

Formelsammlung zur Vorlesung Statistik I/II für Statistiker, Mathematiker und Informatiker (WS 08/09)

Diese Formelsammlung darf in der Klausur verwendet werden.
 Eigene Notizen und Ergänzungen dürfen eingefügt, aber keine
 zusätzlichen Blätter eingeklebt werden.

Inhaltsverzeichnis

1 Eindimensionale Merkmale 2	
1.1 Häufigkeiten und Häufigkeitsverteilungen . . .	2
1.2 Lageparameter	3
1.3 Quantile	4
1.4 Streuungsparameter	5
2 Zweidimensionale Merkmale 6	
2.1 Gemeinsame Häufigkeiten, Randhäufigkeiten, bedingte Häufigkeiten	6
2.2 Assoziation bei nominalen Merkmalen . . .	6
2.3 Korrelationsrechnung für metrische und ordi- nale Merkmale	7
3 Konzentrationsmaße 9	
3.1 Lorenzkurve und Gini-Koeffizient	9
3.2 Konzentrationsrate CR_g	10
3.3 Herfindahl-Index	11
4 Elementare Wahrscheinlichkeitsrechnung 11	
4.1 Rechenregeln für Wahrscheinlichkeiten . . .	11
4.2 Kombinatorik	12
4.3 Bedingte Wahrscheinlichkeiten	13
4.4 Satz von der totalen Wahrscheinlichkeit . .	13
4.5 Formel von Bayes	13
4.6 Unabhängigkeit zweier Ereignisse	13
5 Eindimensionale Zufallsvariablen und ihre Verteilungen 14	
5.1 Definition von diskreten und stetigen Zu- fallsvariablen und Dichten	14
5.2 Die Verteilungsfunktion einer Zufallsvariablen	14
5.3 Zusammenhänge zwischen Dichten und Ver- teilungsfunktionen	15
5.4 Modus, Median und Quantile	15
5.5 Erwartungswert, Varianz und Standardab- weichung	16
5.6 Rechenregeln und Eigenschaften von Erwar- tungswerten	16
5.7 Spezielle diskrete Verteilungen	17
5.8 Spezielle stetige Verteilungen	17
5.9 Die Chi-Quadrat-, Student- und Fisher- Verteilung	19
6 Zweidimensionale Zufallsvariablen und ihre Verteilungen 19	
6.1 Definition zweidimensionaler Zufallsvariablen	19
6.2 Unabhängigkeit, Kovarianz und Korrelation	20
7 Ergänzungen zu Zufallsvariablen 22	
7.1 Grenzwertsätze	22
7.2 Approximation von Verteilungen	23
7.3 Ungleichung von Tschebyschew	23
8 Testen und Schätzen 23	
8.1 Punktschätzung	24
8.2 Intervallschätzung	26
8.3 Spezielle Schätzprobleme	27
8.4 Testen von Hypothesen	27
8.5 Spezielle Testprobleme	28
8.5.1 Einstichproben-Testprobleme	28
8.5.2 Zweistichproben-Mittelwertsvergleiche	30
8.5.3 Weitere Testprobleme	31
9 Regressionsanalyse 37	
9.1 Lineare Einfachregression	37
9.2 Multiple lineare Regression in Summennota- tion	39
9.3 Multiple lineare Regression in Matrixnotation	40

10 Varianzanalyse 44	
10.1 Einfaktorielle Varianzanalyse	44
10.2 Zweifaktorielle Varianzanalyse	44
11 Zeitreihenanalyse 46	
12 Verteilungstabellen 48	
12.1 Standardnormalverteilung	48
12.2 Students t -Verteilung	49
12.3 χ^2 -Verteilung	50
12.4 Poissonverteilung	51
12.5 F-Verteilung	52
12.6 Wilcoxon-Vorzeichen-Rang-Test	57
12.7 Wilcoxon-Rangsummen-Test	58

1 Eindimensionale Merkmale

1.1 Häufigkeiten und Häufigkeitsverteilungen

Bezeichnungen:

- Urliste: x_1, \dots, x_n
- geordnete Urliste: $x_{(1)} \leq \dots \leq x_{(n)}$
- (kodierte) Merkmalsausprägungen: $a_1 < \dots < a_k$
- absolute Häufigkeit der Ausprägung a_j :

$$h_j = h(a_j) = \sum_{i=1}^n 1_{\{x_i=a_j\}}, \quad \text{mit } 1_{\{x_i=a_j\}} = \begin{cases} 1, & \text{falls } x_i = a_j, \\ 0, & \text{sonst,} \end{cases}$$

