



Schriftliche Prüfung im Grundwissen

## **Angewandte Stochastik**

### **Klausur mit Lösungen**

gemäß Prüfungsordnung 4  
der Deutschen Aktuarvereinigung e.V.

am 13. Mai 2022

#### *Hinweise:*

- Als Hilfsmittel sind Seminarunterlagen und Aufgaben in Papierform, handschriftliche Notizen im Rahmen der normalen Schulung sowie ein nicht programmierbarer Taschenrechner zugelassen.
- Die Gesamtpunktzahl beträgt 180 Punkte. Die Klausur ist bestanden, wenn mindestens 90 Punkte erreicht werden.
- Bitte prüfen Sie die Ihnen vorliegende Prüfungsklausur auf Vollständigkeit. Die Klausur mit Lösungen besteht aus 28 Seiten.
- Mit Ausnahme der MC-Fragen sind alle Antworten zu begründen und bei Rechenaufgaben muss der Lösungsweg ersichtlich sein.
- Bei MC-Fragen entscheiden Sie für jede Aussage, ob diese wahr oder falsch ist. Eine falsche Antwort gibt 0 Punkte; Minuspunkte werden nicht vergeben. Schreiben Sie Ihre Antworten auf die Lösungsblätter, die Sie abgeben. Eine Begründung ist nicht erforderlich.

#### *Mitglieder der Prüfungskommission:*

Prof. Torsten Becker, Dr. Richard Herrmann,  
Prof. Christian Heumann, Dr. Stefan Pilz,  
Prof. Viktor Sandor, Dr. Dominik Schäfer

**Aufgabe 1.** [Lebensdauermodelle] [30 Punkte]

Ein Lebensversicherungsunternehmen ermittelt die Rechnungsgrundlagen für eine Todesfallversicherung gegen Einmalbeitrag (zur Vereinfachung: Eintrittsalter: 34, Dauer: 3 Jahre, Unisex-Rechnungsgrundlagen). Das im Lebensversicherungsunternehmen vorliegende statistische Material zur Sterblichkeit ergibt folgende Ergebnisse aus einem geschlossenen Bestand.

Im Zeitraum 1.1.2020 bis 31.12.2021 werden für die Geburtsjahre 1984 bis 1986 folgende Bestände untersucht und Todesfälle festgestellt:

Geburtsjahr	Bestand	Todesfälle
1984	1.100	19
1985	1.400	24
1986	1.900	29

- (a) [11 Punkte] Ermitteln Sie die empirischen Sterbehäufigkeiten  $q_x$  für jedes der Alter  $x = 34, \dots, 36$  nach der Sterbejahrmethode. Falls erforderlich gehen Sie von einer Gleichverteilung der Todeszeitpunkte aus. Erläutern Sie Ihre Vorgehensweise anhand einer Graphik.
- (b) [6 Punkte] Begründen Sie, dass die Anzahl der Todesfälle  $T_x$  im Alter  $x$  binomialverteilt ist und bestimmen Sie Erwartungswert  $E(T_x)$  und Varianz  $Var(T_x)$ .
- (c) [5 Punkte] Man erwartet in den ersten drei Jahren nach Einführung des neuen Tarifs einen Bestand von jeweils 1.600 Neuzugängen. Ermitteln Sie einen einheitlichen relativen Zuschlag bzw. Abschlag auf die Sterbewahrscheinlichkeiten, so dass auf Dauer mit Wahrscheinlichkeit 95 % die Anzahl der tatsächlichen Todesfälle geringer als die Anzahl der erwarteten Todesfälle ist. Verwenden Sie die  $q_x$  aus Teilaufgabe a). Falls Sie in Teilaufgabe a) keine Lösung haben, verwenden Sie  $q_{34} = 0,011, q_{35} = 0,012, q_{36} = 0,014$ .

Die Quantile der Standard-Normalverteilung betragen

$p$	0,900	0,950	0,975	0,990
$u_p$	1,28	1,64	1,96	2,33



### MC-Fragen

Bei den beiden folgenden Teilaufgaben entscheiden Sie für jede Aussage, ob diese wahr oder falsch ist. Eine falsche Antwort gibt 0 Punkte; Minuspunkte werden nicht vergeben. Schreiben Sie Ihre Antworten auf die Lösungsblätter, die Sie abgeben. Eine Begründung ist nicht erforderlich.

(d) [4 Punkte]

Bei der Ermittlung der Sterbehäufigkeiten

- (i) berücksichtigt die Sterbejahrmethode sämtliche Todesfälle des Beobachtungszeitraums
- (ii) berücksichtigt die Verweildauer methode nur Todesfälle von Personen, die den gesamten Beobachtungszeitraum im Bestand waren.

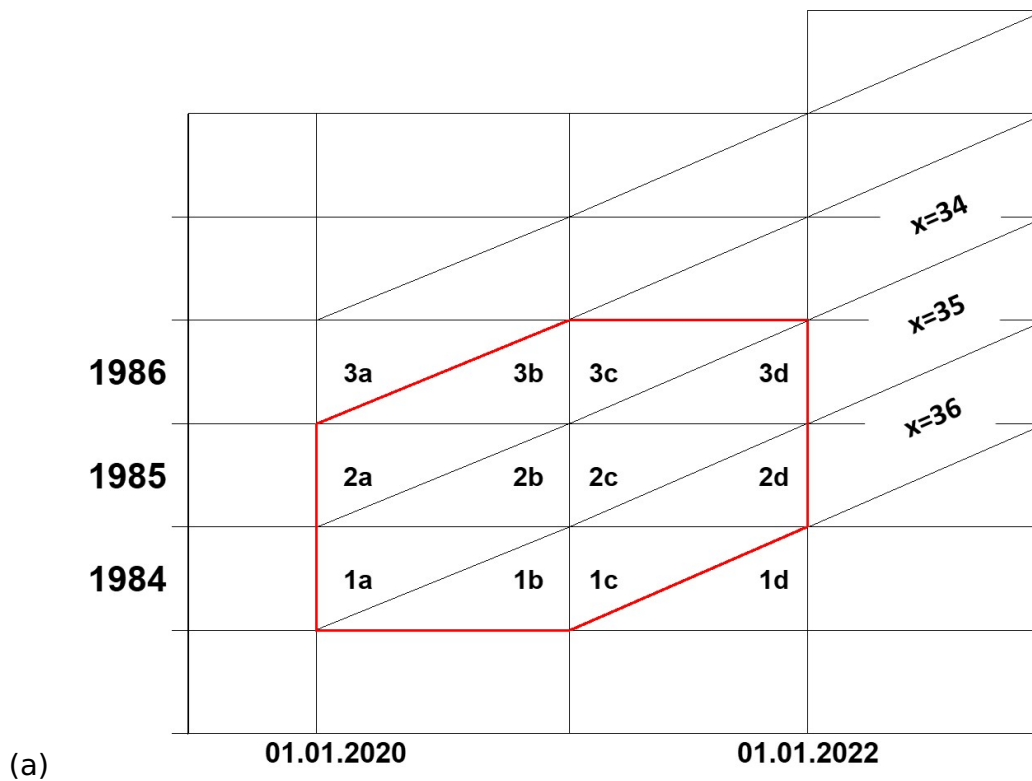
(e) [4 Punkte]

Für geschlossene Personenbestände

- (i) führen die Geburtsjahrmethode und die Sterbejahrmethode immer zu identischen Ergebnissen
- (ii) stimmen Verweildauer methode und Sterbejahrmethode überein, wenn nur die Geburtsjahre ausgewertet werden, deren Todesfälle ausschließlich in dem Beobachtungszeitraum stattfinden können.



