

Statistik Zusammenfassung

Alexander Strobl

July 3, 2014

Contents

1	Grundlagen	2
1.1	Häufigkeitsverteilungen	2
1.2	Verschiedene Diagrammtypen	2
1.3	Maßzahlen zur Beschreibung der Lage einer Verteilung	4
2	Streuung und Konzentration	5
3	Wahrscheinlichkeitstheorie	6
3.1	Rechenbeispiele	6
3.1.1	Paarbildung	6
3.1.2	Sortierung	6
3.1.3	Satz von Bayes - HIV-Test	7
3.1.4	Elemente aus einer Menge nehmen	7
4	Diskrete Zufallsvariablen	7
4.1	Verteilungsfunktionen	8
4.1.1	Bernoulli-Verteilung	8
4.1.2	Binomialverteilung	8
4.1.3	Diskrete Gleichverteilung	8
4.1.4	Geometrische Verteilung	8
4.1.5	Hypergeometrische Verteilung	8
4.1.6	Poisson-Verteilung	9
5	Stetige Zufallsvariablen	9
5.1	Normalverteilung	9
5.1.1	Rechenbeispiel - Normalverteilung	9
6	title	12

List of Figures

1	Begriffserklärungen: Häufigkeitsverteilung	2
2	Beispiel: Häufigkeitsverteilung von Noten	2
3	Histogramm	2
4	Relative Klassenhäufigkeiten	2
5	Häufigkeitsdichte	3
6	klassierte Altersverteilung	3
7	Empirische Verteilungsfunktion	3
8	Approximierende Verteilungsfunktion	4
9	Boxplot	5
10	Begriffserklärungen: Streuung und Konzentration	5
11	Scatterplot	6
12	HIV-Test - Krankheitswahrscheinlichkeit bei pos. Testergebnis	7
13	Verschiedene Lösungsformeln der Kombinatorik	7

1 Grundlagen

1.1 Häufigkeitsverteilungen

Eigenschaft	Beschreibung
Merkmalsträger	Objekt von Interesse bei einer empirischen Untersuchung
Gesamtheit	Menge der relevanten Merkmalsträger; Die Anzahl nennt man Umfang der Gesamtheit
Mikrodaten	Daten, welche ausgewertet werden sollen
Häufigkeitsverteilung	Ausprägungen der einzelnen Merkmalsträger

Figure 1: Begriffserklärungen: Häufigkeitsverteilung

Merkmalsausprägung x_i	absolute Häufigkeiten n_i	relative Häufigkeiten f_i
1	6	0.3 / 30%
2	7	0.35
3	4	0.2
4	2	0.1
5	1	0.05
Σ	20	1

