

# 25 Prognose von Softwarezuverlässigkeit, Softwareversagensfällen und Softwarefehlern

von Michael Grottke

## 25.1 Einleitung

Im Laufe der letzten Jahrzehnte hat Software nicht nur in der Form von Computerprogrammen an Bedeutung gewonnen, sondern zudem als integraler Bestandteil verschiedenster Systeme technischer und/oder kommerzieller Natur fast jeden Lebensbereich erobert. Beispiele für solche Systeme sind medizinische Geräte, Automobile, Zahlungstransaktionssysteme, etc. Diese Entwicklung hat zweifelsohne die Möglichkeiten der betroffenen Maschinen sowie die Bandbreite der mit ihrer Hilfe produzierten Güter oder angebotenen Dienstleistungen erhöht. Allerdings ist der Mensch durch den Siegeszug der Software auch in zunehmendem Maße von ihrem korrekten Verhalten abhängig geworden.

Während die Unterstützung der Analyse-, Design- und Implementierungsphase der Softwareerstellung durch systematische Methoden und Werkzeuge gewisse Irrtümer von Softwareentwicklern verhindern kann, lässt sich eine völlige Programmkorrektheit praktisch nicht garantieren. Vielmehr ist immer damit zu rechnen, dass eine Software *Fehler* (insbesondere falsche oder fehlende Programmzeilen) aufweist, deren Aktivierung bei der Programmausführung ein *Versagen* zur Folge hat, also ein Softwareverhalten, welches von dem eigentlich spezifizierten abweicht. Unter *Softwarezuverlässigkeit* versteht man nun allgemein die Wahrscheinlichkeit dafür, dass in einer definierten Umgebung die Software innerhalb einer bestimmten Nutzungsperiode kein Versagen zeigt [3]. Bei kontinuierlich laufender Software wird die Länge der „Nutzungsperiode“ in der Regel in Form der zeitlichen Dauer der Programmnutzung gemessen [40]. Hingegen ist es bei Transaktionssystemen und Ähnlichem sinnvoll, die „Nutzungsperiode“ an der Anzahl der Programmläufe festzumachen [3]. Da die beobachteten Versagensfälle dynamische Phänomene sind, deren Frequenz z. B. davon abhängt, wie häufig fehlerhafte Codestellen aufgerufen werden, bezieht sich ein geschätzter Zuverlässigkeitswert immer auf eine bestimmte Art der Programmnutzung. Das betont die oben genannte Zuverlässigkeitsdefinition durch den expliziten Hinweis auf die „definierte Umgebung“.

Dieses Kapitel bietet einen Überblick über Modelle zur Prognose von Softwarezuverlässigkeit, Softwareversagensfällen und Softwarefehlern. Die in Abschnitt 25.2 behandelten Modelle nutzen den Versagensverlauf während der Testphase zur Prognose der Zuverlässigkeit im Nutzungsbetrieb oder zur Prognose der bis zum Testende zu erwartenden Versagensfälle. Abschnitt 25.3 enthält knappe Darstellungen einiger weiterer Modellklassen. In Abschnitt 25.3.1 sind Modelle zusammengefasst, die basierend auf ein oder zwei Test-Stichproben die Zuverlässigkeit der Software prognostizieren oder ihren Fehlergehalt abschätzen. Demgegenüber benötigen die in Abschnitt 25.3.2 beschriebenen Modelle keine Beobachtungen aus der Testphase; sie versuchen, aufgrund von Informationen über das Softwareprodukt und seine Erstellung die Anzahl der Softwarefehler vorherzusagen. Dass in dem vorgegebenen Rahmen die Fragestellungen nur angerissen und die Methoden nur skizziert werden können, liegt auf der Hand. Der interessierte Leser sei deshalb auf weitere Überblicksartikel [3], [4], [12], [16], [19],

[28], [36], [46], [48], [55], [56] sowie auf Bücher [9], [32], [39], [41], [44], [49], [54] zum Themenkomplex verwiesen.

## 25.2 Softwarezuverlässigkeitswachstumsmodelle

Die in diesem Abschnitt vorgestellten Modelle betrachten allesamt die Entwicklung der Anzahl der im Laufe der Integrations- oder Systemtestphase der Softwareentwicklung beobachteten Versagensfälle. Während des Testens wird das Versagen der Software zum Anlass genommen, die ursächlichen Fehler im Code aufzuspüren und zu verbessern, sodass die Zuverlässigkeit der Software sich im Zeitablauf verändert und dabei tendenziell zunimmt. Da die hier beschriebenen Modelle versuchen, diesen Effekt abzubilden, werden sie als „Softwarezuverlässigkeitswachstumsmodelle“ (kurz: SZWM) bezeichnet.

Zur Anwendung (d. h. Schätzung) der Modelle benötigt man entweder die Zeitpunkte des Versagens der Software oder die Anzahl der Versagensfälle innerhalb von Zeitintervallen, welche den gesamten Beobachtungszeitraum partitionieren. Hierbei kann das verwendete Zeitmaß im Prinzip auch die Kalenderzeit sein. Allerdings gehen die meisten der im Folgenden diskutierten Modelle davon aus, dass die Intensität der Programmnutzung im Zeitablauf konstant ist. Deshalb sollte ein Maß zugrunde gelegt werden, für welches dies in etwa zutrifft, z. B. die von den Testern zur Durchführung der Testfälle benötigte Zeit oder die CPU-Ausführungszeit.

Seien  $T_1 < T_2 < \dots$  die Zeitpunkte, zu denen das erste, zweite, ... Versagen auftritt. Obwohl Software deterministisch reagiert – auf exakt identische Eingabewerte unter denselben Nebenbedingungen also immer das gleiche Ergebnis produziert – ist es aus den folgenden Gründen dennoch sinnvoll, die Versagenszeitpunkte als zufälligen Prozess zu modellieren [39], S. 29 f.:

1. Die Irrtümer seitens der Softwareentwickler, die zur Einbringung von Fehlern in das Programm führen, sind nicht mit Sicherheit vorhersehbar. Deshalb sind die Positionen der Fehler in der Software unbekannt.
2. Selbst wenn ein bestimmtes Nutzungsprofil mit vorgegebenen Frequenzen für die Aufrufe der einzelnen Funktionsbereiche zugrunde gelegt wird, ist die exakte Sequenz der Nutzereingaben nicht von vornherein festgelegt.

Die Versagensfälle, die sich aufgrund der Aktivierung der Fehler durch die Programmausführung ergeben, können deshalb mithilfe zufälliger Punktprozesse modelliert werden. Bezeichnet man mit  $M(t)$  die Anzahl der Versagensfälle im Intervall  $(0, t]$ , so handelt es sich bei  $\{M(t), t \geq 0\}$  um einen mit dem Punktprozess verbundenen zufälligen Zählprozess. Ein solcher zeichnet sich dadurch aus, dass er nur Null oder ganzzahlige positive Werte annehmen kann und im Zeitablauf nicht abnehmend ist. Verschiedene Zählprozesse unterscheiden sich darin, wie der Zuwachs in der Zahl der beobachteten Ereignisse erfolgen kann. Stellt man sich die Wertausprägungen  $0, 1, 2, \dots$  als mögliche Zustände des Prozesses vor, so läuft die Frage darauf hinaus, wie die Übergangsraten zwischen diesen Zuständen spezifiziert sind.

Eine große Gruppe von Zählprozessen geht davon aus, dass nicht mehr als ein Ereignis gleichzeitig eintreten kann. Somit ist von einem beliebigen Zustand  $j$  direkt nur der

Zustand  $j+1$  zu erreichen. Solange das Zeitintervall  $(t, t+\Delta t]$  kurz ist, ist es plausibel anzunehmen, dass die Wahrscheinlichkeit für den Wechsel von Zustand  $j$  in den Zustand  $j+1$  proportional zur Intervalllänge  $\Delta t$  ist. Im einfachsten Fall ergibt sich die Übergangswahrscheinlichkeit somit als Produkt von  $\Delta t$  mit einer Konstanten  $r$ , die als Übergangsrate bezeichnet wird.

Gaudoin [14], S. 37 ff., und später unabhängig von ihm Chen und Singpurwalla [10] haben gezeigt, dass sich viele SZWM als Spezialfälle des so genannten selbstanregenden Punktprozesses (*self-exciting point process*) darstellen lassen. In diesem komplizierteren Punktprozess ist die Übergangsrate kein konstanter Wert, sondern sie kann sowohl vom aktuellen Zeitpunkt  $t$  als auch von der gesamten Vorgeschichte des Zählprozesses,  $V_t = \{M(t), T_1, T_2, \dots, T_{M(t)}\}$ , abhängen. Formal ist die Übergangsrate aus dem Zustand  $j$  als Grenzwert definiert (vgl. [14], S. 39, und [50], S. 289):

$$r_j(t, V_t) = \lim_{\Delta t \rightarrow 0} \frac{P(M(t + \Delta t) - M(t) = 1 | V_t = \{M(t) = j, T_1, T_2, \dots, T_j\})}{\Delta t} .$$

Die Annahmen des selbstanregenden Punktprozesses im Kontext eines SZWM können informell folgendermaßen beschrieben werden (vgl. [19]):

1. Zum Zeitpunkt  $t = 0$  ist noch kein Versagensfall eingetreten, d. h.  $M(0) = 0$ .
2. Wenn die Software bis zum Zeitpunkt  $t$  bereits  $j$ -mal versagt hat, dann entspricht die Wahrscheinlichkeit für genau ein Softwareversagen in dem Zeitintervall  $(t, t+\Delta t]$  annähernd dem Produkt aus dessen Länge  $\Delta t$  und der Übergangsrate  $r_j(t, V_t)$ .
3. Die Wahrscheinlichkeit für das Auftreten von mehr als einem Versagensfall in dem Intervall  $(t, t+\Delta t]$  geht mit  $\Delta t \rightarrow 0$  schneller gegen Null als die Wahrscheinlichkeit für das Auftreten von genau einem Versagensfall. Praktisch bedeutet dies, dass in einem sehr kurzen Zeitintervall nicht mehr als ein Versagensfall beobachtet werden kann.

Die Struktur dieses allgemeinen selbstanregenden Punktprozesses ist in Abbildung 1 dargestellt. Hierbei entsprechen die Kreise den Zuständen des Zählprozesses  $M(t)$ , von denen jeweils nur ein Übergang in den nächsthöheren Zustand möglich ist. Die Grafik deutet an, dass sich in manchen Modellen maximal  $u_0$  Versagensfälle einstellen können und somit die Zustände  $u_0+1, u_0+2, \dots$  nicht existieren.

Dem Tester gegenüber äußert sich freilich immer nur diejenige Übergangsrate, die mit dem jeweils aktuellen Zustand verbunden ist. Die Versagensrate der Software, bezeichnet als Programmhazardrate  $z(t, V_t)$ , stellt sich im Zeitablauf also als Aneinanderreihung der Übergangsraten der unterschiedlichen Zustände dar.

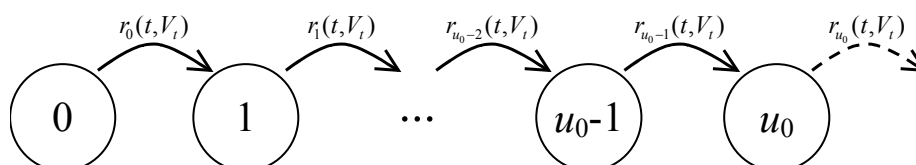


Abb. 1 Struktur der Anzahl der Versagensfälle als selbstanregender Punktprozess

Hierbei kommt zum Zeitpunkt  $t$  die Versagensrate des gegenwärtigen Zustands  $M(t)$  zum Tragen:

$$z(t, V_t) = r_{M(t)}(t, V_t) . \quad (1)$$

Aufgrund ihres stückweisen Aufbaus wird die Programmhazardrate mitunter auch „verkettete Versagensrate“ [10] genannt. Es ist zu beachten, dass die Programmhazardrate zum Zeitpunkt  $t$  eine Funktion der Vorgeschichte bis zu diesem Zeitpunkt,  $V_t$ , ist. Falls diese Vorgeschichte (noch) nicht beobachtet wurde, ist die Programmhazardrate aus zwei Gründen potenziell stochastisch:

- Aus welcher der Übergangsraten  $r_j(t, V_t)$  die Programmhazardrate zum Zeitpunkt  $t$  tatsächlich besteht, hängt vom Zufall ab, nämlich von der Anzahl der bis dahin aufgetretenen Versagensfälle. Nur dann, wenn die Übergangsraten aller Zustände identisch sind, ist dieser Punkt ohne Belang.
- Der Wert einer jeden Übergangsrate  $r_j(t, V_t)$  kann wiederum durch die Prozessvorgeschichte bestimmt sein, die zufällig ist, solange ihre Realisation nicht beobachtet wurde. Jede einzelne Übergangsrate ist nur dann deterministisch, wenn sie lediglich eine Funktion der Zeit  $t$  ist.

Im Folgenden werden wir deshalb die Bezeichnungen  $r_j(t, V_t | V_t)$  und  $z(t, V_t | V_t)$  verwenden, falls diejenige deterministische Übergangs- bzw. Programmhazardrate gemeint ist, die sich unter Kenntnis einer ganz bestimmten Prozessvorgeschichte  $V_t = \{M(t) = m(t), T_1 = t_1, T_2 = t_2, \dots, T_{M(t)} = t_{m(t)}\}$  bis zum Zeitpunkt  $t$  ergibt. Hierbei stellen die Kleinbuchstaben die Realisationen der entsprechenden Zufallsvariablen dar.

