.0123	-5.0351	1.0000
g-Alter 31 bis 50	-0.0084	0.0179	-0.4686	0.3606
g-Alter über 50	-0.0209	0.0383	-0.5460	0.4150
<i>Log-logistisch (Typ I) (Log-likelihood: -1248.61)</i>				
Konstante	-1.3210	0.0731	-18.0728	1.0000
Alter 31 bis 50	-0.0876	0.1199	-0.7306	0.5350
Alter über 50	-2.3385	0.2207	-10.5946	1.0000
ln(p)	0.2925	0.0404	7.2349	1.0000
<i>Log-logistisch (Typ II) (Log-likelihood: -1236.48)</i>				
Konstante	-1.3287	0.0707	-18.7821	1.0000
Alter 31 bis 50	-0.0805	0.1032	-0.7800	0.5646
Alter über 50	-1.9901	0.2399	-8.2950	1.0000
ln(l)	-0.7707	0.1541	-5.0010	1.0000
ln(p)	0.5353	0.0761	7.0368	1.0000
<i>Sichel (Log-likelihood: -1278.72)</i>				
Konstante	-1.7516	0.0865	-20.2569	1.0000
Alter 31 bis 50	-0.0914	0.1031	-0.8865	0.6246
Alter über 50	-2.0712	0.2396	-8.6430	1.0000
ln(g)	1.3489	0.0529	25.4761	1.0000

Das Gompertz-Modell impliziert bei negativen  $\gamma$ -Koeffizienten, daß die Hazardrate mit der Zeit gegen Null tendiert. Das *Gompertz-Makeham-Modell* enthält eine zusätzliche Konstante  $\alpha$ , gegen die die Rate - sofern sie fällt - tendieren würde. Auch diese Konstante kann (muß aber nicht) mit Kovariaten verknüpft werden, so daß sich als vollständige Formulierung ergibt:

$$r(t; \mathbf{X}) = \exp(\alpha_0 + \mathbf{X}\boldsymbol{\alpha}) + \exp(\beta_0 + \mathbf{X}\boldsymbol{\beta}) + \exp((\gamma_0 + \mathbf{X}\boldsymbol{\gamma}) t) \quad (10)$$



Dieses Modell ist bereit relativ komplex; es kann zwar sehr verschiedenartige (wenngleich nur monoton steigende oder fallende) Verläufe modellieren, dies wird jedoch mit dem Nachteil einer sehr aufwendigen Modellsuche sowie im übrigen auch mit Problemen bei der Durchführung der Modellschätzung erkauft. Solche Schwierigkeiten deuten dann allerdings im allgemeinen darauf hin, daß das Modell den Daten nicht sehr gut angemessen ist. Auch bei unseren - eigentlich sehr einfachen - Beispieldaten traten solche Probleme auf (der Schätzalgorithmus konvergierte erst nach 75 Iterationen und erbrachte Schätzungen mit extrem hohen Standardfehlern), was als Indikator dafür gelten könnte, daß der unterstellte monotone Verlauf nur eine grobe Annäherung an den tatsächlichen Verlauf ist. (Daher verzichte ich auf eine Wiedergabe der Ergebnisse in Darstellung 4).

Im Gegensatz zu Weibull- und Gompertz-(Makeham-)Modellen sind die beiden folgenden Modelle in der Lage, Hazardfunktionen zu modellieren, die zunächst steigen und erst dann fallen. Relativ häufig angewendet wird die *log-logistische Verteilung* mit folgender Formulierung (A: 292 f.; BHM: 39 f. und 240 f.; DM: 153):<sup>22</sup>

$$r(t; \mathbf{X}) = \frac{\exp(\beta_0 + \mathbf{X}\boldsymbol{\beta})^p p t^{p-1}}{1 + (\exp(\beta_0 + \mathbf{X}\boldsymbol{\beta}) t)^p} \quad (11)$$

Hier ergibt sich bei einem Koeffizienten  $p$  von  $> 1$  (bzw. von  $\ln(p) > 0$ ) eine zunächst steigende und dann fallende Hazardrate, während bei  $p < 1$  bzw.  $\ln(p) < 0$  die Hazardrate von Anfang an fällt. Im vorliegenden Falle ergibt die Modellschätzung tatsächlich den nach den explorativen Analysen in Teil I, Abschnitt 3 zu erwartenden Verlauf einer zunächst steigenden und anschließend fallenden Hazardrate.

Ein noch flexibleres log-logistisches Modell mit einem weiteren Parameter wurde von Schneider (1991) und Brüderl (vgl. Brüderl/Diekmann 1995) entwickelt. In der Formulierung von Brüderl - diejenige von Schneider unterscheidet sich nur geringfügig - wird folgende Rate geschätzt:

$$r(t; \mathbf{X}) = \exp(\beta_0 + \mathbf{X}\boldsymbol{\beta}) \frac{p(\lambda t)^{p-1}}{1 + (\lambda t)^p} \quad (12)$$

Mit dieser Formulierung - zur Unterscheidung von der vorhergehenden Formulierung spreche ich vom log-logistischen Modell Typ II<sup>23</sup> - können noch steiler ansteigende Hazardfunktionen modelliert werden als mit dem zuerst genannten Modell vom Typ I. Dementsprechend weist das Modell vom Typ II eine noch größere Log-Likelihood auf als das Modell vom Typ I, d.h., es ist den Daten vermutlich noch besser angemessen. Nach der Höhe der Log-Likelihood handelt es sich hierbei um das erklärungskräftigste Modell überhaupt; es übertrifft sogar die weiter oben

dargestellten komplexen Exponentialmodelle, obwohl diese noch mehr Parameter zur Modellierung der Hazardfunktion verwenden.

Ebenfalls eine zunächst steigende und dann fallende Rate kann mit Hilfe der *Sichel-Verteilung* modelliert werden (vgl. A: 231; DM: 152 f.):

$$r(t; \mathbf{X}) = \exp(\beta_0 + \mathbf{X}\boldsymbol{\beta})t \exp(-t / \gamma) \quad (13)$$

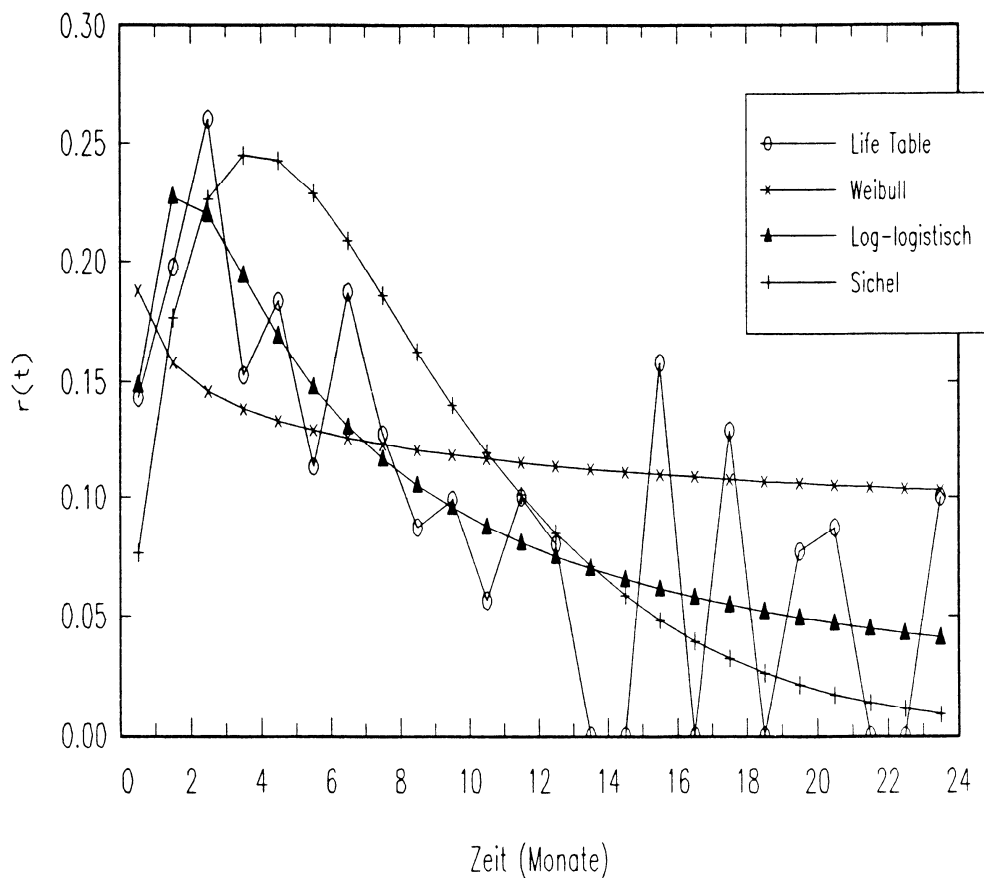
Dieses Modell läßt sich relativ gut inhaltlich interpretieren: Angenommen wird ein Maximum der Rate bei  $t = \gamma$  (in unserem Beispiel:  $\exp(1.3489) = 3.85$ ) und ein Wendepunkt bei  $t = 2\gamma$ .<sup>24</sup> Eine wesentliche Eigenschaft unterscheidet es von dem vorgenannten Modell: Es geht davon aus, daß ein Teil der Untersuchungspopulation niemals den Ausgangszustand verläßt, was in dem Kontext, in dem dieses Modell entwickelt wurde, der Analyse von Heiratsdauern, durchaus realistisch ist, da - trotz ansteigender Scheidungsziffern - die Mehrzahl der Ehen nicht mit einer Scheidung endet. Im vorliegenden Fall ist diese Annahme allerdings möglicherweise unangemessen, und so zeigt dieses Modell eine deutlich geringere Log-Likelihood als die meisten anderen Modelle.