## Lösung



Geburtsjahr	Alter x	$L_x$ am 1.1.2020	$T_x$ [1.1.2020, 1.1.2022)	Bereich	Anteil
1984	35	} 1100	19	1a	25%
	36		19	1b und 1c	50%
1985	34	} 1400	24	2a	25%
	35		24	2b und 2c	50%
	36		24	2d	25%
1986	34	} 1900	29	3b und 3c	50%
	35		29	3d	25%

Beobachtungszeitraum B = [1.1.2020, 1.1.2022)

$$\begin{aligned}
 q_{34} &= \frac{|T_{34}(B, 1985 \cup 1986)|}{\frac{1}{2}|L_{34}(B, 1985)| + |L_{34}(B, 1986)|} \\
 &= \frac{\frac{1}{4} \cdot 24 + \frac{1}{2} \cdot 29}{\frac{1}{2} \cdot 1400 + 1900} \\
 &= \frac{20,5}{2.600} \\
 &= 0,00788
 \end{aligned}$$



$$\begin{aligned} q_{35} &= \frac{|T_{35}(B, 1984 \cup 1985 \cup 1986)|}{\frac{1}{2}|L_{35}(B, 1984)| + |L_{35}(B, 1985)| + \frac{1}{2}|L_{35}(B, 1986)|} \\ &= \frac{\frac{1}{4} \cdot 19 + \frac{1}{2} \cdot 24 + \frac{1}{4} \cdot 29}{\frac{1}{2} \cdot 1100 + 1400 + \frac{1}{2} \cdot 1900} \\ &= \frac{24}{2.900} \\ &= 0,00828 \end{aligned}$$

$$\begin{aligned} q_{36} &= \frac{|T_{36}(B, 1984 \cup 1985)|}{|L_{36}(B, 1984)| + \frac{1}{2}|L_{36}(B, 1985)|} \\ &= \frac{\frac{1}{2} \cdot 19 + \frac{1}{4} \cdot 24}{1100 + \frac{1}{2} \cdot 1400} \\ &= \frac{15,5}{1.800} \\ &= 0,00861 \end{aligned}$$

(b) Bezeichne

$L_x$  die Anzahl der Lebenden im Alter  $x$

$T_x$  die Anzahl der Toten im Alter  $x$

Die Anzahl der Toten  $T_x$  lässt sich auch darstellen als Summe von unabhängigen Bernoulli-verteilten Zufallsvariablen  $X_i$ ,  $i=1, \dots, L_x$ , d.h.

$$T_x = \sum_{i=1}^{L_x} X_i \text{ mit } X_i \sim B(1; q_x).$$

Daher gilt  $T_x \sim \text{Bin}(L_x; q_x)$  mit  $E(T_x) = L_x q_x$  und  $\text{Var}(T_x) = L_x q_x (1 - q_x)$ .

(c) Für den Lebensversicherer besteht in den Altern 34 bis 36 ein Todesfallrisiko, da die Todesfallsumme immer höher als die mögliche Deckungsrückstellung ist. Um bzgl. des Schwankungsrisikos die geforderte Sicherheit von  $1 - \alpha = 0,95$  zu erreichen, muss daher für jedes Alter  $x$  der gleiche Zuschlag  $s^\alpha \geq 0$  auf die Sterbewahrscheinlichkeit  $q_x$  ermittelt werden, so dass die Wahrscheinlichkeit, dass die Anzahl der Toten kleiner gleich der erwarteten Anzahl der Toten ist, 95% beträgt.

Für alle Alter  $x = 34, \dots, 36$  muss also gelten:

$$P\left(\sum_{x=34}^{36} T_x \leq \sum_{x=34}^{36} L_x \cdot q_x \cdot (1 + s^\alpha)\right) \stackrel{!}{=} 0,95 \quad (*)$$

mit

$s^\alpha$  relativer Zuschlag auf die Sterbewahrscheinlichkeiten

Wegen (b) ist Gleichung (\*) äquivalent zu

$$P\left(\frac{\sum_{x=34}^{36} T_x - E\left(\sum_{x=34}^{36} T_x\right)}{\sqrt{\text{Var}\left(\sum_{x=34}^{36} T_x\right)}} \leq \frac{\sum_{x=34}^{36} L_x \cdot q_x \cdot (1 + s^\alpha) - \sum_{x=34}^{36} L_x \cdot q_x}{\sqrt{\sum_{x=34}^{36} L_x \cdot q_x \cdot (1 - q_x)}}\right) \stackrel{!}{=} 0,95$$



Da die Zufallsvariablen  $T_x = \sum_{i=34}^{36} X_i$  näherungsweise normal verteilt sind, ist die Zufallsvariable

$$\frac{\sum_{x=34}^{36} T_x - E\left(\sum_{x=34}^{36} T_x\right)}{\sqrt{\text{Var}\left(\sum_{x=34}^{36} T_x\right)}}$$

näherungsweise standardnormalverteilt und es gilt

$$\frac{\sum_{x=34}^{36} L_x \cdot q_x \cdot (1 + s^\alpha) - \sum_{x=34}^{36} L_x \cdot q_x}{\sqrt{\sum_{x=34}^{36} L_x \cdot q_x \cdot (1 - q_x)}} = \frac{\sum_{x=34}^{36} L_x \cdot q_x \cdot s^\alpha}{\sqrt{\sum_{x=34}^{36} L_x \cdot q_x \cdot (1 - q_x)}} = u_{0,95} \approx 1,64$$

$$\Leftrightarrow s^\alpha = 1,64 \cdot \frac{\sqrt{\sum_{x=34}^{36} L_x \cdot q_x \cdot (1 - q_x)}}{\sum_{x=34}^{36} L_x \cdot q_x} \approx 0,2699$$

Bei Vorgabe der  $q_x$  gem. Teilaufgabe c) gilt  $s^\alpha \approx 0,2187$

- (d) (i) wahr  
(ii) falsch
- (e) (i) falsch  
(ii) wahr



**Aufgabe 2.** [Deskriptive Statistik] [30 Punkte]

(a) [6 Punkte] Bei den folgenden MC-Fragen entscheiden Sie für jede Aussage, ob diese wahr oder falsch ist. Eine falsche Antwort ergibt 0 Punkte; Minuspunkte werden nicht vergeben. Schreiben Sie Ihre Antworten auf die Lösungsblätter, die Sie abgeben. Eine Begründung ist nicht erforderlich.

- A Gegeben seien zwei Stichproben  $x_1, \dots, x_n$  und  $y_1, \dots, y_m$  mit  $n \neq m$ . Dann sind die beiden Stichproben unabhängig.
- B Gegeben seien zwei Stichproben  $x_1, \dots, x_n$  und  $y_1, \dots, y_m$  mit  $n = m$  mit paarweise verschiedenen  $x_i$ ,  $i = 1, \dots, n$ . Dann besteht der Q-Q-Plot für zwei Stichproben aus  $n$  Punkten.
- C Gegeben sei der Q-Q-Plot für zwei Stichproben  $x_1, \dots, x_n$  und  $y_1, \dots, y_m$  mit  $n < m$ . Dann kann man den Boxplot von  $y_1, \dots, y_m$  zeichnen.

In (b)-(g) werden die beigefügten Abbildungen verwendet. Sie beschreiben die bivariate Stichprobe  $(x_1, y_1), \dots, (x_n, y_n)$  bzw. die daraus entstehenden einzelnen Stichproben  $x := (x_1, \dots, x_n)$ ,  $y := (y_1, \dots, y_n)$ .

- (b) [4 Punkte] Gegeben seien die Plots in Abbildung 1. Es handelt sich um ein Streudiagramm von  $(x_1, y_1), \dots, (x_n, y_n)$  und einen Zweistichproben Q-Q-Plot von  $x = (x_1, \dots, x_n)$  und  $y = (y_1, \dots, y_n)$ . Welche der beiden Abbildungen ist das Streudiagramm, welches der Q-Q-Plot? Begründen Sie Ihre Antwort. Erläutern Sie auch bei beiden Abbildungen die Werte auf den Achsen.
- (c) [4 Punkte] Erläutern Sie die Fragestellungen, die man mit Hilfe dieser beiden Graphiken in Abbildung 1 untersuchen kann.
- (d) [2 Punkte] Wie beurteilen Sie die Annahme, dass die Stichproben  $x$  und  $y$  in Abbildung 1 unabhängig sind? Begründen Sie Ihre Antwort qualitativ.
- (e) [8 Punkte] Die Abbildung 2 enthält den einfachen Boxplot von  $x$  oder  $y$ . Erstellen Sie einen einfachen Box-Plot für den nicht dargestellten Datensatz.
- (f) [2 Punkte] In der Abbildung 3 sind die Normal Q-Q Plots der Stichproben  $x$  bzw.  $y$  gegeben. Wie beurteilen Sie die Annahme der univariaten Normalverteilung für die einzelnen Stichproben  $x$  und  $y$ ? Begründen Sie Ihre Antwort qualitativ.
- (g) [4 Punkte] Wie beurteilen Sie die Annahme, dass die beiden Stichproben  $x$  und  $y$  identisch verteilt sind? Begründen Sie Ihre Antwort qualitativ, nennen Sie mindestens zwei Gründe Ihrer Überlegung.

