

Lösungen – Weitere Aufgaben zur Prüfungsvorbereitung

12.1 Zufallsvariable

Aufgabe: Geben Sie zwei reelle Zufallsvariablen X und Y , für welche gilt $E[X] > E[Y]$ und $E[X \cdot Y] \neq E[X] \cdot E[Y]$.

Lösung: Zum Beispiel:

	Y	
X	0	1
1	$\frac{1}{8}$	$\frac{1}{4}$
2	$\frac{7}{16}$	$\frac{3}{16}$

$$E[X] = \frac{3}{2}; E[Y] = \frac{9}{16}; E[X] \cdot E[Y] = \frac{27}{32}; E[X \cdot Y] = \frac{3}{4}.$$

12.2 Zufallswahl

Aufgabe: Im folgenden Algorithmus mit Eingabe $k > 0$ wählt die Funktion $random(a, b)$ einen zufälligen Integer in $[a, b]$. Modellieren Sie die Ausgabe des Algorithmus als eine Zufallsvariable Y , beschreiben Sie diese, und begründen Sie Ihre Antwort.

```

x ← 0
y ← 2k
while x < y - 1 do
  r ← random(x, y - 1)
  m ← x +  $\frac{y-x}{2}$ 
  if m ≤ r then
    x ← m
  else
    y ← m
return y

```

Lösung: Für $i = 1, \dots, k$ seien x_i und y_i die Werte der Variablen x und y jeweils am Ende einer Wiederholung und $x_0 = 0$ und $y_0 = 2^k$ deren Anfangswerte. Die Zufallsvariable R_i wird bestimmt durch den Aufruf von $random()$ in Wiederholung i ; die ZV $B_i \in \{0, 1\}$ modelliert den Ausgang des Experiments in Wiederholung i , wobei das Ereignis $B_i = 0$ dem Fall $R_i < m$ entspricht und $B_i = 1$ für $R_i \geq m$ steht. Unter Verwendung von $x + \frac{y-x}{2} = \frac{x+y}{2}$ folgt

$$\begin{aligned}
P[B_i = 0] &= P[R_i < m] = \frac{m - x}{y - x} = \dots = \frac{1}{2} \\
P[B_i = 1] &= P[R_i \geq m] = \frac{y - m}{y - x} = \dots = \frac{1}{2}
\end{aligned}$$

Die Variable x entwickelt sich nach der Rekursion

$$x_{i+1} = \begin{cases} x_i & \text{falls } B_i = 0 \\ \frac{x_i + y_i}{2} & \text{falls } B_i = 1 \end{cases}$$

oder vereinfacht als $x_{i+1} = x_i + B_i \left(\frac{x_i + y_i}{2} - x_i \right) = x_i + B_i \left(\frac{y_i - x_i}{2} \right)$. Dual dazu ist

$$y_{i+1} = \begin{cases} \frac{x_i + y_i}{2} & \text{falls } B_i = 0 \\ y_i & \text{falls } B_i = 1 \end{cases}$$

und $y_{i+1} = y_i + (1 - B_i) \frac{x_i - y_i}{2}$. Mit $y_0 - x_0 = 2^k$ folgt aus der Rekursion ebenfalls, dass $y_i - x_i = 2^{k-i}$ und damit $x_{i+1} = x_i + B_i 2^{k-i-1}$. Der Algorithmus terminiert nach k Wiederholungen, da $y_k - x_k = 1$.

Daraus folgt weiter, dass $x_k = \sum_{i=1}^k B_i 2^{k-i-1}$. Anders ausgedrückt, x_k ist uniform verteilt in $[0, 2^k - 1]$ und y_k ist uniform verteilt in $[1, 2^k]$.

12.3 Bitte anschnallen

Aufgabe: In ein Flugzeug mit r Sitzreihen steigen g Passagiere auf Geschäftsreise und f Passagiere auf dem Weg in die Ferien. Jeder Passagier wählt zufällig und uniform eine Reihe und nimmt dort Platz, wobei die Geschäftsreisenden sich auf die vorderen $2/3$ des Flugzeugs beschränken und die Ferienpassagiere nur in den hinteren $2/3$ aller Reihen Platz nehmen (r ist durch 3 teilbar).