Wenn wir jetzt wiederum einige der erörterten Modelle graphisch vergleichen<sup>25</sup>, so wird man auch aufgrund der Verläufe der Hazardraten (*Darstellung 5*) geneigt sein, dem log-logistischen Modell (Typ II) den Vorzug zu geben, welches auch den höchsten Wert der Log-Likelihood aufweist. Zusätzlich möchte ich die von den Modellen vorhergesagten Survivorfunktionen abbilden (*Darstellung 6*). Zunächst bestätigt sich die weiter oben getroffene Feststellung, daß die sehr unterschiedlichen Hazardraten, die von den verschiedenen Modellen geschätzt werden, zu recht ähnlichen Survivorfunktionen führen. Bei genauem Hinsehen zeigt sich aber doch, daß jedenfalls im mittleren Bereich das Weibull-Modell den Verlauf der Survivorfunktion tendenziell unter- und das Sichelmodell diesen Verlauf tendenziell überschätzt, während das log-logistische Modell insgesamt die größte Nähe zu dem »datennahen« Life-Table-Schätzer aufweist.

Trotzdem möchte ich noch einmal betonen, daß die »Ergebnisse« dieser Modelle teilweise durch die Modellwahl *vorausgesetzt* sind: Auch wenn die Hazardrate tatsächlich ansteigt, führt z.B. ein log-logistisches oder log-normales Modell zu dem »Ergebnis«, daß die Rate nach einem initialen Anstieg absinkt! Umgekehrt führen Weibull- und Gompertz-Modell immer zu dem »Ergebnis« monoton steigender oder fallender Raten, auch wenn kein monotoner Verlauf vorliegt.

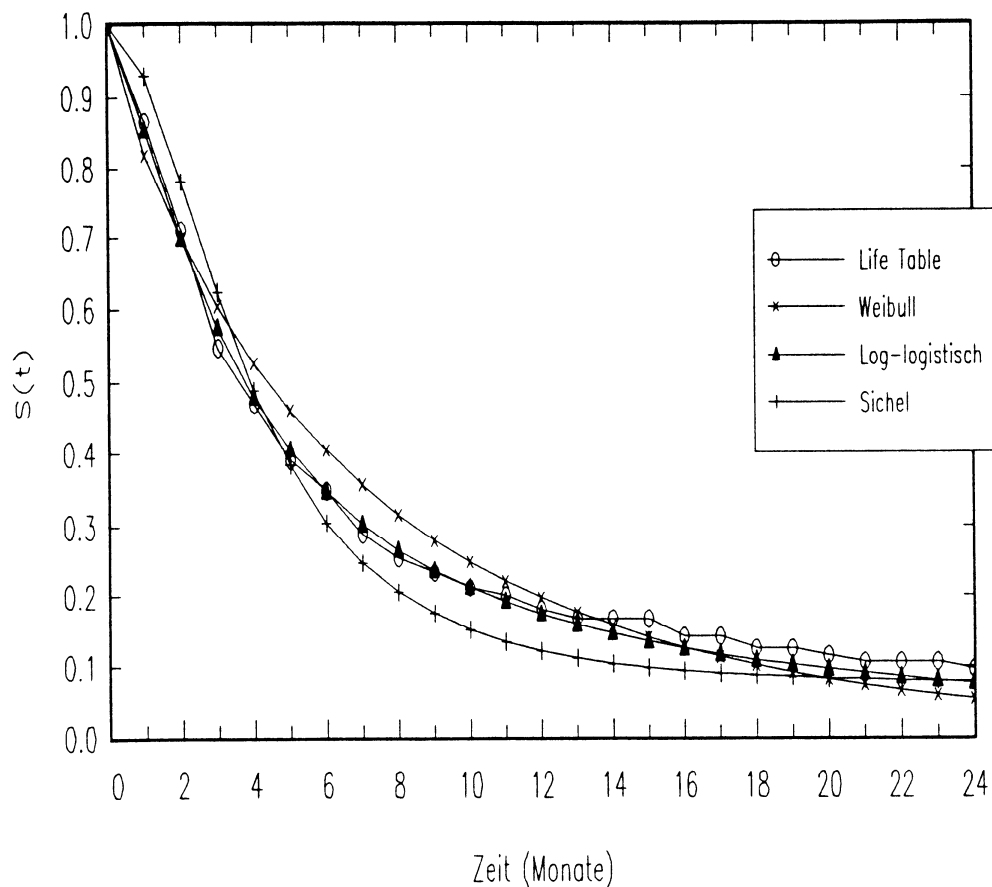
Für weitere Überprüfungen, welches aus der Vielzahl möglicher Modelle den Daten am besten angemessen ist, werden in der Literatur insbesondere Residuenanalysen vorgeschlagen. Deren Darstellung würde hier zu weit führen, und es sei wieder auf die Lehrbuchliteratur verwiesen (A: 272 f.; BHM v.a. S. 189, 217, 237, 246).

**Darstellung 5:** Hazardfunktionen für die Altersgruppe bis 30 (Schätzungen nach Life-Table, Weibull-, log-logistischem (Typ II) und Sichelmodell)



Zum Schluß ist darauf hinzuweisen, daß hier nur die gebräuchlichsten Verfahren parametrischer Analyse dargestellt werden konnten. So befindet sich unter den hier diskutierten Verfahren keines, welches eine zunächst fallende und dann steigende Hazardrate modellieren kann. Daher sind Versuche zu erwähen, sehr allgemeine Modelle zu formulieren, welche sowohl die hier vorgestellten Modelle als Spezialfall enthalten als auch weitere Verläufe zu modellieren geeignet sind. Zu nennen sind hier Modelle mit einer verallgemeinerten Gamma-Verteilung oder solche mit einer Box-Cox-Transformation, zu denen in der Lehrbuchliteratur jedoch leider nur kurze Hinweise vorliegen (A: 234 f.).<sup>26</sup> Auch hier gilt die oben geäußerte Warnung: Einerseits erlaubt die hohe Flexibilität der Modelle die Prüfung einer Vielzahl von Hypothesen; andererseits legt gerade dies nahe, mit der Interpretation von Ergebnissen, die nicht aufgrund gezielter Hypothesen, sondern durch exploratives »Herumprobieren« zustande gekommen sind, vorsichtig zu sein.

**Darstellung 6:** Survivorfunktionen für die Altersgruppe bis 30 (Schätzungen nach Life-Table, Weibull-, log-logistischem (Typ II) und Sichelmodell)



## 5. Modelle für diskrete Verweildauern

In diesem Abschnitt möchte ich kurz auf Modelle eingehen, die angewendet werden können, wenn *diskrete Verweildauern vorliegen*. Grundsätzlich - vgl. die Ausführungen in Teil I, Abschnitt 2 - ist damit gemeint, daß der interessierende Zustandswechsel nicht jederzeit, sondern nur zu fixen Zeitpunkten stattfinden kann (Wahlen, Versetzungen in der Schule). In der Praxis wird jedoch vielfach auch vorgeschlagen, solche Modelle anzuwenden, wenn es sich um stark gruppierte oder aggregierte Dauern handelt (Hamerle/Tutz 1989). So haben z.B. Licht/Steiner (1991) die hier herangezogenen Arbeitslosigkeitsdaten aus dem SOEP mit einem Modell für diskrete Dauern untersucht, weil es sich bei den monatlichen Messungen um gruppierte Daten handelt.<sup>27</sup>

Zwei relativ einfache und doch ziemlich flexible Modelle basieren auf der *logistischen Verteilung* sowie auf der *komplementären Log-Log-Verteilung* (auch um-

gekehrt als doppelte Exponentialverteilung bezeichnet). Die erste Verteilung ist aus der Analyse binärer abhängiger Variablen mit dem Verfahren der logistischen Regression inzwischen hinlänglich vertraut (Urban 1993). Als einfaches Modell erhalten wir auch hier ein Modell ohne Zeitabhängigkeit:

*Logistisches Modell für diskrete Zeit ohne Zeitabhängigkeit:*

$$r(t; \mathbf{X}) = \frac{\exp(\beta_0 + \mathbf{X}\boldsymbol{\beta})}{1 + \exp(\beta_0 + \mathbf{X}\boldsymbol{\beta})} \quad (14)$$

Zeitabhängigkeit des Prozesses kann ganz analog zum Exponentialmodell durch Hinzufügen eines Polynom-Terms für die Zeit modelliert werden:

*Logistisches Modell für diskrete Zeit mit Zeitabhängigkeit:*

$$r(t; \mathbf{X}) = \frac{\exp(\beta_0 + \mathbf{X}\boldsymbol{\beta} + \beta_{t_1}t + \beta_{t_2}t^2 + \dots + \beta_m t^m)}{1 + \exp(\beta_0 + \mathbf{X}\boldsymbol{\beta} + \beta_{t_1}t + \beta_{t_2}t^2 + \dots + \beta_m t^m)} \quad (15)$$

