

Probeklausur Frühjahr 2009

Statistik für Human- und Sozialwissenschaftler

Zugelassene Hilfsmittel: Taschenrechner aller Art, Fremdsprachenwörterbücher.

Verlangt und gewertet werden **vier der folgenden fünf** Aufgaben. Lösungsschritte und Teilergebnisse sind ausreichend zu begründen. Eine Angabe des Endergebnisses allein genügt nicht.

Aufgabe 1

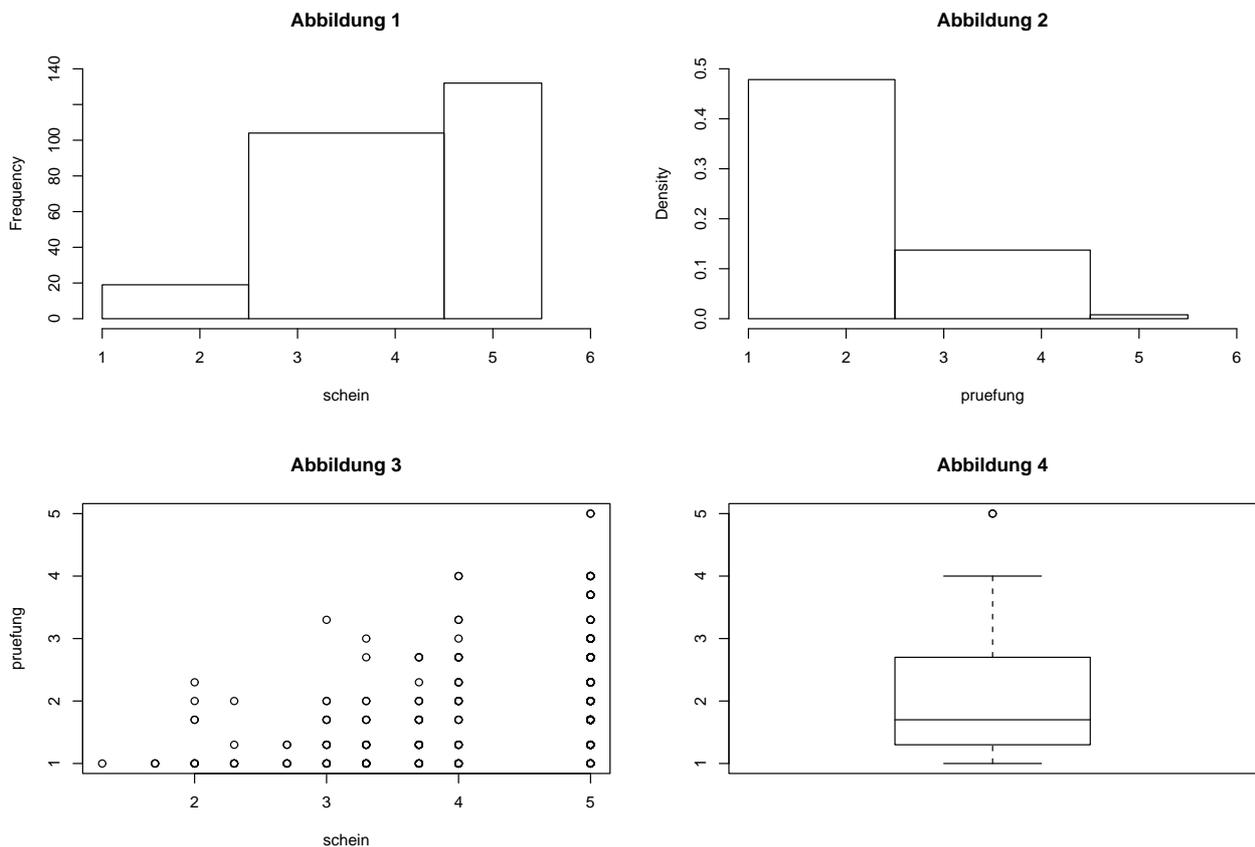
- a) Nach einer Studie des Hochschul-Information-Systems (HIS) brechen an Universitäten ungefähr 30 Prozent der Studenten und an Fachhochschulen ungefähr 22 Prozent der Studenten ihr Studium vorzeitig ab.
- a₁) Warum vergleicht man hier prozentuale Anteile und nicht die absoluten Zahlen der Abbrecher ?
 - a₂) Was würden Sie vermuten: Handelt es sich bei dieser Studie um eine kontrollierte Studie mit Randomisierung, um eine kontrollierte Studie ohne Randomisierung oder um eine Beobachtungsstudie ? Begründen Sie ihre Antwort.
 - a₃) An Fachhochschulen ist der Studienverlauf stärker strukturiert und der Praxisbezug intensiver als an Universitäten. Kann man aus der obigen Studie schließen, dass dies in kausalem Zusammenhang steht mit der geringeren Abbrecherquote ?
- b) Beschreiben Sie **kurz**, was man bei einer Umfrage unter einer Verzerrung durch Auswahl (sampling bias) versteht.
- c) Bei Umfragen steht man häufig vor dem Problem, dass ein Teil der Befragten die Antwort verweigert. Warum kann man das Problem nicht dadurch lösen, dass man solange Leute zur Befragung (zufällig) auswählt, bis man eine gewisse Mindestanzahl an Leuten gefunden hat, die bereit sind, an der Umfrage mitzuwirken ?