Mit der Spezifikation der Übergangsraten  $r_j(t, V_t)$  bzw. der Programmhazardrate  $z(t, V_t)$  legt ein SZWM indirekt zugleich weitere Größen fest, z. B. die so genannte Mittelwertfunktion  $\mu(t) = E(M(t))$ , welche die erwartete Anzahl an Versagensfällen bis zum Zeitpunkt  $t$  angibt, sowie deren Ableitung nach  $t$ , die als Versagensintensität  $\lambda(t)$  bezeichnet wird. Während für  $\mu(t)$  kein allgemeingültiger Zusammenhang mit der Programmhazardrate formuliert werden kann, ist dies für die Versagensintensität möglich. Diese entspricht jeweils der Programmhazardrate, die ohne Kenntnis der Prozessvorgeschichte  $V_t$  zum Zeitpunkt  $t$  erwartet wird,

$$\lambda(t) = E(z(t, V_t)) . \quad (2)$$

Hängt die Programmhazardrate ohnehin nicht von  $V_t$ , sondern nur von der Zeit  $t$  ab, so ist sie mit der Versagensintensität identisch. Die grafische Darstellung von Mittelwertfunktion und Versagensintensität vermittelt einen guten Eindruck davon, wie sich aufgrund der Modellannahmen die Versagensauftritte erwartungsgemäß über die Beobachtungszeit verteilen.

Falls man davon ausgehen kann, dass sich die Programmhazardrate weiter wie vom Modell spezifiziert entwickeln wird, so lässt sich mit ihrer Hilfe auch die zukünftige Zuverlässigkeit der Software bestimmen. Bedingt auf die beobachtete Vorgeschichte  $V_t$  beträgt die Wahrscheinlichkeit für kein Versagen im Zeitintervall  $(t, t+\Delta t]$

$$R(\Delta t | t, V_t) = \exp\left(-\int_t^{t+\Delta t} z(x, V_t | V_t) dx\right) = \exp\left(-\int_t^{t+\Delta t} r_{m(t)}(x, V_t | V_t) dx\right) . \quad (3)$$

Die im Folgenden dargestellten SZWM sind alle Spezialfälle des selbstanregenden Punktprozesses. Ihre Gruppierung erfolgt unter dem Gesichtspunkt, welche Teile der Vorgeschichte  $V_t$  Einfluss auf die Programmhazardrate haben.

### 25.2.1 Markovprozess-Modelle

Die während des Testens aufgetretenen Versagensfälle geben den Anstoß zur Beseitigung der diese verursachenden Fehler. Eine idealisierte Modellierung des Testprozesses geht davon aus, dass jede Fehlerkorrektur unmittelbar nach dem beobachteten Versagen erfolgt und erst im Anschluss daran die Testausführung und die Zeitnahme fortgesetzt werden. (Obwohl dies in den allerwenigsten Fällen der Realität entspricht, kann man sich dieser modellhaften Situation annähern, indem man für jeden Fehler nur das erste von ihm bewirkte Versagen berücksichtigt und somit im weiteren Testverlauf so tut, als ob der Fehler nicht mehr im Code vorhanden wäre.) Unter diesen Voraussetzungen ist es plausibel anzunehmen, dass die Versagensrate eines jeden Zustands von eben diesem Zustand selbst abhängt, da auf jeden Versagenseintritt eine sofortige Änderung des Fehlergehalts der Software folgt. Bei den in diesem Abschnitt besprochenen Modellen ist der aktuelle Prozesszustand  $M(t)$  die *einzig*e Information der Prozessvorgeschichte  $V_t$ , die Einfluss auf die Übergangsraten hat. Es handelt sich deshalb um zeitstetige Markovprozess-Modelle. (Mit zeitdiskreten Markovprozess-Modellen befasst sich Kapitel 16 dieses Buches.) Manche der Modelle gehen zudem davon aus, dass die Übergangsraten auch von der momentanen Beobachtungszeit  $t$  selbst (die nicht Teil der Prozessvorgeschichte ist) abhängen.

Eine besondere Unterklasse bilden die Binomialmodelle (s. Musa und andere [41], Kapitel 10.3 und 11.1.1, sowie Shantikumar [47]). Diese Modelle nehmen an, dass sich zu Testbeginn eine bestimmte Zahl von Fehlern,  $u_0$ , in der Software befindet. Zudem unterstellen sie, dass in Bezug auf die Tendenz, ein Versagen zu verursachen, alle Softwarefehler zu jedem Zeitpunkt im Durchschnitt gleich gefährlich sind. (Wie hier beziehen wir im Folgenden die „Gefährlichkeit“ eines Fehlers einzig auf seine Äußerungsrate und nicht auf den im Falle seiner Aktivierung verursachten Schaden.) Jeder Fehler weist also die gleiche Hazardrate  $z_a(t)$  auf, die zwar von der Zeit, nicht aber von der Prozessvorgeschichte abhängen kann. Somit ist für jeden einzelnen Fehler die Wahrscheinlichkeit dafür, zum Zeitpunkt  $t$  oder früher zu einem Versagen geführt zu haben, identisch

$$F_a(t) = 1 - \exp\left(-\int_0^t z_a(x) dx\right).$$

Da die Binomialmodelle zudem davon ausgehen, dass die Fehlerkorrektur perfekt erfolgt, ergibt sich die Programmhazardrate zu jedem Zeitpunkt als Produkt der Anzahl der noch verbliebenen Fehler,  $u_0 - M(t)$ , mit dem Beitrag eines einzelnen Fehlers,  $z_a(t)$ . Sie hängt deshalb von der Prozessvorgeschichte  $V_t$  ausschließlich über  $M(t)$  ab:

$$z(t, V_t) = z(t, M(t)) = r_{M(t)}(t, M(t)) = (u_0 - M(t))z_a(t). \quad (4)$$

Da nach dem  $u_0$ -ten Versagen auch der letzte Fehler behoben wird, beträgt die Übergangsrate aus dem Zustand  $u_0$  gleich Null. Deshalb kann der Zählprozess  $M(t)$  die Zustände  $u_0+1, \dots$  nicht einnehmen.

Auch die Mittelwertfunktion lässt sich für die Binomialmodelle in intuitiv eingängiger Weise darstellen. Sie ist die Wahrscheinlichkeit eines jeden Fehlers, bereits ein Versagen verursacht zu haben, multipliziert mit der Anzahl der ursprünglichen Softwarefehler,

$$\mu(t) = u_0 F_a(t) = u_0 \left[ 1 - \exp\left(-\int_0^t z_a(x) dx\right) \right]. \quad (5)$$

Wurde also z. B. zu einem bestimmten Zeitpunkt jeder Fehler mit 50-prozentiger Wahrscheinlichkeit entdeckt, so kann man erwarten, dass dies im Schnitt tatsächlich bei der Hälfte der ursprünglichen Fehler eingetreten ist.

### Jelinski-Moranda-Modell

Das einfachste Binomialmodell wurde 1972 von Jelinski und Moranda [22] vorgeschlagen – als eines der ersten SZWM überhaupt. Es geht nicht nur davon aus, dass alle Fehler zu jedem Zeitpunkt gleich gefährlich sind, sondern unterstellt zudem eine konstante fehlerbezogene Hazardrate  $z_a(t)$  von  $\phi$ . Aus der für Binomialmodelle allgemein geltenden Gleichung (4) folgt deshalb:

$$z(t, V_t) = z(t, M(t)) = r_{M(t)}(t, M(t)) = (u_0 - M(t))\phi. \quad (6)$$

Unter steter Kenntnis der Prozessvorgeschichte  $V_t$  nimmt die Programmhazardrate damit einen treppenförmigen Verlauf an: Sofort nach einem Versagensfall sinkt sie exakt um den Betrag  $\phi$  (die Hazardrate des soeben entdeckten und korrigierten Fehlers) und ist dann bis zum nächsten Versagen ein konstanter Wert. Im linken Teil von Abbildung 2 ist diese typische Entwicklung von  $z(t, V_t | V_t)$  dargestellt.

Für die Mittelwertfunktion ergibt sich aus Gleichung (5)

$$\mu(t) = u_0 \left[ 1 - \exp\left(-\int_0^t \phi dx\right) \right] = u_0 [1 - \exp(-\phi t)]. \quad (7)$$

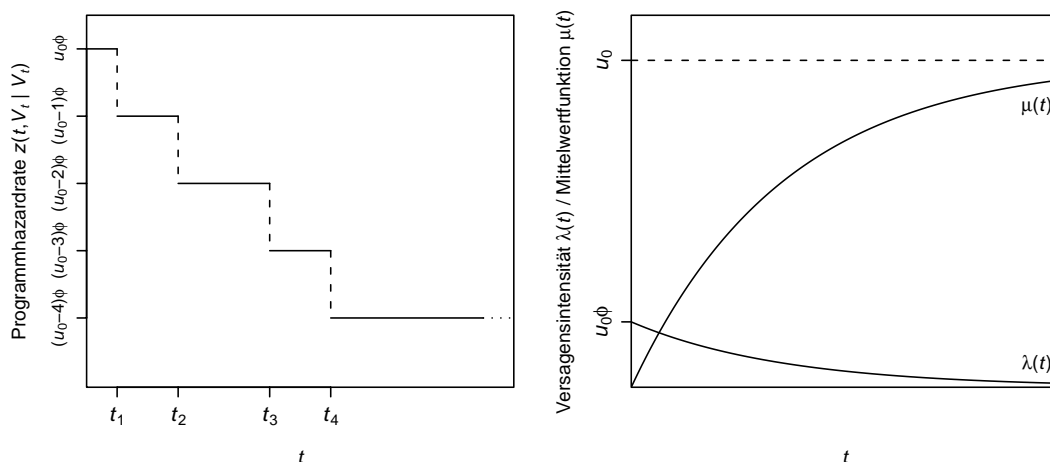


Abb. 2 Typische Entwicklung von Programmhazardrate, Versagensintensität und Mittelwertfunktion gemäß dem Jelinski-Moranda-Modell

Da sich dieser Ausdruck mit wachsendem  $t$  immer mehr an  $u_0$  annähert, ist aufgrund des Modells zu erwarten, dass nach einer ausreichenden Testdauer alle Softwarefehler gefunden und behoben sein werden.

Eine für das Jelinski-Moranda-Modell typische Mittelwertfunktion und ihre Ableitung, die Versagensintensität

$$\lambda(t) = u_0 \phi \exp(-\phi t),$$

zeigt das rechte Diagramm von Abbildung 2. Die Grafik macht noch einmal deutlich, dass die Versagensintensität als eine von der konkreten Prozessvorgeschichte unbeeinflusste *erwartete* Programmhazardrate eine stetige Funktion der Zeit ist, die kontinuierlich abnimmt.

Wurde die Software bis zum Zeitpunkt  $t$  getestet, wobei sie  $m(t)$ -mal versagte, so beträgt ihre Zuverlässigkeit im Intervall  $(t, t+\Delta t]$  unter Anwendung von Gleichung (3)

$$R(\Delta t | t, V_t) = \exp\left(-\int_t^{t+\Delta t} z(x, V_t | V_t) dx\right) = \exp(-(u_0 - m(t))\phi \Delta t).$$

Offensichtlich ist der Zuverlässigkeitswert (als Wahrscheinlichkeit einer versagensfreien Programmnutzung) um so geringer, je mehr Fehler sich noch in der Software befinden, je größer die Hazardrate eines einzelnen Fehlers ist und je länger die Zeitspanne  $\Delta t$  ist, in der die Software verwendet wird.

Um für ein konkretes Softwareprodukt die erwartete Anzahl von Versagensfällen oder die Zuverlässigkeit prognostizieren zu können, müssen die beiden Parameter  $u_0$  und  $\phi$  geschätzt werden. Die Anwendung der Maximum-Likelihood-Methode führt zu einem System aus zwei nichtlinearen Gleichungen [12], die leicht simultan gelöst werden können. Allerdings hat sich gezeigt, dass die so gewonnenen Schätzer unschöne Eigenschaften besitzen [13], [31]. Littlewood und Sofer entwickeln deshalb eine Bayes-Erweiterung des Jelinski-Moranda-Modells, deren Prognosequalität allerdings nur wenig besser ist (s. [1]). Dies spricht dafür, dass die Probleme nicht nur vom Schätzverfahren, sondern auch von den simplen Modellannahmen selbst verursacht werden. Insbesondere die beiden folgenden Einwände werden oft erhoben:

1. Es ist unrealistisch anzunehmen, dass alle Fehler die gleiche Hazardrate besitzen. Tatsächlich ist es plausibler davon auszugehen, dass diejenigen Fehler, die früh im Test ein Versagen verursachen, im Hinblick auf ihre Äußerungsrate auch gefährlicher waren. Deshalb sollten die Sprünge, welche die Programmhazardrate nach den Versagensfällen macht, im Zeitablauf tendenziell abnehmen.
2. Obwohl nach jedem Versagensfall der Versuch unternommen wird, den für ihn kausalen Fehler zu korrigieren, gibt es in Wirklichkeit keine Gewähr dafür, dass die Bereinigung gelingt und keine neuen Fehler in die Software eingebracht werden. Deshalb könnte sich nach einem Versagenseintritt die Programmqualität sogar verschlechtern. Ein SZWM sollte diese Möglichkeit zulassen.

Viele der komplexeren Modelle entstanden als Antwort auf diese Kritikpunkte.