Dieses Modell hat den Vorzug, daß es grundsätzlich mit jedem Programm realisiert werden kann, welches ein logistisches Regressionsmodell schätzen kann (vgl. Allison 1982; Yamaguchi 1991, Kap. 2). Alternativ dazu wird auch ein Modell vorgeschlagen, welches auf der doppelten Exponentialverteilung bzw. deren Komplementärwert beruht; daher wird es vielfach als komplementäres Log-Log-Modell bezeichnet. Auch hier lassen sich ein einfaches Modell ohne Zeitabhängigkeit und ein Modell mit Polynom-Term für die Zeitabhängigkeit formulieren:

*Komplementäres Log-log-Modell für diskrete Zeit ohne Zeitabhängigkeit*

$$r(t; \mathbf{X}) = 1 - \exp\{-\exp(\beta_0 + \mathbf{X}\boldsymbol{\beta})\} \quad (16)$$

*Komplementäres Log-log-Modell für diskrete Zeit mit Zeitabhängigkeit*

$$r(t; \mathbf{X}) = 1 - \exp\{-\exp(\beta_0 + \mathbf{X}\boldsymbol{\beta} + \beta_{t_1}t + \beta_{t_2}t^2 + \dots + \beta_m t^m)\} \quad (17)$$

In unserem Beispiel führen beide Modelle zu praktisch identischen Ergebnissen. Daher werden nur die Ergebnisse des logistischen Modells vorgestellt (*Darstellung 7*). Die Ergebnisse ähneln sehr stark denjenigen, die mit dem Exponentialmodell für stetige Zeiten erhalten wurden. Im einfachen, zeitkonstanten Modell ist auch hier der Parameter für die Altersgruppe der 31- bis 50jährigen beinahe auf dem 5-Prozent-Niveau signifikant.

**Darstellung 7:** Ergebnisse logistischer Hazardratenmodelle für diskrete Zeit  
(erweiterte Beispieldaten aus Teil I, Darstellung 3)

Variable	Coeff	Error	T-Stat	Signif
<i>Zeitkonstantes Modell</i> (Log-likelihood: -1313.86)				
Konstante	-1.8945	0.0685	-27.6468	1.0000
Alter 31 bis 50	-0.2130	0.1094	-1.9461	0.9484
Alter über 50	-2.3404	0.2410	-9.7105	1.0000
<i>Polynom 1. Grades</i> (Log-likelihood: -1264.08)				
Konstante	-1.4773	0.0812	-18.1974	1.0000
Alter 31 bis 50	-0.0521	0.1113	-0.4679	0.3601
Alter über 50	-2.0359	0.2433	-8.3670	1.0000
$\beta_{t1}$	-0.0716	0.0092	-7.8140	1.0000
<i>Polynom 4. Grades</i> (Log-likelihood: -1258.07)				
Konstante	-1.8674189	0.1576898	-11.8423561	1.0000
Alter 31 bis 50	-0.0536346	0.1114339	-0.4813137	0.3697
Alter über 50	-2.0370082	0.2438730	-8.3527419	1.0000
$\beta_{t1}$	0.1665113	0.0768530	2.1665921	0.9697
$\beta_{t2}$	-0.0302462	0.0098949	-3.0567520	0.9978
$\beta_{t3}$	0.0011732	0.0004203	2.7915498	0.9948
$\beta_{t4}$	-0.0000136	0.0000055	-2.4630662	0.9862

Log-likelihood des Grundmodells nur mit Konstante: -1403.38

Allerdings zeigt dieses Modell wiederum eine deutlich schlechtere Modellanpassung als die zeitabhängigen Modelle, unter denen hier die beiden mit Polynom 1. und 4. Grades dargestellt werden, da diejenigen mit Polynom 2. und 3. Grades keine signifikanten Koeffizienten für die Polynomglieder und auch keinen Erklärungszuwachs erbrachten.

Einige weitere Modelle wie etwa ein Cox-Modell für gruppierte Zeiten sowie ein von Aranda-Ordaz vorgeschlagenes, relativ flexibles Modell werden in der einschlägigen Spezialliteratur diskutiert (Hamerle/Tutz 1989: 31 ff.).

Insgesamt läßt sich festhalten, daß vermutlich bei den hier untersuchten Daten die Meßeinheit für die Verweildauern in Form von Monaten klein genug ist, um sie auch mit Modellen für stetige Zeiten analysieren zu können. Die Tatsache, daß Modelle für diskrete Verweildauern, soweit sie vergleichbar sind, zu den gleichen Schlußfolgerungen führen - zunächst kurzfristig steigende, dann sinkende Hazardrate, kann aber als zusätzliche Abstützung der Ergebnisse gewertet werden. Liegen andere Datenstrukturen vor, insbesondere also noch größere Meßintervalle oder tatsächlich diskrete Ereignisse, sind die hier genannten Modelle eine wichtige alternative Auswertungsmöglichkeit.

## 6. Zeitabhängige (zeitveränderliche) Kovariaten

In diesem Abschnitt möchte ich noch kurz auf die Möglichkeit eingehen, mit den hier besprochenen Verfahren die Einflüsse von Kovariaten zu untersuchen, die sich selbst während des untersuchten Prozesses ändern. Grundsätzlich ist hierin eine der wichtigsten Anwendungsmöglichkeiten zu sehen, die mit anderen statistischen Verfahren nicht in der gleichen Weise erreicht werden kann. Eine ausführliche Erörterung würde aber den Rahmen dieser Arbeit sprengen, ich will daher nur einige Beispiele ansprechen und kurz auf die praktische Durchführung eingehen.

Zunächst ist auf die wichtige Unterscheidung zwischen *extern* und *intern* zeitabhängigen Kovariaten einzugehen. *Interne* Zeitabhängigkeit meint den Sachverhalt, daß sich Kovariaten in Abhängigkeit von der Dauer des untersuchten Prozesses selbst verändern. So könnten Arbeitslose mit zunehmender Arbeitslosigkeitsdauer immer größere Entmutigungseffekte zeigen.<sup>28</sup> Solche Effekte sind für die Datenauswertung problematisch, weil sie schwer von der Zeitabhängigkeit des Prozesses selbst zu trennen sind, so daß man häufig genötigt ist, die Zeitabhängigkeit einfach festzustellen und entweder ihr aufgrund theoretischer Überlegungen eine substantielle Interpretation zu geben oder zu versuchen, diese substantielle Interpretation durch andere Daten oder Analysemethoden zu erhärten (vgl. zum Beispiel der Arbeitslosigkeit Winter-Ebmer 1992).<sup>29</sup>

*Externe Zeitabhängigkeit* von Kovariaten liegt vor, wenn sich Werte von erklärenden Variablen zwar *während*, aber (mutmaßlich) *nicht infolge* des untersuchten Prozesses ändern. Ein Beispiel hierfür wäre die *Arbeitsmarktlage*, gemessen in regionalen Arbeitslosigkeitsquoten (vgl. Hujer/Schneider 1987, 1992). Diese ändert sich (praktisch) unabhängig vom Arbeitslosigkeitsschicksal des einzelnen Arbeitslosen. Im konkreten Beispiel ist insbesondere an Saisonalitätseffekte zu denken, d.h., an die Verschlechterung der Arbeitsmarktsituation im Winter und ihre Verbesserung im Frühjahr. Aber auch *Individualdaten* lassen sich als extern zeitveränderlich auffassen, z.B. Daten zur Familiensituation. Es ist bekannt, daß Heirat und Geburt von Kindern sich auf das Erwerbsverhalten insbesondere von Frauen auswirken, und es kann nach bisherigen Analysen davon ausgegangen werden, daß dies auch für die Arbeitslosigkeit zutrifft (Ludwig-Mayerhofer 1990). Man kann sich solche Kovariaten so vorstellen: Es gibt, ähnlich wie in unseren Beispieldaten, verschiedene Gruppen von Arbeitslosen, z.B. verheiratete und ledige Personen, mit je unterschiedlichen Hazardraten. Wenn eine Person während der Arbeitslosigkeit heiratet (oder umgekehrt sich trennt oder scheiden läßt), so heißt das praktisch, daß sie während der Arbeitslosigkeit aus der einen in die andere Gruppe wechselt, d.h., die Risikomenge der einen Gruppe verläßt und - während des Prozesses, also nicht zum Zeitpunkt Null, sondern zum Zeitpunkt des Wechsels! - derjenigen der anderen Gruppen zugeschlagen wird.<sup>30</sup>

Allgemein läßt sich also sagen: Bei der Einbeziehung zeitveränderlicher (oder zeitabhängiger) Kovariaten muß - ganz ähnlich wie bei dem im Vordergrund der

Analyse stehenden Prozeß selbst - bekannt sein, zu welchem Zeitpunkt ein Wechsel in der erklärenden Variablen eintritt, und natürlich, welche Werte diese Variable zu den verschiedenen Zeitpunkten hat. D.h. die zeitabhängigen Kovariaten können ohne weiteres mehrmals oder häufig ihre Werte wechseln, wie es z.B. bei einer Untersuchung der Arbeitslosigkeit der Fall wäre, die sich über mehrere Jahre erstreckt und monatliche Arbeitslosenquoten als zeitveränderliche Kovariate einbezieht. In anderen Fällen geht es möglicherweise nur um einmalige Veränderungen, etwa den Zeitpunkt, zu dem eine Person eine Ausbildung abgeschlossen hat. In solchen Fällen genügt es im allgemeinen, diesen Zeitpunkt als Variable in den Datensatz aufzunehmen und Fälle, bei denen die betreffende Änderung nicht eintritt, mit einem geeigneten Wert zu kodieren.