Lösung

- a₁) Weil mit absoluten Werten keine Informationen über den Abbrecheranteil bekommen kann. Dieser hängt von der Studentenzahl der jeweiligen Hochschule ab.
 - a₂) Es handelt sich um eine Beobachtungsstudie, da kein Einfluss auf die Studiengruppe ausgeübt wird.
 - a₃) Man kann nur eine Gleichzeitigkeit beobachten. Einen kausalen Zusammenhang kann man aus obiger Studie nicht schließen, da es sich hier um eine Beobachtungsstudie handelt.
- b) Unter "sampling bias", also der Verzerrung durch Auswahl versteht man die nicht repräsentative Auswahl einer Stichprobe. Wählt man z.B. eine Stichprobe aller Busfahrer/Innen aus und wählt dabei nur Männer aus, so ist die Stichprobe nicht repräsentativ.

- c) Weil man mit dieser Strategie noch immer keine Personen der Gruppe “Antwortverweigerer” befragen kann, welche in der endgültigen Stichprobe fehlt. Das “non-respond bias” bleibt also erhalten.

Aufgabe 2

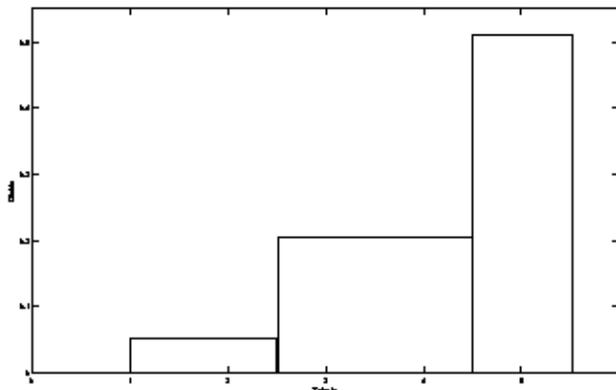


- a) Das Säulendiagramm in Abbildung 1 beschreibt die Noten von $n = 255$ Studenten in der Scheinklausur zur Vorlesung “Statistik I für WiWi”.
- a₁) Inwieweit ist die Darstellung in diesem Säulendiagramm irreführend ?
 - a₂) Stellen Sie die Daten in einem Histogramm so dar, dass die Flächeninhalte der einzelnen Balken (aufgrund von Problemen beim Ablesen der Werte in Abbildung 1 eventuell nur ungefähr) proportional zur Anzahl der Datenpunkte in den zugrundeliegenden Intervallen sind.
- b) Das Histogramm in Abbildung 2 beschreibt die Noten von $n = 255$ Studenten in der Diplom-Vorprüfung zur Vorlesung “Statistik I für WiWi”. Bestimmen Sie mit Hilfe dieses Histogramms (approximativ) die Anzahl der Studenten, die in der Prüfung eine Note besser als 2.5 hatten (d.h., die als Note eine 1.0, 1.3, 1.7, 2.0 oder 2.3 hatten).

- c) Im Streudiagramm in Abbildung 3 sind die Noten in der Scheinklausur gegen die Noten in der Diplom-Vorprüfung abgetragen. Ist die Korrelation zwischen diesen Noten größer oder kleiner als Null ? Begründen Sie ihre Antwort.
- d) Der Boxplot in Abbildung 4 beschreibt die Noten in der Diplom-Vorprüfung zur Vorlesung "Statistik I für WiWi". Wie groß ist der Median und wie groß ist der Interquartilabstand dieser Noten ?