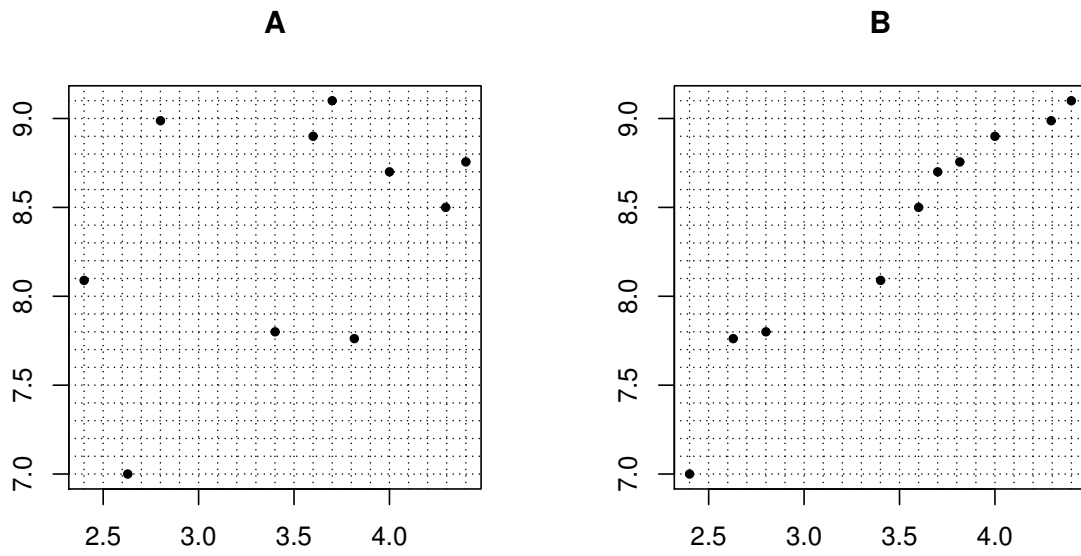


Abbildung 1: Zwei Plots zu Aufgabe 2 (b)-(d)



### Einfacher Boxplot von x oder y

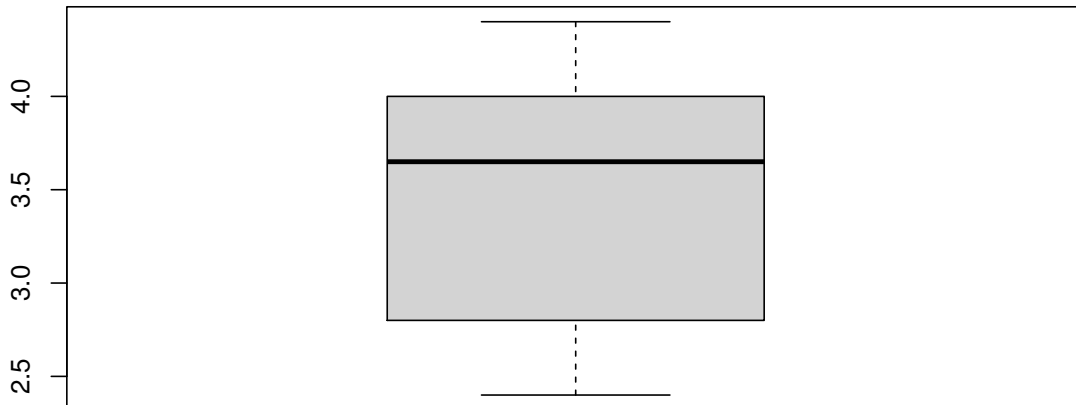


Abbildung 2: Boxplot für x oder y zu Aufgabe 2 (e)

### Lösung

(a) [6 Punkte = je 2] Begründungen sind nicht erforderlich.

A: Falsch, die Größe der Stichproben geben keine Information über Abhängigkeiten.

B: Richtig, es werden die Punktepaare  $(x_{(k)}, y_{(k)})$ ,  $k = 1, \dots, n$  gezeichnet.

C: Falsch, aus den empirischen Quantilen von  $y_1, \dots, y_m$  können die Daten nicht rekonstruiert werden.

(b) [4 Punkte = je 2 für die Beschreibung und 2 für die richtige Zuordnung] Die Abbildung A ist ein Streudiagramm, es werden die Punkte  $(x_i, y_i)$ ,  $i = 1, \dots, n$  eingezeichnet. Die Abbildung B ist ein Q-Q-Plot, es werden die Punkte  $(x_{(i)}, y_{(i)})$ ,  $i = 1, \dots, n$  eingezeichnet. Das sind die Ordnungsstatistiken von x und y, also sind die x-Werte aufsteigend sortiert, ebenso die y-Werte.

(c) [4 Punkte = je 2] Q-Q-Plot: sind die Verteilungen von x und y gleich?  
 Streudiagramm: sind Zusammenhänge von x, y zu erkennen?

(d) [2 Punkte] Es sind keine Zusammenhänge erkennbar, weder lineare noch nicht lineare. Die Punktepaare scheinen regellos zu schwanken.

(e) [8 Punkte=4 für die Werte, 4 für den Plot] Dargestellt ist der Boxplot von x. Die fünf für den Boxplot von y benötigten Werte kann man ablesen. Maximum und Minimum sind

$$y_{max} = 9, 1 \text{ und } y_{min} = 7.$$

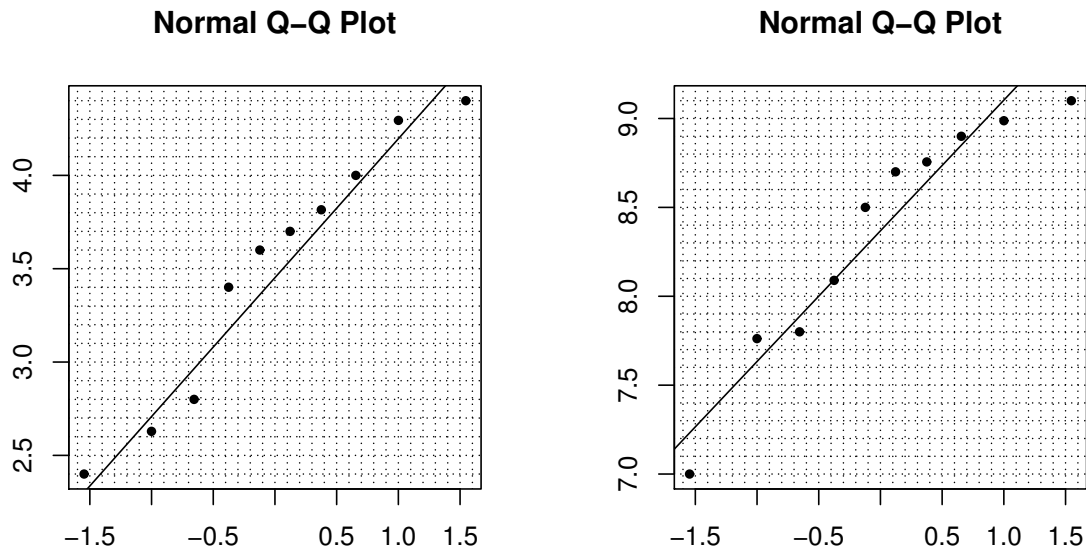


Abbildung 3: Zwei Q-Q-Normalplots zu Aufgabe 2 (f)

Das obere und untere Quartil ergeben sich beispielsweise aus

$$y_p = y_{(\lfloor np \rfloor + 1)} \text{ falls } np \notin \mathbb{N},$$

da  $n = 10$  gilt, also

$$y_{0,25} = y_{(3)} = 7,8 \text{ und } y_{0,75} = y_{(8)} = 8,9.$$

Für den Median wählen wir

$$\frac{1}{2}(y_{(5)} + y_{(6)}) = \frac{1}{2}(8,5 + 8,7) = 8,6.$$

Damit ergibt sich der Boxplot in Abb. 4.