Figure 2: Beispiel: Häufigkeitsverteilung von Noten

1.2 Verschiedene Diagrammtypen

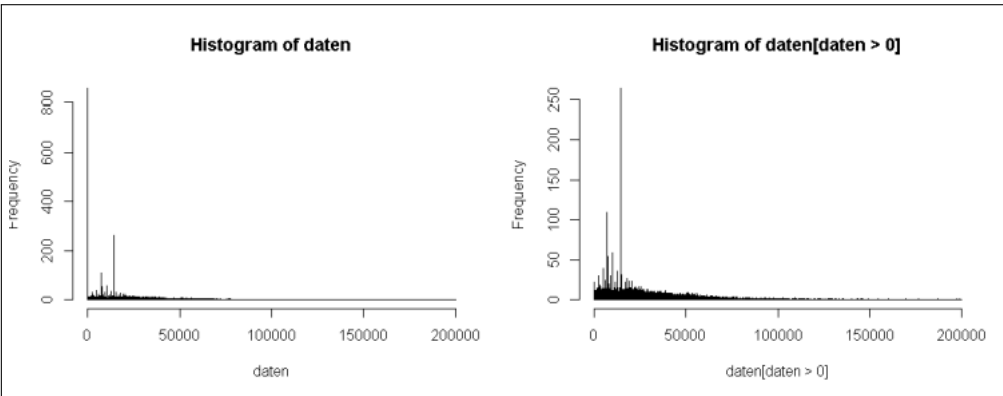


Figure 3: Histogramm

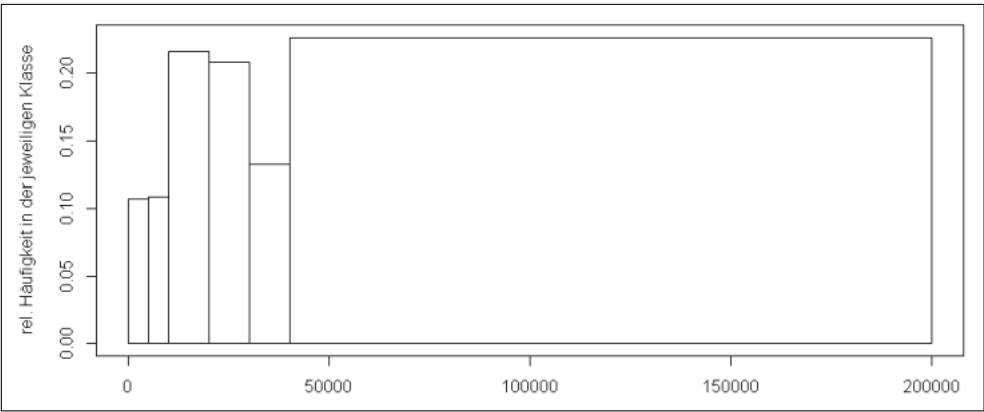


Figure 4: Relative Klassenhäufigkeiten

Wird allerdings nicht mehr verwendet, sondern stattdessen:

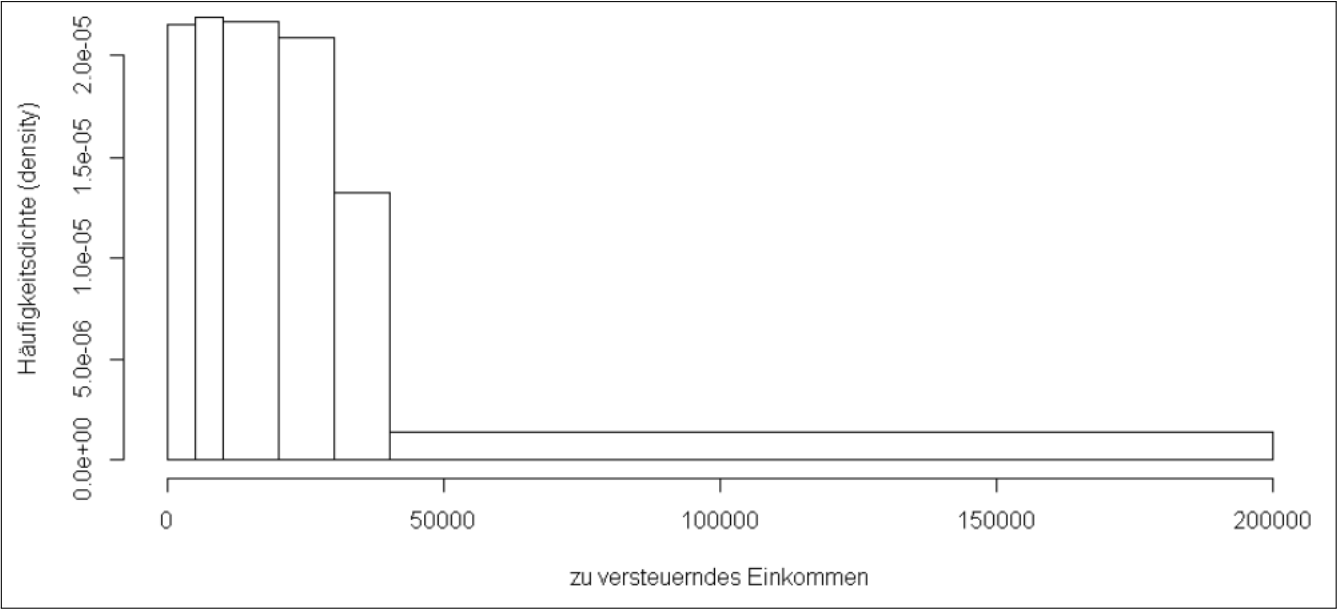


Figure 5: Häufigkeitsdichte

Klasse i x_i	Klassenobergrenze x_i^o	absolute Häufigkeiten n_i	relative Häufigkeiten f_i	emp. Verteilungsfunktion an der Klassenobergrenze $F(x_i^o)$
1	29	7	0.01165	1.17 %
2	39	59	0.09817	10.98 %
3	49	127	0.21131	32.11 %
4	54	120	0.19967	52.08 %
5	59	146	0.24293	76.37 %
6	64	112	0.18636	95.01 %
7	73	30	0.04992	100.00 %
\sum		601	1	