Wie lassen sich solche Variablen mit den hier untersuchten Modellen analysieren? Im Rahmen des *semiparametrischen Cox-Modells* ist das Verfahren grundsätzlich relativ einfach. Es genügt hier, wie soeben geschildert, die entsprechenden Variablen in einer Form im Datensatz verfügbar zu haben, die eine Beziehung auf die Dauer des untersuchten Prozesses ermöglicht. Wie das konkret geschieht, muß in den einschlägigen Programmhandbüchern und der Lehrbuchliteratur nachgesehen werden (BHM: 155 ff.). Allerdings wird, sofern es sich nicht nur um wenige Kovariaten handelt, die Rechenzeit u.U. überproportional lang, so daß sich bei vielen zeitveränderlichen Kovariaten auch im Cox-Modell ein Vorgehen empfiehlt, das bei den *parametrischen Modellen unumgänglich ist*. Dieses Vorgehen besteht darin, die untersuchten Episoden jeweils an der Stelle, an der die zeitveränderlichen Kovariaten ihren Wert ändern, in Teilepisoden zu zerlegen (»Episodensplitting«). Konkret: Wenn ein Individuum z.B. eine Arbeitslosigkeitsdauer von 10 Monaten aufweist und nach 5 Monaten heiratet, so müssen aus der entsprechenden Episode zwei Episoden gemacht werden: Die erste mit einer Dauer von 0 bis 5, die zweite mit einer Dauer von 5 bis 10. Die betreffende zeitabhängige Kovariate müßte in den beiden Teilepisoden Werte aufweisen, die den Zustandswechsel wiedergeben (also z.B. 0 für ledig und 1 für verheiratet), die zeitkonstanten Kovariaten müßten für beide Teilepisoden identisch sein. Die erste Teilepisode wäre als zensiert zu betrachten, da die Person ja nach 5 Monaten noch arbeitslos ist. Ein solches Episodensplitting ist entgegen der Darstellung von BHM (193 ff., v.a. 196) auch problemlos mit SPSS oder vergleichbaren Programmen durchführbar.<sup>31</sup> Zu beachten ist, daß nicht alle Programme in der Lage sind, entsprechende Datensätze auszuwerten; Voraussetzung ist, daß sie andere Anfangszeiten als »0« zulassen.<sup>32</sup> - Die Analyse mit parametrischen Modellen hat natürlich den Vorteil, daß hier zusätzlich die Zeitabhängigkeit der Hazardrate selbst modelliert werden kann.

Abschließend ist darauf hinzuweisen, daß man auf diese Weise auch die *wechselseitige Beeinflussung* von Prozessen untersuchen kann. Wenn wir hier das Beispiel angesprochen haben, daß die Geburt von Kindern die Arbeitslosigkeitsverläufe verändern kann, so könnte man auch der Frage nachgehen, ob umgekehrt die Arbeitslosigkeit von Frauen die Neigung beeinflusst, Kinder zu bekommen - wobei



die Wirkungsrichtung auch vom Familienkontext abhängen könnte, d.h., bei einer günstigen ökonomischen Situation könnte die Arbeitslosigkeit möglicherweise als Gelegenheit gesehen werden, einen Kinderwunsch zu realisieren, während umgekehrt eine ohnehin ungünstige Situation durch die eigene Arbeitslosigkeit noch verschlechtert würde und dadurch eher einem Kinderwunsch entgegenwirken könnte. Entsprechende Analysen für den Zusammenhang von Kindern und Unterbrechungen der Erwerbstätigkeit hat z.B. Huinink (1991, 1992) vorgestellt. Grundsätzlich sind solche Analysen wiederum, wie Huinink in den genannten Arbeiten verdeutlicht, vor allem dann unproblematisch, wenn *keine Dauerabhängigkeit* vorliegt, d.h., wenn z.B. zwar die Arbeitslosigkeit an sich, aber nicht ihre Dauer, sich auf die Geburt von Kindern auswirkt (und umgekehrt). Auch hier steht aber die Entwicklung von statistischen Modellen erst in den Anfängen.

## 7. Mehrere Zielzustände

Bislang wurden nur Übergänge in einen einzigen Zielzustand untersucht, in unserem Beispiel von der Arbeitslosigkeit in die Vollzeitbeschäftigung. Hier soll kurz darauf eingegangen werden, wie zu verfahren ist, wenn Übergänge in mehrere Zielzustände möglich sind (vgl. A: 77 ff.; BHM: 59 ff., 78 ff., 133 ff., 164 ff.; DM: 51 f., 174 ff.).

Entscheidend ist hier die Annahme, daß die verschiedenen Risiken voneinander *unabhängig* sein sollten. Diese Annahme zeigt eine deutliche Analogie zur Forderung nach der Unabhängigkeit von Zensurierungen und Ereignissen, und tatsächlich besteht ein unmittelbarer Bezug dazu. Denn, wie schon in Teil I kurz angedeutet, verfährt man bei der Analyse mehrerer Zielzustände so, daß die Übergänge in die einzelnen Zielzustände separat analysiert werden,<sup>33</sup> wobei jeweils sämtliche anderen Zielzustände als rechtszensierte Daten betrachtet werden (so wie wir auch in unserer Beispielsanalyse immer mit den Übergängen in andere Zustände verfahren sind). Den Analysen für die einzelnen Zielzustände können dabei ohne weiteres verschiedene Modelle des Verlaufs der Hazardfunktion zugrundegelegt werden, also z.B. ein zeitkonstantes Exponentialmodell für den Übergang in Zustand A, ein Weibull-Modell für den Übergang in Zustand B, usw.

Die Annahme der Unabhängigkeit der verschiedenen Zielzustände bzw. Übergänge ist nicht immer plausibel. Auch in unserem Beispiel müssen wir davon ausgehen, daß z.B. Personen mit geringen Chancen einer (Wieder-)Beschäftigung eine höhere Wahrscheinlichkeit aufweisen, aus dem Arbeitslosenbestand auszuschneiden und den Arbeitsmarkt ganz zu verlassen. Das kann dazu führen, daß die Beschäftigungschancen u.U. sogar überschätzt werden. Insofern müssen unsere Ergebnisse mit einer gewissen Vorsicht betrachtet werden.

Eine ausführliche Darstellung der Probleme bei voneinander abhängigen Risiken findet sich bei Klein (1988). Die Diskussion über Möglichkeiten der Schätzung

von Modellen bei abhängigen Risiken hat bislang noch keine leicht anwendbaren Verfahren erbracht (Hinweise hierzu etwa bei Schneider/Hujer 1992, Fn. 8).

## 8. Karrieren und wiederholte Episoden

Bislang wurden sehr einfache Modelle geschätzt, und das Versprechen einer Analyse von »Karrieren«, wie es in der Einleitung zu Teil I angedeutet wurde, läßt sich damit sicher nur begrenzt einlösen. In diesem Abschnitt sollen kurz einige Möglichkeiten angesprochen werden, komplexere Abfolgen von Episoden zu untersuchen, die sich gegebenenfalls im Sinne von Karrieren interpretieren lassen. D.h., es geht darum, *mehrere Episoden* im Lebenslauf von Individuen in geeigneter Weise aufeinander zu beziehen.

Eine erste Analysemöglichkeit könnte darin bestehen, *verschiedenartige Zustände bzw. Übergänge* zu untersuchen, aber jeweils zu fragen, ob Merkmale des vorangegangenen Zustandes bzw. Prozesses (gegebenenfalls auch mehrerer Zustände oder Prozesse) den aktuellen Prozeß beeinflussen. So könnte man fragen, ob z.B. die Arbeitslosigkeitsdauer durch Merkmale des vorherigen Zustandes beeinflusst wird, ob sich also Merkmale aus der Beschäftigung (oder einem anderen Zustand) *vor* der Arbeitslosigkeit *in* dieser auswirken. Umgekehrt ließe sich fragen, ob die Art oder Dauer von Beschäftigungsverhältnissen davon abhängt, ob oder wie lange Individuen vorher arbeitslos waren. So zeigte sich in einer Untersuchung an Berufsanfängern (Ludwig-Mayerhofer 1992b), daß eine ganz lange Arbeitslosigkeit (über 12 Monate) vor der ersten Beschäftigung zu einem geringerem Anfangseinkommen führt. Weiterhin wurde gefragt, ob Arbeitslosigkeit *vor* der ersten Beschäftigung auch das Risiko *erneuter* Arbeitslosigkeit erhöht. Tatsächlich ist ein solcher Effekt, jedenfalls als direkter, nicht oder kaum zu beobachten; der entscheidende Effekt hinsichtlich des Arbeitslosigkeitsrisikos ist vielmehr das Einkommen aus dem aktuellen Beschäftigungsverhältnis, wodurch immerhin ein indirekter Effekt einer früheren Arbeitslosigkeit denkbar ist.