- (f) [2 Punkte] Die beiden Punktwolken liegen um die Gerade, die Annahme der Normalverteilung erscheint plausibel.
- (g) [4 Punkte=je 2 für jede der beiden Aspekte]  $x$  und  $y$  sind vermutlich nicht identisch verteilt, weil die Erwartungswerte sich unterscheiden dürften. Das sieht man in den beiden Q-Q-Plots in Abbildung 3, die Achsenabschnitte der Anpassungsgeraden sind ca. 3,45 bzw. 8,35. Auch die beiden Boxplots stützen die Annahme gleicher Verteilungen nicht, da

$$x_{(10)} = 4,4 < y_{(1)} = 7.$$

Die obere Antenne von  $x$  liegt unter der unteren Antenne von  $y$ , vgl. Abb. 5.

### Einfacher Boxplot von y

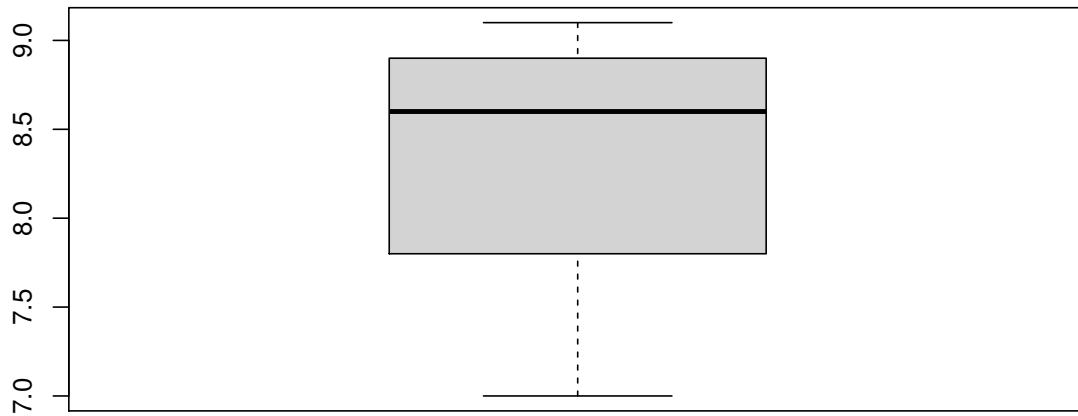


Abbildung 4: Boxplot von y

### Einfacher Boxplot von x und y

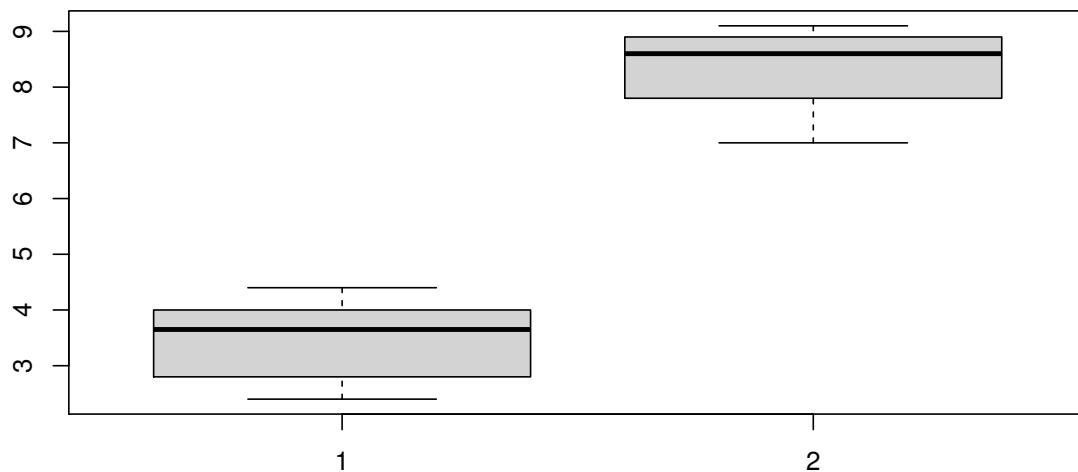


Abbildung 5: Boxplot von x und y

**Aufgabe 3.** [Induktive Statistik] [36 Punkte]

*Hinweis: Bei allen Ergebnissen genügen 2 Nachkommastellen, als Dezimaltrennzeichen in der Programmausgabe unten wird der Punkt verwendet.*

Ein Versicherer möchte seinen Bestandskunden eine neue Reiseversicherung mit Covid-Schutz anbieten. Aufgrund vorhandener historischer Bestandsdaten soll mittels eines statistischen Modells herausgefunden werden, welchen Einfluss bestimmte Kundenmerkmale auf den Abschluss einer Reiseversicherung hatten (die Zielvariable ist also Reiseversicherung). Es werden folgende Merkmale betrachtet:

- Reiseversicherung: hat der Kunde schon mal eine Reiseversicherung abgeschlossen? (Ausprägungen: 1 für ja, 0 für nein)
- Alterkat: Alter des Kunden in 4 Kategorien (1,2,3,4), die für (geordnete) Altersintervalle stehen.
- Vielflieger (Ausprägungen: 1 für Vielflieger, 0 sonst)
- Monatliches Einkommen (*in Tausend Euro*)

Die statistische Software *R* liefert folgende Ausgabe für das geschätzte Modell:

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )	
(Intercept)	-3.583428	0.186148	-19.250	< 2e-16	***
Alterkat2	0.006244	0.152893	0.041	0.967	
Alterkat3	1.050184	0.128501	8.173	3.02e-16	***
Alterkat4	0.952924	0.218791	4.355	1.33e-05	***
Vielflieger	?	0.114483	5.982	2.20e-09	***
MonatlichesEinkommen	0.567195	0.039416	14.390	< 2e-16	***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 2590.5 on 1986 degrees of freedom  
Residual deviance: 2132.1 on 1981 degrees of freedom  
AIC: 2144.1

- (a) [2 Punkte] Welche Verteilung wählen Sie für die Zielvariable *Y*? (kurze Begründung).



- (b) [2 Punkte] Die Anzahl der Bestandskunden im Regressionsmodell sei  $n$ . Wie groß ist  $n$  in obigem Modell? (kurze Begründung).
- (c) [2 Punkte] Berechnen Sie die Punktschätzung für das Merkmal Vielflieger.
- (d) [4 Punkte] Für das obige Modell wurde die logistische Linkfunktion verwendet.
- (i) [2 Punkte] Geben Sie das Modell unter Verwendung dieser Linkfunktion und der geschätzten Koeffizienten an.
  - (ii) [1 Punkt] Wie lautet das Modell, wenn statt der Linkfunktion die Responsefunktion verwendet werden soll?
  - (iii) [1 Punkt] Stellen Sie die Likelihood als Funktion der Erwartungswerte  $\mu_i = P(Y_i = 1)$  der Zufallsvariablen  $Y_i$ ,  $i = 1, \dots, n$ , dar.
- (e) [2 Punkte] Für das Merkmal Alterkat wurde die Dummykodierung verwendet. Welche Kategorie ist die Referenzkategorie? Welcher Regressionskoeffizient ergibt sich für die Referenzkategorie?
- (f) [4 Punkte]
- (i) [2 Punkte] Wie lautet der gesamte kodierte Vektor der Kovariablen (inklusive Intercept) für einen Kunden in Alterskategorie 3, der kein Vielflieger ist und 2500 Euro im Monat verdient?
  - (ii) [2 Punkte] Wie hoch ist die geschätzte Wahrscheinlichkeit, dass dieser Kunde eine Reiseversicherung abschließt?
- (g) [3 Punkte] Interpretieren Sie die geschätzten Koeffizienten für das Merkmal Alterkat. Verwenden Sie Odds Ratios für den Vergleich mit der Referenzkategorie.
- (h) [2 Punkte] Interpretieren Sie den geschätzten Koeffizienten für das Merkmal MonatlichesEinkommen.
- (i) [3 Punkte] Geben Sie die maximale Wahrscheinlichkeit für den Abschluss einer Reiseversicherung für einen Kunden mit 2000 Euro monatlichem Einkommen an. *Hinweis: wenn Sie den Koeffizienten für Vielflieger nicht berechnen können, wählen Sie 0.75 als Ersatzwert.*
- (j) [3 Punkte] Das Merkmal MonatlichesEinkommen ist als metrisches Merkmal linear in das Modell aufgenommen worden. Nennen Sie drei weitere Möglichkeiten.



(k) [9 Punkte] Das obige Modell wird erweitert um die Interaktion von Alterkat und Vielflieger. Der AIC-Wert dieses Modells ist 2088.40. Führen Sie einen Likelihood-Quotienten-Test für die Hypothesen

$H_0$ : alle Koeffizienten der Interaktion von Alterkat und Vielflieger sind 0

versus

$H_1$ : mindestens ein Koeffizient der Interaktion ist ungleich 0

durch.