### Moranda-Modell

So berücksichtigt z. B. Moranda [38] in seinem so genannten „geometrischen Modell“ den ersten Einwand, indem er unterstellt, dass die Programmhazardrate nach jedem Softwareversagen auf das  $k$ -fache des vorherigen Wertes sinkt, wobei  $k$  zwischen Null und Eins liegt. Zu Testbeginn beträgt die Programmhazardrate  $D$ . Es gilt also:

$$z(t, V_t) = z(t, M(t)) = r_{M(t)}(t, M(t)) = k \cdot r_{M(t)-1}(t, M(t)-1) = Dk^{M(t)}.$$

Wie die typische Entwicklung der Programmhazardrate in Abhängigkeit von der Vorgeschichte in der linken Grafik von Abbildung 3 zeigt, nimmt unter diesen Annahmen der Umfang der Reduzierung der Programmhazardrate mit jeder Fehlerkorrektur ab. Die Ermittlung von Mittelwertfunktion und Versagensintensität gestaltet sich für dieses Modell schwierig. Musa und andere [41], S. 572, leiten die Näherungen

$$\mu(t) \approx -\frac{1}{\ln(k)} \ln\left(1 - \frac{D}{k} \ln(k)t\right) \quad \text{und} \quad \lambda(t) \approx \frac{D}{k - D \ln(k)t}$$

her. Die exakten aber unhandlichen Ausdrücke, deren Werte sich für bestimmte Wertkombinationen für  $D$  und  $k$  deutlich von den Approximationen unterscheiden können, werden von Boland und Singh [5] angegeben.

Auf der rechten Seite von Abbildung 3 sind typische Verläufe der angenäherten Mittelwertfunktion und Versagensintensität angetragen. Im Vergleich mit dem Jelinski-Moranda-Modell wird deutlich, dass die Mittelwertfunktion im Zeitablauf nicht gegen einen festen Wert strebt, sondern die erwartete Anzahl der Versagensfälle beliebig groß werden kann, wenn die Testphase entsprechend lange dauert. Während es unplausibel erscheint, eine unendliche Anzahl von ursprünglichen Softwarefehlern anzunehmen, liegt eine mögliche Interpretation dieser Mittelwertfunktion darin, dass die verminderte Wirkung der Korrekturen auf die Programmhazardrate nicht nur durch die geringere Gefährlichkeit der später entdeckten Fehler verursacht wird. Vielmehr werden zudem neue Fehler in die Software eingebracht, was in diesem Modell jedoch in keinem Fall zu einem Anwachsen der Programmhazardrate führt.

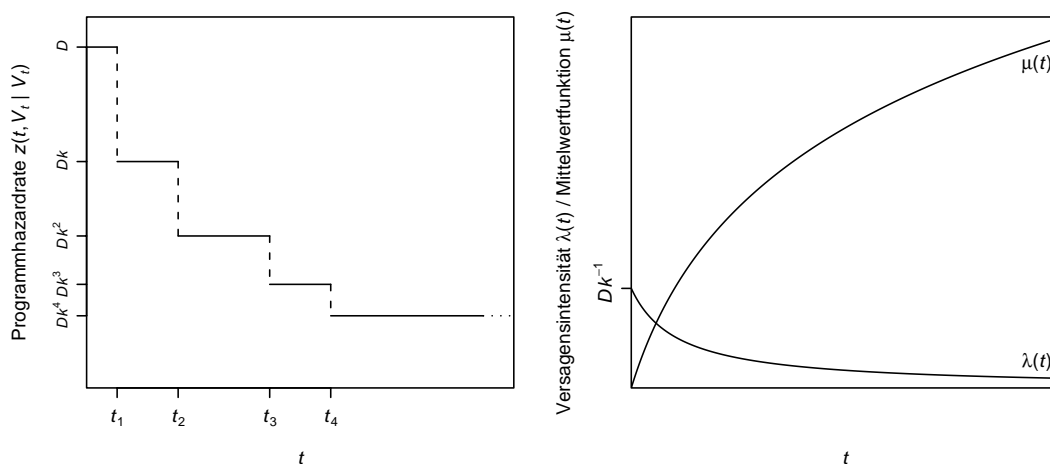


Abb. 3 Typische Entwicklung von Programmhazardrate, approximierter Versagensintensität und approximierter Mittelwertfunktion gemäß dem Moranda-Modell

Nachdem bis zum Zeitpunkt  $t$  insgesamt  $m(t)$  Versagensfälle beobachtet wurden, errechnet sich die Zuverlässigkeit des Programms gemäß Gleichung (3) als

$$R(\Delta t | t, V_t) = \exp\left(-\int_t^{t+\Delta t} z(x, V_t | V_t) dx\right) = \exp(-Dk^{m(t)} \Delta t).$$

Auch im Moranda-Modell resultiert die Parameterschätzung nach der Maximum-Likelihood-Methode in einem System aus zwei gemeinsam zu lösenden Gleichungen [12].

**Littlewood-Modell**

Der nachlassende Effekt jeder weiteren Fehlerkorrektur auf die Programmhazardrate lässt sich auch im Rahmen eines Binomialmodells darstellen. An das Jelinski-Moranda-Modell anknüpfend geht Littlewood [29] zwar davon aus, dass jeder Fehler eine über die Zeit hinweg konstante Hazardrate  $\phi$  besitzt. Diese hat jedoch nicht für alle Fehler den gleichen fixen Wert. Stattdessen handelt es sich bei den Hazardraten der verschiedenen Fehler um Zufallszüge aus einer Gamma( $\alpha, \beta$ )-Verteilung. Für die Wahrscheinlichkeit, dass ein bestimmter Fehler bis zum Zeitpunkt  $t$  noch kein Versagen ausgelöst hat, folgt aus diesen Annahmen die Verteilungsfunktion einer Paretoverteilung,

$$F_a(t) = 1 - \left(\frac{\beta}{\beta + t}\right)^\alpha, \tag{8}$$

deren Hazardrate

$$z_a(t) = \frac{\alpha}{\beta + t} \tag{9}$$

beträgt. Diese im Zeitablauf abnehmende Hazardrate scheint im Widerspruch zu der postulierten konstanten Hazardrate  $\phi$  zu stehen. Sie erklärt sich dadurch, dass ein Fehler mit hoher konstanter Hazardrate erwartungsgemäß früher zu einem Versagen führt. Im Umkehrschluss ist die Gefährlichkeit eines Fehlers tendenziell um so geringer, je länger er bereits unentdeckt geblieben ist. Genau dies spiegelt sich in (9) wider.

Da es sich bei dem Littlewood-Modell um ein Binomialmodell handelt, ist die Programmhazardrate zum Zeitpunkt  $t$  das Produkt aus der Anzahl der noch verbliebenen Fehler und der augenblicklichen fehlerbezogenen Hazardrate,

$$z(t, V_t) = z(t, M(t)) = r_{M(t)}(t, M(t)) = (u_0 - M(t))z_a(t) = (u_0 - M(t)) \frac{\alpha}{\beta + t}.$$

Die Darstellung einer typischen Entwicklung der Programmhazardrate in Abbildung 4 zeigt, dass diese anders als bei den bisher vorgestellten Modellen sich nicht nur zu den Versagens- und Fehlerkorrekturzeitpunkten sprunghaft vermindert, sondern zudem *zwischen* ihnen kontinuierlich absinkt. Natürlich ist diese Eigenschaft eine Folge der fehlerbezogenen Hazardrate; die Interpretation fällt denn auch analog aus: Je länger die Software bereits getestet wurde, desto mehr unterstützt dies die subjektive Einschätzung, dass die Gefährlichkeit der noch nicht behobenen Fehler gering und die Programmqualität somit hoch ist.

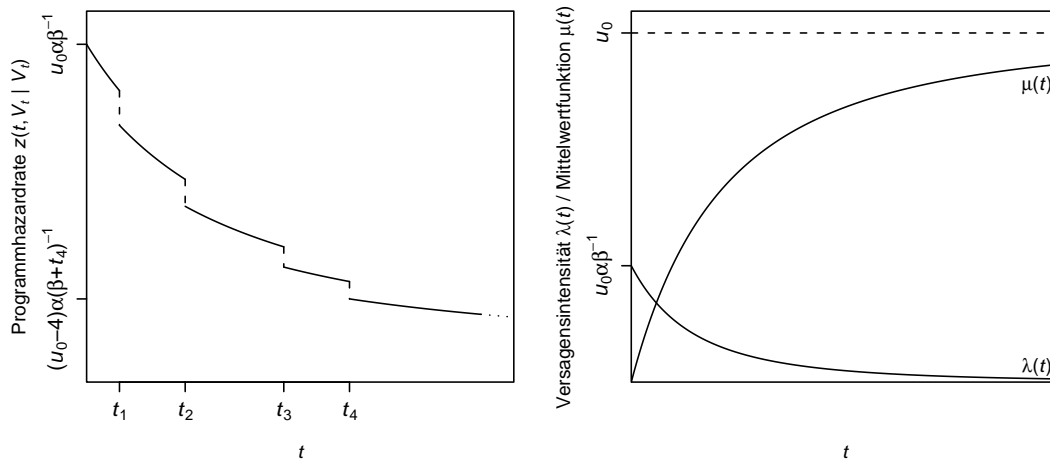


Abb. 4 Typische Entwicklung von erwarteter Programmhazardrate, Versagensintensität und Mittelwertfunktion gemäß dem Littlewood-Modell

Die sich aus der für Binomialmodelle geltende Gleichung (5) in Verbindung mit der Verteilungsfunktion (8) ergebende Mittelwertfunktion

$$\mu(t) = u_0 F_a(t) = u_0 \left[ 1 - \left( \frac{\beta}{\beta + t} \right)^\alpha \right]$$

und die mit ihr verbundene Versagensintensität

$$\lambda(t) = u_0 \frac{\alpha}{\beta + t} \left( \frac{\beta}{\beta + t} \right)^\alpha$$

sind im rechten Teil von Abbildung 4 angetragen. Waren zum Zeitpunkt  $t$  bereits  $m(t)$  Versagensfälle aufgetreten, beträgt die Zuverlässigkeit im Intervall  $(t, t + \Delta t]$  dem Littlewood-Modell gemäß

$$R(\Delta t | t, V_t) = \exp\left(-\int_t^{t+\Delta t} z(x, V_t | V_t) dx\right) = \left( \frac{\beta + t}{\beta + t + \Delta t} \right)^{(u_0 - m(t))\alpha}.$$

Die Anwendung von Gleichung (3) impliziert hierbei die Erwartung, dass sich die Programmhazardrate in der Zukunft bis zu einem Versagenseintritt gemäß der Übergangsrate des aktuellen Zustands  $m(t)$  verändern und somit kontinuierlich absinken wird. Diese Annahme ist selbst dann vernünftig, wenn die Prognose am Ende der Testphase erfolgt und sich auf die Nutzungsphase bezieht, in der typischerweise kaum Verbesserungen an der Software vorgenommen werden; denn die stetige Verkleinerung der Übergangsrate ist in diesem Modell – wie gezeigt – nur subjektiv begründet und nicht durch eine erwartete Fehlerkorrektur bedingt.

Die drei Parameter des Littlewood-Modells ( $u_0$ ,  $\alpha$  und  $\beta$ ) können wiederum mittels der Maximum-Likelihood-Methode geschätzt werden [29].

### 25.2.2 Ein Semi-Markovprozess-Modell: Littlewood-Verrall-Modell

Bei den im letzten Abschnitt vorgestellten Modellen hängen die einzelnen Übergangsraten  $r_j(t, V_t)$  ausschließlich vom jeweiligen Zustand  $M(t)$  und ggf. vom aktuellen Zeitpunkt  $t$  ab, weshalb sie zu den Markovprozess-Modellen zählen. Es ist jedoch durchaus denkbar, dass weitere Teile der Prozessvorgeschichte  $V_t$  die Übergangsraten beeinflussen, z. B. der Zeitpunkt des letzten Versagens,  $T_{M(t)}$ .

In einem Semi-Markovprozess-Modell sind die Übergangsraten Funktionen der Zeitspanne *seit* der letzten Versagensbeobachtung, der Differenz  $t - T_{M(t)}$ . ( $T_0$  wird hierbei gleich Null gesetzt und ist kein Versagenszeitpunkt, sondern der Testbeginn.) Der Name des Modells rührt daher, dass das künftige Verhalten des Zählprozesses zur Zeit  $T_0 \equiv 0$  und zu den Versagenszeitpunkten  $T_1, T_2, \dots$  lediglich vom Zustand  $M(t)$  bestimmt wird, so wie dies in einem Markovprozess-Modell der Fall ist.

Das bedeutendste SZWM aus der Klasse der Semi-Markovprozess-Modelle wurde bereits 1973 von Littlewood und Verrall [30] entwickelt. Die beiden Autoren gehen davon aus, dass die Übergangsrate aus dem  $j$ -ten Zustand konstant  $\phi_j$  beträgt. Allerdings ist  $\phi_j$  kein bekannter, deterministischer Wert. Vielmehr handelt es sich um die Realisation einer Zufallsvariablen  $\Phi_j$ , die einer Gammaverteilung mit den Parametern  $\alpha$  und  $\psi(j)$  folgt. Da der Parameter  $\psi(j)$  eine Funktion des Zustands  $j$  ist, beeinflusst dieser die Verteilung der Übergangsrate. Hierbei wird die Funktion  $\psi(j)$  so gewählt, dass für jedes beliebige  $x$  gilt:

$$P(\Phi_j \leq x) \geq P(\Phi_{j-1} \leq x). \tag{10}$$

Unabhängig von  $x$  ist es also immer wahrscheinlicher, dass  $\Phi_j$  diesen Wert unterschreitet, als dass dies auf die Übergangsrate aus dem vorherigen Zustand,  $\Phi_{j-1}$ , zutrifft. In einem ganz spezifischen Sinne verbessert sich somit die Softwarequalität erwartungsgemäß nach jeder Fehlerkorrektur. Man kann zeigen, dass die Gültigkeit von Gleichung (10) dann gewährleistet ist, wenn  $\psi(j)$  eine streng monoton steigende Funktion von  $j$  ist. Im Folgenden werden wir diejenige Funktion betrachten, die in der Diskussion und Anwendung des Littlewood-Verrall-Modells die größte Aufmerksamkeit erfahren hat,  $\psi(j) = \beta_0 + \beta_1(j+1)$ .