- Wie viele Passagiere sitzen im Mitteldrittel durchschnittlich in einer Reihe?
- Angenommen g und h sind viel kleiner als \sqrt{r} . Geben Sie eine Abschätzung für die Wahrscheinlichkeit, dass im Mitteldrittel nirgends zwei oder mehr Passagiere in derselben Reihe sitzen.
- Unter derselben Annahme wie vorher, geben Sie eine Abschätzung für die Wahrscheinlichkeit, dass alle Passagiere im Flugzeug allein in einer Reihe sitzen.

Hinweis: Führen Sie Indikator-ZV ein; daraus konstruieren Sie eine ZV für die gesuchte Anzahl doppelt besetzter Reihen; auf diese wenden Sie dann zur Abschätzung die Markov-Ungleichung an.

Lösung: Sei $m = r/3$ die Grösse eines Abschnitts.

- Die erwartete Passagierzahl in den m Reihen in der Mitte ist

$$\left(\frac{g}{2} + \frac{f}{2} \right) / m = \frac{1}{2} \frac{g + f}{m}$$

- Wir geben eine Abschätzung mittels Indikator-ZV. Sei \mathcal{M} die Menge der m Reihen des Mitteldrittels. Die Auswahl einer Reihe einer Geschäftsreisenden i ist eine Zufallsvariable G_i mit uniformer Verteilung über alle $2m$ vorderen Reihen. Analog dazu bezeichnet eine Zufallsvariable F_k die Wahl des Ferienreisenden k über die hinteren $2m$ Reihen.

Wir definieren folgende Gruppen von Indikator-ZV für die Reihen im Mitteldrittel \mathcal{M} :

- M_{ij}^{GG} für das Ereignis, dass Geschäftsreisende i und j in \mathcal{M} in die gleiche Reihe

sitzen;

$$\begin{aligned}
E[M_{ij}^{GG}] &= P[(G_i = G_j) \cap G_i \in \mathcal{M} \cap G_j \in \mathcal{M}] \\
&= P[G_i = G_j | G_i \in \mathcal{M} \cap G_j \in \mathcal{M}] \cdot P[G_i \in \mathcal{M}] \cdot P[G_j \in \mathcal{M}] \\
&= \frac{1}{m} \cdot \frac{1}{2} \cdot \frac{1}{2} \\
&= \frac{1}{4m}
\end{aligned}$$

- M_{ik}^{GF} für das Ereignis, dass die Geschäftsreisende i und der Ferienreisende k in \mathcal{M} in die gleiche Reihe sitzen; analog zu M_{ij}^{GG} gilt

$$E[M_{ik}^{GF}] = P[(G_i = F_k) \cap G_i \in \mathcal{M} \cap F_k \in \mathcal{M}] = \frac{1}{m} \cdot \frac{1}{2} \cdot \frac{1}{2} = \frac{1}{4m}$$

- M_{kl}^{FF} für das Ereignis, dass Ferienreisende k und l in \mathcal{M} in die gleiche Reihe sitzen; analog zu oben gilt

$$E[M_{kl}^{FF}] = \frac{1}{4m}.$$

Sei X die Anzahl doppelt besetzter Reihen in \mathcal{M} :

$$X = \sum_{(i,j)} M_{ij}^{GG} + \sum_{i=1}^g \sum_{k=1}^f M_{ik}^{GF} + \sum_{(k,l)} M_{kl}^{FF}.$$

Wegen der Linearität des Erwartungswerts gilt

$$\begin{aligned}
E[X] &= E\left[\sum_{(i,j)} M_{ij}^{GG} + \sum_{i=1}^g \sum_{k=1}^f M_{ik}^{GF} + \sum_{(k,l)} M_{kl}^{FF}\right] \\
&= \sum_{(i,j)} E[M_{ij}^{GG}] + \sum_{i=1}^g \sum_{k=1}^f E[M_{ik}^{GF}] + \sum_{(k,l)} E[M_{kl}^{FF}] \\
&= \frac{g(g-1)}{2} \frac{1}{4m} + gf \frac{1}{4m} + \frac{f(f-1)}{2} \frac{1}{4m} \\
&\leq \frac{1}{4m} \frac{g^2 + 2gf + f^2}{2} \\
&= \frac{(g+f)^2}{8m}.
\end{aligned}$$