(i) [2 Punkte] Formulieren sie  $H_0$  und  $H_1$  konkret mit der korrekten Anzahl an Koeffizienten.

(ii) [4 Punkte] Berechnen Sie die Teststatistik.

(iii) [3 Punkte] Geben sie die Testentscheidung an (kurze Begründung) und interpretieren Sie das Ergebnis.

*Hinweis:* Kritische Werte der  $\chi^2_{df}$ -Verteilung mit  $df$  Freiheitsgraden sind

$$\chi^2_1 = 3.84 \quad \chi^2_2 = 5.99 \quad \chi^2_3 = 7.81 \quad \chi^2_4 = 9.49 .$$

### Lösung Aufgabe 3

(a) Bernoulli-Verteilung, da die Zielvariable binär ist.

(b)  $n = 1987$ . Begründung: Null deviance bzw. Residual deviance sind auf Basis des Modells nur mit Intercept (1 Parameter) bzw. des angegebenen Modells (hier: 6 Parameter incl. Intercept) berechnet.  $n$  errechnet sich als angegebene Freiheitsgrade (degrees of freedom) plus Anzahl der Freiheitsgrade.

(c)  $z = \hat{\beta} / \hat{\sigma}_{\hat{\beta}}$ , d.h.  $\hat{\beta} = z \cdot \hat{\sigma}_{\hat{\beta}}$ :

$$\hat{\beta}_{\text{Vielflieger}} = 0.114483 \cdot 5.982 = 0.6848.$$

(d) Linkfunktion ist der logit-Link.

(i) Hinweis: für  $\beta$  werden die Punktschätzungen eingesetzt, d.h. eigentlich stehen hier  $\hat{\mu}_i$  und  $\hat{\beta}$ , das "Dach" wird der Einfachheit halber weggelassen.

$$\begin{aligned} \log\left(\frac{\mu_i}{1-\mu_i}\right) &= x_i^T \hat{\beta} \\ &= -3.58 + 0.006 * \text{Alterkat2} + 1.05 * \text{Alterkat3} + 0.95 * \text{Alterkat4} + 0.68 * \text{Vielflieger} + 0.57 * \text{MonatlichesEinkommen} \end{aligned}$$



(ii) Die Responsefunktion ist die Umkehrfunktion

$$\mu_i = P(Y_i = 1|x_i) = \frac{\exp(x_i^T \beta)}{1 + \exp(x_i^T \beta)}$$

(iii) Die Likelihood lautet

$$L = \prod_{i=1}^{1987} \mu_i^{y_i} (1 - \mu_i)^{1-y_i}$$

(Wird auch voll gewertet, wenn nur  $n$  statt 1987 angegeben wird).

(e) Alterkat1, also die Alterskategorie 1 ist die Referenzkategorie. Diese hat praktisch den Koeffizienten 0, da alle Dummies 0 sind.

(f) (i)  $x_i = [1, (0, 1, 0), 0, 2.5]$ . Achtung: Angabe des monatlichen Einkommens ist in Tausend Euro.

(ii) Prädiktor:

$$x_i' \beta = -3.58 + 0.006 * 0 + 1.05 * 1 + 0.95 * 0 + 0.68 * 0 + 0.57 * 2.5 = -1.105$$

$$P(Y_i = 1) = \exp(-1.105) / [1 + \exp(-1.105)] = 0.25$$

(g) Odds für Abschluss der Reiseversicherung erhöht sich in allen Alterskategorien (2,3,4) im Vergleich zur Referenzkategorie, da alle Koeffizienten positiv sind. Odds Ratio =  $\exp(\beta)$ . Beispiel für Odds Ratio: Odds (Alterskategorie 3) / Odds (Alterskategorie 1, Referenz) =  $\exp(1.05) = 2.86$ , d.h. Odds für Abschluss der Reiseversicherung ist um Faktor 2.86 höher in Alterskategorie 3 im Vergleich zu Alterskategorie 1.

(h) Erhöht sich das monatliche Einkommen um 1 (also um 1000 Euro!), dann erhöht sich der Odds um den Faktor  $\exp(0.57) = 1.77$ .

(i) Maximale Wahrscheinlichkeit ist bei Alterskategorie 3, Vielflieger (1): Prädiktor:

$$x_i' \beta = -3.58 + 0.006 * 0 + 1.05 * 1 + 0.95 * 0 + 0.68 * 1 + 0.57 * 2.0 = -0.71$$

$$P(Y_i = 1) = \exp(-0.71) / [1 + \exp(-0.71)] = 0.33$$

(j) Linear + quadratisch, glatte Funktion (GAM), Bildung von Kategorien.

(k) (i) Likelihood-Quotienten-Test.  $H_0 : \beta_{IA1} = \beta_{IA2} = \beta_{IA3} = 0$  versus  $H_1$  : mindestens einer dieser 3 Koeffizienten ist ungleich 0.



- (ii) LQT benötigt die Differenz der beiden log-likelihoods (Modell mit IA - Modell ohne IA), multipliziert mit dem Faktor 2. Beim AIC gilt: „Kleiner ist besser“.  $AIC(\text{mit Interaktion}) = 2088.4$ ,  $AIC(\text{ohne IA}) = 2144.1$ .

$$AIC = -2 \loglik + 2 \cdot \text{Anzahl geschätzter Parameter}$$

Modell ohne IA hat 6 Parameter, Modell mit IA hat 9 Parameter (da Alterkat 3 Dummies hat und Vielflieger 1 Dummy,  $3 \cdot 1 = 3$ ). Damit ist die LQT-Statistik gleich 2\* Differenz der logliks:

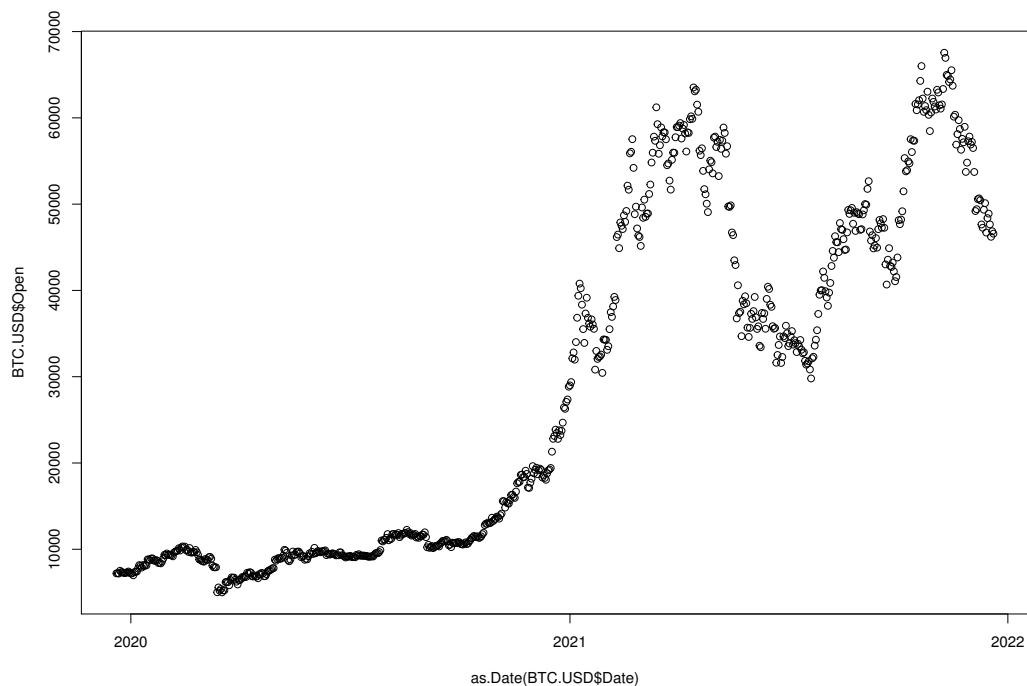
$$LQT = -(2088.4 - 2 * 9) - (-(2144 - 2 * 6)) = 61.6 .$$

- (iii) Bei 3 Freiheitsgraden ist der kritische Wert  $\chi_3^2 = 7.81$ . Damit wird  $H_0$  zugunsten von  $H_1$  abgelehnt, d.h. die IA ist statistisch signifikant.