Ähnlich wie im Littlewood-Modell für die fehlerbezogene Hazardrate ergibt sich hier für die Übergangsrate eines jeden Zustands ein Ausdruck, der im Zeitablauf absinkt,

$$r_j(t, V_t) = r_j(t, j, T_j) = \frac{\alpha}{\beta_0 + \beta_1(j+1) + (t - T_j)},$$

wobei jeweils die Zeitspanne seit dem Betreten des aktuellen Zustands,  $t - T_j$ , von Bedeutung ist. Auch die Erklärung dieser Eigenschaft ist analog zu derjenigen, die uns aus dem Littlewood-Modell bekannt ist: Je mehr Zeit bereits seit dem letzten beobachteten Versagen verstrichen ist, desto plausibler ist es, dass die augenblicklich geltende Übergangsrate einen niedrigen Wert besitzt. Die Programmhazardrate ist aus solchen kontinuierlich fallenden Übergangsraten stückweise zusammengesetzt:

$$z(t, V_t) = z(t, M(t), T_{M(t)}) = \frac{\alpha}{\beta_0 + \beta_1(M(t)+1) + (t - T_{M(t)})}. \tag{11}$$

Der im linken Diagramm der Abbildung 5 dargestellte beispielhafte Verlauf lässt allerdings ein Phänomen erkennen, welches in keinem der bislang diskutierten Modelle präsent war: Unter bestimmten Umständen kann die Programmhazardrate direkt nach einem Versagensfall höher sein als kurz zuvor. Wie Gleichung (11) zeigt, ist dies mathematisch gesehen nach dem  $j$ -ten Versagen genau dann der Fall, wenn die Zeitspanne seit dem  $(j - 1)$ -ten Versagen so lange war, dass  $t_j - t_{j-1}$  größer ist als  $\beta_1$ . Das durch das fehlerhafte Verhalten der Software zerstörte Vertrauen in deren Zuverlässigkeit überwiegt dann die durch die Fehlerkorrektur erwartete Verbesserung. Mitunter wird diese Modelleigenschaft auch dahingehend interpretiert, dass das Littlewood-Verrall-Modell die Möglichkeit zusätzlicher im Rahmen einer Korrekturmaßnahme eingebrachter Fehler berücksichtigt und somit auf den zweiten wichtigen Einwand gegen das Jelinski-Moranda-Modell eingeht (so z. B. von Gaudoin [14], S. 62 f.).

Mittelwertfunktion und Versagensintensität lassen sich für dieses Modell wiederum nicht so einfach herleiten wie für die Binomialmodelle. Musa und andere [41], S. 295 f., ermitteln die Näherungen

$$\mu(t) \approx \frac{1}{\beta_1} \left( \sqrt{\beta_0^2 + 2\alpha\beta_1 t} - \beta_0 \right) \quad \text{und} \quad \lambda(t) \approx \frac{\alpha}{\sqrt{\beta_0^2 + 2\alpha\beta_1 t}}.$$

Die approximierte Mittelwertfunktion und ihr typischer Verlauf in Abbildung 5 zeigen, dass bei einer unendlich langen Testdauer auch unendlich viele Versagensfälle zu erwarten sind. Auch deshalb ist das Modell potenziell für Situationen geeignet, in denen bei der Fehlerverbesserung neue Fehler entstehen.

Wurde die Software bis zum Zeitpunkt  $t$  getestet, wobei der letzte der  $m(t)$  beobachteten Versagensfälle zum Zeitpunkt  $t_{m(t)}$  aufgetreten war, dann beträgt unter Anwendung von Gleichung (3) die Zuverlässigkeit im Intervall  $(t, t+\Delta t]$

$$R(\Delta t | t, V_t) = \exp\left(-\int_t^{t+\Delta t} z(x, V_t | V_t) dx\right) = \left( \frac{\beta_0 + \beta_1(m(t) + 1) + t - t_{m(t)}}{\beta_0 + \beta_1(m(t) + 1) + t - t_{m(t)} + \Delta t} \right)^\alpha.$$

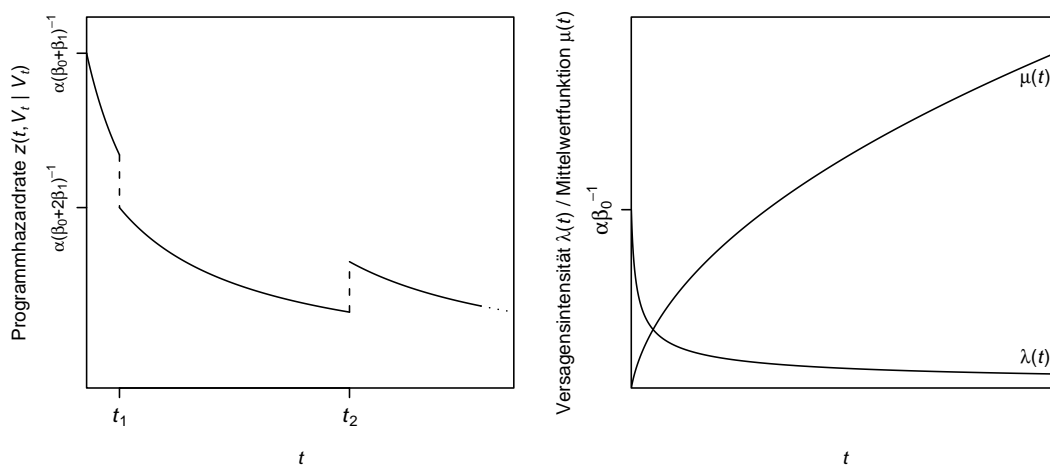


Abb. 5 Typische Entwicklung von erwarteter Programmhazardrate, approximierter Versagensintensität und approximierter Mittelwertfunktion gemäß dem Littlewood-Verrall-Modell

Littlewood und Verrall [30] berücksichtigen die Unsicherheit bzgl. des Parameters  $\alpha$ , indem sie für ihn eine so genannte uninformative a priori Verteilung unterstellen. Die Parameter der Funktion  $\psi(j)$  maximieren sie dagegen durch Optimierung einer Metrik, welche die Güte der Anpassung misst. Diesem hybriden Vorgehen stellen Mazzuchi und Soyer [37] ein geschlossenes Bayes-Verfahren gegenüber, welches für jeden der drei Modellparameter eine a priori Verteilung spezifiziert.

### 25.2.3 Nichthomogene Poissonprozess-Modelle

Im Vergleich zu den in Abschnitt 25.2.1 besprochenen Modellen werden im Littlewood-Verrall-Modell die Übergangsraten von einer zusätzlichen Information beeinflusst, nämlich dem Zeitpunkt des letzten Zustandswechsels. Demgegenüber hängen die Übergangsraten bei den nun vorgestellten nichthomogenen Poissonprozess-Modellen (kurz: NHPP-Modellen) *überhaupt nicht* von der Vorgeschichte ab; sie sind also lediglich Funktionen der Zeit  $t$ ,  $r_j(t)$ . Da noch nicht einmal der aktuelle Zustand selbst einen Einfluss ausübt, müssen die einzelnen Übergangsraten  $r_j(t)$  zudem identisch sein. Die Programmhazardrate ergibt sich durch Aneinanderreihung von Teilstücken der Übergangsraten, s. Gleichung (1). Wenn diese Übergangsraten sich wie hier aber gar nicht unterscheiden, dann ist auch die Programmhazardrate nur eine Funktion der Zeit und identisch mit den Übergangsraten. Es gilt demnach:

$$z(t) = r_i(t) = r_j(t) \text{ für alle } i, j = 0, 1, 2, \dots$$

Dass die Programmhazardrate nach dem Auftreten eines Versagens und der Korrektur des ihn verursachenden Fehlers nicht schlagartig einen anderen Wert annimmt, scheint wenig realistisch zu sein. Eine mögliche Begründung dieser Annahme, die Ascher und Feingold [2], S. 51, als „minimale Reparatur“ bezeichnen, liegt darin, dass ein Computerprogramm sehr viele relativ unbedeutende Fehler enthält; die Bereinigung eines einzelnen dieser Fehler hat somit kaum einen Einfluss auf die Programmhazardrate [14], S. 64. Aber auch die kontinuierliche Veränderung der Programmhazardrate zwischen den Versagensauftritten muss erklärt werden. Zwei Ansätze kommen in Betracht:

1. Wie im Littlewood- und im Littlewood-Verrall-Modell sind diese Modifikationen ausschließlich subjektiver Natur: Das erhöhte Vertrauen in die Qualität der Software kann die Programmhazardrate zwischen zwei Versagensfällen sinken lassen; die erwartete Zunahme der Fähigkeit der Tester, Fehler aufzuspüren, mag dagegen für ihren Anstieg sorgen.
2. Die Fehlerkorrekturen finden nicht zwangsweise sofort nach dem beobachteten Softwareversagen, sondern u. U. erst in der Folgezeit statt. Somit kann es auch zwischen den Versagensauftritten zu leichten Änderungen der Programmhazardrate kommen.

Anders als in dem zu Beginn des Abschnitts 25.2 dargestellten allgemeinen Fall ist die Programmhazardrate völlig unabhängig von der zufälligen Prozessvorgeschichte und somit auch nicht stochastisch. Ihr Erwartungswert, die Versagensintensität, fällt deshalb mit der Programmhazardrate selbst zusammen:

$$\lambda(t) = E(z(t, V_t)) = E(z(t)) = z(t). \quad (12)$$

Zur vollständigen Spezifikation eines NHPP-Modells genügt es deshalb, die Versagensintensität oder ihr Integral, die Mittelwertfunktion  $\mu(t)$ , anzugeben. Die Bezeichnung der Modellklasse rührt daher, dass die Anzahl der Versagensfälle im Intervall  $(t, t+\Delta t]$  Poisson-verteilt ist mit Erwartungswert  $\mu(t+\Delta t) - \mu(t)$ . Das Modell ist insofern nichthomogen, als bei vorgegebener Intervalllänge der Erwartungswert von der Lage des Intervalls abhängt, d. h. von seinem Startzeitpunkt  $t$ .

Aus Gleichung (3) in Verbindung mit Gleichung (12) ergibt sich für die Zuverlässigkeit der Software in der Zeitspanne  $(t, t+\Delta t]$

$$R(\Delta t | t, V_t) = \exp\left(-\int_t^{t+\Delta t} \lambda(x) dx\right) = \exp\{-[\mu(t+\Delta t) - \mu(t)]\}. \quad (13)$$

Dies ist identisch mit der Wahrscheinlichkeit dafür, dass eine Poisson-verteilte Zufallsvariable mit Erwartungswert  $\mu(t+\Delta t) - \mu(t)$  den Wert Null annimmt. Gleichung (13) wird generell für die Zuverlässigkeitsprognose auch in der Nutzungsphase eingesetzt. Dies impliziert die Annahme, dass sich die innerhalb der Testphase erwarteten weiteren Veränderungen der Programmhazardrate auch nach dem Software-Release fortschreiben lassen, z. B. weil sie rein subjektiver Natur sind. Yang und Xie [59] stellen diesem Ansatz eine Berechnung der Zuverlässigkeit im Nutzungsbetrieb gegenüber, bei der die Programmhazardrate konstant belassen wird.

### **Goel-Okumoto-Modell**

Das einfachste NHPP-Modell erhält man, wenn man für die Versagensintensität (und damit zugleich für die Programmhazardrate) die Form der Versagensintensität im Jelinski-Moranda-Modell wählt,

$$\lambda(t) = z(t) = N\phi \exp(-\phi t),$$

womit die Mittelwertfunktion

$$\mu(t) = N[1 - \exp(-\phi t)]$$

beträgt. Der linke Teil von Abbildung 6 zeigt typische Verläufe der beiden Funktionen. Anders als im Jelinski-Moranda-Modell ist in diesem von Goel und Okumoto [15] eingeführten Modell die Anzahl der ursprünglich in der Software enthaltenen Fehler wie auch die Anzahl der bei unendlichem Testaufwand auftretenden Versagensfälle kein fixer Wert  $u_0$ . Vielmehr sind beide Größen Poisson-verteilt, und ihr Erwartungswert entspricht dem Parameter  $N$ . Somit kann der Zählprozess  $M(t)$  also durchaus die Zustände  $N+1, N+2, \dots$  annehmen; es können also mehr als  $N$  Versagensfälle auftreten.

Unter der Annahme, dass man die Programmhazardrate des Modells in die Zukunft extrapolieren kann, ergibt sich für die Zuverlässigkeit im Intervall  $(t, t+\Delta t]$ :

$$R(\Delta t | t, V_t) = \exp\left(-\int_t^{t+\Delta t} \lambda(x) dx\right) = \exp\{N[\exp(-\phi(t+\Delta t)) - \exp(-\phi t)]\}.$$

Aufgrund der unterschiedlichen Verteilungsannahmen führt die Parameterschätzung mittels der Maximum-Likelihood-Methode zu anderen Schätzern als für das Jelinski-Moranda-Modell, s. [15].