- relative Häufigkeit der Ausprägung a_j : $f_j = f(a_j) = h_j/n$,

Häufigkeitsfunktion (-verteilung):

$$\text{absolut: } h(x) = \begin{cases} h_j, & x = a_j, \quad j = 1, \dots, k \\ 0 & \text{sonst} \end{cases}$$

$$\text{relativ: } f(x) = \begin{cases} f_j, & x = a_j, \quad j = 1, \dots, k \\ 0 & \text{sonst} \end{cases}$$

Empirische Verteilungsfunktion (kumulierte relative Häufigkeitsverteilung):

$$F(x) = \sum_{a_j \leq x} f(a_j)$$

Klassenbildung (Gruppierung):

- k Klassen der Form $[c_0, c_1), [c_1, c_2), \dots, [c_{k-1}, c_k)$
- Klassenbreite: $d_j = c_j - c_{j-1} \quad j = 1, \dots, k$
- Klassenmitte: $m_j = (c_j + c_{j-1})/2$
- absolute Häufigkeit der Klasse j : $h_j = \sum_{a_i \in [c_{j-1}, c_j]} h(a_i)$
- relative Häufigkeit der Klasse j : $f_j = h_j/n$

Histogramm (flächentreue Häufigkeitsverteilung):

„Blockhöhe (abzutragende Höhe)“ = $\tilde{f}(x) = \begin{cases} f_j/d_j & x \in [c_{j-1}, c_j) \\ 0 & \text{sonst} \end{cases}$ für $j = 1, \dots, k$.

1.2 Lageparameter**Modus:**

x_{mod} : Ausprägung mit der größten Häufigkeit

- nichtklassierte Daten: $x_{mod} = \{a_j \mid h_j = \max_{a_i} h_i\}$
- klassierte Daten:

$$x_{mod} = c_{j-1} + \frac{\tilde{f}_j - \tilde{f}_{j-1}}{2\tilde{f}_j - \tilde{f}_{j-1} - \tilde{f}_{j+1}}(c_j - c_{j-1}), \quad \tilde{f}_j := \tilde{f}(x), \quad x \in [c_{j-1}, c_j]$$

mit j : Modalklasse

Median:

- nichtklassierte Daten:

$$x_{med} = \begin{cases} x_{(\frac{n+1}{2})} & n \text{ ungerade} \\ \frac{1}{2} (x_{(\frac{n}{2}+1)} + x_{(\frac{n}{2})}) & n \text{ gerade} \end{cases}$$

- klassierte Daten:

$$x_{med} = c_{j-1} + \frac{0.5 - F(c_{j-1})}{F(c_j) - F(c_{j-1})}(c_j - c_{j-1}), \quad \text{mit Medianklasse } j$$

Arithmetisches Mittel (Durchschnittswert):

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

Harmonisches Mittel:

$$\bar{x}_H = \frac{\sum_{i=1}^n g_i}{\sum_{i=1}^n \frac{g_i}{x_i}}, \quad g_i : \text{Gewicht der } i\text{-ten Beobachtung}$$

Geometrisches Mittel:

$$x_G = \sqrt[n]{x_1 \cdot \dots \cdot x_n}$$

mittleres Entwicklungstempo:

$$i_G = \sqrt[n]{i_1 \cdot \dots \cdot i_n}, \quad i_t = \frac{x_t}{x_{t-1}}, \quad t = 1, \dots, n$$

Gewichtetes arithmetisches Mittel:

$$\bar{x}_w = \sum_{i=1}^n w_i x_i \quad \text{wobei} \quad \sum_{i=1}^n w_i = 1 \quad \text{und} \quad 0 \leq w_i \leq 1 \quad \text{für alle } i$$

Spezialfall: arithmetisches Mittel für $w_i = 1/n$.