**Aufgabe 4.** [Zeitreihenanalyse] [24 Punkte]

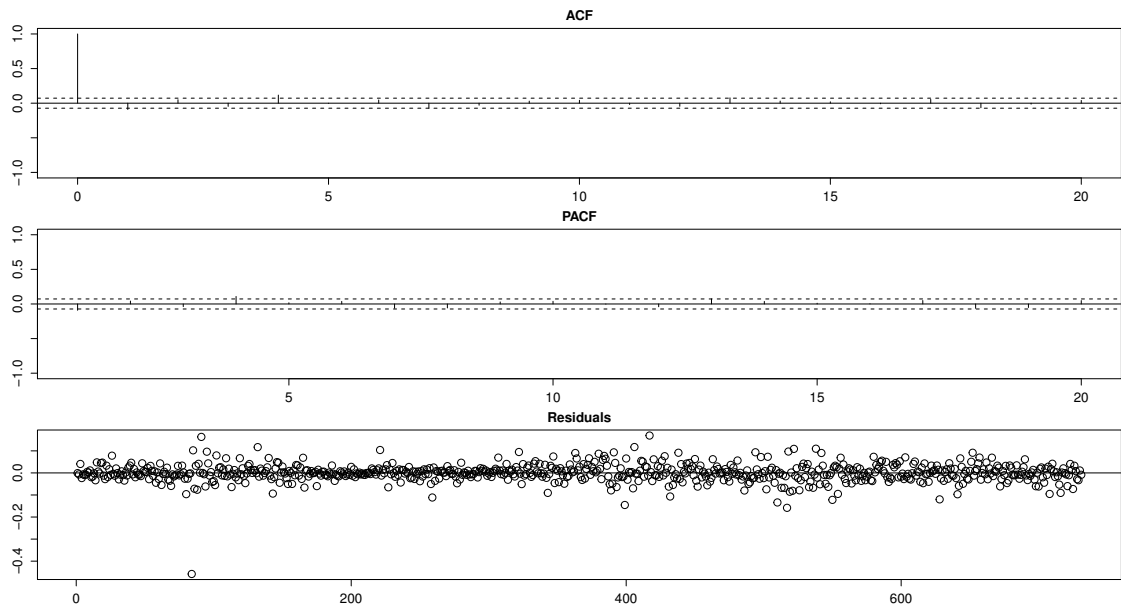
- (a) [3 Punkte] Nennen Sie drei Annahmen für schwach stationäre Prozesse.
- (b) [6 Punkte] Bei den folgenden MC-Fragen entscheiden Sie für jede Aussage, ob diese wahr oder falsch ist. Eine falsche Antwort ergibt 0 Punkte; Minuspunkte werden nicht vergeben. Schreiben Sie Ihre Antworten auf die Lösungsblätter, die Sie abgeben. Eine Begründung ist nicht erforderlich.
- A) Der random walk  $y_t = y_{t-1} + u_t$ ,  $u_t$  iid  $N(0,1)$ , ist stationär.
  - B) Bei AR(1) Modellen klingt die Autokorrelation als Funktion des lags linear ab.
  - C) Die partielle Autokorrelation ist ein Spezialfall der Autokorrelation.
- (c) (i) [3 Punkte] Die folgende Grafik enthält den täglichen Bitcoin - Kurs (in US Dollar, Dezember 2019 - Dezember 2021). Nennen Sie drei Charakteristika, die geeignet sind den Verlauf der Zeitreihe zu beschreiben!



- (ii) [1 Punkt] Nennen Sie einen Grund warum die Zeitreihe nicht stationär ist.
- (iii) [4 Punkte] Welche Transformationen bzw. Bereinigungen schlagen Sie vor um eine stationäre Zeitreihe zu erreichen? (mit kurzer Begründung)
- (iv) [3 Punkte] Für die bereinigte Zeitreihe wurde die Autokorrelationsfunktion bzw. die partielle Autokorrelationsfunktion bzw. die Residuen berechnet



und visualisiert. Interpretieren Sie die Grafiken (jeweils 1 Punkt).



- (v) [2 Punkte] Sie sollen für die bereinigte Zeitreihe ein Modell schätzen. Machen Sie einen Vorschlag und begründen ihn kurz!
- (vi) [2 Punkte] Wie beurteilen Sie den Einsatz von diesem Zeitreihenmodell für die kurzfristige Prognose (mit kurzer Begründung)?

## Lösung

(a) [3 Punkte]

- (i) Konstanter Erwartungswert.
- (ii) Konstante Varianz.
- (iii) Die Autokorrelationsfunktion hängt nur von dem Lag ab, nicht aber von den jeweiligen Zeitpunkten.

(b) [6 Punkte]

- A) falsch, da die white noise Komponente einen positiven Erwartungswert hat.
  - B) falsch; Gegenbeispiel, negativer Parameter  $\alpha_1$  führt zu Oszillation.
  - C) falsch; es handelt sich um verwandte Konzepte der Abhängigkeit: bedingte versus marginale Betrachtung.
- (c) (i) [3 Punkte] Die Werte der Zeitreihe sind positiv. Es gibt einen Trend. Die Varianz wächst über die Zeit. Die Aussage, dass die Zeitreihe keine saisonalen Komponenten hat wird auch akzeptiert.
- (ii) [1 Punkte] Die Zeitreihe ist nicht stationär; Gründe: Trend, wachsende Volatilität (ein Grund genügt).
- (iii) [4 Punkte] Transformationen
- i. Trendbereinigung durch Differenzen erster Ordnung.
  - ii. log-Transformation mit dem Ziel die Varianz zu stabilisieren.
- (iv) [3 Punkte] Interpretation
- i. Autokorrelationsfunktion: Die Werte sind für alle Lags  $> 0$  nahe dem Wert Null.
  - ii. Partielle Autokorrelationsfunktion: Die Werte sind für alle Lags nahe dem Wert Null.
  - iii. Residuen-Plot: Die Residuen streuen um den Wert Null; es gibt einen Wert mit -0.4



- (v) [2 Punkte] Aufgrund der Werte der (partiellen) Autokorrelationsfunktion empfiehlt sich ein white noise Modell.
- (vi) [2 Punkte] Da es keine nennenswerten (partiellen) Autokorrelationen gibt, sind die Zeitreihenmodelle nicht hilfreich für die kurzfristige Prognose.

**Aufgabe 5.** [Credibility-Theorie, 30 Punkte]

(a) [2 Punkte] Bei den folgenden Fragen entscheiden Sie für jede Aussage, ob diese wahr oder falsch ist. Eine falsche Antwort ergibt 0 Punkte, Minuspunkte werden nicht vergeben. Schreiben Sie Ihre Antworten auf die Lösungsblätter, die Sie abgeben. Eine Begründung ist nicht erforderlich.

- A Im Bayes'schen Credibility-Modell stellt die a-posteriori-Verteilung die Verteilung des Strukturparameters bedingt auf die Beobachtungen der Schadenhöhen dar.
- B Die allgemeine Credibility-Prämie lässt sich immer in der Form  $H^* = z\bar{X} + (1 - z)E(X)$  darstellen, mit dem empirischen Mittel  $\bar{X}$  und dem Erwartungswert  $E(X)$  des Schadens, sowie dem Credibility-Faktor  $z$ .

(b) [22 Punkte] Der Jahresgesamtschaden  $X$  in einer bestimmten Region werde durch eine Paretoverteilung mit Dichte

$$f_{X|\Theta=\vartheta}(x) = \vartheta x^{-(\vartheta+1)} \quad (\text{für } x \geq 1)$$

und Parameter  $\vartheta > 1$  beschrieben. Der Parameter der Paretoverteilung wird als Strukturparameter  $\Theta$  in einem Bayes'schen Credibility-Modell betrachtet.

In Regionen mit hohem Unwetterpotenzial gelte der Parameter  $\vartheta = 3$ . In Regionen mit niedrigem Unwetterpotenzial gelte dagegen  $\vartheta = 7$ . Zunächst geht der Versicherer davon aus, dass die betrachtete Region mit einer Wahrscheinlichkeit  $p = 25\%$  hohes Unwetterpotenzial besitze.

In den letzten 5 Jahren wurden in der betrachteten Region die Jahresgesamtschäden

$x_1$	$x_2$	$x_3$	$x_4$	$x_5$
1,4	1,6	1,1	1,5	2,2

beobachtet.