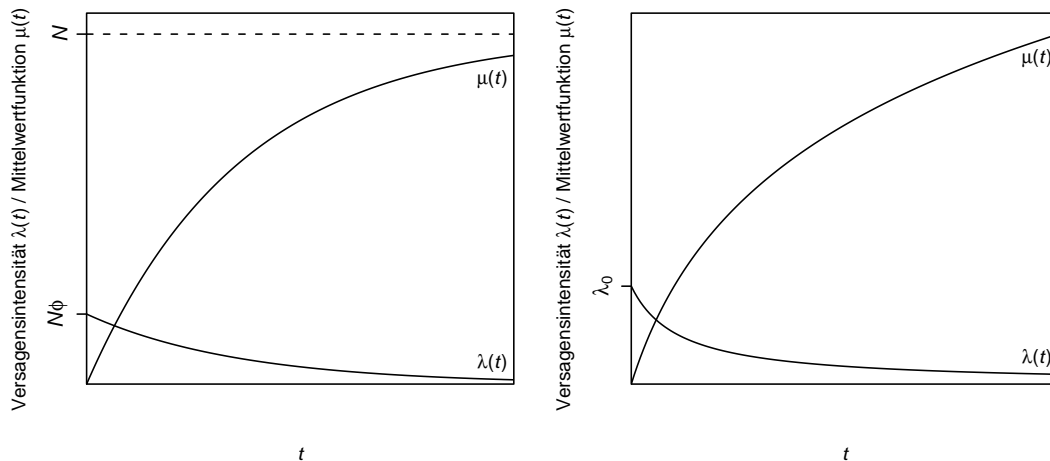


Abb. 6 Typische Entwicklung von Versagensintensität und Mittelwertfunktion gemäß dem Goel-Okumoto-Modell (links) und dem Musa-Okumoto-Modell (rechts)

**Musa-Okumoto-Modell**

Ähnlich wie Moranda gehen Musa und Okumoto [40] von einer Versagensintensität aus, die zu Testbeginn schneller abnimmt als in einer späteren Testphase.

Genauer unterstellen sie, dass die Versagensintensität mit der *erwarteten* Anzahl der aufgetretenen Versagensfälle exponentiell absinkt:

$$\lambda(t) = \lambda_0 \exp(-\theta\mu(t)).$$

Berücksichtigt man, dass die Versagensintensität die Ableitung der Mittelwertfunktion ist, so erhält man eine Differentialgleichung, deren Auflösung zu der Mittelwertfunktion

$$\mu(t) = \frac{1}{\theta} \ln(\lambda_0\theta t + 1)$$

führt. Wie im Moranda-Modell strebt auch hier die Mittelwertfunktion nicht gegen einen festen Wert. Die Anzahl der bei unendlichem Testaufwand erwarteten Versagensbeobachtungen ist also ebenfalls unendlich. Im rechten Diagramm der Abbildung 6 ist der Verlauf dieser Mittelwertfunktion und ihrer Ableitung, der Versagensintensität

$$\lambda(t) = \frac{\lambda_0}{\lambda_0\theta t + 1},$$

beispielhaft dargestellt. Bei Fortschreiben der Versagensintensität folgt für die Zuverlässigkeit im Intervall  $(t, t+\Delta t]$ :

$$R(\Delta t | t, V_t) = \exp\left(-\int_t^{t+\Delta t} \lambda(x) dx\right) = \left(\frac{\lambda_0\theta t + 1}{\lambda_0\theta(t + \Delta t) + 1}\right)^{1/\theta}.$$

Aus den beobachteten Versagensdaten kann man zunächst durch Maximierung der bedingten Likelihoodfunktion (unter der Bedingung, dass bis zum Beobachtungsende  $t_e$  insgesamt  $m(t_e)$  Versagensfälle aufgetreten sind) einen Schätzer für das Produkt  $\lambda_0\theta$  gewinnen und danach den Parameter  $\theta$  separat schätzen. Aufgrund der Invarianzei-

genschaft der Maximum-Likelihood-Schätzung ergibt sich der Schätzer für  $\lambda_0$  aus dem Quotienten dieser beiden Größen. Details finden sich in [40] und in [41], S. 326 und 347.

### Goel-Okumoto-Modell mit Weibull-Testaufwand

Wie bereits bemerkt, gehen fast alle SZWM davon aus, dass die Belastung, dem ein Programm ausgesetzt ist, im Zeitablauf konstant ist. Insbesondere dann, wenn es sich bei dem verwendeten Zeitmaß  $t$  um die Kalenderzeit handelt, ist jedoch damit zu rechnen, dass die Intensität der Programmnutzung variiert. Oftmals werden zu Beginn einer Testphase nur wenige Tester eingesetzt – z. B. weil Teile der Software gerade noch programmiert werden –, und erst in der Folgezeit wird die Anzahl der Tester und damit der Testaufwand deutlich erhöht, um gegen Ende der Testphase, wenn sich die Anzahl der je Zeiteinheit gefundenen Fehler stark verringert hat, wieder zurückgefahren zu werden. Yamada und andere [58] erweitern das Goel-Okumoto-Modell, indem sie die Verteilung des Testaufwands über die Zeit mit einer Weibullfunktion beschreiben. Dieser Ansatz führt zu der Mittelwertfunktion

$$\mu(t) = N \left[ 1 - \exp\left(-\phi\alpha\left(1 - \exp(-\beta t^\gamma)\right)\right)\right]$$

und der mit ihr verbundenen Versagensintensität

$$\lambda(t) = N\phi\alpha\beta\gamma t^{\gamma-1} \exp\left(-\phi\alpha\left(1 - \exp(-\beta t^\gamma)\right) - \beta t^\gamma\right),$$

deren typischer Verlauf im linken Teil von Abbildung 7 dargestellt ist. Offensichtlich bewirkt eine zunächst wachsende und dann fallende Testintensität eine im Zeitablauf S-förmige Mittelwertfunktion. Die Weibull-Verteilung ist aber flexibel genug, um auch einen kontinuierlich fallenden Testaufwand modellieren zu können; die Steigung der Mittelwertfunktion nimmt dann stetig ab.

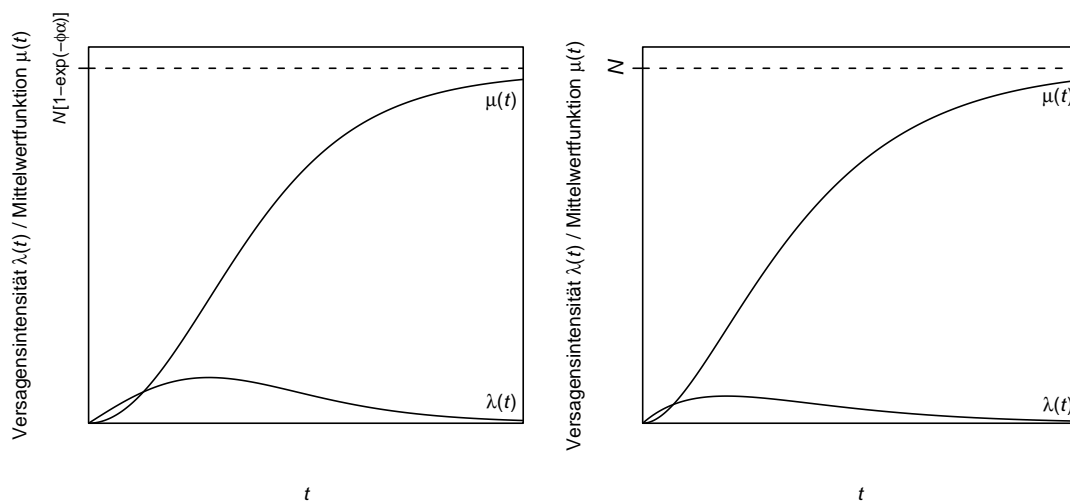


Abb. 7 Typische Entwicklung von Versagensintensität und Mittelwertfunktion gemäß dem Goel-Okumoto-Modell mit Weibull-Testaufwand (links) und dem verzögert S-förmigen Modell (rechts)

Falls die Entwicklung sowohl des Testaufwands als auch des Programmverhaltens sich in der Zukunft ohne Strukturbruch fortsetzen wird, so beträgt die Wahrscheinlichkeit dafür, dass im Intervall  $(t, t+\Delta t]$  kein Softwareversagen auftritt,

$$R(\Delta t | t, V_t) = \exp \left[ N \exp \left( -\phi \alpha \left( 1 - \exp(-\beta(t + \Delta t)^\gamma) \right) \right) - N \exp \left( -\phi \alpha \left( 1 - \exp(-\beta t^\gamma) \right) \right) \right].$$

Alleine auf Versagensbeobachtungen basierend ist die getrennte Schätzung sämtlicher Modellparameter nicht möglich. Die Parameter  $\alpha$  und  $\phi$  können nicht identifiziert werden, das Produkt aus ihnen dagegen schon. Yamada und andere [58] gehen davon aus, dass zusätzlich Daten zur Entwicklung des Testaufwands vorliegen, und schlagen ein zweistufiges Schätzverfahren vor.

### **Verzögert S-förmiges Modell**

Schon vor der expliziten Einbeziehung eines variierenden Testaufwands in SZWM waren S-förmige Modelle betrachtet worden. Yamada und andere [57] untersuchen ein NHPP-Modell mit der Mittelwertfunktion

$$\mu(t) = N \left[ 1 - (1 + \phi t) \exp(-\phi t) \right] \quad (14)$$

und der mit ihr verbundenen Versagensintensität

$$\lambda(t) = N \phi^2 t \exp(-\phi t) . \quad (15)$$

Ohba [43] bezeichnet es als „verzögert S-förmiges Modell“, da man es auch zur Modellierung der Verzögerung zwischen der Versagensbeobachtung und der Fehlerisolierung (d. h. der Bestätigung der Versagens-Reproduzierbarkeit) verwenden kann. Besteht zwischen der Hazardrate des Prozesses  $M(t)$ , der die Versagensfälle zählt, und dem augenblicklichen erwarteten Fehlergehalt der Software ein proportionales Verhältnis und ist die Hazardrate des Prozesses  $G(t)$ , welcher die Anzahl der isolierten Fehler zählt, ihrerseits ein Vielfaches von der erwarteten Anzahl der beobachteten aber noch nicht reproduzierten Versagensfälle, dann hat der Erwartungswert von  $G(t)$  die durch Gleichung (14) gegebene Form.

Allerdings wird das Modell durchaus auch auf reine *Versagensbeobachtungen* angewandt, wenn diese einen S-förmigen Verlauf aufweisen. Insbesondere erfreut sich das Modell – wie auch andere S-förmige Modelle – in Japan einer großen Beliebtheit zur Anpassung an in Kalenderzeit gemessene Versagensdaten, die aufgrund der im letzten Abschnitt angesprochenen Effekte oftmals eine S-Form besitzen [19].

Kann man davon ausgehen, dass die Versagensintensität auch zukünftig durch Gleichung (15) adäquat beschrieben wird, so beträgt die Zuverlässigkeit im Zeitintervall  $(t, t+\Delta t]$

$$R(\Delta t | t, V_t) = \exp \left\{ N \left[ (1 + \phi(t + \Delta t)) \exp(-\phi(t + \Delta t)) - (1 + \phi t) \exp(-\phi t) \right] \right\}.$$

Die Schätzung der beiden Parameter  $N$  und  $\phi$  erfolgt wiederum mittels der Maximum-Likelihood-Methode [43], [57].

#### 25.2.4 Weitere Ansätze zur Modellvereinheitlichung

Der selbstanregende Punktprozess bildet einen sehr weiten Rahmen, in den sich viele existierende SZWM eingliedern lassen. Er ist aber keineswegs die einzig mögliche Sichtweise zur Vereinheitlichung von Modellen. So haben bereits Langberg und Singpurwalla [27] gezeigt, dass sich sowohl das Goel-Okumoto-Modell als auch das Littlewood-Verrall-Modell gewinnen lässt, indem man das Jelinski-Moranda-Modell in einen Bayes-Kontext einbettet und für seine Parameter spezifische (mitunter degenerierte) a priori Verteilungen unterstellt.

Kuo und Yang [25] weisen nach, dass sich für diejenigen NHPP-Modelle, für welche auch bei unendlichem Testaufwand nur eine endliche Anzahl an Versagensfällen  $N$  erwartet wird, die Versagenszeitpunkte als die Ordnungsstatistiken von  $v$  unabhängig und identisch verteilten Beobachtungen auffassen lassen. Hierbei ist  $v$  eine Poisson-verteilte Zufallsvariable mit Erwartungswert  $N$ . Für diejenigen NHPP-Modelle, bei denen die erwartete Zahl an Versagensfällen nicht beschränkt ist, stellen die Versagenszeitpunkte dagegen so genannte Rekorde dar.

Ein völlig anderer Ansatz zur Modellvereinheitlichung betrachtet nur die Mittelwertfunktionen und führt diese auf verschiedene Einflussfaktoren zurück [17], [18], S. 14 ff. Bei den treibenden Größen, die jeweils miteinander in Beziehung stehen, handelt es sich um die Kalenderzeit, den kumulierten Testaufwand, die Anzahl der ausgeführten Testfälle und die Codeabdeckung. Eine Differentialgleichung, die all diese Faktoren berücksichtigt, enthält als Spezialfälle unter anderem das Jelinski-Moranda-Modell, das Goel-Okumoto-Modell, das Goel-Okumoto-Modell mit Weibull-Testaufwand und das verzögert S-förmige Modell. Der Modellrahmen hilft, die Modellannahmen den einzelnen Beziehungen zuzuordnen und auf ihre Realitätsnähe zu überprüfen. Zudem kann er als Ausgangspunkt für die Konstruktion neuer Modelle dienen.

#### 25.2.5 Systematisches und nutzungsprofilorientiertes Testen

Alle bislang vorgestellten SZWM gehen implizit davon aus, dass die Software nutzungsprofilorientiert getestet wird. Während des Testens soll ein Programm also in etwa so bedient werden, wie man es von den späteren Nutzern typischerweise erwartet. Grundsätzlich bedeutet dies, dass sich die Gewichtungen der einzelnen Funktionalitäten nach deren (geschätzten) Nutzungsfrequenzen richten, die Eingabewerte aus adäquaten Verteilungen gezogen und die spezifizierten Testfälle in einer zufälligen Reihenfolge ausgeführt werden sollten [39], S. 165 ff.