Arithmetisches Mittel bei Schichtenbildung:

$$\bar{x} = \frac{1}{n} (n_1 \bar{x}_1 + \dots + n_r \bar{x}_r) = \frac{1}{n} \sum_{j=1}^r n_j \bar{x}_j$$

1.3 Quantile

Jeder Wert x_p , mit $0 < p < 1$, für den mindestens ein Anteil p der Daten $\leq x_p$ und mindestens ein Anteil $1 - p$ der Daten $\geq x_p$ ist, heißt p -Quantil:

$$\frac{\text{Anzahl}(x\text{-Werte} \leq x_p)}{n} \geq p \quad \text{und} \quad \frac{\text{Anzahl}(x\text{-Werte} \geq x_p)}{n} \geq 1 - p.$$

Äquivalent dazu ist: x_p ist der kleinste x -Wert, für den $F(x) \geq p$ gilt, d.h.

$$F(x) < p \quad \text{für } x < x_p \quad \text{und} \quad F(x_p) \geq p.$$

Damit gilt für das p -Quantil

- bei nichtklassierten Daten:

$$x_p = \begin{cases} x_{(\lceil np \rceil + 1)} & \text{wenn } np \text{ nicht ganzzahlig,} \\ x_{(np)}, x_{(np+1)} & \text{wenn } np \text{ ganzzahlig.} \end{cases}$$

- bei klassierten Daten:

$$x_p = c_{j-1} + \frac{p - F(c_{j-1})}{F(c_j) - F(c_{j-1})}(c_j - c_{j-1}), \quad p \in (0, 1), \quad \text{mit } j \text{ Quantilkategorie}$$

Speziell: $x_{0.5}$ = Median, $x_{0.25}$ = unteres Quartil, $x_{0.75}$ = oberes Quartil

1.4 Streuungsparameter

Spannweite:

$$SP = x_{(n)} - x_{(1)} = x_{max} - x_{min}$$

Quantilsabstand:

$$x_{1-p} - x_p$$

Interquartilsabstand:

$$d_Q = x_{0.75} - x_{0.25}$$

Empirische Varianz (mittlere quadratische Abweichung):

$$\text{Urliste: } \tilde{s}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \left(\frac{1}{n} \sum_{i=1}^n x_i^2 \right) - \bar{x}^2$$

$$\text{Häufigkeitsdaten: } \tilde{s}^2 = \frac{1}{n} \sum_{j=1}^k (a_j - \bar{x})^2 h(a_j) \quad \text{bzw.} \quad \tilde{s}^2 = \sum_{j=1}^k (a_j - \bar{x})^2 f(a_j)$$

Empirische Standardabweichung:

$$\tilde{s} = +\sqrt{\tilde{s}^2}$$

Stichprobenvarianz:

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

Variationskoeffizient:

$$v = \frac{\tilde{s}}{\bar{x}} \quad (\bar{x} > 0)$$

Streuungszerlegung:

Für r disjunkte statistische Massen E_1, \dots, E_r , deren jeweilige arithmetische Mittel bzw. mittlere quadratische Abweichungen mit $\bar{x}_1, \dots, \bar{x}_r$ bzw. $\tilde{s}_1^2, \dots, \tilde{s}_r^2$ bezeichnet sind, berechnet sich die mittlere quadratische Abweichung für die Gesamtmasse folgendermaßen:

$$\tilde{s}_{Ges}^2 = \frac{1}{n} \sum_{j=1}^r n_j \tilde{s}_j^2 + \frac{1}{n} \sum_{j=1}^r n_j (\bar{x}_j - \bar{x}_{Ges})^2$$

wobei $n_j = |E_j|$ und $\bar{x}_{Ges} = \frac{1}{n} \sum_{j=1}^r n_j \bar{x}_j$.

2 Zweidimensionale Merkmale

2.1 Gemeinsame Häufigkeiten, Randhäufigkeiten, bedingte Häufigkeiten

Bezeichnungen:

- Urliste: $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$
- Merkmalsausprägungen: a_1, a_2, \dots, a_k für X bzw. b_1, b_2, \dots, b_m für Y

Gemeinsame Häufigkeiten:

$$h_{ij} = h(a_i, b_j) = \sum_{\ell=1}^n 1_{\{x_\ell=a_i\}} 1_{\{y_\ell=b_j\}} \quad \text{absolute Häufigkeiten}$$

$$f_{ij} = f(a_i, b_j) = \frac{h_{ij}}{n} \quad \text{relative Häufigkeiten}$$

Randhäufigkeiten:

$$h_{i\bullet} = \sum_{j=1}^m h_{ij}, \quad h_{\bullet j} = \sum_{i=1}^k h_{ij} \quad (\text{absolut})$$

$$f_{i\bullet} = \frac{h_{i\bullet}}{n} = f_X(a_i), \quad f_{\bullet j} = \frac{h_{\bullet j}}{n} = f_Y(b_j) \quad (\text{relativ})$$