Figure 6: klassierte Altersverteilung

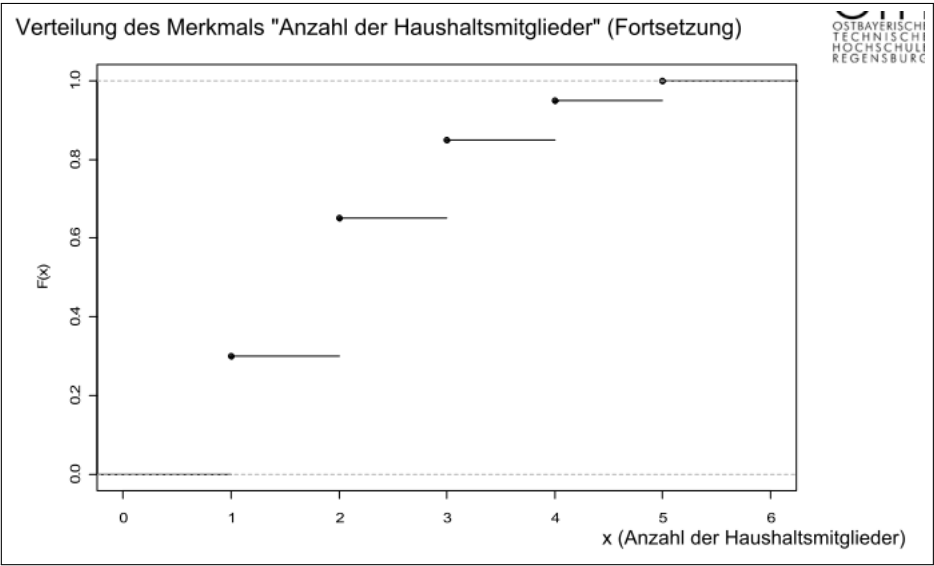


Figure 7: Empirische Verteilungsfunktion

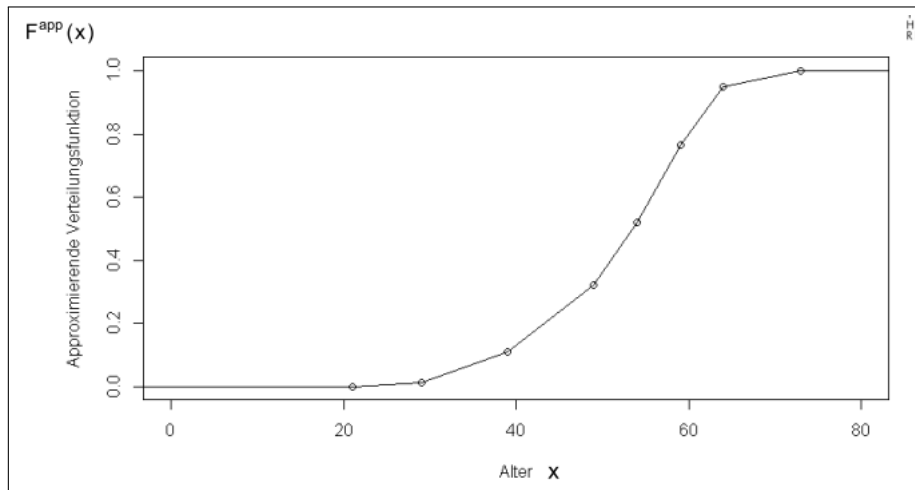


Figure 8: Approximierende Verteilungsfunktion

1.3 Maßzahlen zur Beschreibung der Lage einer Verteilung

Eigenschaft	Formeln	Beschreibung
Modus	$\max(x_i)$	Merkmalsausprägung die am Häufigsten vorkommt
Median	$x_{0.5}$	$50\% \leq x \ \&\& \ 50\% \geq x$
Arithmetisches Mittel \bar{x}	$= \frac{\text{Summe aller Merkmalswerte}}{\text{Anzahl aller Merkmalswerte}}$	Ist ein Spezialfall des Erwartungswertes, mit gleicher Wahrscheinlichkeit für alle Elemente
	$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$	dabei sind die x_i die Merkmalswerte, und n ist die Anzahl der Merkmalswerte, d.h. die Anzahl der Merkmalsträger)
	$\bar{x} = \frac{1}{n} \sum_{i=1}^k x_i n_i$	dabei ist x_i die i-te Merkmalsausprägung, n_i die absolute Häufigkeit der Ausprägung x_i und k ist die Anzahl der verschiedenen Merkmalsausprägungen
	$\bar{x} = \sum_{i=1}^k x_i f_i$	dabei ist x_i die i-te Merkmalsausprägung, f_i die relative Häufigkeit der Ausprägung x_i und k ist die Anzahl der verschiedenen Merkmalsausprägungen
	$\bar{x} = \frac{1}{n} \sum_{i=1}^r \bar{x}_i n_i$	dabei bezeichnet x_i das arithmetische Mittel in der i-ten Schicht, n_i den Umfang der i-ten Schicht und r die Zahl der Schichten
	$\bar{x} = \sum_{i=1}^r \bar{x}_i f_i$	dabei bezeichnet x_i das arithmetische Mittel in der i-ten Schicht, f_i den Anteil der Merkmalsträger in der i-ten Schicht und r die Zahl der Schichten
Quantile	$x_p, x_{1-p} \rightarrow x_{0.3}, x_{0.7}$	Spezialfall: Quartile $x_{0.25}, x_{0.75}$