Weshalb diese Methodik eine wichtige Voraussetzung für die Prognose der Zuverlässigkeit im Nutzungsbetrieb ist, liegt auf der Hand: Falls beim Testen die Software in einer völlig anderen Weise verwendet wird, ändert sich mit ihrer Veröffentlichung der Prozess, welcher die Versagensfälle generiert. Es ist dann nicht sinnvoll, die Hazardrate des SZWM über das Ende der Testphase hinaus zu extrapolieren. Das nutzungsprofilorientierte Testen kann diesen Strukturbruch verhindern oder zumindest sein Ausmaß verringern.

Allerdings wird diese Teststrategie vielfach für ineffizient und nicht praktikabel gehalten. Der Großteil der softwareproduzierenden Unternehmen setzt so genannte systematische Testtechniken ein. Diese Ansätze versuchen, ausgehend von Informationen

über die Funktionalität der Software oder ihre Implementierung Testfälle zu generieren, die möglichst viele unterschiedliche und insbesondere fehleranfällige Bereiche der Software ausführen. (Für eine detailliertere Diskussion des nutzungsprofilorientierten und des systematischen Testens sowie der jeweiligen Vor- und Nachteile s. [18], S. 6 ff.)

Zwar können ohne genaue Kenntnis der Unterschiede zwischen dem Testprofil und dem Nutzungsprofil die während des systematischen Testens gesammelten Daten nicht zur Prognose der Zuverlässigkeit im Feld verwendet werden. Allerdings ist es auf ihrer Grundlage z. B. möglich, die Anzahl der weiteren Versagensfälle bis zum Testende vorherzusagen, solange es bis dahin nicht zu Strukturbrüchen kommt. Für den Testmanager und das Programmiererteam, welches sich um die Fehlerkorrektur zu kümmern hat, stellen auch diese Informationen wertvolle Planungsgrößen dar.

Die Anwendung eines klassischen SZWM auf Versagensdaten, welche dem systematischen Testen entstammen, führt jedoch nicht unbedingt zu vertrauenswürdigen Ergebnissen. Der Grund hierfür liegt darin, dass die Modelle für das nutzungsprofilorientierte Testen geschaffen wurden und sich dies mitunter in den Strukturen der unterstellten Programmhazardrate und der von ihr abgeleiteten Größen widerspiegelt. So ergibt sich z. B. die Mittelwertfunktion des Jelinski-Moranda- und des Goel-Okumoto-Modells, wenn man ein Ziehen von Codekonstrukten mit Zurücklegen unterstellt [45]; dieser Aufbau ähnelt stark dem nutzungsprofilorientierten Testen mit einem homogenen Nutzungsprofil. Ausgehend von dem im letzten Abschnitt erwähnten Modellrahmen, der aus sukzessiven Beziehungen zwischen treibenden Faktoren besteht, wird in [18], S. 37 ff., ein Modell für die Entwicklung der Anzahl der Versagensfälle während des systematischen Testens hergeleitet.

### 25.2.6 Evaluierung und Verbesserung der Modellgüte

Aufgrund der Vielzahl der existierenden Modelle scheint für jeden Fall ein adäquates Modell bereitzustehen. Das übergroße Angebot hat aber auch Nachteile, ist es doch Ausdruck der Tatsache, dass keines der Modelle für jeden Datensatz gute Ergebnisse liefert. Schlimmer noch: Da jedes Modell nur einen kleinen Teil der mannigfaltigen technischen und sozialen Einflussfaktoren des versagensverursachenden Prozesses abbilden kann, ist es nicht möglich, im Vorfeld der Datenerhebung mit Sicherheit zu entscheiden, welches der Modelle am besten zu den Versagensbeobachtungen passen wird [6].

Um so wichtiger ist es, für einen vorhandenen Datensatz die Güte verschiedener Modelle zu vergleichen. Hierbei gibt es eine Reihe von Kriterien, die unterschiedliche Aspekte der Modellqualität operationalisieren und erfassen. Das grundsätzliche Vorgehen zur Berechnung dieser Maße ist dabei immer gleich: Beginnend mit den ersten z. B. fünf Datenpunkten des gesamten Datensatzes werden die Parameter eines Modells geschätzt und zur Prognose einer bestimmten Größe (z. B. des nächsten Versagenszeitpunkts) bzw. deren Verteilung verwendet. Unter Hinzunahme jeweils eines weiteren Datenpunktes zu dem gestutzten Datensatz wird diese Prozedur sukzessive wiederholt. Man simuliert also die begleitende Anwendung des Modells während des gesamten bisherigen Projektverlaufs. Aus dem Vergleich der einzelnen Prognosen untereinander bzw. mit den tatsächlichen Beobachtungen errechnet sich schließlich

das Gütekriterium für das jeweilige Modell. Folgende konkrete Maße werden oftmals betrachtet:

1. Absolute relative Prognosefehler [11], S. 3 f., [18], S. 78 f.:

Zur Beurteilung der Qualität der *kurzfristigen* Prognose wird in jedem Schritt die geschätzte Anzahl der Versagensfälle zum nächsten Versagenszeitpunkt mit dem tatsächlichen Wert verglichen und der Absolutbetrag der relativen Abweichung bestimmt. Der kurzfristige absolute relative Prognosefehler ergibt sich dann als Mittelwert all dieser Größen. Um das *langfristige* Verhalten eines Modells quantifizieren zu können, ist der maximale Prognosehorizont zu wählen, für den eine Gegenüberstellung mit der Realität möglich ist. Deshalb wird für den so genannten mittleren absoluten relativen Prognosefehler für jeden gestutzten Datensatz die prognostizierte Anzahl von Versagensfällen bis zum Ende des Beobachtungszeitraums mit dem tatsächlichen Wert verglichen. Selbstredend ist ein Modell um so besser, je geringer seine Prognosefehler ausfallen.

2. Variabilitätsmaße [1], [18], S. 80 f.:

Um als Planungsgrundlage dienen zu können, dürfen sich die von einem Modell gelieferten Qualitätseinschätzungen bei der Hinzunahme einer weiteren Beobachtung nicht zu stark verändern. Zur Beurteilung des Ausmaßes dieser unerwünschten Variabilität berechnet man aus der Sequenz der Prognosen (z. B. der Zeitspanne bis zum nächsten Softwareversagen) für je zwei aufeinanderfolgende Werte den Absolutbetrag der relativen Abweichung. Das so genannte Variabilitätsmaß ergibt sich dann als Summe dieser Abweichungsgrößen; je kleiner sein Wert ist, desto besser.

3. Präquenzielle Likelihoodfunktion [1]:

Mit der Schätzung eines (zeitbasierten) SZWM aufgrund der bisherigen Beobachtungen wird indirekt zugleich die Verteilung der Zeit bis zum nächsten Versagensfall prognostiziert. Falls das Modell adäquat ist, sollte man erwarten, dass die später tatsächlich eintretende Realisation aus einem Bereich der Verteilung stammt, welcher eine große Eintrittswahrscheinlichkeit aufweist. Bei einer stetigen Zufallsvariablen sollte die Dichtefunktion an der Stelle dieser Beobachtung tendenziell einen hohen Wert annehmen. Evaluiert man in der sequenziellen Modellanwendung jede der Prognosedichten an der jeweils eingetroffenen Realisation und multipliziert die so erhaltenen Größen, dann ergibt sich ein Maß für die Plausibilität des Modells anhand des gesamten Datensatzes, welches als präquenzielle Likelihoodfunktion (*prequential likelihood*) bezeichnet wird. Zum Vergleich zweier Modelle bildet man den Quotienten der beiden präquenziellen Likelihoodfunktionen. Tendiert dieses präquenzielle Likelihoodverhältnis mit zunehmender Datensatzlänge gegen unendlich, dann ist das Modell, dessen präquenzielle Likelihoodfunktion im Zähler steht, dem anderen vorzuziehen; tendiert es gegen Null, so trifft das Gegenteil zu.

4.  $u$ -Plot und  $y$ -Plot [1], [6]:

Obwohl die tatsächlich beobachteten Wartezeiten bis zum nächsten Versagensfall überwiegend aus denjenigen Bereichen der Prognosedichten stammen sollten,

welche eine hohe Wahrscheinlichkeitsmasse umfassen (wie von der präquenziellen Likelihoodfunktion betont), sind durchaus auch – einige wenige – Realisationen aus den Rändern der Verteilungen zu erwarten. So dürften etwa 5% der Werte kleiner als diejenigen Schranken sein, die mit fünfprozentiger Wahrscheinlichkeit unterschritten werden. Beim so genannten  $u$ -Plot handelt sich um ein grafisches Instrument, mit dessen Hilfe überprüft werden kann, ob die Beobachtungen in diesem Sinne zu der Gestalt der von einem Modell prognostizierten Verteilungsfunktionen passen. Hierbei kann nicht nur die Stärke der Abweichung quantifiziert und mit derjenigen, die mit einem anderen Modell verbunden ist, verglichen werden. Es zeigt sich zudem, ob die tatsächlichen Wartezeiten bis zum nächsten Softwareversagen tendenziell in den oberen (unteren) Rändern der prognostizierten Verteilungen liegen und das Modell somit die Zuverlässigkeit der Software systematisch unterschätzt (überschätzt). Falls allerdings für eine Hälfte der Daten die Prognosen zu optimistisch sind, während sie für die andere Hälfte zu pessimistisch ausfallen, dann können sich diese gegensätzlichen Abweichungen des Modells von der Wirklichkeit ausgleichen, sodass sie im  $u$ -Plot nicht zu entdecken sind. Der auf den Werten des  $u$ -Plots aufbauende  $y$ -Plot kann solche Trends identifizieren.

Hat man mithilfe eines Gütemaßes erkannt, dass ein SZWM bei der Prognose systematische Fehler macht, ist es möglich, dieses Wissen zur Verbesserung der Vorhersagen zu nutzen. Dies ähnelt dem Vorgehen eines Schützen, der bei seinen bisherigen Schüssen immer links am Ziel vorbeigeschossen hat: Die Lage rekapitulierend wird er beim nächsten Mal weiter nach rechts zielen. Brocklehurst und andere [7] schlagen eine solche „Rekalibrierung“ von Prognosen vor, die auf den Ergebnissen des (geglätteten)  $u$ -Plots basiert. Sie zeigen, dass diese Technik die Prognosen verschiedener Modelle aneinander angleicht und dass zudem dem präquenziellen Likelihoodverhältnis gemäß die rekalierten „Modelle“ ihren ursprünglichen Varianten vorzuziehen sind.

Einen deutlich einfacheren Ansatz wählen Lyu und Nikora [33]. Sie raten dazu, eine Zuverlässigkeitsprognose als arithmetisches Mittel der Prognosen mehrerer SZWM zu berechnen. Insbesondere empfehlen sie die Mittelung der Prognosen des Goel-Okumoto-Modells (welches generell als zu optimistisch gilt), des Littlewood-Verrall-Modells (welches zu pessimistischen Prognosen tendiert) und des Musa-Okumoto-Modells (dessen Verzerrungsrichtung stärker variiert). Bei Erweiterungen der Methodik können die Gewichte für die einzelnen Modelle unterschiedlich ausfallen und sich sogar dynamisch nach der relativen Güte der Anpassung des jeweiligen Modells an die Daten richten [34]. In dieser adaptiven Variante handelt es sich bei den Gewichten um so genannte Bayesfaktoren; diese sind eng mit den präquenziellen Likelihoodfunktionen der verschiedenen Modelle verbunden [49], S. 148 ff.

Ein grundsätzlicher Nachteil der Rekalibrierung und der Mittelung von Prognosen liegt darin, dass die Annahmen der ursprünglichen Modelle und die Interpretierbarkeit einzelner Modellparameter (z. B. des Parameters  $u_0$  als der Anzahl der zu Beginn vorhandenen Softwarefehler im Jelinski-Moranda-Modell) verloren gehen. Die Gesamtheit aus Modell(en), Schätz- und Prognoseverfahren wird vollends zur Blackbox.

### 25.3 Weitere Modellklassen

Zwar haben die SZWM in der Literatur über die Schätzung und Prognose der Qualität und Zuverlässigkeit von Software die größte Aufmerksamkeit erfahren. Es handelt sich bei ihnen aber keineswegs um die einzige existierende Modellklasse. Der Vollständigkeit halber sollen in diesem Abschnitt einige weitere Modellansätze skizziert werden. Völlig ausklammern wollen wir aus Platzgründen Zertifizierungstests zur Überprüfung der Zuverlässigkeit von fertigen Softwareprodukten (s. [41], S. 201 ff.), welche auf der statistischen Theorie des sequenziellen Testens [53] beruhen.

#### 25.3.1 Stichprobenmodelle

Anders als die SZWM in Abschnitt 25.2 versuchen die in diesem Abschnitt besprochenen Modelle nicht, die Entwicklung der Zuverlässigkeit oder der Anzahl der verbliebenen Softwarefehler während einer Testphase mit Fehlerkorrektur nachzuvollziehen und vorherzusagen. Vielmehr dienen sie dazu, den aktuellen Fehlergehalt oder die Zuverlässigkeit einer Software zu bestimmen. Insofern mag man sie eher als Schätz- denn als Prognosemodelle bezeichnen. Allerdings ist zu beachten, dass Zuverlässigkeitswerte als Wahrscheinlichkeiten für einen Versagenseintritt bei zukünftiger Nutzung immer auch Vorhersagen sind, selbst wenn das Softwareprodukt unverändert bleibt.