Lösung

- a) a₁) Die graphische Darstellung ist irreführend, da die Klassen, bzw. die Länge der Intervalle nicht alle gleich lang sind. Vergleicht man beispielsweise den Flächeninhalt der mittleren mit der rechten Klasse, so entsteht der Eindruck, dass die mittlere Klasse fast anderthalb mal so viele Datenpunkte enthält wie die rechte Klasse und das ist falsch!
- a₂) Anzahl der Studenten in der Klasse 1, d.h. im Intervall $I_1 (= [1, 2.5)) \approx 20$, in $I_2 \approx 105$ und in $I_3 \approx 130$. Somit ergibt sich folgendes Histogramm:



- b) Beim Histogramm gibt der Flächeninhalt (FI) einer Klasse j den prozentualen Anteil der Datenpunkte (PAD) im zugrunde liegenden Intervall (I_j) an. Somit lässt sich die Anzahl der Datenpunkte in Klasse j wie folgt berechnen:

1. Möglichkeit:

$$\begin{aligned} \text{PAD in } I_j &= \text{FI von } I_j \\ \Rightarrow \text{Anzahl der Datenpunkte in } I_j &= \text{GD} \times \text{PAD in } I_j \end{aligned}$$

mit GD = Gesamtzahl der Datenpunkte

2. Möglichkeit

$$\begin{aligned} \text{Höhe von } I_j &= \frac{n_j}{n \cdot \lambda(I_j)} \\ \Rightarrow n_j &= n \cdot \lambda(I_j) \cdot \text{Höhe von } I_j \end{aligned}$$

mit n_j = Anzahl der Datenpunkte im j -ten Intervall und $\lambda(I_j)$ = Länge des j -ten Intervalls. Hier: $n_j \approx 182$

c) Korrelation ist größer Null, da die zugehörige Regressionsgerade steigend ist.

d) Median ≈ 1.8
IQR ≈ 1.45

Aufgabe 3

Die reelle Zufallsvariable X nehme die Werte 1, 2 und 6 mit den Wahrscheinlichkeiten $1/2$, $1/6$ bzw. $1/3$ an. Die Zufallsvariable Y sei stetig verteilt mit Dichte

$$f : \mathbb{R} \rightarrow \mathbb{R}_+, \quad f(x) = \begin{cases} 20 \cdot x^3(1-x) & \text{für } 0 \leq x \leq 1, \\ 0 & \text{für } x < 0 \text{ oder } x > 1. \end{cases}$$

X und Y seien unabhängig.

a) Bestimmen Sie $\mathbf{E}X$ und $V(X)$.

b) Bestimmen Sie $\mathbf{E}Y$ und $V(Y)$.

c) Bestimmen Sie $\mathbf{E}(X+Y)$ und $V(X+Y)$. An welcher Stelle benötigen Sie hier die Unabhängigkeit von X und Y ?

Lösung

a) Wir berechnen EX mit

$$\begin{aligned} EX &= \sum_{i=1}^3 x_i \cdot P[X = x_i] \\ &= 1 \cdot \frac{1}{2} + 2 \cdot \frac{1}{6} + 6 \cdot \frac{1}{3} \\ &= \frac{3 + 2 + 12}{6} \\ &= \frac{17}{6} \end{aligned}$$

und

$$\begin{aligned} EX^2 &= \sum_{i=1}^3 x_i^2 \cdot P[X = x_i] \\ &= \frac{3 + 4 + 72}{6} \\ &= \frac{79}{6} \end{aligned}$$

womit

$$\begin{aligned} VX &= EX^2 - (EX)^2 \\ &= \frac{79}{6} - \frac{289}{36} \\ &= \frac{185}{36} \end{aligned}$$

gilt.

b) Wir berechnen EY mit

$$\begin{aligned} EY &= \int_{-\infty}^{\infty} x \cdot f(x) dx \\ &= \int_0^1 20 \cdot x^4(1-x) dx \\ &= 20 \left(\int_0^1 x^4 dx - \int_0^1 x^5 dx \right) \\ &= 20 \left(\frac{1}{5} - \frac{1}{6} \right) \\ &= \frac{2}{3} \end{aligned}$$

und EY^2 mit

$$\begin{aligned} EY^2 &= \int_{-\infty}^{\infty} x^2 \cdot f(x) dx \\ &= 20 \left(\int_0^1 x^5 dx - \int_0^1 x^6 dx \right) \\ &= 20 \left(\frac{1}{6} - \frac{1}{7} \right) \\ &= \frac{10}{21} \end{aligned}$$