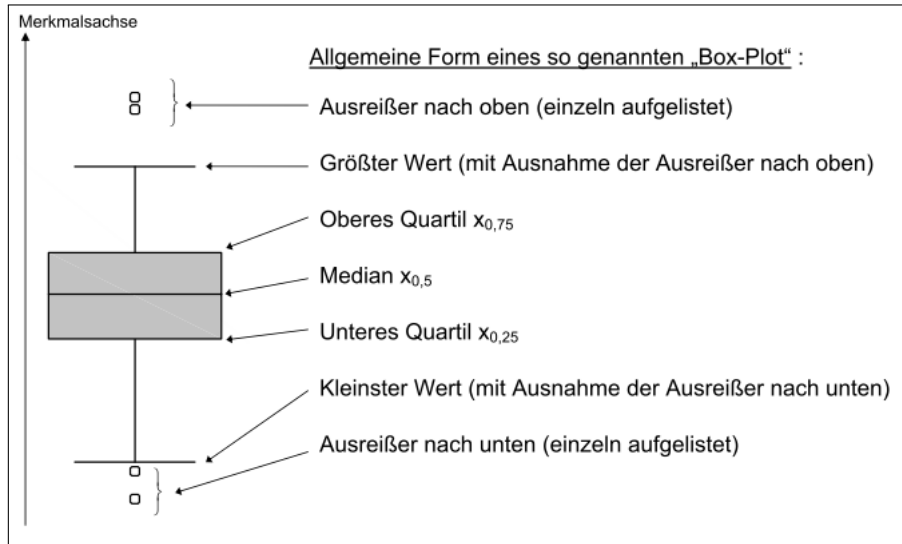


Figure 9: Boxplot

2 Streuung und Konzentration

Eigenschaft	Formel	Beschreibung
Spannweite	$\Delta = \max - \min$	Differenz zwischen max. und min. Merkmalswerten
Quartilsabstand	$\Delta = x_{0,75} - x_{0,25}$	Differenz zwischen den beiden Quartilen
Varianz	s^2	Mittlere quad. Abweichung vom Mittelwert, invariant gegenüber Verschiebungen
	$= \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$	$= \frac{1}{n} \sum_{i=1}^k (x_i - \bar{x})^2 n_i = \sum_{i=1}^k (x_i - \bar{x})^2 f_i$
	$= \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2$	$= \frac{1}{n} \sum_{i=1}^k x_i^2 n_i - \bar{x}^2 = \sum_{i=1}^k x_i^2 f_i - \bar{x}^2$
	$= \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$	Bei Stichproben ein n weniger nehmen
Standardabweichung	$s = \sqrt{s^2}$	Mittlere Abweichung vom Mittelwert, invariant gegenüber Verschiebungen
Standardisierte Merkmale	$\bar{x} = 0$ und $s^2 = 1$	Merkmalsverteilung gilt als Standardisiert
Kovarianz	$s_{XY} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$	
	$= \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y}$	
Korrelationskoeffizient	r	nach Bravais-Pearson
	$\frac{s_{XY}}{s_X s_Y} = \frac{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2}}$	$r = 1 \parallel r = -1$: steigende / fallende Gerade $r = 0$: kein lin. Zusammenhang
Regressionsgerade	$\hat{y} = \hat{\alpha} + \hat{\beta}x$ $\hat{\alpha} = \bar{y} - \hat{\beta}\bar{x}$	$\hat{\beta} = \frac{s_{XY}}{s_X}$
Simpsons Paradoxon		Widersprüchliche Ergebnisse bei genauerer Betrachtung