Bedingte relative Häufigkeiten:

$$f_X(a_i|b_j) = \frac{f(a_i, b_j)}{f_Y(b_j)} = \frac{h_{ij}}{h_{\bullet j}}, \quad f_Y(b_j|a_i) = \frac{f(a_i, b_j)}{f_X(a_i)} = \frac{h_{ij}}{h_{i\bullet}}$$

2.2 Assoziation bei nominalen Merkmalen

 χ^2 -Koeffizient:

$$\chi^2 = \sum_{i=1}^k \sum_{j=1}^m \frac{(h_{ij} - \tilde{h}_{ij})^2}{\tilde{h}_{ij}} = n \sum_{i=1}^k \sum_{j=1}^m \frac{(f_{ij} - \tilde{f}_{ij})^2}{\tilde{f}_{ij}}$$

$$\text{wobei } \tilde{h}_{ij} = \frac{h_{i\bullet} h_{\bullet j}}{n}, \quad \tilde{f}_{ij} = f_{i\bullet} f_{\bullet j}.$$

Kontingenzkoeffizient:

$$K = \sqrt{\frac{\chi^2}{\chi^2 + n}}$$

korrigierter Kontingenzkoeffizient:

$$K^* = \frac{K}{K_{max}}$$

mit: $K_{max} = \sqrt{\frac{M-1}{M}}$, wobei $M = \min\{k, m\}$.

Chance (odds) bzw. Risiko:

$$\gamma(j_1, j_2|i) = \gamma_{Y|X}(b_{j_1}, b_{j_2}|a_i) = \frac{h_{ij_1}}{h_{ij_2}}$$

Relative Chance (odds ratio):

$$\gamma(j_1, j_2|i_1, i_2) = \gamma_{Y|X}(b_{j_1}, b_{j_2}|a_{i_1}, a_{i_2}) = \frac{\frac{h_{i_1 j_1}}{h_{i_1 j_2}}}{\frac{h_{i_2 j_1}}{h_{i_2 j_2}}} = \frac{h_{i_1 j_1} h_{i_2 j_2}}{h_{i_1 j_2} h_{i_2 j_1}}$$

Spezialfall: Vierfeldertafel

χ^2 -Koeffizient:

a	b	a+b
c	d	c+d
a+c	b+d	n

$$\chi^2 = \frac{n(ad - bc)^2}{(a+b)(a+c)(b+d)(c+d)}$$

Kreuzproduktverhältnis (Odds-Ratio, relative Chance):

$$\begin{array}{|c|c|} \hline h_{11} & h_{12} \\ \hline h_{21} & h_{22} \\ \hline \end{array} \quad \gamma = \frac{h_{11}/h_{12}}{h_{21}/h_{22}} = \frac{h_{11}h_{22}}{h_{12}h_{21}}$$

2.3 Korrelationsrechnung für metrische und ordinale Merkmale

Bravais-Pearson-Korrelationskoeffizient

$$r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}} = \frac{\sum_{i=1}^n x_i y_i - n\bar{x}\bar{y}}{\sqrt{\left(\sum_{i=1}^n x_i^2 - n\bar{x}^2\right) \left(\sum_{i=1}^n y_i^2 - n\bar{y}^2\right)}}$$

Rangkorrelationskoeffizient von Spearman

$$r_{SP} = \frac{\sum_{i=1}^n (rg(x_i) - \bar{rg}_X)(rg(y_i) - \bar{rg}_Y)}{\sqrt{\sum_{i=1}^n (rg(x_i) - \bar{rg}_X)^2 \sum_{i=1}^n (rg(y_i) - \bar{rg}_Y)^2}}$$

Alternative Darstellungsform (bei Abwesenheit von Bindungen)

$$r_{SP} = 1 - \frac{6 \sum_{i=1}^n d_i^2}{n(n^2 - 1)}$$

mit $d_i = rg(x_i) - rg(y_i)$

korrigierter Rangkorrelationskoeffizient (bei Vorliegen von Bindungen)

$$r_{SP}^* = \frac{n(n^2 - 1) - \frac{1}{2} \sum_{i=1}^k s_i(s_i^2 - 1) - \frac{1}{2} \sum_{j=1}^m r_j(r_j^2 - 1) - 6 \sum_{i=1}^n d_i^2}{\sqrt{n(n^2 - 1) - \sum_{i=1}^k s_i(s_i^2 - 1)} \sqrt{n(n^2 - 1) - \sum_{j=1}^m r_j(r_j^2 - 1)}}$$

mit $s_i = \sum_{\ell=1}^n \mathbf{1}_{\{x_\ell = a_i\}}$, $i = 1, \dots, k$, und $r_j = \sum_{\ell=1}^n \mathbf{1}_{\{x_\ell = b_j\}}$, $j = 1, \dots, m$