Zu ihrer Schätzung benötigen die Modelle keine Informationen über die Entwicklung der Anzahl der Versagensfälle im Zeitablauf. Die Daten der sukzessive durchgeführten Tests können gruppiert vorliegen, entweder in Form einer globalen Stichprobe oder getrennt in zwei Stichproben.

#### **Nelson-Modell**

Dieses Modell [52], S. 217 ff., gründet sich auf derjenigen Zuverlässigkeitsdefinition, welche die Länge der „Nutzungsperiode“ anhand der Anzahl der Programmläufe misst (s. oben Abschnitt 25.1). Genauer bezeichnet es als Zuverlässigkeit  $R$  die Wahrscheinlichkeit dafür, dass im Rahmen *eines* Laufs kein Versagen auftritt. Unter einem Programmlauf wird hierbei die Ausführung der Software mit einer bestimmten Kombination von Eingabewerten für die Inputvariablen verstanden. Diese entstammt der sehr großen aber endlichen Menge aller möglichen Wertekombinationen. Kam es während des Testens bei  $m$  von insgesamt  $n$  Programmläufen zu einem Softwareversagen, dann lautet die Zuverlässigkeitsschätzung nach dem Nelson-Modell

$$\hat{R} = 1 - \frac{m}{n}.$$

Damit dieser Wert auch ein unverzerrter Schätzer für die versagensfreie Programmausführung im normalen Nutzungsbetrieb sein kann, müssen natürlich die Auswahlwahrscheinlichkeiten der Inputkombinationen denjenigen entsprechen, welche auch nach der Softwareveröffentlichung vorherrschen; kurz: Es muss nutzungsprofilorientiert getestet werden. Obwohl der Schätzer unter dieser Voraussetzung unverzerrt ist [52], S. 222 ff., ist eine große Zahl an Testläufen nötig, um seine Varianz gering zu halten und somit ein hohes Vertrauen in die Punktschätzung legen zu können [3].

Dass das Modell an das nutzungsprofilorientierte Testen gebunden ist und deshalb nicht während des Testens gemäß systematischer Strategien verwendet werden kann, wird als weiterer Nachteil gesehen [3].

### **Brown-Lipow-Modell**

Eine Lösung des letztgenannten Problems bietet der Ansatz von Brown und Lipow [8]. Zu seiner Anwendung ist es nicht nötig, dass dem Nutzungsprofil gemäß getestet wird; dieses Profil muss aber in folgender Weise explizit spezifiziert sein: Die große Menge der möglichen Eingabekombinationen sei in Teilmengen  $Z_1, Z_2, \dots, Z_K$  aufgespalten, die überschneidungsfrei sind und gemeinsam die gesamte Menge ergeben. Bei diesen Teilmengen kann es sich beispielsweise um Äquivalenzklassen handeln, deren Elemente bei ihrer Eingabe erwartungsgemäß jeweils die gleiche Reaktion der Software bewirken. Das Nutzungsprofil muss dann in Form der Auftrittswahrscheinlichkeiten all dieser Teilmengen bei normaler Programmnutzung,  $P(Z_1), P(Z_2), \dots, P(Z_K)$ , bekannt sein. Als Ergebnis der Testdurchführung ist für jedes  $Z_j$  zum einen die Anzahl  $n_j$  der Programmläufe festzuhalten, deren Eingabekombinationen zu dieser Teilmenge gehören. Des Weiteren muss jeweils gezählt werden, wie viele der  $n_j$  Läufe zu einem Versagen führen; für  $Z_j$  sei dieser Wert mit  $m_j$  bezeichnet. Die geschätzte augenblickliche Zuverlässigkeit der Software bei Bedienung dem Nutzungsprofil entsprechend beträgt dann

$$\hat{R} = 1 - \sum_{j=1}^K \frac{m_j}{n_j} P(Z_j).$$

Nelson [42] überträgt diese Formel auf eine Situation, in der die Läufe aller spezifizierten Testfälle zu *keinem* Softwareversagen führen (z. B. weil bereits zuvor alle Testfälle durchgeführt und die durch sie aufgedeckten Fehler bereinigt wurden). Die Zuverlässigkeitsschätzung errechnet er als

$$\hat{R} = 1 - \sum_{j=1}^K \varepsilon_j P(Z_j),$$

wobei  $\varepsilon_j$  die Wahrscheinlichkeit dafür bezeichnet, dass eine beliebige aus  $Z_j$  gewählte Eingabekombination zu einem Versagen führt. Für die Bestimmung dieser  $\varepsilon_j$ -Werte gibt Nelson heuristische Regeln an, welche unter anderem die Anzahl der Testfälle berücksichtigen, die sich auf die Teilmenge  $Z_j$  beziehen.

### **Mills-Modell**

Ist man nicht an der Zuverlässigkeit, sondern lediglich an der Anzahl der Programmfehler interessiert, so kann man sich so genannte Capture-Recapture-Modelle zunutze machen, statistische Modelle, welche ursprünglich zur Schätzung der Größe von Populationen (z. B. der Anzahl der Fische in einem Teich) verwendet wurden [23], S. 248 ff. In einer spezifischen Form [44], S. 81 ff., wurden sie erstmals von Mills auf das Gebiet des Softwaretestens übertragen. In ein Programm, welches eine unbekannte Anzahl von Fehlern  $u_0$  aufweist, werden bewusst  $u_1$  Fehler eingebaut. Es sei angenommen, dass alle Fehler in etwa gleich leicht entdeckt werden können und insbesondere die  $u_1$  „gesäten“ Fehler nicht leichter oder schwerer zu finden sind als die  $u_0$

von Anfang an vorhandenen. Zudem liege keine Interaktion zwischen den verschiedenen Fehlern vor.

Wird nun eine Reihe von Testfällen durchgeführt, wobei insgesamt  $f$  Fehler gefunden werden, dann ist die Anzahl derjenigen unter ihnen, bei denen es sich um gesäte Fehler handelt, zufällig. Unter den oben genannten Voraussetzungen folgt diese Zufallsvariable, die mit  $F_1$  bezeichnet sei, einer hypergeometrischen Verteilung; d. h. die Wahrscheinlichkeit dafür, dass sie den Wert  $f_1$  annimmt und sich also  $f_1$  gesäte Fehler unter allen entdeckten befinden, beträgt

$$P(F_1 = f_1; u_0, u_1, f) = \frac{\binom{u_0}{f-f_1} \binom{u_1}{f_1}}{\binom{u_0+u_1}{f}}.$$

Wurden tatsächlich  $f_1$  der gesäten Fehler wiedergefunden, so ergibt sich mittels der Maximum-Likelihood-Methode folgende Schätzung für den Parameter  $u_0$ , die Anzahl der ursprünglichen Fehler [9], S. 108:

$$\hat{u}_0 = \left\lfloor \frac{u_1(f-f_1)}{f_1} \right\rfloor.$$

Hierbei bezeichnet  $\lfloor x \rfloor$  die größte ganze Zahl, die kleiner oder gleich  $x$  ist. Der Schätzwert für  $u_0$  entspricht also in etwa derjenigen Größe, die man erhält, wenn man die Anzahl der gefundenen ursprünglichen Fehler  $(f-f_1)$  durch die Entdeckungsquote bei den gesäten Fehlern  $(f_1/u_1)$  dividiert. Offensichtlich wird dieser Ansatz  $u_0$  tendenziell unterschätzen, wenn die bewusst eingefügten Fehler leichter zu finden sind als die ursprünglichen und damit erwartungsgemäß über eine höhere Entdeckungsquote verfügen.

### **Basin-Modell**

Ein etwas veränderter Aufbau des Experiments, der kein Einbringen weiterer Fehler erfordert, wird mit Basin in Verbindung gebracht [9], S. 113. Das Fangen und Wiederfangen, welches in dem Begriff „Capture-Recapture-Modell“ zum Ausdruck kommt, wird in der Form des unabhängigen Testens der Software durch zwei Personen realisiert. Von den insgesamt  $u_0$  Fehlern habe der erste Tester  $f_1$  und der zweite Tester  $f_2$  entdeckt. Falls die Schwierigkeit des Auffindens für jeden Fehler gleich groß ist und zudem nicht von der Person des Testers abhängt, so beträgt die Wahrscheinlichkeit dafür, dass  $w$  der  $f_2$  vom zweiten Tester aufgespürten Fehler bereits von seinem Kollegen entdeckt worden waren,

$$P(W = w; u_0, f_1, f_2) = \frac{\binom{f_1}{w} \binom{u_0-f_1}{f_2-w}}{\binom{u_0}{f_2}}.$$

Wiederum tritt also die hypergeometrische Verteilung in Erscheinung.

Unter Verwendung des für  $w$  tatsächlich beobachteten Wertes lautet der Maximum-Likelihood-Schätzer für die Gesamtzahl an Fehlern in der Software

$$\hat{u}_0 = \left\lfloor \frac{f_1 f_2}{w} \right\rfloor.$$

Wurden im zweiten „Fang“ also  $(w / f_1) = y\%$  der zuvor gefundenen Fehler wiederentdeckt, so kann man davon ausgehen, dass die  $f_2$  Fehler selbst etwa  $y\%$  der Gesamtfehlerzahl ausmachen.

Auf Probleme der Verwendung von Capture-Recapture-Modellen zur Schätzung des Fehlergehalts einer Software weist Isoda [21] hin.

### 25.3.2 Modelle zur Prognose von Softwarefehlern

Die bisher behandelten Modelle haben gemein, dass sie Daten aus der Ausführung des betrachteten Programms verwenden, um Rückschlüsse über dessen Qualität zu ziehen. In diesem Abschnitt sind einige Ansätze zusammengefasst, die aufgrund anderer (typischerweise bereits vor der Testphase verfügbarer) Informationen versuchen, den Fehlergehalt der Software vorherzusagen. Gegenüber den anderen Modellen unterscheiden sich diese Querschnittsmodelle grundlegend darin, dass zur Schätzung ihrer Parameter nicht nur die Daten eines einzigen Projekts verwendet werden. Vielmehr greift man entweder auf in der Literatur publizierte Erfahrungs- und Schätzwerte zurück, oder man schätzt die Parameter basierend auf einer Sammlung von früheren Projekten des eigenen Unternehmens. In jedem Fall unterstellt man, dass die Zusammenhänge, die anhand der für die Modellspezifikation genutzten Projekte identifiziert wurden, auch für die zukünftigen Projekte gültig sind. Die Prognoseergebnisse sind mit um so größerer Vorsicht zu genießen, je stärkere Zweifel an der Vergleichbarkeit der Projekte bestehen.

#### **Multiplikative Modelle**

Um multiplikative Modelle handelt es sich z. B. bei dem Modell des Rome Laboratory der Air Force (s. [12] und [26], S. 7-4 ff.) sowie bei dem Ansatz von Malaiya und Denton [35]. In ihnen ergibt sich die prognostizierte Fehlerdichte – die Anzahl der Fehler je 1000 Sourcecode-Zeilen – als Produkt einer Reihe von Faktoren, deren Werte in Abhängigkeit von den Gegebenheiten der Software und des gesamten Entwicklungsprojekts bestimmt werden. Beide Modelle verfügen über einen Faktor, welcher eine Basis-Fehlerdichte repräsentiert. Im Modell des Rome Laboratory wird der Wert dieses Faktors anhand einer Checkliste ermittelt, welche die Schwierigkeit der Entwicklung der vorliegenden Art von Software beurteilt; bei Malaiya und Denton beruht er auf der von dem betrachteten Unternehmen im Durchschnitt erreichten Fehlerdichte. Die weiteren Faktoren berücksichtigen Aspekte der verwendeten Entwicklungsmethoden, der institutionalisierten Entwicklungs- und Testprozesse, der Eignung der Mitarbeiter und der Struktur des implementierten Codes. All diese Faktoren weisen einen Wertebereich um die Zahl Eins auf. Je nach Ausprägung der einzelnen Aspekte wird also die Basis-Fehlerdichte aufgebläht oder verringert. So nimmt z. B. der Programmiererteam-Faktor in Malaiyas und Dentons Modell bei einer durchschnittlichen Leistungsfähigkeit den Wert

Eins und bei einem hohen bzw. niedrigen Leistungsniveau die Werte 0,4 resp. 2,5 an. Der Vorteil dieses multiplikativen Aufbaus besteht darin, dass bei Fehlen einzelner Informationen die jeweiligen Faktoren weggelassen (und dabei implizit auf ihren Grundwert Eins gesetzt) werden können. Da es möglich ist, die Faktoren des Rome Laboratory – Modells den Entwicklungsphasen Analyse, Design und Implementierung/ Test zuzuordnen, kann man somit für jeden Entwicklungsstand ein Submodell aufstellen, welches eine Teilmenge der Faktoren umfasst [26], S. 7-4 ff.

### **Lineare Regressionsmodelle**

Lineare Regressionsmodelle versuchen, eine abhängige Variable  $y$  auf eine Linearkombination von erklärenden Variablen  $x_1, \dots, x_k$  zurückzuführen:

$$y = \alpha_0 + \alpha_1 x_1 + \dots + \alpha_k x_k + \eta .$$

Hierbei steht  $\eta$  für die zufällige Abweichung von dem linearen Zusammenhang, für welche unter anderem ein Erwartungswert von Null unterstellt wird.

In der konkreten Anwendung der Fehlerprognose handelt es sich bei der abhängigen Variablen um die Anzahl der Fehler im Programmcode oder um die Fehlerdichte. Mitunter wird jedoch auch auf den natürlichen Logarithmus der Fehlerdichte zurückgegriffen; ein Vorteil dieses Vorgehens liegt darin, dass der Logarithmus im Gegensatz zur Fehlerdichte selbst nicht auf Werte größer oder gleich Null beschränkt ist.