- (b1) [3 Punkte] Berechnen Sie ausgehend von der gegebenen Dichte die Funktion  $H(\vartheta) = E(X|\Theta = \vartheta)$ .
- (b2) [2 Punkte] Geben Sie die Zähldichte  $f_\Theta(\vartheta) = P(\Theta = \vartheta)$  der a-priori-Verteilung an.

- (b3) [7 Punkte] Zeigen Sie, dass die gemeinsame Dichte  $g(x_1, \dots, x_n, \vartheta)$  der Schäden und des Strukturparameters durch

$$g(x_1, \dots, x_n, \vartheta) = \begin{cases} 0,25 \cdot 3^5 \cdot \exp(-4 \cdot \sum_{i=1}^5 \ln(x_i)) & \text{für } \vartheta = 3 \\ 0,75 \cdot 7^5 \cdot \exp(-8 \cdot \sum_{i=1}^5 \ln(x_i)) & \text{für } \vartheta = 7 \end{cases}$$

gegeben ist und berechnen Sie deren Wert für  $\vartheta = 3$  und  $\vartheta = 7$  (Ergebnis als Prozentwerte mit zwei Nachkommastellen angeben).

- (b4) [5 Punkte] Berechnen Sie die Werte der Zähldichte  $f_{\Theta|X_1=x_1, \dots, X_n=x_n}(\vartheta) = P(\Theta = \vartheta | X_1 = x_1, \dots, X_n = x_n)$  der a-posteriori-Verteilung des Strukturparameters  $\Theta$  unter den gegebenen Schadenbeobachtungen (Ergebnis als Prozentwerte mit zwei Nachkommastellen angeben).

Spricht diese eher für eine Region mit hohem Unwetterpotenzial?

- (b5) [5 Punkte] Berechnen Sie den Wert der allgemeinen Credibility-Prämie  $H^*$  für die betrachtete Region.

(Hinweis: Falls Sie (b3) nicht gelöst haben, können Sie davon ausgehen, dass es sich unter der a-posteriori-Verteilung mit Wahrscheinlichkeit 90% um eine Region mit hohem Unwetterpotenzial handelt.)

- (c) [6 Punkte] In 7 Regionen wurden in den vergangenen 5 Jahren folgende Beobachtungen von  $X$  gemacht und die zugehörigen empirischen Mittelwerte und Varianzen pro Zeile bzw. Spalte ermittelt:

	Jahr 1	Jahr 2	Jahr 3	Jahr 4	Jahr 5	Mittelwert	Varianz
Region 1	1,4	1,6	1,1	1,5	2,2	1,5600	0,1630
Region 2	1,3	1,2	2,7	1,7	1,5	1,6800	0,3620
Region 3	1,4	1,1	1,0	1,0	1,2	1,1400	0,0280
Region 4	1,1	1,2	1,2	1,1	1,2	1,1600	0,0030
Region 5	1,1	1,0	1,2	1,9	1,6	1,3600	0,1430
Region 6	1,0	1,0	1,5	1,7	1,9	1,4200	0,1670
Region 7	1,1	1,2	1,2	1,1	1,2	1,1600	0,0030
					Mittelwert	1,3543	0,1241
					Varianz	0,0457	

Berechnen Sie mit diesen Werten unter den Annahmen des *Bühlmann-Modells* einen Schätzer für den Credibility-Faktor sowie die Credibility-Prämie von Region 1.

**Lösung:**

(a)

A richtig

B falsch (dies gilt nur, wenn die a-priori-Verteilung und die Schadenhöhenverteilung konjugierte Verteilungen sind)

(b)

(b1) Es ergibt sich

$$\begin{aligned} H(\vartheta) = E(X|\Theta = \vartheta) &= \int_1^{\infty} x \cdot f_{X|\Theta=\vartheta}(x) dx = \int_1^{\infty} x \cdot \vartheta x^{-(\vartheta+1)} dx \\ &= \vartheta \int_1^{\infty} x^{-\vartheta} dx = \vartheta \cdot \left[ \frac{x^{-\vartheta+1}}{-\vartheta+1} \right]_{x=1}^{\infty} = \frac{\vartheta}{\vartheta-1} \end{aligned}$$

(b2) Die a-priori-Verteilung ist auf  $\{3; 7\}$  konzentriert, mit der Zähldichte  $f_{\Theta}(3) = 25\%$  und  $f_{\Theta}(7) = 75\%$ .

(b3) Die gemeinsame Dichte ergibt sich gemäß

$$\begin{aligned} g(x_1, \dots, x_n, \vartheta) &= f_{\Theta}(\vartheta) \cdot \prod_{i=1}^n f_{X|\Theta=\vartheta}(x_i) = f_{\Theta}(\vartheta) \cdot \prod_{i=1}^n \vartheta x_i^{-(\vartheta+1)} \\ &= f_{\Theta}(\vartheta) \cdot \vartheta^n \prod_{i=1}^n \exp(-(\vartheta+1) \ln(x_i)) \\ &= f_{\Theta}(\vartheta) \cdot \vartheta^n \exp\left(-(\vartheta+1) \sum_{i=1}^n \ln(x_i)\right) \end{aligned}$$

und mit  $n = 5$  und den Werten aus (b1) ergibt sich die Behauptung. Mit  $\sum_{i=1}^5 \ln(x_i) = 2,0957$  erhält man

$$g(x_1, \dots, x_n, 3) = 0,25 \cdot 3^5 \exp(-4 \cdot 2,0957) = 1,39\%$$

$$g(x_1, \dots, x_n, 7) = 0,75 \cdot 7^5 \exp(-8 \cdot 2,0957) = 0,07\%$$

(b4) Die a-posteriori-Dichte ist

$$f_{\Theta|X_1=x_1, \dots, X_n=x_n}(\vartheta) = \frac{g(x_1, \dots, x_n, \vartheta)}{g(x_1, \dots, x_n, 3) + g(x_1, \dots, x_n, 7)}$$

so dass



$$f_{\Theta|X_1=x_1, \dots, X_n=x_n}(3) = \frac{1,39\%}{1,39\% + 0,07\%} = 95,21\%$$

$$f_{\Theta|X_1=x_1, \dots, X_n=x_n}(7) = \frac{0,07\%}{1,39\% + 0,07\%} = 4,79\%$$

Mit hoher Wahrscheinlichkeit liegt also hohes Unwetterpotenzial vor.

(b5) Für die allgemeine Credibility-Prämie gilt

$$\begin{aligned} H^* &= \sum_{\vartheta} H(\vartheta) f_{\Theta|X_1=x_1, \dots, X_n=x_n}(\vartheta) \\ &= H(3) \cdot f_{\Theta|X_1=x_1, \dots, X_n=x_n}(3) + H(7) \cdot f_{\Theta|X_1=x_1, \dots, X_n=x_n}(7) \\ &= \frac{3}{3-1} \cdot 95,21\% + \frac{7}{7-1} \cdot 4,79\% = 1,4840 \end{aligned}$$

bzw. 1,4667 wenn man mit dem Hinweis rechnet.

(c) Der Credibility-Faktor im Bühlmann-Modell ist

$$\hat{z} = 1 - \frac{\hat{E}(\sigma^2(\Theta))}{n \cdot \widehat{Var}(\bar{X})} = 1 - \frac{0,1241}{5 \cdot 0,0457} = 0,4569.$$

Damit errechnet man die Credibility-Prämie für Region 1 zu

$$H_1^{**} = \hat{z} \cdot \bar{X}_1 + (1 - \hat{z}) \cdot \hat{E}(\mu(\Theta)) = 0,4569 \cdot 1,56 + 0,5431 \cdot 1,3543 = 1,4483.$$



**Aufgabe 6.** [Stochastische Prozesse und deren Simulation] [30 Punkte]

(a) [6 Punkte] Wir betrachten die stochastische DGL

$$dX_t = \frac{1 - X_t}{1 - t} dt + dW_t, \quad 0 \leq t < 1 \quad (*)$$

mit Startwert  $X_0 \equiv 0$ . Zeigen Sie mit Hilfe des Lemmas von Ito, dass

$$X_t = t + (1 - t) \int_0^t \frac{1}{1 - u} dW_u.$$

Hinweis: Verwenden Sie den Prozess  $Y_t := \frac{X_t}{1-t}$ .

(b) [6 Punkte] Geben Sie die Verteilung von  $X_t$  für  $0 \leq t < 1$  an.