und insgesamt

$$\begin{aligned} VY &= EY^2 - (EY)^2 \\ &= \frac{10}{21} - \left(\frac{2}{3} \right)^2 \\ &= \frac{2}{63} \end{aligned}$$

c) Es gilt $E(X+Y) = EX + EY = \frac{17}{6} + \frac{2}{3} = \frac{21}{6}$ und $V(X+Y) = V(X) + V(Y) = \frac{185}{36} + \frac{2}{63} \approx 5,17$, wofür die Unabhängigkeit nur für die Varianz benötigt wird.

Aufgabe 4

Ein Glücksrad bleibt nach dem Drehen rein zufällig auf einem von insgesamt 50 Feldern stehen. Bleibt es auf einem der 10 blau gefärbten Felder stehen, so wird ein Gewinn von 5 Euro ausgezahlt. Bleibt es auf einem der 5 grün gefärbten Felder stehen, so wird ein Gewinn von 10 Euro ausgezahlt. Und bleibt es auf dem *einzigsten* roten Feld stehen, so wird ein Gewinn von 100 Euro ausgezahlt. Auf den übrigen 34 weiß gefärbten Feldern wird kein Gewinn ausgezahlt.

- Wie groß ist der Gewinn “im Mittel”, und wie groß ist die “mittlere quadratische Abweichung” zwischen dem zufälligen Gewinn und dem Gewinn “im Mittel” ?
- Für einmaliges Drehen verlangt der Besitzer des Glücksrads einen Einsatz von 5 Euro. Damit beträgt sein Verdienst bei einmaligen Drehen $Y = 5 - X$, wobei X der ausgezahlte Gewinn ist. Wie groß ist sein Verdienst “im Mittel”, und wie groß ist die “mittlere quadratische Abweichung” zwischen dem zufälligen Verdienst und dem Verdienst “im Mittel”?
- Der Besitzer betreibt sein Glücksrad einen Monat lang auf einem Jahrmarkt. In dieser Zeit drehen dabei $n = 6000$ Personen (unbeeinflusst voneinander) am Glücksrad. Bestimmen Sie mit Hilfe des Zentralen Grenzwertsatzes eine untere Schranke für den Verdienst des Besitzers in diesem Monat, die (ungefähr) mit Wahrscheinlichkeit 0.95 überschritten wird.

Hinweise:

Seien Y_1, \dots, Y_n unabhängige Zufallsvariablen mit $\mathbf{P}_{Y_i} = \mathbf{P}_Y$ ($i = 1, \dots, n$), wobei Y die in b) eingeführte ZV ist. Bestimmen Sie mit Hilfe des Zentralen Grenzwertsatzes $x \in \mathbb{R}$ so, dass gilt:

$$\mathbf{P}\left[\sum_{i=1}^n Y_i \geq x\right] \approx 0.95.$$

Verwenden Sie dabei, dass für die Verteilungsfunktion Φ der $N(0, 1)$ -Verteilung gilt: $\Phi(1.65) \approx 0.95$ und $\Phi(-1.65) \approx 0.05$.

Nach dem Zentralen Grenzwertsatz gilt

$$\mathbf{P}\left[\frac{1}{\sqrt{n} \cdot \sqrt{V(Y_1)}} \sum_{i=1}^n (Y_i - \mathbf{E}Y_1) \leq x\right] \approx \Phi(x).$$

Lösung

- Als zugrundeliegenden Wahrscheinlichkeitsraum kann man z.B. (Ω, P) betrachten mit

$$\Omega = \{\text{weiß, blau, grün, rot}\}$$

und

$$P(\{\text{weiß}\}) = \frac{17}{25}, P(\{\text{blau}\}) = \frac{1}{5}, P(\{\text{grün}\}) = \frac{1}{10}, P(\{\text{rot}\}) = \frac{1}{50}.$$

Als nächstes definieren wir die Zufallsvariable X , die den Gewinn in Euro bei dem beschriebenen Glücksradspiel beschreibt. Man erhält

ω	weiß	blau	grün	rot
$X(\omega)$	0	5	10	100

oder anders ausgedrückt

x	0	5	10	100
$P(X = x)$	$\frac{17}{25}$	$\frac{1}{5}$	$\frac{1}{10}$	$\frac{1}{50}$