Wir wenden nun die Markov-Ungleichung an:

$$P[X \geq 1] \leq \frac{E[X]}{1}.$$

Sofern $g, f \ll \sqrt{r}$ wie vorgegeben, so gilt auch $g, f \ll \sqrt{m}$. Wegen $E[X] \leq \frac{(g+f)^2}{8m} \ll 1$, so ist auch die WSK einer doppelt besetzten Reihe beschränkt als $P[X \geq 1] \leq \frac{(g+f)^2}{8m}$.

- c) Analog zur letzten Teilaufgabe definiert man eine Menge von Indikator-ZV V_{ij}^{GG} für zwei Geschäftsreisende in der gleichen Reihe *vorne* und ZV H_{kl}^{FF} für zwei Ferienreisende in

der gleichen Reihe *hinten*. Wie vorher folgt $E[V_{ij}^{GG}] = \frac{1}{4m}$ und $E[H_{kl}^{FF}] = \frac{1}{4m}$. Die Anzahl doppelt besetzter Reihen *irgendwo* im Flugzeug ist nun

$$Y = \sum_{(i,j)} V_{ij}^{GG} + X + \sum_{(k,l)} H_{kl}^{FF};$$

mit den Resultaten von oben folgt $E[Y] \leq \frac{g^2+(g+f)^2+f^2}{8m}$ und $P[Y \geq 1] \leq E[Y]$. Anders gesagt, mit WSK mindestens $1 - E[Y]$ sitzen alle allein in einer Reihe.

12.4 Warteschlange

Aufgabe: In einem Supermarkt mit c Kassen wählen k Kunden unabhängig und uniform verteilt je eine Kasse und stellen sich an.

- Wie viele Kunden stehen durchschnittlich an jeder Kasse an?
- Geben Sie eine Abschätzung unter Verwendung der Chebyshev-Ungleichung für die Wahrscheinlichkeit, dass an der ersten Kasse x oder mehr Kunden anstehen.
- Basierend auf der letzten Teilaufgabe, geben Sie eine Abschätzung für die Wahrscheinlichkeit dafür, dass in irgendeiner Warteschlange x oder mehr Kunden stehen.

Hinweis: Bei (b) bestimmen Sie den Erwartungswert und die Varianz der ZV für die Anzahl Kunden. (Chebyshev-Ungleichung und Varianz sind quasi unzertrennbar!) Für (c) benutzen Sie den Union-Bound.

Lösung:

- Die k Kunden verteilen sich gleichmässig auf alle Kassen, d.h., der Erwartungswert der Kunden für jede Kasse ist k/c .
- Sei X_1 eine Zufallsvariable, welche die Anzahl Kunden an Kasse 1 modelliert. Es gilt $E[X_1] = k/c$. Um die Varianz von X_1 abzuschätzen, definiert man für $j = 1, \dots, k$ eine Indikator-ZV X_{1j} für das Ereignis Kunde j steht bei Kasse 1 an. Damit wird $X_1 = \sum_{j=1}^k X_{1j}$. Es gelten:

$$\begin{aligned} E[X_{1j}] &= \frac{1}{c} \\ \text{Var}[X_{1j}] &= \frac{1}{c} - \frac{1}{c^2} \\ \text{Var}[X_1] &= \sum_{j=1}^k \text{Var}[X_{1j}] = \frac{k}{c} - \frac{k}{c^2} \leq \frac{k}{c}. \end{aligned}$$

Angenommen $x \geq E[X_j] = k/c$ folgt daraus nach Chebyshev

$$P[X_1 \geq x] = P\left[X_1 \geq \frac{k}{c} + \left(x - \frac{k}{c}\right)\right] \leq P\left[\left|X_1 - \frac{k}{c}\right| \geq x - \frac{k}{c}\right] \leq \frac{\text{Var}[X_1]}{(x - k/c)^2}$$

und

$$P[X_1 \geq x] \leq \frac{k/c}{(x - k/c)^2} = \frac{1}{cx^2/k - 2x + k/c}.$$

- Diese WSK ergibt sich aus dem Union-Bound:

$$\begin{aligned} P[\exists i : X_i \geq x] &= P[\cup_{i=1}^c X_i \geq x] \leq \sum_{i=1}^c P[X_i \geq x] \\ &\leq cP[X_1 \geq x] \leq \frac{c}{cx^2/k - 2x + k/c} = \frac{1}{x^2/k - 2x/c + k/c^2}, \end{aligned}$$

wobei die Abschätzung für $P[X_1 \geq x]$ aus der vorigen Teilaufgabe stammt.