*Hinweis für die Aufgabenteile (c)-(e):* Man kann zeigen, dass für den obigen Prozess  $(X_t)_{0 \leq t < 1}$  gilt  $\lim_{t \rightarrow 1} X_t \equiv 1$ . Man kann also den Prozess auf  $[0; 1]$  sinnvoll fortsetzen durch  $X_1 := 1$ . Dies müssen Sie nicht beweisen.

(c) Es soll eine Simulation des Prozesses  $(X_t)_{0 \leq t \leq 1}$  durchgeführt werden. Dazu wird der Zeitbereich  $[0; 1]$  durch die Punkte  $t_k := \frac{k}{N}$  ( $N \in \mathbb{N}$ ,  $k = 0, \dots, N$ ) diskretisiert. Ein bekanntes Verfahren liefert näherungsweise Realisierungen eines Pfades des Prozesses durch Werte  $\hat{x}_{t_k}$ , die der Rekursion

$$\hat{x}_{t_k} = \hat{x}_{t_{k-1}} + a(t_{k-1}, \hat{x}_{t_{k-1}}) \cdot \frac{1}{N} + b(t_{k-1}, \hat{x}_{t_{k-1}}) \cdot \frac{1}{\sqrt{N}} \cdot z_k \quad (*)$$

mit  $\hat{x}_0 = 0$  genügen ( $k = 1, \dots, N$ ).

(i) [1 Punkt] Um welches Verfahren handelt es sich?

(ii) [2 Punkte] Geben Sie explizit  $a(t, x)$  und  $b(t, x)$  an.

(iii) [1 Punkt] Worum handelt es sich bei den Zahlen  $z_k$ ?

(d) [6 Punkte] Simulieren Sie auf Grundlage der Rekursion (\*) für  $N = 2$  die Werte  $\hat{x}_{0,5}$  und  $\hat{x}_1$ . Zur Bestimmung der Werte  $z_1, z_2$  verwenden Sie eine Methode, die auf  $\mathcal{U}(0; 1)$ -verteilten Zufallszahlen basiert, hier speziell  $u_1 = 0,134$  und  $u_2 = 0,789$ . Erläutern Sie, welche Methode Sie verwenden.

(e) [2 Punkte] Wie interpretieren Sie Ihr Ergebnis für  $\hat{x}_1$ ?

(f) [6 Punkte] Bei den folgenden MC-Fragen entscheiden Sie für jede Aussage, ob diese wahr oder falsch ist. Eine falsche Antwort ergibt 0 Punkte; Minuspunkte werden nicht vergeben. Schreiben Sie Ihre Antworten auf die Lösungsblätter, die Sie abgeben. Eine Begründung ist nicht erforderlich.



- (i) Die asymptotische Verteilung eines homogenen Markovprozesses ist ein Linkseigenvektor der Fundamentalmatrix des Prozesses.
- (ii) Ist  $(X_t)_{t \geq 0}$  ein stochastischer Prozess mit  $X_0 \equiv 0$  sowie stationären Zuwächsen und  $X_t \sim \mathcal{LN}(0, t)$  für alle  $t > 0$ , dann gilt für  $s < t$ :  $X_t - X_s \sim \mathcal{LN}(0, t - s)$ .
- (iii) Sind  $x_1, \dots, x_{2n}$  Simulationen einer Zufallsvariablen  $X$  mit  $\text{Var}(X) < \infty$ , dann liegt  $\frac{1}{2n}(x_1 + \dots + x_{2n})$  näher an  $E(X)$  als  $\frac{1}{n}(x_1 + \dots + x_n)$ .



## Lösungen

(a) Für  $dX_t = D_t dt + V_t dW_t$  und  $Y_t = f(t, X_t)$  lautet Ito's Lemma

$$dY_t = \left[ \frac{\partial f}{\partial t}(t, X_t) + \frac{\partial f}{\partial x}(t, X_t)D_t + \frac{1}{2} \frac{\partial^2 f}{\partial x^2}(t, X_t)V_t^2 \right] dt + \frac{\partial f}{\partial x}(t, X_t)V_t dW_t.$$

Mit  $D_t = \frac{1-X_t}{1-t}$  und  $V_t = 1$  sowie  $f(t, x) = \frac{x}{1-t}$  folgt

$$dY_t = \left[ \frac{X_t}{(1-t)^2} + \frac{1}{1-t} \cdot \frac{1-X_t}{1-t} + 0 \right] dt + \frac{1}{1-t} dW_t = \frac{1}{(1-t)^2} dt + \frac{1}{1-t} dW_t.$$

Daraus folgt mit  $Y_0 \equiv 0$

$$Y_t = Y_0 + \int_0^t \frac{1}{(1-u)^2} du + \int_0^t \frac{1}{1-u} dW_u = \frac{1}{1-t} - 1 + \int_0^t \frac{1}{1-u} dW_u$$

und damit

$$X_t = (1-t)Y_t = 1 - (1-t) + (1-t) \int_0^t \frac{1}{1-u} dW_u = t + (1-t) \int_0^t \frac{1}{1-u} dW_u.$$

(b) Wir verwenden: Ist  $h$  eine stetig differenzierbare Funktion, dann gilt

$$\int_0^t h(u) dW_u \sim \mathcal{N}(0, \sigma^2)$$

mit  $\sigma^2 = \int_0^t h^2(u) du$ . Da

$$\int_0^t \frac{1}{(1-u)^2} du = \frac{1}{1-t} - 1$$

folgt

$$X_t \sim \mathcal{N}\left(t, (1-t)^2 \left(\frac{1}{1-t} - 1\right)\right) = \mathcal{N}(t, t(1-t)).$$

(ci) Es handelt sich um das Euler-Verfahren.

(cii) Es gilt

$$a(t, x) = \frac{1-x}{1-t}, \quad b(t, x) = 1.$$

(ciii) Die  $z_k$  sind iid standardnormalverteilte Zufallszahlen.

(d) Die  $z_k$  können entweder mit der Inversionsmethode oder der Box-Muller-Methode erzeugt werden (eine Methode reicht):



- Die Inversionsmethode liefert (unter Zuhilfenahme eines Taschenrechners, der  $\Phi^{-1}$  auswerten kann)

$$z_1 = \Phi^{-1}(0,134) \approx -1,11, \quad z_2 = \Phi^{-1}(0,789) \approx 0,8$$

Das ergibt

$$\hat{x}_{0,5} = \hat{x}_0 + a(0, \hat{x}_0) \cdot \frac{1}{2} + b(0, \hat{x}_0) \cdot \frac{1}{\sqrt{2}} \cdot z_1 = -0,28$$

$$\hat{x}_1 = \hat{x}_{0,5} + a(0,5, \hat{x}_{0,5}) \cdot \frac{1}{2} + b(0,5, \hat{x}_{0,5}) \cdot \frac{1}{\sqrt{2}} \cdot z_2 = 1,57$$

- Box-Muller ergibt

$$z_1 = \sqrt{-2 \ln(0,134)} \cos(2\pi \cdot 0,789) \approx 0,49$$

$$z_2 = \sqrt{-2 \ln(0,134)} \sin(2\pi \cdot 0,789) \approx -1,95$$

Mit diesen Werten erhält man

$$\hat{x}_{0,5} = \hat{x}_0 + a(0, \hat{x}_0) \cdot \frac{1}{2} + b(0, \hat{x}_0) \cdot \frac{1}{\sqrt{2}} \cdot z_1 = 0,84$$

$$\hat{x}_1 = \hat{x}_{0,5} + a(0,5, \hat{x}_{0,5}) \cdot \frac{1}{2} + b(0,5, \hat{x}_{0,5}) \cdot \frac{1}{\sqrt{2}} \cdot z_2 = -0,38$$

(e)  $\hat{x}_1$  ist ungleich 1. Es handelt sich hierbei um den Diskretisierungsfehler, der bei der Euler-Methode immer auftritt, da nicht die tatsächliche Verteilung der  $X_{t_k}$  simuliert wird (daher auch das Dach über den  $x$ ).

(fi) Richtig, nämlich zum Eigenwert 0.

(fii) Richtig: Die Stationarität der Zuwächse bedeutet nämlich

$$X_t - X_s \sim X_{t-s} - X_0 = X_{t-s}.$$

(fiii) Falsch, die Folge der arithmetischen Mittel muss ja gar nicht gegen  $E(X)$  konvergieren. Und wenn sie es tut, dann muss die Konvergenz nicht monoton sein.