Der Gewinn im Mittel ist dann

gerade der Erwartungswert der Zufallsvariable X , bei der mittleren quadratischen Abweichung handelt es sich um die Varianz.

$$\begin{aligned}
 E(X) &= 0 \cdot P(X = 0) + 5 \cdot P(X = 5) + 10 \cdot P(X = 10) + 100 \cdot P(X = 100) \\
 &= 5 \cdot \frac{1}{5} + 10 \cdot \frac{1}{10} + 100 \cdot \frac{1}{50} = 4 \\
 E(X^2) &= 0^2 \cdot P(X = 0) + 5^2 \cdot P(X = 5) + 10^2 \cdot P(X = 10) + 100^2 \cdot P(X = 100) \\
 &= 25 \cdot \frac{1}{5} + 100 \cdot \frac{1}{10} + 100^2 \cdot \frac{1}{50} = 215 \\
 V(X) &= E(X^2) - (EX)^2 = 215 - 4^2 = 199
 \end{aligned}$$

- b) Die Zufallsvariable X aus der Aufgabenstellung entspricht der Zufallsvariable X aus der Lösung von Aufgabenteil a). Gesucht sind Erwartungswert und Varianz von Y .

$$\begin{aligned}
 E(Y) &= E(5 - X) = E(5) - E(X) = 5 - 4 = 1 \\
 V(Y) &= V(5 - X) = V(-X) = V(X) = 199
 \end{aligned}$$

- c) Seien Y_1, \dots, Y_n unabhängig und identisch verteilte Zufallsvariablen mit der selben Verteilung wie Y aus Aufgabenteil b). Nach der Aufgabenstellung ist $n = 6000$. Der Gewinn, den der Glückradbetreiber in dem Monat macht wird beschrieben durch

$$\sum_{i=1}^n Y_i.$$

Wir suchen jetzt eine Zahl x mit

$$P\left(\sum_{i=1}^n Y_i \geq x\right) = P\left(\sum_{i=1}^n Y_i > x\right) \approx 0.95.$$

Dies ist äquivalent mit

$$1 - P\left(\sum_{i=1}^n Y_i \leq x\right) \approx 0.95.$$

Wir wollen nun $\sum_{i=1}^n Y_i$ so transformieren, dass wir den zentralen Grenzwertsatz anwenden können.

$$\begin{aligned}
 &\sum_{i=1}^n Y_i \leq x \\
 \Leftrightarrow &\sum_{i=1}^n Y_i - n \cdot E(Y_1) \leq x - n \cdot E(Y_1) \\
 \Leftrightarrow &\sum_{i=1}^n (Y_i - E(Y_1)) \leq x - n \cdot E(Y_1) \\
 \Leftrightarrow &\frac{1}{\sqrt{n} \cdot \sqrt{V(Y_1)}} \left(\sum_{i=1}^n (Y_i - E(Y_1))\right) \leq \frac{1}{\sqrt{n} \cdot \sqrt{V(Y_1)}} (x - n \cdot E(Y_1))
 \end{aligned}$$

Nach dem zentralen Grenzwertsatz, ist die Zufallsvariable

$$Z = \frac{1}{\sqrt{n} \cdot \sqrt{V(Y_1)}} \left(\sum_{i=1}^n (Y_i - E(Y_1))\right)$$

näherungsweise Normalverteilt mit Erwartungswert 0 und Varianz 1. Es gilt

$$\begin{aligned}
 0.95 &= 1 - P\left(\sum_{i=1}^n Y_i \leq x\right) \\
 &= 1 - P\left(\frac{1}{\sqrt{n} \cdot \sqrt{V(Y_1)}} \left(\sum_{i=1}^n (Y_i - E(Y_1))\right) \leq \frac{1}{\sqrt{n} \cdot \sqrt{V(Y_1)}} (x - n \cdot E(Y_1))\right) \\
 &= 1 - P\left(Z \leq \frac{1}{\sqrt{n} \cdot \sqrt{V(Y_1)}} (x - n \cdot E(Y_1))\right).
 \end{aligned}$$

Dies ist äquivalent zu

$$0.05 = P\left(Z \leq \frac{1}{\sqrt{n} \cdot \sqrt{V(Y_1)}} (x - n \cdot E(Y_1))\right) = P(Z \leq -1.65).$$