12.5 Quiz

Aufgabe: In einem Quiz gibt es rote, grüne, blaue und weisse Fragen. Die Gewinnchancen unterscheiden sich je nach Farbe einer Frage und sind:

Farbe	rot	grün	blau	weiss
Chance	0.4	0.3	0.2	0.1

Das Quiz wird n Mal gespielt und Fragen aller Farben kommen gleich häufig vor. Für jeden Gewinn gibt es einen Preis. Die Quizmaster errechnen die erwartete Anzahl Gewinne, erhöhen diese als Puffer um 20% und kaufen so viele Preise ein. Berechnen Sie eine Schranke abhängig von n für die Wahrscheinlichkeit, dass die bereitgestellten Preise trotzdem nicht ausreichen.

Lösung: Die Anzahl Gewinne X ist die Summe aus n unabhängigen Poisson-Experimenten (d.h., Bernoulli-Experimenten mit unterschiedlichen WSK). Die durchschnittliche Gewinnchance pro Frage ist 0.25 und deshalb $E[X] = n/4$. Aus einer Chernoff-Ungleichung folgt

$$P[X > (1 + \frac{1}{5})n/4] \leq e^{-n/4 \cdot (1/5)^2/3} = e^{-n/300}.$$

12.6 Vergleichen durch Hashing

Aufgabe: Anna und Barbara möchten zwei grosse Datensätze $D_A, D_B \in \mathcal{D}$ miteinander vergleichen, wobei Anna D_A kennt und Barbara D_B . Anna und Barbara können miteinander kommunizieren, aber die Datensätze sind zu gross, um sie auszutauschen.

Die beiden nehmen dazu eine ideale Hashfunktion $H : \mathcal{D} \rightarrow [1, n]$. Für die Anwendung von H wird der Datensatz zuerst in eine kanonische Ordnung gebracht. In der Analyse wird die Hashfunktion so modelliert, dass H für jeden einzelnen Datensatz eine uniform verteilte Zahl in $[1, n]$ ausgibt. In diesem Modell wirkt die Hashfunktion wie ein Orakel, welches aber wichtige Eigenschaften konkreter Hashfunktionen hat. Wenn ein Wert $d \in \mathcal{D}$ zum ersten Mal angefragt wird (von einem beliebigen Teilnehmer), so wählt H eine zufällige Zahl h in $[1, n]$; dies widerspiegelt, dass die Hashfunktion einen scheinbar zufälligen Wert ausgibt. Jeder weitere Aufruf von $H(d)$, egal ob von demselben Teilnehmer oder einem anderen, ergibt wiederum das gleiche h , entsprechend einer realen, deterministischen Hashfunktion.

- Beschreiben Sie einen Algorithmus für den Vergleich. Das Ergebnis soll möglichst sicher bestimmt werden und höchstens mit Wahrscheinlichkeit ϵ falsch sein.
- In welcher Art kann sich dieser Algorithmus irren?
- Wie viele Aufrufe von H sind nötig mit $n = 10$, damit Anna und Barbara sich höchstens mit Wahrscheinlichkeit 2^{-10} irren?

Hinweis: Parameter A aus der ursprünglichen Aufgabe ist hier n .

Lösung: Die Idee hinter dem Algorithmus ist, dass Anna $h = H(D_A)$ berechnet, den kurzen Hashwert h an Barbara sendet, Barbara selbst $H(D_B)$ ausrechnet und danach testet, ob $h \stackrel{?}{=} H(D_B)$. Falls $D_A = D_B$, so liefert diese Methode immer das richtige Ergebnis. Falls jedoch $D_A \neq D_B$, so ist es trotzdem möglich, dass $H(D_A) = H(D_B)$ aufgrund der Zufälligkeit, welche wir für H im Modell annehmen. In diesem Monte-Carlo-Algorithmus passiert der einseitige Fehler mit WSK $1/n$ (Frage (b)).