Bei solchen Analysen ist grundsätzlich nicht anders vorzugehen als bisher geschildert, da jeweils nur ein Zustandswechsel untersucht wird und die Merkmale aus vorangegangenen Episoden oder Zuständen als Kovariaten eingesetzt werden. Anders ist dies, wenn man *wiederholte* Episoden der *gleichen Art* untersucht, wenn also z.B. in einer Stichprobe - wie auch in den Daten aus dem SOEP - Individuen mit mehreren (hier: Arbeitslosigkeits-)Episoden enthalten sind. Würde man diese mit einem der gängigen Modelle untersuchen, ohne das mehrfache Auftreten ein und derselben Person zu berücksichtigen, wäre die Annahme der Unabhängigkeit der einzelnen Beobachtungen voneinander verletzt. Es ist also grundsätzlich geboten, zwischen verschiedenen Episoden ein und desselben Individuums zu unterscheiden, zumal sich die Dauer und/oder die Einflüsse von Kovariaten zwischen verschiedenen Episoden unterscheiden können. Zudem wird es auch in diesem Fall

ratsam sein, Aspekte des Verlaufs vor der jeweiligen Episode (der »Vorgeschichte«) einzubeziehen, insbesondere bei den wiederholten Arbeitslosigkeitsepisoden Merkmale der früheren Arbeitslosigkeitsepisoden (vgl. grundsätzlich Hamerle 1989 sowie BHM: 62 ff., Diekmann/Mitter 1993: 56 f.).

Im Rahmen dieser Arbeit ist natürlich keine umfassende Diskussion möglich. Ich will aber die jedenfalls grundsätzlich vorhandene Fruchtbarkeit einer solchen Analyse anhand der Ergebnisse eines einfachen Modells zeigen. Dieses geht in zwei Hinsichten über die bisherigen Auswertungen hinaus: Erstens werden bis zu vier Arbeitslosigkeitsepisoden jeder Person analysiert, und zweitens werden die Einflüsse dreier zusätzlicher Variablen aus der »Vorgeschichte« untersucht: der Dauer der jeweils vorangegangenen Arbeitslosigkeitsepisode, der Zeit zwischen dem Ende der vorangegangenen und dem Beginn der gegenwärtigen Arbeitslosigkeitsepisode, und des Zustands, in den die vorangegangene Arbeitslosigkeitsepisode mündete, hier nur dichotomisiert nach Voll- und Teilzeitbeschäftigung vs. alle anderen Zustände (Haushalt, Ausbildung, Bundeswehr usw.). Zugrundegelegt wird ein log-logistisches Modell vom Typ I (*Darstellung 8*).<sup>34</sup> Selbstverständlich ließen sich noch komplexere Einflüsse heranziehen (z.B. bei der dritten Episode die Dauer der ersten *und* der zweiten Episode, usw.).

Vorbehaltlich der Tatsache, daß man die Ergebnisse wegen des Fehlens weiterer relevanter Variablen und angesichts der vor allem für die dritte und vierte Episode schon recht kleinen Fallzahlen nicht überbewerten sollte, zeigt sich, daß es zum einen durchaus Unterschiede zwischen den verschiedenen Episoden und zum anderen auch Einflüsse der Vorgeschichte geben könnte. So ist nach den Ergebnissen bei der zweiten bis vierten Episode auch eine Verschlechterung der Wiederbeschäftigungschancen für die 31- bis 50jährigen festzustellen (über 50jährige haben so wenige dritte und vierte Arbeitslosigkeitsepisoden, daß die Schätzungen sehr unzuverlässig werden). Ferner zeigt sich - teilweise nur als Tendenz -, daß mit längerer Dauer der vorangegangenen Arbeitslosigkeitsepisode auch die Hazardrate für die gegenwärtige Episode zurückgeht und daß die Personen, die schon früher nach der Arbeitslosigkeit direkt in ein Beschäftigungsverhältnis übergingen, bei einer späteren Arbeitslosigkeit wiederum bessere Wiederbeschäftigungschancen haben. Nur der Einfluß der Dauer seit der letzten Arbeitslosigkeit ist inkonsistent.<sup>35</sup>

Die Ergebnisse sind grundsätzlich die gleichen, die man erhalten würde, wenn man jeweils separate Modelle für die erste, zweite, dritte und vierte Arbeitslosigkeitsepisode schätzen würde (so sind die Ergebnisse für die erste Episode identisch mit jenen aus *Darstellung 4*). Der Vorteil der simultanen Schätzung in einem Modell ist darin zu sehen, daß geprüft werden kann, ob sich die Parameter für die einzelnen Episoden signifikant voneinander unterscheiden. So erbringt eine Modellschätzung unter der Annahme, daß die vier Modellkonstanten für die Basisrate miteinander identisch sind, ein Modell mit einer Log-Likelihood von 1919,61.

**Darstellung 8:** *Ergebnisse logistischer Hazardratenmodelle (Typ I) für kontinuierliche Zeiten, wiederholte Arbeitslosigkeitsepisoden*

Variable	Coeff	Error	T-Stat	Signif
<i>1. Episode</i>				
Konstante	-1.3210	0.0731	-18.0728	1.0000
Alter 31 bis 50	-0.0876	0.1199	-0.7306	0.5350
Alter über 50	-2.3385	0.2207	-10.5946	1.0000
Konstante	0.2925	0.0404	7.2349	1.0000
<i>2. Episode (214 Episoden mit 158 Übergängen)</i>				
Konstante	-1.1948	0.2720	-4.3929	1.0000
Alter 31 bis 50	-0.6784	0.1874	-3.6204	0.9997
Alter über 50	-1.1617	0.3588	-3.2377	0.9988
Dauer d. vorh. Arbeitslosigk.	-0.0266	0.0205	-1.2996	0.8063
Zeit seit vorher. Arbeitslosigk.	-0.0076	0.0082	-0.9249	0.6450
Zielzustand d. vorher. Arblk.	0.5627	0.2330	2.4151	0.9843
ln(p)	0.3861	0.0651	5.9331	1.0000
<i>3. Episode (96 Episoden mit 68 Übergängen)</i>				
Konstante	-1.1297	0.6889	-1.6399	0.8990
Alter 31 bis 50	-0.3633	0.2118	-1.7157	0.9138
Alter über 50	-0.3835	0.3521	-1.0891	0.7239
Dauer d. vorh. Arbeitslosigk.	-0.0872	0.0243	-3.5934	0.9997
Zeit seit vorher. Arbeitslosigk.	-0.0229	0.0121	-1.8993	0.9425
Zielzustand d. vorher. Arblk.	0.6626	0.6815	0.9723	0.6691
ln(p)	0.6880	0.1000	6.8781	1.0000
<i>4. Episode (43 Episoden mit 28 Übergängen)</i>				
Konstante	-2.2347	0.5389	-4.1465	1.0000
Alter 31 bis 50	-0.6374	0.3228	-1.9749	0.9517
Alter über 50	-0.5415	0.4276	-1.2665	0.7947
Dauer d. vorh. Arbeitslosigk.	-0.1081	0.0391	-2.7646	0.9943
Zeit seit vorher. Arbeitslosigk.	0.0837	0.0231	3.6278	0.9997
Zielzustand d. vorher. Arblk.	1.0436	0.4963	2.1028	0.9645
ln(p)	0.7917	0.1567	5.0514	1.0000

Log likelihood: -1917.96

Log-likelihood des Null-Modells (Exponentialmodell mit 4 Konstanten): -2124.65

Der Likelihood-Ratio-Test ergibt also im Vergleich zum oben dargestellten Modell einen  $\chi^2$ -Wert von 3,30, der bei 3 Freiheitsgraden (es wird im Vergleich zu vorher nur mehr eine einzige Regressionskonstante geschätzt) nicht signifikant auf dem 5-Prozent-Niveau ist. D.h., wir können - auf der Basis unseres unvollständigen Modells - nicht annehmen, daß wiederholte Arbeitslosigkeitsepisoden *ceteris paribus* länger sind als frühere.

Wie an diesem kleinen Beispiel schon deutlich geworden sein könnte, lassen sich durch die Verknüpfung wiederholter Episoden mit Daten aus anderen Prozes

sen oder Zuständen unter Umständen sehr komplexe Prozesse modellieren. Die Forschungspraxis beschränkt sich bislang allerdings ganz weitgehend auf Ein-Episoden-Modelle, und die Analyse komplexer Verlaufsmuster wurde noch kaum in Angriff genommen. Hier stehen für künftige Untersuchungen noch erhebliche Potentiale offen, bei denen allerdings auch ein beträchtlicher Datenerhebungsaufwand erforderlich ist.

## 9. Abschließende Bemerkungen

Ich hoffe, in dieser Arbeit aufgezeigt zu haben, daß die Verweildaueranalyse für Längsschnittuntersuchungen ganz erhebliche und wichtige Analysemöglichkeiten bietet, die weit über die früher verfügbaren Verfahren hinausgehen. Gleichzeitig hoffe ich, die Grundlagen dieser Modelle soweit erläutert zu haben, daß damit ein Verständnis einschlägiger Forschungsergebnisse ebenso möglich ist wie - in Verbindung mit der Lehrbuchliteratur - ein Einstieg in eigene Datenauswertungen.