Also ist

$$-1.65 = \frac{1}{\sqrt{n} \cdot \sqrt{V(Y_1)}} (x - n \cdot E(Y_1))$$

und damit

$$x = -1.65 \cdot \sqrt{6000} \cdot \sqrt{199} + 6000 \approx 4197.04.$$

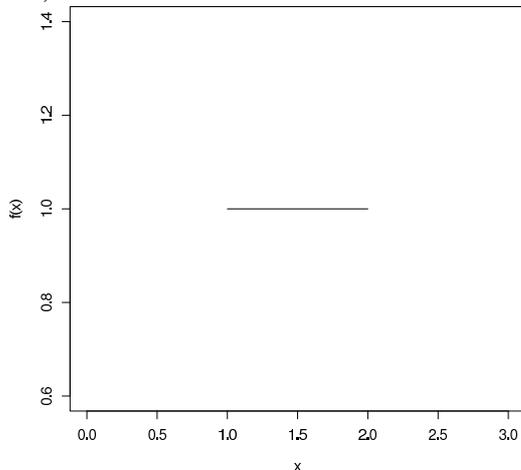
Der berechnete Wert von x ist die gesuchte untere Schranke für den Gewinn.

Aufgabe 5

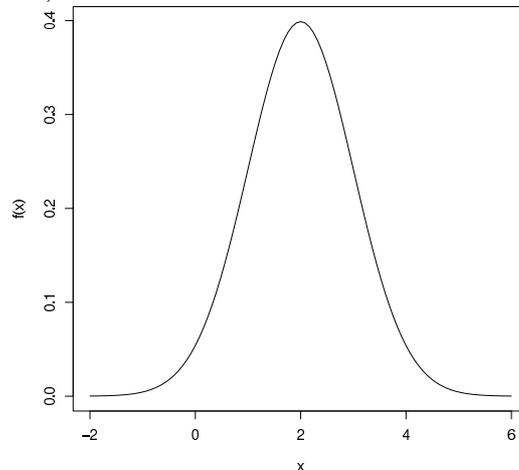
- a) Angenommen, die Lotto-Zahlen beim Lotto 6 aus 49 werden unbeeinflusst voneinander jede Woche rein zufällig aus den Zahlen von 1 bis 49 ausgewählt. Ist es dann wahrscheinlicher, dass nächste Woche die Zahlen 1, 9, 10, 15, 32, 37 gezogen werden, oder dass nächste Woche die gleichen Zahlen wie in dieser Woche gezogen werden?

- b) Welche der beiden Dichten

b₁)



b₂)



gehört zu einer Normalverteilung? Was können Sie über den Erwartungswert dieser Normalverteilung aussagen?

- c) Erläutern Sie jeweils kurz (und evtl. auch anschaulich) die Aussagen des

- c₁) empirischen Gesetzes der großen Zahlen,
- c₂) starken Gesetzes der großen Zahlen,
- c₃) zentralen Grenzwertsatzes.

Lösung

- a) Da wir es hier mit einem Laplace Wahrscheinlichkeitsraum zu tun haben und die Ziehungen unbeeinflusst voneinander durchgeführt werden, haben beide Ereignisse die gleiche Wahrscheinlichkeit $\frac{1}{\binom{49}{6}}$.
- b) Die Dichte b_2). Der Erwartungswert lässt sich ablesen und ist 2. Im Graph liest man einfach den x -Wert ab zu dem der größte Funktionswert gehört.
- c₁) Führt man ein Zufallsexperiment unbeeinflusst voneinander immer wieder durch, so nähert sich die relative Häufigkeit des Auftretens eines festen Ereignisses A einer festen Zahl $P(A) \in [0, 1]$ an.
- c₂) Sind die auf dem selben Wahrscheinlichkeitsraum definierten reellen Zufallsvariablen X_1, X_2, \dots unabhängig und identisch verteilt, und existiert EX_1 , so gilt:

$$\frac{1}{n} \sum_{i=1}^n X_i \rightarrow EX_1 \text{ f.s.}$$

- c₃) Sind X_1, X_2, \dots unabhängige und identisch verteilte reelle Zufallsvariablen mit $EX_1^2 < \infty$ so ist für n groß

$$\frac{\sum_{i=1}^n X_i - E(\sum_{i=1}^n X_i)}{\sqrt{V(\sum_{i=1}^n X_i)}} = \frac{\sqrt{n}}{\sqrt{V(X_1)}} \cdot \left(\frac{1}{n} \sum_{i=1}^n X_i - EX_1 \right)$$

standardnormalverteilt.