Figure 10: Begriffserklärungen: Streuung und Konzentration

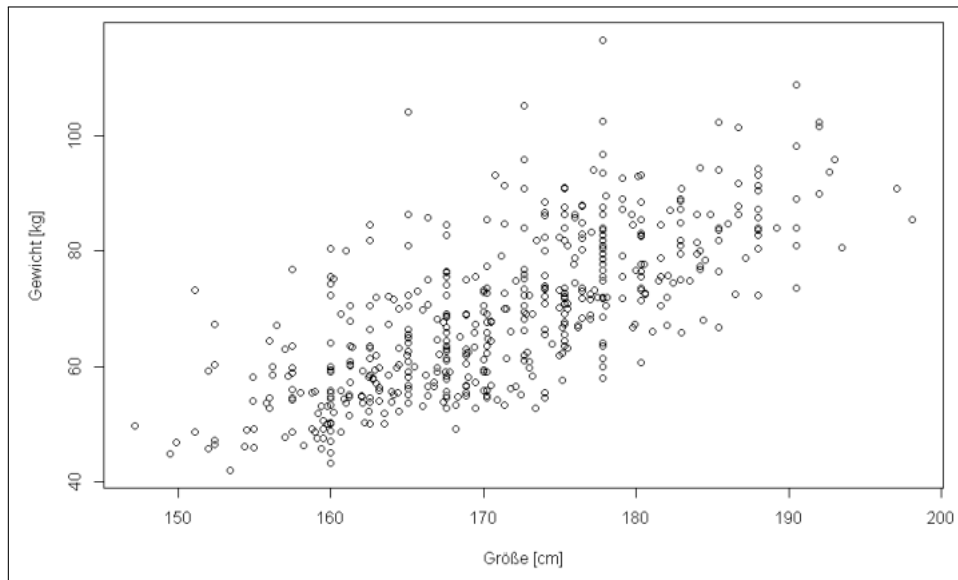


Figure 11: Scatterplot

3 Wahrscheinlichkeitstheorie

Eigenschaft	Formel	Beschreibung
Wahrscheinlichkeitsraum	Ω	Alle Ergebnisse, die stattfinden können
LaPlace-Wahrscheinlichkeit	$P(A) = \frac{\text{Anz. Elem in } A}{\text{Anz. Elem in } \Omega} = \frac{\#A}{\#\Omega}$	Alles gleich wahrscheinlich
Bedingte Wahrscheinlichkeiten	$P(A_1 A_2) = \frac{P(A_1 \cap A_2)}{P(A_2)}$	Wahrscheinlichkeit für A_1 nachdem A_2 eingetreten
Axiome von Kolmogorov	$P(A \cup B) = P(A) + P(B)$	Nur bei sich gegenseitig ausschliessenden Ereignissen
	$P(A \cup B C)$ $= P(A C) + P(B C) - P(A \cap B C)$	
Abhängigkeit	$P(A \cup B) = P(A) + P(B) - P(A)P(B)$	
Unabhängigkeit	$P(A \cap B) = P(A)P(B)$	
Totale Wahrscheinlichkeit	$P(B) = \sum_{i=1}^m P(B \cap A_i) = P(A_i)P(B A_i)$	
Satz von Bayes	$P(A_j B) = \frac{P(B A_j)P(A_j)}{\sum_{i=1}^m P(B A_i)P(A_i)}$	

Table 1: Begriffserklärungen: Wahrscheinlichkeitstheorie

3.1 Rechenbeispiele

3.1.1 Paarbildung

12 Männer, 10 Frauen: Wie viele Paare können gebildet werden?

$$\#A\#B = 12 \times 10 = 120$$

3.1.2 Sortierung

n Elemente sollen sortiert werden: Wie viele Möglichkeiten?

n!

3.1.3 Satz von Bayes - HIV-Test

$P(HIV^+) = \text{Prävalenz} = \text{Anteil Kranke an der Bevölkerung}$

$$P(HIV^+ | T^+) = \frac{P(T^+ | HIV^+)P(HIV^+)}{P(T^+ | HIV^+)P(HIV^+) + P(T^+ | HIV^-)P(HIV^-)}$$

Figure 12: HIV-Test - Krankheitswahrscheinlichkeit bei pos. Testergebnis

3.1.4 Elemente aus einer Menge nehmen

Anzahl der Möglichkeiten, n aus N Elementen auszuwählen	mit Zurücklegen	ohne Zurücklegen
mit Berücksichtigung der Reihenfolge	N^n	$\frac{N!}{(N-n)!}$
ohne Berücksichtigung der Reihenfolge	$\binom{N+n-1}{n}$	$\binom{N}{n}$