Ich will abschließend noch einmal daran erinnern, daß der Einsatz der hier besprochenen Verfahren stets in Abhängigkeit von der konkreten Fragestellung und den verfügbaren Daten erfolgen muß. Verschiedene inhaltlich durchaus wichtige Probleme sind noch nicht oder jedenfalls nicht definitiv gelöst, etwa in der Analyse voneinander abhängiger konkurrierender Risiken oder in der simultanen Analyse von sich wechselseitig beeinflussenden Prozessen. Allerdings kann das kein Einwand gegen die Anwendung der hier diskutierten Verfahren sein, im Gegenteil. Es muß vielmehr hervorgehoben werden, daß erst die Beschäftigung mit diesen Verfahren die zugrundeliegenden Probleme ins Bewußtsein gehoben hat, die andernfalls vielleicht mit gänzlich unzulänglichen statistischen Mitteln angegangen worden wären, mit der Folge von gravierenden Methodenartefakten. Notwendig ist also, die Verfahren in dem Bewußtsein einzusetzen, daß zwar manche Probleme noch offen sind, daß dem aber nicht durch Abstinenz, sondern durch möglichst sachgerechten Einsatz der Verfahren und die daraus resultierende Akkumulation von Erfahrungen abgeholfen werden kann.

Selbstverständlich soll auch nicht der Eindruck erweckt werden, daß keine anderen Verfahren für Längsschnittuntersuchungen sinnvoll seien. Methoden zur Panelanalyse für metrische Variablen sind längst etabliert (Arminger/Müller 1990). Auch in der Panelanalyse von diskreten (bislang allerdings nur binären oder ordinalen) Variablen wurden in den letzten Jahren erhebliche Fortschritte erzielt (Andreß 1992b; Maddala 1987; Petersen 1993; Schneider/Hujer 1992).

\* \* \* \* \*

## Anhang

### Programme zur Verlaufsdatenanalyse.

Wie schon in Teil I erwähnt, stellt TDA das umfassendste Programm zur Analyse von Verlaufsdaten dar. Unter anderem ist die simultane Analyse wiederholter Episoden nur in diesem Programm als Standardprozedur verfügbar, sieht man von einzelnen schwer zugänglichen und sehr speziellen Programmen ab. Auch Panelmodelle für diskrete Daten sind in TDA implementiert. Einzig SAS bietet unter den verbreiteten Programmen eine annähernd große Vielfalt an Verfahren. In BMDP sind immerhin die gängigsten Modelle für kontinuierliche Verweildauern implementiert, während SPSS zur Zeit nur Prozeduren für nicht- und semiparametrische Verfahren enthält, und auch dies nur in der Windows-Version.

Unter den weniger verbreiteten Programmen sind noch zu nennen: LIMDEP (Greene 1992) und PARAT (Schneider 1991), die verschiedene Modelle für kontinuierliche Verweildauern schätzen können, GLAMOUR mit Verfahren für diskrete Verweildauern (Tutz/Georg 1991) sowie EGRET (Statistics and Epidemiology Research Corporation 1991), das insbesondere für Mediziner und Epidemiologen von Interesse ist, allerdings neben nicht- und semiparametrischen Verfahren nur wenige Modelle für kontinuierliche Dauern enthält. Mit dem Programm GLIM können sehr viele Modelle durch benutzerdefinierte Routinen geschätzt werden.

Nur am Großrechner verfügbar ist das Programm RATE (Tuma 1980), welches einmal eines der wichtigsten Programme war (vgl. die Beispiele in BHM), inzwischen aber durch einige andere Programme überholt ist.

Bei allen Programmen ist vor allem bei den multivariaten Verfahren sehr genau darauf zu achten, in welcher Art und Weise die verschiedenen Modelle formuliert (»parametrisiert«) werden. Beispielsweise verwenden BMDP und SAS vielfach Parametrisierungen, die stark von den hier (in Anlehnung vor allem an die deutschsprachige Lehrbuchliteratur) vorgestellten Formulierungen abweichen. Die angegebenen Modellparameter haben teilweise genau das umgekehrte Vorzeichen, so daß ein Effekt mit positivem Vorzeichen tatsächlich einer Verlängerung der Verweildauer, also einer Verringerung der Hazardrate entspricht. Auch die Parameter für die Zeitbezogenheit der Übergangsraten müssen zum Teil umgerechnet werden (die Formulierungen entsprechen offenbar überwiegend denjenigen bei Kalbfleisch/Prentice 1980, vor allem 24 ff.).

### Anmerkungen

- 1 Es gibt allerdings – soviel im Vorgriff – Techniken, die »Basisrate« dennoch zu schätzen; dies stellt jedoch einen zweiten Schritt zusätzlich zur Schätzung der Modellparameter dar.
- 2 Ein Beispiel, wo die »Hazardkomponenten« explizit vorgestellt werden, findet sich bei Kiefer (1988).

- 3 Zur Klarstellung sei darauf hingewiesen, daß sich das Alter in den später vorgestellten multivariaten Modellen auch als metrische Variable behandeln ließe. Für die hier diskutierten nicht-parametrischen Verfahren ist es aber immer erforderlich, metrische Variable zu gruppieren. Allerdings ist dies im Sinne einer explorativen Datenanalyse durchaus nicht nur ein notwendiges Übel, sondern sogar sehr empfehlenswert, um etwaigen nicht-linearen Zusammenhängen auf die Spur zu kommen.
- 4 Die von SPSS/PC<sup>+</sup> (Version 4) ausgegebene Lee-Desu-Statistik führt praktisch zu identischen Ergebnissen wie die Gehan/Breslow-Statistik.
- 5 Fettgedruckte Buchstaben verweisen darauf, daß sich bei den Größen um Matrizen bzw. Vektoren handelt.
- 6 Anschauliche Beispiele aus dem Bereich der ML-Schätzung für kategoriale Daten finden sich bei Maier/Weiss (1990: 80 ff.) und Urban (1993: 53 ff.). Konkrete Beispiele für die Verweildaueranalyse finden sich vor allem bei A: 191 ff. zu ML- und 241 ff. zu PL-Schätzungen.
- 7 Auf die naheliegende Frage, welcher Stichprobenumfang als »hinreichend« gelten kann, findet sich in der einschlägigen Literatur keine klare Antwort. Aus verschiedenen Simulationsstudien, allerdings mit meist einfachen Datensätzen mit nur einer oder zwei Kovariaten, läßt sich ersehen, daß bereits Stichprobenumfänge von 50 bis 100 Fällen zu einigermaßen zuverlässigen Ergebnissen führen (vgl. A: 203 sowie Tuma 1982). Ganz grundsätzlich läßt sich sagen (was aber für alle multivariaten Verfahren gilt), daß natürlich bei kleinen Stichprobenumfängen die Zahl der geprüften Kovariaten ebenfalls nur sehr klein sein sollte. Studien, von denen gelegentlich berichtet wird, in denen die Zahl der geprüften Variablen kaum kleiner oder gar größer ist als der Stichprobenumfang, können nur als unsinnig bezeichnet werden. – Ebenso relevant sind Fragen der Teststärke, also der Vermeidung von  $\beta$ -Fehlern. Hierzu gibt es leider noch weniger Aufschluß. (Zu nicht-parametrischen Verfahren vgl. die Arbeit von Freedman 1982.)
- 8 Allerdings legt das die Gefahr nahe, diese Signifikanzniveaus einfach zu übernehmen. Tatsächlich beziehen sich diese immer auf zweiseitige Signifikanztests. Liegen gerichtete Hypothesen vor, müßte man sich an den Werten für einseitige Signifikanztests orientieren (also z.B. 1,645 für ein Signifikanzniveau von 0,05).
- 9 Der Ausdruck rührt daher, daß Formel 16 identisch ist mit der Formel

$$-2 \ln \frac{L_1}{L_0} \quad (3),$$

wobei  $L_0$  und  $L_1$  der (nicht logarithmierten) Likelihood entspricht. Grundsätzlich geben aber alle Programme die Log-Likelihood aus.

- 10 Dies entspricht also etwa dem globalen F-Test eines linearen Regressionsmodells.
- 11 Grundsätzlich ist dies auch mit der Wald-Statistik möglich, allerdings ist diese in den meisten Programmen nur als Test implementiert, ob einzelne Koeffizienten von Null verschieden sind.
- 12 Zu beachten ist, daß für die folgenden Beispiele jeweils der gesamte Datensatz herangezogen wurde, also auch die Verweildauern, die mehr als 12 Monate betragen. Daher führt ein Nachrechnen anhand von Darstellung 3 aus Teil I zu anderen Ergebnissen!
- 13 Für die Schätzungen wurde das Programm TDA verwendet. Dieses gibt den (auf 4 Stellen gerundeten) Wert  $1 - p$  des Signifikanzniveaus aus, d.h., der Wert von 1,0000 in der Spalte „Signif“ bedeutet, daß das Signifikanzniveau unter 0,00005 liegt.
- 14 Beim Nachrechnen ergeben sich hier wie anderswo wegen Rundungen leichte Abweichungen.
- 15 Diese Schlußfolgerung gilt natürlich nur im Rahmen unsers einfachen Beispiels. Es könnte sein, daß sich durch Einbeziehen weiterer wichtiger Variablen die Ergebnisse auch hinsichtlich des Alters ändern!