Falls $1/n \leq \epsilon$ erfüllt ein solcher Vergleich schon die Vorgabe. Falls nicht, so wiederholt man den Vergleich, um an Sicherheit zu gewinnen. Wie in verschiedenen Beispielen gezeigt, sollten die Wiederholungen für eine einfache Analyse unabhängig voneinander sein. Dann ist die WSK, dass der Fehler in k Wiederholungen auftritt, höchstens n^{-k} . Damit der Fehler kleiner als 2^{-10} wird, müssen mit $n = 10$ also vier Wiederholungen durchgeführt werden (Frage (c)).

Anna und Barbara erreichen unabhängige Wiederholungen, indem sie H jeweils auf einen neuen Wert anwenden, welcher jedoch verschieden von den Werten vorher sein muss. $H(D_A)$ und $H(D_B)$ allein würden ja in jeder Wiederholung gleich bleiben, da die modellierte zufällige Hashfunktion für schon angefragte Werte deterministisch ist! Sie können in Runden vorgehen und in Runde r den Wert $H(r\|D_A)$ von Anna zu Barbara senden. (Die Notation $\|$ bezeichnet einfach die Konkatenation von zwei Werten als Bitstring.)

Algorithmus Vergleich(D_A, D_B, ϵ)

```

 $f \leftarrow 1$       (Fehler-WSK)
 $r \leftarrow 1$     (Runde)
while  $f > \epsilon$  do
     $h \leftarrow H(r\|D_A)$       (dies macht nur Anna)
    Anna sendet  $h$  and Barbara
    if  $h \neq H(r\|D_B)$  then    (dies macht nur Barbara)
        return VERSCHIEDEN
     $f \leftarrow f/n$ 
     $r \leftarrow r + 1$ 
return MÖGLICHERWEISE GLEICH

```

12.7 Entropie

Aufgabe: Sei $Y = f(X)$ für eine Zufallsvariable X und eine deterministische Funktion f . Was können Sie allgemein aussagen über Grössen oder Beziehungen zwischen $H(X)$, $H(Y)$ und $H(X|Y)$, $H(Y|X)$?

Lösung: Allgemein gelten $H(X) \geq H(X|Y) \geq 0$ und $H(Y) \geq H(Y|X) \geq 0$. Da $Y = f(X)$ ist $H(Y|X) = 0$; falls f bijektiv ist, dann auch $H(X|Y) = 0$ und sonst $H(X|Y) > 0$.

12.8 Entropie zweier Zufallsvariablen

Aufgabe: Berechnen Sie folgende Grössen für die Verteilung P_{XY} von $X, Y \in \{0, 1\}^2$ unten:

- $H(X)$, $H(Y)$, $H(XY)$;
- $H(X|Y)$, $H(Y|X)$;
- $H(X) - H(X|Y)$, $H(Y) - H(Y|X)$.

	Y	
	0	1
X	$\frac{1}{8}$	$\frac{1}{4}$
	$\frac{7}{16}$	$\frac{3}{16}$

(Ohne Taschenrechner benutzen Sie hierzu Ausdrücke wie $\log(5)$ oder $h(\frac{1}{4})$; Logarithmen sind zur Basis 2.)

Lösung: $P_X(0) = 3/8$ und $P_Y(0) = 9/16$.

$$H(X) = h\left(\frac{3}{8}\right);$$

$$H(Y) = h\left(\frac{9}{16}\right);$$

$$H(X|Y) = P_Y(0)H(X|Y=0) + P_Y(1)H(X|Y=1) = \frac{9}{16}h\left(\frac{2}{9}\right) + \frac{7}{16}h\left(\frac{3}{7}\right);$$

$$H(Y|X) = P_X(0)H(Y|X=0) + P_X(1)H(Y|X=1) = \frac{3}{8}h\left(\frac{1}{3}\right) + \frac{5}{8}h\left(\frac{3}{10}\right);$$

$$H(XY) = H(X) + H(Y|X) = h\left(\frac{3}{8}\right) + \frac{3}{8}h\left(\frac{1}{3}\right) + \frac{5}{8}h\left(\frac{3}{10}\right);$$

$$H(X) - H(X|Y) = h\left(\frac{3}{8}\right) - \frac{9}{16}h\left(\frac{2}{9}\right) - \frac{7}{16}h\left(\frac{3}{7}\right);$$

$$H(Y) - H(Y|X) = h\left(\frac{9}{16}\right) - \frac{3}{8}h\left(\frac{1}{3}\right) - \frac{5}{8}h\left(\frac{3}{10}\right).$$