Takahashi und Kamayachi [51] definieren neun quantitative Indikatoren, die potenziell einen Einfluss auf die Anzahl der Programmfehler haben. Bei denjenigen drei Variablen, die für ihren Datensatz aus 30 Projekten die größte Erklärungskraft für die Gesamtfehlerzahl zeigen, handelt es sich um die Häufigkeit von Änderungen der Programmspezifikation (gemessen in Seiten der Änderungswünsche), die durchschnittliche Programmiererfahrung der Entwickler (in Jahren) und den Umfang der Designdokumente (in Seiten). Für diese Variablen stellen die Autoren ein lineares Regressionsmodell mit der Anzahl der Programmfehler als exogene Variable auf.

Zhang und Pham [60] erweitern Takahashis und Kamayachis Liste der Einflussfaktoren deutlich. Mittels eines Fragebogens erheben sie bei verschiedenen Gruppen von Mitarbeitern in Softwareunternehmen (z. B. Managern, Programmierern und Testern) die subjektiv empfundene Bedeutung dieser Faktoren für die Zuverlässigkeit der entwickelten Software, ohne allerdings anhand echter Projektdaten die Erwartungen zu verifizieren oder ein Regressionsmodell zu schätzen.

In [18] wird eine Auswahl der von Zhang und Pham zusammengetragenen Einflussfaktoren weiter operationalisiert und damit objektiv messbar gemacht. Zudem enthält ein ausführlicher Fragebogen [18], S. 223 ff., Fragen und detaillierte Szenarien, mithilfe derer die Reife von Softwareentwicklungsprozessen in Anlehnung an den zukünftigen Standard ISO/IEC 15504 [20] – das so genannte SPICE-Modell – bestimmt werden kann. Die Analyse der dreizehn verfügbaren Projektdatensätze führt zu einem linearen Regressionsmodell, bei dem die Fehlerdichte durch eine selektive Reifegradbewertung, das Verhältnis zwischen der tatsächlichen und der geplanten Entwicklungsdauer und den Anteil der nach der Spezifikationsphase geänderten Anforderungen erklärt wird.

Insbesondere dann, wenn es sich bei den exogenen Faktoren um Maße der Programmkomplexität handelt, welche erwartungsgemäß stark miteinander verbunden sind, können sich die geschätzten Regressionskoeffizienten bei der Aufnahme weiterer erklärender Variablen deutlich verändern. Um dies und weitere Probleme der so genannten Multikollinearität in den Griff zu bekommen, schlagen Khoshgoftaar und Munson [24] die Anwendung der Faktorenanalyse zur Gewinnung von orthogonalen (voneinander völlig unabhängigen) Faktoren vor.

Eine Übersicht über weitere – nicht notwendigerweise lineare – Regressionsmodelle für Fehlerdaten findet sich bei Cai [9], S. 47 ff.

## 25.4 Abschließende Bemerkung

Dieses Kapitel gibt einen knappen Überblick über verschiedene Ansätze zur Prognose von Softwarezuverlässigkeit, Softwareversagensfällen und Softwarefehlern. Da im Rahmen eines Softwareentwicklungsprojekts das korrekte Softwareverhalten zumeist nur eines der zu beachtenden Kriterien darstellt (neben dem Funktionsumfang, der Entwicklungszeit, den Lebenszykluskosten, usw.), kann es sinnvoll sein, die hier diskutierten Modelle als Elemente umfassenderer Optimierungsprobleme einzusetzen. So betrachten z. B. Pham [44], S. 159 ff., und Yamada [56] Ansätze zur Bestimmung derjenigen Testdauer, welche die Gesamtkosten der Testdurchführung und der Gewährleistung minimiert.

## 25.5 Literatur

- [1] Abdel-Ghaly, A.A., Chan, P.Y. und Littlewood, B., Evaluation of competing software reliability predictions, *IEEE Transactions on Software Engineering* 12 (1986), S. 950 ff.
- [2] Ascher, H. und Feingold, H., *Repairable systems reliability – Modeling, inference, misconceptions and their causes*, New York 1984.
- [3] Bastani, F.B. und Ramamoorthy, C.V., Software reliability, in: Krishnaiah, P.R. und Rao, C.R. (Hrsg.), *Handbook of statistics*, Vol. 7, Amsterdam 1988, S. 7 ff.
- [4] Belli, F., Grochtmann, M. und Jack, O., Erprobte Modelle zur Quantifizierung der Software-Zuverlässigkeit, *Informatik Spektrum* 21 (1998), S. 131 ff.
- [5] Boland, P.J. und Singh, H., A birth-process approach to Moranda's geometric software-reliability model, *IEEE Transactions on Reliability* 52 (2003), S. 168 ff.
- [6] Brocklehurst, S. und Littlewood, B., Techniques for prediction analysis and recalibration, in: Lyu, M.R. (Hrsg.), *Handbook of software reliability engineering*, New York 1996, S. 119 ff.
- [7] Brocklehurst, S., Chan, P.Y. und Littlewood, B., Recalibrating software reliability models, *IEEE Transactions on Software Engineering* 16 (1990), S. 458 ff.
- [8] Brown, J.R. und Lipow M., Testing for software reliability, *Proceedings of the International Conference on Reliable Software*, New York 1975, S. 518 ff.
- [9] Cai, K.-Y., *Software defect and operational profile modeling*, Boston 1998.
- [10] Chen, Y. und Singpurwalla, N.D., Unification of software reliability models by self-exciting point processes, *Advances in Applied Probability* 29 (1997), S. 337 ff.
- [11] Denton, J.A., *Accurate software reliability estimation*, Master's thesis, Colorado State University, Fort Collins 1999.
- [12] Farr, W., Software reliability modeling survey, in: Lyu, M.R. (Hrsg.), *Handbook of software reliability engineering*, New York 1996, S. 71 ff.

- [13] Forman, E.H. und Singpurwalla, N.D., An empirical stopping rule for debugging and testing computer software, *Journal of the American Statistical Association* 72 (1997), S. 750 ff.
- [14] Gaudoin, O., Outils statistiques pour l'évaluation de la fiabilité des logiciels, Thèse de doctorat, Université de Joseph Fourier – Grenoble 1, Grenoble 1990.
- [15] Goel, A.L. und Okumoto, K., Time-dependent error-detection model for software reliability and other performance measures, *IEEE Transactions on Reliability* 28 (1979), S. 206 ff.
- [16] Gokhale, S.S., Marinou, P.N. und Trivedi, K.S., Important milestones in software reliability modeling, *Proceedings of the Eighth International Conference on Software Engineering and Knowledge Engineering*, Skokie 1996, S. 345 ff.
- [17] Grottke, M., A vector Markov model for structural coverage growth and the number of failure occurrences, *Proceedings of the Thirteenth IEEE International Symposium on Software Reliability Engineering*, Los Alamitos 2002, S. 304 ff.
- [18] Grottke, M., Modeling software failures during systematic testing – The influence of environmental factors, Aachen 2003.
- [19] Iannino, A., Software reliability theory, in: Marciniak, J. (Hrsg.), *Encyclopedia of software engineering*, New York 1994, S. 1223 ff.
- [20] ISO/IEC JTC 1/SC 7/WG 10, Information technology – Software process assessment – Part 2: A reference model for processes and process capability, Technical report ISO/IEC TR 15504-2, Genf 1998.
- [21] Isoda, S., A criticism on the capture-and-recapture method for software reliability assurance, *The Journal of Systems and Software* 43 (1998), S. 3 ff.
- [22] Jelinski, Z. und Moranda, P., Software reliability research, in: Freiburger, W. (Hrsg.), *Statistical computer performance evaluation*, New York 1972, S. 465 ff.
- [23] Johnson, N.L. und Kotz, S., *Urn models and their application*, New York 1977.
- [24] Khoshgoftaar, T.M. und Munson, J.C., Predicting software development errors using software complexity metrics, *IEEE Journal on Selected Areas in Communications* 8 (1990), S. 253 ff.
- [25] Kuo, L. und Yang, T.Y., Bayesian computation for nonhomogeneous Poisson processes in software reliability, *Journal of the American Statistical Association* 91 (1996), S. 763 ff.
- [26] Lakey, P.B. und Neufelder, A.M., *System and software reliability assurance guidebook*, Rome Laboratory, Rome 1997.  
Verfügbar unter <http://www.softrel.com/notebook.zip> (Abruf am 22.12.2003).
- [27] Langberg, N. und Singpurwalla, N.D., A unification of some software reliability models, *SIAM Journal of Scientific and Statistical Computing* 6 (1985), S. 781 ff.
- [28] Ledoux, J., Software reliability modeling, in: Pham, H. (Hrsg.), *Handbook of reliability engineering*, London 2003, S. 213 ff.
- [29] Littlewood, B., Stochastic reliability growth: A model for fault-removal in computer-programs and hardware-design, *IEEE Transactions on Reliability* 30 (1981), S. 313 ff.
- [30] Littlewood, B. und Verrall, J.L., A Bayesian reliability growth model for computer software, *Journal of the Royal Statistical Society, series C* 22 (1973), S. 332 ff.
- [31] Littlewood, B. und Verrall, J.L., Likelihood function of a debugging model for computer software reliability, *IEEE Transactions on Reliability* 30 (1981), S. 145 ff.
- [32] Lyu, M.R. (Hrsg.), *Handbook of software reliability engineering*, New York 1996.
- [33] Lyu, M.R. und Nikora, A., A heuristic approach for software reliability prediction: The equally-weighted linear combination model, *Proceedings of the 1991 IEEE International Symposium on Software Reliability Engineering*, Los Alamitos 1991, S. 172 ff.
- [34] Lyu, M.R. und Nikora, A., CASRE – A computer-aided software reliability estimation tool, *Proceedings of the 1992 IEEE Computer-Aided Software Engineering Workshop*, Los Alamitos 1992, S. 264 ff.

- [35] Malaiya, Y.K. und Denton, J.A., What do the software reliability growth model parameters represent?, Technical report CS-97-115, Computer Science Department, Colorado State University, Fort Collins 1997.
- [36] Mazzuchi, T.A. und Singpurwalla, N.D., Software reliability models, in: Krishnaiah, P.R. und Rao, C.R. (Hrsg.), Handbook of statistics, Vol. 7, Amsterdam 1988, S. 73 ff.
- [37] Mazzuchi, T.A. und Soyer, R., A Bayes empirical-Bayes model for software reliability, IEEE Transactions on Reliability 37 (1988), S. 248 ff.
- [38] Moranda, P.B., Event-altered rate models for general reliability analysis, IEEE Transactions on Reliability 28 (1979), S. 376 ff.
- [39] Musa, J.D., Software reliability engineering, New York 1999.
- [40] Musa, J.D. und Okumoto, K., A logarithmic Poisson execution time model for software reliability measurement, Proceedings of the Seventh International Conference on Software Engineering, Piscataway 1984, S. 230 ff.
- [41] Musa, J.D., Iannino, A. und Okumoto, K., Software reliability: Measurement, prediction, application, New York 1987.
- [42] Nelson, E., Estimating software reliability from test data, Microelectronics and Reliability 17 (1978), S. 67 ff.
- [43] Ohba, M., Software reliability analysis models, IBM Journal of Research and Development 28 (1984), S. 428 ff.
- [44] Pham, H., Software reliability, Singapore 2000.
- [45] Piwowarski, P., Ohba, M. und Caruso, J., Coverage measurement experience during function test, in: Proceedings of the Fifteenth International Conference on Software Engineering, Los Alamitos 1993, S. 287 ff.
- [46] Schick, G.J. und Wolverton, R.W., An analysis of competing software reliability models, IEEE Transactions on Software Engineering 4 (1978), S. 104 ff.
- [47] Shantikumar, J.G., A general software reliability model for performance prediction, Microelectronics and Reliability 21 (1981), S. 671 ff.
- [48] Singpurwalla, N.D. und Wilson, S.P., Software reliability modeling, International Statistical Review 62 (1994), S. 289 ff.
- [49] Singpurwalla, N.D. und Wilson, S.P., Statistical methods in software engineering: Reliability and risk, New York 1999.
- [50] Snyder, D.L. und Miller, M.I., Random point processes in time and space, New York 1991.
- [51] Takahashi, M. und Kamayachi, Y., An empirical study of a model for program error prediction, Proceedings of the Eighth International Conference on Software Engineering, Los Alamitos 1985, S. 330 ff.
- [52] Thayer, T.A., Lipow, M. und Nelson, E.C., Software reliability, Amsterdam 1978.
- [53] Wald, A., Sequential analysis, New York 1947.
- [54] Xie, M., Software reliability modelling, Singapore 1991.
- [55] Xie, M. und Hong, G.Y., Software reliability modeling, estimation and analysis, in: Balakrishnan, N. und Rao, C.R. (Hrsg.), Handbook of statistics, Vol. 20, Amsterdam 2001, S. 707 ff.
- [56] Yamada, S., Software reliability models, in: Osaki, S. (Hrsg.), Stochastic models in reliability and maintenance, Berlin 2002, S. 253 ff.
- [57] Yamada, S., Ohba, M. und Osaki, S., S-shaped reliability growth modeling for software error detection, IEEE Transactions on Reliability 32 (1983), S. 475 ff.
- [58] Yamada, S., Hishitani, J. und Osaki, S., Software-reliability growth with a Weibull test-effort: A model & application, IEEE Transactions on Reliability 42 (1993), S. 100 ff.
- [59] Yang, B. und Xie, M., A study of operational and testing reliability in software reliability analysis, Reliability Engineering and System Safety 70 (2000), S. 323 ff.
- [60] Zhang, X. und Pham, H., An analysis of factors affecting software reliability, The Journal of Systems and Software 50 (2000), S. 43 ff.