Figure 13: Verschiedene Lösungsformeln der Kombinatorik

4 Diskrete Zufallsvariablen

Eigenschaft	Formel	Beschreibung
Verteilungsfunktion	$F_x(x) = P(X \leq x)$	definiert die Wahrscheinlichkeit der Zufallsvariable X, dass X höchstens den Wert x annimmt
Unabhängigkeit	$P(X_1 = x_1, X_2 = x_2) = P(X_1 = x_1) * P(X_2 = x_2)$	gilt ebenfalls für andere Operationen wie z.B. \leq
Erwartungswert	$E(X) = \mu_x = \mu$ $= \sum_{i=1}^k x_i p_i = \sum_{i=1}^k x_i * P(X = x_i)$ $E(Y) = E(g(X)) = \sum_i g(x_i) p_i$ $E(X + Y) = E(X) + E(Y)$ $E\left(\sum_{i=1}^n X_i\right) = \sum_{i=1}^n E(X_i)$	Ist der Mittelwert von X Weitere Rechenregeln: Wenn g(x) eine reelle Funktion und $Y = g(X)$
Varianz	$Var(X) = E((X - \mu_x)^2)$	
	$= \sum_{i=1}^k (x_i - \mu_x)^2 * p_i$	
	$= E(X^2) - E(X)^2$	
	$Var(aX + b) = Var(aX)$	

Table 2: Begriffserklärungen: Diskrete Zufallsvariablen

4.1 Verteilungsfunktionen

4.1.1 Bernoulli-Verteilung

Wird verwendet, wenn es nur 2 Möglichkeiten als Ausgang gibt (Sieg / Niederlage, Gewinn / Verlust usw.), z.B. Münzwurf, Los ziehen.

$$f(x) = \begin{cases} 1 - \pi & \text{für } x = 0 \\ \pi & \text{für } x = 1 \\ 0 & \text{sonst} \end{cases}$$

$$E(X) = \sum_{i=1}^2 x_i p_i = 0 * P(X = 0) + 1 * P(X = 1) = \pi$$

$$\text{Var}(X) = \sum_{i=1}^2 (x_i - \mu_x)^2 p_i = (0 - \pi)^2 (1 - \pi) + (1 - \pi)^2 \pi = \pi(1 - \pi)$$

4.1.2 Binomialverteilung

X sind die Anzahl der Treffer bei n Versuchen, z.B. Multiple Choice

$$X_i = \begin{cases} 1 & \text{falls im } i\text{ten Versuch ein Treffer erzielt wird} \\ 0 & \text{sonst} \end{cases}$$

und es gilt $X = \sum_{i=1}^n X_i$

$$f(x|n, \pi) = \begin{cases} \binom{n}{x} \pi^x (1 - \pi)^{n-x} & \text{für } x = 0, 1, 2, \dots, n \\ 0 & \text{sonst} \end{cases}$$

$$E(X) = n * \pi$$

$$\text{Var}(X) = n * \pi * (1 - \pi)$$

4.1.3 Diskrete Gleichverteilung

$$f(x|n) = \begin{cases} \frac{1}{n} & \text{für } x = 1, 2, \dots, n \\ 0 & \text{sonst} \end{cases}$$

$$E(X) = \frac{n+1}{2}$$

$$\text{Var}(X) = \frac{n^2 - 1}{12}$$

4.1.4 Geometrische Verteilung

Ein Bernoulli-Experiment (0/1-Verteilung) wird immer wieder wiederholt. Die Zufallsvariable X ist die Versuchsnummer des ersten Treffers.

$$f(n|\pi) = \begin{cases} (1 - \pi)^{n-1} * \pi & \text{für } n = 1, 2, 3, \dots \\ 0 & \text{sonst} \end{cases}$$

$$E(X) = \frac{1}{\pi}$$

$$\text{Var}(X) = \frac{1 - \pi}{\pi^2}$$

4.1.5 Hypergeometrische Verteilung

M aus N Elementen haben eine bestimmte Eigenschaft (Loggewinn).