12.9 Bernoulli-Verteilung erzeugen

Aufgabe: Gegeben $p \in [0, 1]$ soll eine Bernoulli-Zufallsvariable X erzeugt werden mit $P_X(1) = p$. Zur Verfügung steht eine Quelle von uniform verteilten und unabhängigen Zufallsbits Z_1, Z_2, \dots .

Hinweis: Stellen Sie p im Binärsystem dar. Für rationale p siehe auch Übung 6, aber hier geht es um $p \in \mathbb{R}$.

Lösung: Im Binärsystem geschrieben ist $p = 0.p_1p_2p_3\dots$. Sei U die Zufallsvariable $U = 0.Z_1Z_2Z_3\dots$, ebenfalls als binärer Bruch interpretiert. U ist uniform verteilt im halb-offenen Intervall $[0, 1)$. Der folgende Algorithmus gibt $X = 1$ aus, falls $U < p$ und 0 sonst:

$i \leftarrow 1$

while TRUE **do**

if $B_i < p_i$ **then**

return 1

else if $B_i > p_i$ **then**

return 0

$i \leftarrow i + 1$

Da $P[B_i = p_i] = \frac{1}{2}$ für jede Eingabe p , ist die erwartete Anzahl benötigter Zufallsbits genau 2:

$$E[\text{Wert von } i \text{ bei Terminierung}] = \sum_{i=1} i 2^{-i} = 2.$$

12.10 Code

Aufgabe: Betrachten Sie die folgenden Codes und geben Sie jeweils an, ob ein Code (1) präfixfrei und/oder (2) eindeutig decodierbar ist:

- (a) $[0, 01, 001]$
- (b) $[00, 01, 100, 101, 11]$
- (c) $[0, 00, 000, 0000]$

Lösung:

- a) 0 ist ein Präfix von 01; da der Code aber suffixfrei ist (d.h., der Code aus den umgekehrten Codewörtern ist präfixfrei), so ist er eindeutig decodierbar.
- b) Der Code ist präfixfrei und deshalb auch eindeutig decodierbar.
- c) Der Code ist nicht eindeutig decodierbar (und deshalb auch nicht präfixfrei).

12.11 Shannon- und Huffman-Code

Aufgabe: Gegeben eine Zufallsvariable X mit der Verteilung

x	a	b	c	d
$P_X(x)$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{4}$	$\frac{1}{12}$

- a) Finden Sie einen Huffman-Code für X .
- b) Zeigen Sie, dass zwei optimale Codes gibt, nämlich einen Code C_1 mit den Codewort-Längen $[1, 2, 3, 3]$ und einen Code C_2 mit den Längen $[2, 2, 2, 2]$.
- c) Ist C_1 oder C_2 ein Shannon-Code? Warum?

Lösung:

- a) Zwei mögliche Huffman-Codes sind:

x	a	b	c	d
$C(x)$	00	01	10	11

x	a	b	c	d
$C(x)$	0	10	110	111

- b) Siehe (a). Wenn man die Knoten für die Symbole c und d zusammengefasst hat, entsteht ein neuer Knoten (e) mit Gewicht $\frac{1}{3}$. Je nachdem, ob man e mit **a oder b** verbindet, oder ob man **a oder b** zuerst verbindet, resultieren die zwei Codes.
- c) In einem binären Shannon-Code ist die Codewortlänge für Symbol x_i immer

$$w_i = \lceil -\log P_X(x_i) \rceil.$$

Für die Quelle X wäre das ein Code mit Profil $[2, 2, 2, 4]$. Deshalb ist weder C_1 noch C_2 ein Shannon-Code.