- 16 Inzwischen können - mit dem Programm TDA - Modelle mit unbeobachteter Heterogenität nicht für das Exponentialmodell (A: 266 ff.; BHM: 251 ff.), sondern auch für zahlreiche andere Modelle geschätzt werden.
- 17 Das Exponentialmodell hat allerdings nicht deshalb seinen Namen, sondern weil es unterstellt, daß  $S(t)$  exponentiell mit der Zeit abnimmt.
- 18 Es ist noch einmal daran zu erinnern, daß die hier vorgestellten Ergebnisse mit dem Datensatz über die gesamte Beobachtungsdauer berechnet wurden, also nicht nur mit den Daten für die ersten 12 Monate, die in Darstellung 3 enthalten sind. Zieht man nur diese Daten heran, so erhält man teilweise konträre Ergebnisse. Insbesondere würden diejenigen Verfahren, die nur monoton steigende oder fallende Verläufe modellieren (vgl. die folgende Darstellung von Weibull- und Gompertz-Verteilung), nicht zu den hier berichteten Ergebnissen einer fallenden, sondern zu einer steigenden Hazardrate führen. Die Ergebnisse werden also nicht unerheblich von der Untersuchungsdauer beeinflußt! Dieser Sachverhalt könnte auch die unterschiedlichen Ergebnisse von Klein (1990) und Ludwig-Mayerhofer (1992a) erklären, da in der erstgenannten Arbeit nur vier, in der zweiten jedoch sechs Wellen des Sozio-ökonomischen Panels herangezogen wurden (vgl. hierzu auch Hujer/Schneider 1992: 328).
- 19 Zu diesen beiden Modellen gibt es in den deutschen Lehrbüchern leider nur eine kurze Fundstelle (A: 231 ff.). Das in BHM (205 ff.) vorgestellte Modell mit periodisierter Verweildauer ist eine noch umfassendere Ausweitung; es werden nicht nur periodenspezifische Konstanten, sondern auch periodenspezifische Koeffizienten für die Kovariateneinflüsse geschätzt.
- 20 Aus Gründen der Übersichtlichkeit ist die Hazardrate für das Piecewise Constant Exponentialmodell nicht aufgeführt.
- 21 In den Lehrbüchern werden die einzelnen Bestandteile - man ist geneigt zu sagen: wie könnte es auch anders sein - unterschiedlich bezeichnet: Der hier nach DM mit „ $p$ “ bezeichnete Parameter heißt bei BHM „ $\alpha$ “, und bei A „ $\gamma$ “. Zu beachten ist, daß die bei A: 233 angegebene Formel 5.58 mißverständlich ist, da es dort den Anschein hat, als wäre  $\gamma$  ein Index zu  $\lambda$ !
- 22 Dieser Verteilung sehr ähnlich ist die log-normale Verteilung, auf die ich nicht weiter eingehen will (ebenfalls nur kurze Darstellungen bei A: 74, 292 f.; BHM: 54 f.).
- 23 Diese Bezeichnungen lehnen sich an das Programm TDA an, in dem auch diese beiden Modelle implementiert sind.
- 24 Es ist auch denkbar, Kovariaten-Einflüsse hinsichtlich des  $\gamma$ -Terms zu schätzen, was zu einem relativ komplexen Modell führt. Daher verzichte ich auf eine Darstellung, zumal die Einflüsse im vorliegenden Fall wiederum nicht signifikant waren.
- 25 Die Schätzungen des Gompertz-Modells sind in den folgenden Abbildungen nicht enthalten, weil sie sich kaum von denen des Weibull-Modells unterscheiden. Das log-logistische Modell Typ I verläuft ähnlich dem Modell Typ II, aber in den ersten Monaten deutlich flacher.
- 26 Allerdings können diese Modelle mit dem Programm TDA geschätzt werden; das Manual des Programms enthält auch Erläuterungen zu den Modellen.
- 27 Man könnte sogar argumentieren, daß es sich um diskrete Daten im Wortsinn handelt, weil die meisten Entlassungen und Einstellungen jeweils zu Monatsende bzw. Monatsanfang erfolgen.
- 28 Ein weiteres, oft in Zusammenhang mit Arbeitslosigkeit diskutiertes Beispiel ist der Bezug von Arbeitslosenunterstützung: Der Bezug von Arbeitslosenhilfe tritt sehr häufig erst nach Ablauf des Arbeitslosengeldes ein, d.h., Arbeitslosenhilfe ist fast per definition mit längeren Arbeitslosigkeitsdauern verbunden. Ob und wie diese Problematik zu lösen ist, ist in der Literatur umstritten (vgl. unter anderem Hunt 1992).
- 29 Für eine ausführliche, aber wiederum nicht einfache Diskussion sei auf die Arbeiten von Petersen (1986a, b) verwiesen.
- 30 Man kann sich dies sogar an den Beispieldaten vorstellen: Da die Personen während des Arbeitslosigkeitsprozesses selbst älter werden, können sie auch von der einen Altersgruppe in die andere wechseln. Bei dieser Konzeption ist allerdings schon wieder fraglich, ob es sich tatsächlich um eine externe zeitveränderliche Kovariate handelt.



- 31 Eine ausführliche Beschreibung des Vorgehens mit SPSS/PC<sup>+</sup> bzw. SPSS for Windows, das aber auch auf die meisten anderen Programme übertragbar sein dürfte, findet sich in Brüderl/Ludwig-Mayerhofer 1994). Die Lehrbuchliteratur enthält leider keine geeigneten Darstellungen; das Beispiel von BHM (195 f.) ist - vorsichtig formuliert - mißverständlich.
- 32 Das trifft z.B. für BMDP oder SPSS nicht zu. Besonders hinzuweisen ist wieder auf das Programm TDA, das eine eingebaute Funktion zum Episodensplitting enthält, die allerdings nicht ganz einfach zu handhaben und nicht sehr gut dokumentiert ist.
- 33 Mit dem Programm TDA ist es möglich, mehrere Übergänge in einem einzigen Modell simultan zu schätzen. Der Vorteil ist, daß es auf diese Weise auch möglich ist, Unterschiede zwischen den Parametern für die verschiedenen Modelle zu testen. Der Nachteil ist, daß in diesem Fall allen Übergängen die gleiche funktionale Form der Hazardrate zugrundegelegt wird.
- 34 Ein Modell vom Typ II konnte nicht geschätzt werden, da der Schätzalgorithmus - trotz Vergabe von Startwerten - nicht konvergierte. Dies sei wiederum als Hinweis darauf gegeben, daß mit zunehmender Komplexität der Modelle auch die Schwierigkeiten mit ihrer Handhabung zunehmen.
- 35 Da die auf die Dauer in der Vorgeschichte bezogenen Variablen eine ziemlich schiefe Verteilung aufweisen, wäre hier auch daran zu denken, Transformationen dieser Variablen (z.B. den Logarithmus) zu verwenden. An der je nach Episode unterschiedlichen Wirkung der letztgenannten Variablen würde dies im konkreten Fall jedoch nichts Grundsätzliches ändern.

## Literatur

- Allison, P. D., 1982: Discrete-Time Methods for the Analysis of Event Histories. S. 61-98 in: Leinhardt, S. (Hrsg.), *Sociological Methodology* 1982. San Francisco: Jossey-Bass.
- Andreß, H.-J., 1988: Spezifikationsfehler und unbeobachtete Heterogenität in Regressionsmodellen für Übergangsraten. *Kölner Zeitschrift für Soziologie und Sozialpsychologie* 40: 93-116.
- Andreß, H.-J., 1992a: Verlaufsdatenanalyse (Historical Social Research/Historische Sozialforschung, Supplement/Beiheft No. 5). Köln: Zentrum für Historische Sozialforschung.
- Andreß, H.-J., 1992b: Logistische Regressionsmodelle für Paneldaten. Analyse dichotomer Variablen im Zeitverlauf unter besonderer Berücksichtigung unbeobachteter Heterogenität. S. 35-66 in: Andreß, H.-J./Huinink, J./Meinken, H./Rumianek, D./Sodeur, W./Sturm, G. (Hrsg.), *Theorie, Daten, Methoden. Neue Modelle und Verfahrensweisen in den Sozialwissenschaften*. München: R. Oldenbourg.
- Arminger, G., 1984: Modelltheoretische und methodische Probleme bei der Analyse von Paneldaten mit qualitativen Variablen. *Vierteljahreshefte zur Wirtschaftsforschung*: 470-479.
- Arminger, G./Müller, F., 1990: *Lineare Modelle zur Analyse von Paneldaten*. Opladen: Westdeutscher Verlag.
- Blossfeld, H. P./Hammerle, A./Mayer, K. U., 1986: *Ereignisanalyse*. Frankfurt/New York: Campus.
- Brüderl, J./Diekmann, A., 1995: The Log-Logistic Rate Model. Two Generalizations with an Application to Demographic Data. *Sociological Methods & Research* 1995 (im Erscheinen).
- Brüderl, J./Ludwig-Mayerhofer, W., 1994: Aufbereitung von Verlaufsdaten mit zeitveränderlichen Kovariaten mit SPSS. *ZA-Information* 34: 79-105.