Wenn n Elemente mit Zurücklegen gezogen werden: $\pi = \frac{M}{N}$

Wenn n Elemente ohne Zurücklegen gezogen werden, ist die Anzahl X nicht mehr exakt binomialverteilt.

$$f(n|n, N, M) = \begin{cases} \frac{\binom{M}{x} * \binom{N-M}{n-x}}{\binom{N}{n}} & \text{für } x = \max(0, n - (N-M)), \dots, \min(n, M) \\ 0 & \text{sonst} \end{cases}$$

$$E(X) = n * \frac{M}{N}$$

$$\text{Var}(X) = n * \frac{M}{N} * \left(1 - \frac{M}{N}\right) * \frac{N-n}{N-1}$$

4.1.6 Poisson-Verteilung

Muss in der Klausur erwähnt werden, kann nicht selbst entschieden werden.

z.B. Anrufe in best. Zeitabschnitt in einer Hotline oder Anzahl der Personen, die in einem best. Zeitabschnitt einen Schalter besuchen

$$f(x|\lambda) = \begin{cases} \frac{\lambda^x}{x!} e^{-\lambda} & \text{für } x = 0, 1, 2, 3, \dots \\ 0 & \text{sonst} \end{cases}$$

$$E(X) = \lambda$$

$$\text{Var}(X) = \lambda$$

$$(n+1)! = (n+1) * n! \leq (n+1) * n^n = n^{n+1} + n^n \leq (n+1)^{n+1}$$

5 Stetige Zufallsvariablen

Eigenschaft	Formel	Beschreibung
Definition	$F_x(x) = \int_{-\infty}^x f(t)dt$	Verteilungsfunktion, $f(t)$ = Dichtefunktion
	$P(X = x) = 0$	Wahrscheinlichkeit für einen Wert gleich x ist immer 0
	$P(x_1 \leq X \leq x_2) = F_x(x_2) - F_x(x_1)$	$F'_x(x) = f_x(x)$ Die Dichtefunktion ist die Ableitung der Verteilungsfunktion
Erwartungswert μ_x	$E(X) = \int_{-\infty}^{+\infty} x * f(x)dx$	Die Dichtefunktion $f(x)$ wird nie verändert! $E(\frac{1}{X}) = \int \frac{1}{X} * f(x)dx$
Rechenregeln	$E(Y) = E(g(X)) = \int_a^b g(x) * f(x)dx$	$g(x)$ ist eine reelle Funktion
	$E(aX + b) = a * E(X) + b$	lineare Transformation
	$E(X + Y) = E(X) + E(Y)$	
Modus	$F_x(x_p) = p$	Die Wahrscheinlichkeit, dass X höchstens den Wert x_p annimmt, ist mind. $p/100\%$
Varianz	$\text{Var}(X) = \int_{-\infty}^{+\infty} (x - \mu_x)^2 * f(x)dx$	Standardabweichung ist $\sqrt{\text{Var}(X)}$
Rechenregeln	vgl. Diskrete Zufallsvariablen	

Table 3: Begriffserklärungen: Stetige Zufallsvariablen

5.1 Normalverteilung

Gegeben ist $\mu = E(X)$, sowie $\sigma^2 = \text{Var}(X)$, z.B.: $N(E(X), \text{Var}(x))$

Dies kann auf die Standardnormalverteilung $\mu = 0$, sowie $\sigma = 1$ zurückgerechnet werden:

$$P(z = \frac{x-\mu}{\sigma}) = z_p$$

5.1.1 Rechenbeispiel - Normalverteilung

$N(4000, 10^6)$

Wie hoch ist die Wahrscheinlichkeit für weniger als 3000?

$$\begin{aligned} P(Z \leq \frac{3000-4000}{10^3}) &= z_p \\ P(Z \leq -1) &= 1 - P(Z \geq 1) \\ &= 1 - 0.84 \\ &= 0.16 \end{aligned}$$