- Cox, D. R., 1972: Regression Models and Life-tables (with discussion). *Journal of the Royal Statistical Society, Series B*, 34: 187-220.
- Diekmann, A./Klein, T., 1991: Bestimmungsgründe des Ehescheidungsrisikos. *Kölner Zeitschrift für Soziologie und Sozialpsychologie* 43: 271-290.
- Diekmann, A./Mitter, P., 1984: *Methoden zur Analyse von Zeitverläufen*. Stuttgart: Teubner.
- Diekmann, A./Mitter, P., 1993: Methoden der Ereignisanalyse in der Bevölkerungssoziologie: Stand und Probleme. S. 20-65 in: Diekmann, A./Weick, S. (Hrsg.), *Der Familienzyklus als sozialer Prozeß. Bevölkerungssoziologische Untersuchungen mit den Methoden der Ereignisanalyse. (Sozialwissenschaftliche Schriften, Heft 26)*. Berlin: Duncker & Humblot.
- Freedman, L. S., 1982: Tables of the Number of Patients Required in Clinical Trials Using the Logrank Test. *Statistics in Medicine* 1: 121-129.
- Galler, H. P./Pötter, U., 1987: Unobserved Heterogeneity in Models of Unemployment Duration. S. 628-650 in: Mayer, K. U./Tuma, N. B. (Hrsg.), *Applications of Event History Analysis in Life Course Research. (Materialien aus der Bildungsforschung, Vol. 30)*. Berlin: Max-Planck-Institut für Bildungsforschung.
- Galler, H. P./Pötter, U., 1992: Zur Robustheit von Schätzmodellen für Ereignisdaten. S. 379-405 in: Hujer, R./Schneider, H./Zapf, W. (Hrsg.), *Herausforderungen an den Wohlfahrtsstaat im strukturellen Wandel*. Frankfurt/New York: Campus.
- Greene, W. E., 1992: *LIMDEP, Version 6.0*. New York: Econometric Software.
- Hamerle, A., 1989: Multiple-spell Regression Models for Duration Data. *Applied Statistics* 38: 127-138.
- Hamerle, A./Tutz, G., 1989: *Diskrete Modelle zur Analyse von Verweildauern und Überlebenszeiten*. Frankfurt/New York: Campus.
- Huinink, J., 1991: The Analysis of Interdependent Social Processes - The Example of Life-Course Analysis. S. 601-615 in: Albrecht, G./Otto, H. U. (Hrsg.), *Social Prevention and the Social Sciences. Theoretical Controversies, Research Problems, and Evaluation Strategies*. Berlin/New York: de Gruyter.
- Huinink, J., 1992: Die Analyse interdependenter Lebensverlaufsprozesse. Zum Zusammenhang von Familienbildung und Erwerbstätigkeit bei Frauen. S. 343-366 in: Andreß, H.-J./Huinink, J./Meinken, H./Rumianek, D./Sodeur, W./Sturm, G. (Hrsg.), *Theorie, Daten, Methoden. Neue Modelle und Verfahrensweisen in den Sozialwissenschaften*. München: R. Oldenbourg.
- Hujer, R./Schneider, H., 1987: Ökonometrische Ansätze zur Analyse von Paneldaten: Schätzung und Vergleich von Übergangsratenmodellen. S. 219-242 in: Krupp, H.-J./Hanefeld, U. (Hrsg.), *Lebenslagen im Wandel: Analysen 1987*. Frankfurt/New York: Campus.
- Hujer, R./Schneider, H., 1992: Strukturelle und institutionelle Determinanten der Arbeitslosigkeit aus mikroanalytischer Sicht. S. 315-341 in: Hujer, R./Schneider, H./Zapf, W. (Hrsg.), *Herausforderungen an den Wohlfahrtsstaat im strukturellen Wandel*. Frankfurt/New York: Campus.

- Hunt, J. 1992: The Effect of Unemployment Compensation on Unemployment Duration in Germany. Diskussionspapier Nr. 50, Deutsches Institut für Wirtschaftsforschung (DIW), Berlin.
- Kalbfleisch, J. D./Prentice, R. L., 1980: The Statistical Analysis of Failure Time Data. New York: Wiley.
- Kiefer, N. M., 1988: Economic Duration Data and Hazard Functions. *Journal of Economic Literature* 26: 646-679.
- Klein, T., 1988: Zur Abhängigkeit zwischen konkurrierenden Mortalitätsrisiken. *Allgemeines Statistisches Archiv* 72: 248-258.
- Klein, T., 1990: Arbeitslosigkeit und Wiederbeschäftigung im Erwerbsverlauf. *Kölner Zeitschrift für Soziologie und Sozialpsychologie* 42: 688-705.
- Klein, T., 1992: Zur Zeitabhängigkeit der Wiederbeschäftigungsrate Arbeitsloser. *Kölner Zeitschrift für Soziologie und Sozialpsychologie* 44: 134-138.
- Licht, G./Steiner, V., 1991: Abgang aus der Arbeitslosigkeit, Individualeffekte und Hysteresis. Eine Panelanalyse für die Bundesrepublik Deutschland. S. 182-205 in: Helberger, C./Bellmann, L./Blaschke, D. (Hrsg.), *Erwerbstätigkeit und Arbeitslosigkeit. Analysen auf der Grundlage des Sozioökonomischen Panels*. (BeitrAB 144). Nürnberg: IAB.
- Liebersohn, S., 1985: *Making It Count*. Berkeley: University of California Press.
- Ludwig-Mayerhofer, W., 1990: Arbeitslosigkeit im Erwerbsverlauf. *Zeitschrift für Soziologie* 19: 345-359.
- Ludwig-Mayerhofer, W., 1992a: Fakt und Artefakt in der Analyse von Arbeitslosigkeitsverläufen. *Kölner Zeitschrift für Soziologie und Sozialpsychologie* 44: 124-133.
- Ludwig-Mayerhofer, W., 1992b: Arbeitslosigkeit an der "zweiten Schwelle" - was sind die Folgen? Eine Analyse anhand des Sozioökonomischen Panels (1984-1988). S. 172-189 in: Kaiser, M./Görlitz, O. (Hrsg.), *Berufliche Bildung im Umbruch - Zur Diskussion der Übergänge in Bildung und Beruf im vereinten Deutschland* (BeitrAB 153.2). Nürnberg: IAB.
- Maddala, G. S., 1987: Limited Dependent Variable Models Using Panel Data. *Journal of Human Resources* 22: 307-338.
- Maier, G./Weiss, P., 1990: *Modelle diskreter Entscheidungen*. Wien: Springer.
- Petersen, T., 1986a: Estimating Fully Parametric Hazard Rate Models with Time-dependent Covariates. *Sociological Methods & Research* 14: 219-246.
- Petersen, T., 1986b: Fitting Parametric Survival Models With Time-Dependent Covariates. *Applied Statistics* 35: 281-288.
- Petersen, T., 1993: Recent Advances in Longitudinal Methodology. *Annual Review of Sociology* 19: 425-454.
- Prentice, R. L., 1978: Linear Rank Tests with Right Censored Data. *Biometrika* 65: 167-179.
- Prentice, R. L./Marek, P., 1979: A Qualitative Discrepancy between Censored Data Rank Tests. *Biometrics* 35: 861-867.
- Rohwer, G., 1994: TDA Working Papers. Bremen, Ms.

- Schneider, H., 1991: Verweildaueranalyse mit GAUSS. Frankfurt/New York: Campus.
- Schneider, H./Hujer, R., 1992: Die Analyse struktureller Wandlungsprozesse mit Hilfe von Panel-daten. S. 39-69 in: Hujer, R./Schneider, H./Zapf, W. (Hrsg.), Herausforderungen an den Wohlfahrtsstaat im strukturellen Wandel. Frankfurt/New York: Campus.
- Schnell, R., 1994: Graphisch gestützte Datenanalyse. München, Wien: Oldenbourg.
- Statistics and Epidemiology Research Corporation, 1991: EGRET Reference Manual and Manual Addendum. Seattle: Statistics and Epidemiology Research Corporation.
- Tarone, R. E./Ware, J., 1977: On Distribution-free Tests for Equality of Survival Distributions. *Biometrika* 64: 156-160.
- Teachman, J. D., 1983: Analyzing Social Processes: Life Tables and Proportional Hazards Models. *Social Science Research* 12: 263-301.
- Tuma, N. B., 1980: Invoking RATE. Menlo Park: SRI International.
- Tuma, N. B., 1982: Nonparametric and Partially Parametric Approaches to Event-History Analysis. S. 1-60 in: Leinhardt, S. (Hrsg.), *Sociological Methodology 1982*. San Francisco: Jossey-Bass.
- Tutz, G./Georg, W., 1991: Diskrete Hazardraten-Modelle in der Shell-Jugendstudie 1985: Eine Anwendung des Programms GLAMOUR. *ZA-Information* 29: 81-93.
- Urban, D., 1993: Logit-Analyse. Statistische Verfahren zur Analyse von Modellen mit qualitativen Response-Variablen. Stuttgart: Gustav Fischer.
- Winter-Ebmer, R., 1992: Persistenz von Arbeitslosigkeit. Frankfurt/New York: Campus.
- Yamaguchi, K., 1991: Event History Analysis. Newbury Park: Sage.
- Ziegler, R./Brüderl, J./Diekmann, A., 1988: Stellensuchdauer und Anfangseinkommen bei Hochschulabsolventen. *Zeitschrift für Wirtschafts- und Sozialwissenschaften* 108: 247-270.