diskrete Verteilungen

Verteilungsname	Wahrscheinlichkeitsgewicht/ Zähldichte	Erwartungs- wert $E(X)$	Varianz $\text{Var}(X)$	Anwendung
Bernoulli-Verteilung Parameter $0 < p < 1$	$P(X = 1) = p,$ $P(X = 0) = 1 - p$	p	$p \cdot (1 - p)$	$X = 1 = \text{Erfolg}, X = 0 = \text{Misserfolge}$ z.B. beim einmaligen Werfen eines Würfels eine 6 geworfen (=Erfolg), hier $p = \frac{1}{6}$.
Binomialverteilung Parameter $0 < p < 1$	$P(X = k) = \binom{n}{k} \cdot p^k \cdot (1 - p)^{n-k}$ für $k \in \{0, 1, 2, 3, \dots, n\}$	$n \cdot p$	$n \cdot p \cdot (1 - p)$	$X = \text{Anzahl der Erfolge bei } n \text{ identischen Bernoulli-Experimenten}$ z.B. $X = \text{Anzahl geworfener 6en}$ beim n -maligen Wurf eines fairen Würfels (hier $p = \frac{1}{6}$). z.B. $X = \text{Anzahl gezogener roter Kugeln}$, beim Zie- hen mit Zurücklegen von n Kugeln aus einer Urne M roten und $N - M$ sonstigen Kugeln, wobei $p = \frac{M}{N}$
Diskrete Gleichverteilung auf $\{1, 2, 3, \dots, n\}$	$P(X = k) = \frac{1}{n}$ für $k \in \{1, 2, 3, \dots, n\}$	$\frac{n+1}{2}$	$\frac{n^2-1}{12}$	z.B. ein Wurf mit einem Würfel beschreibt X die geworfene Augenzahl, hier $n = 6$.
Geometrische Verteilung Parameter $0 < p < 1$	$P(X = k) = (1 - p)^{k-1} \cdot p$ für $k \in \{1, 2, 3, \dots\}$	$\frac{1}{p}$	$\frac{1-p}{p^2}$	X beschreibt die Wartezeit auf den ersten Er- folg, beim fortgesetzten Ausführen eines Bernoulli- Experimentes z.B. beim Würfeln warten auf die erste 6, d.h. $X = k$ bedeutet die erste 6 wurde im k -ten Wurf geworfen.
Hypergeometrische Vert. N Anzahl Kugeln in der Urne M Anzahl roter Kugeln n Anzahl zu ziehende Kugeln k Anzahl roter Kugeln unter den gezogenen Kugeln	$P(X = k) = \frac{\binom{M}{k} \cdot \binom{N-M}{n-k}}{\binom{N}{n}}$ für $k \in$ $\{\max(0, n-(N-M)), \dots, \min(n, M)\}$	$n \cdot \frac{M}{N}$	$n \cdot \frac{M}{N} \cdot \left(1 - \frac{M}{N}\right) \cdot \frac{N-n}{N-1}$	$X = \text{Anzahl gezogener roter Kugeln}$, beim Ziehen ohne Zurücklegen von n Kugeln aus einer Urne M roten und $N - M$ sonstigen Kugeln
Poisson-Verteilung Parameter $\lambda > 0$	$P(X = k) = \frac{\lambda^k}{k!} \cdot e^{-\lambda}$ für $k \in \{0, 1, 2, 3, \dots\}$	λ	λ	Anzahl Ereignisse in einem vorgegebenen Zeitinter- vall z.B. Anzahl radioaktiver Zerfälle, Anzahl Blitz- schläge auf einer gegebenen Fläche,...

stetige Verteilungen

Verteilungsname	Dichte	Verteilungsfunktion	Median	$E[X]$	$\text{Var}(X)$	Anwendung
stetige Gleichverteilung auf $[a, b]$	$f(x) = \begin{cases} \frac{1}{b-a} & \text{für } x \in [a, b] \\ 0 & \text{sonst} \end{cases}$	$F(x) = \begin{cases} 0 & \text{für } x < a \\ \frac{x-a}{b-a} & \text{für } x \in [a, b] \\ 1 & \text{für } x > b \end{cases}$	$\frac{a+b}{2}$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$	stetiges Analogon zur diskreten Gleichverteilung z.B. jede reelle Zahl aus dem Intervall $[a, b]$ wird mit gleicher Wahrscheinlichkeit gewählt.
Exponentialverteilung Parameter $\alpha > 0$	$f(x) = \begin{cases} \alpha \cdot e^{-\alpha \cdot x} & \text{für } x \geq 0 \\ 0 & \text{sonst} \end{cases}$	$F(x) = \begin{cases} 0 & \text{für } x < 0 \\ 1 - e^{-\alpha \cdot x} & \text{für } x \geq 0 \end{cases}$	$\frac{\ln(2)}{\alpha}$	$\frac{1}{\alpha}$	$\frac{1}{\alpha^2}$	stetiges Analogon zur geometrischen Verteilung Warten auf das erste/nächste Eintreffen eines Ereignisses z.B. Warten auf den Ausfall einer Glühbirne
Normalverteilung Parameter $\mu \in \mathbb{R}$ und $\sigma > 0$	$f(x) = \frac{1}{\sigma \cdot \sqrt{2\pi}} \cdot \exp\left(-\frac{(x-\mu)^2}{2 \cdot \sigma^2}\right)$	$F(x)$ kann nicht als Funktion hingschrieben werden, vgl. Tabelle	μ	μ	σ^2	Wenn auf etwas viele verschiedene zufällige Einflussfaktoren einwirken, ist das Ergebnis in etwa normalverteilt, z.B. die Körpergröße von Männern (Ernährung, Veranlagung,...) Wird auch zur Approximation von Binomial- und Poissonverteilungen verwendet