Kendall's τ ohne Bindungen

- Sei (x_i, y_i) das Beobachtungstupel des i -ten Merkmalsträgers. Wir betrachten nun ein Paar von Beobachtungstupeln (x_i, y_i) und (x_j, y_j) . Sei o.B.d.A. $x_i < x_j$. Das Paar heißt **konkordant**, wenn auch $y_i < y_j$. Das Paar heißt **diskordant**, wenn $y_i > y_j$.

Insgesamt gibt es $\binom{n}{2} = \frac{n(n-1)}{2}$ Paare, die man überprüfen muss. Dabei gelten die folgenden Bezeichnungen:

- N_c = Anzahl der konkordanten Paare
- N_d = Anzahl der diskordanten Paare

$$\tau = \frac{N_c - N_d}{\frac{1}{2}n(n-1)}$$

Kendall's τ mit Bindungen

- Ist ein Paar weder konkordant noch diskordant, ist es ein **Tie** bzw. eine **Bindung**, d.h. entweder gilt $x_i = x_j$ oder $y_i = y_j$ oder beides:
- N_c und N_d wie im Fall ohne Bindungen
- T_x = Anzahl der Paare mit $x_i = x_j$ aber $y_i \neq y_j$ (x-Ties)
- T_y = Anzahl der Paare mit $y_i = y_j$ aber $x_i \neq x_j$ (y-Ties)
- T_{xy} = Anzahl der Paare mit $x_i = x_j$ und zugleich $y_i = y_j$ (spielen keine Rolle für die Berechnung)
- Da für jedes zu überprüfende Paar nur einer der oberen fünf Fälle in Frage kommt, gilt also: $\binom{n}{2} = \frac{n(n-1)}{2} = N_c + N_d + T_x + T_y + T_{xy}$

$$\tau = \frac{N_c - N_d}{\sqrt{(N_c + N_d + T_x)(N_c + N_d + T_y)}}$$

3 Konzentrationsmaße**3.1 Lorenzkurve und Gini-Koeffizient****Lorenzkurve:**

Für die geordnete Urliste $x_{(1)} \leq \dots \leq x_{(n)}$ ergibt sich die *Lorenzkurve* als Streckenzug durch die Punkte

$$(0, 0), (u_1, v_1), \dots, (u_n, v_n) = (1, 1)$$

mit

$$u_j = j/n \quad \text{und} \quad v_j = \frac{\sum_{i=1}^j x_{(i)}}{\sum_{i=1}^n x_{(i)}}.$$

Bei Häufigkeitsdaten mit den Klassenmitten a_1, \dots, a_k :

$$\tilde{u}_j = \sum_{i=1}^j h_i/n = \sum_{i=1}^j f_i \quad \text{und} \quad \tilde{v}_j = \frac{\sum_{i=1}^j f_i a_i}{\sum_{i=1}^k f_i a_i} = \frac{\sum_{i=1}^j h_i a_i}{n\bar{x}} \quad j = 1, \dots, k.$$

Gini-Koeffizient:

$$G = \frac{\text{Fläche zwischen Diagonale und Lorenzkurve}}{\text{Fläche zwischen Diagonale und } u\text{-Achse}}$$

$$= 2 \cdot \text{Fläche zwischen Diagonale und Lorenzkurve}$$

Geordnete Urliste:

$$G = \frac{2 \sum_{i=1}^n i x_{(i)}}{n \sum_{i=1}^n x_{(i)}} - \frac{n+1}{n}.$$

Häufigkeitsdaten mit den Klassenmitten $a_1 < \dots < a_k$:

$$G = 1 - \frac{1}{n} \sum_{j=1}^k h_j (\tilde{v}_{j-1} + \tilde{v}_j)$$

bzw.

$$G = \frac{\sum_{i=1}^k (\tilde{u}_{i-1} + \tilde{u}_i) h_i a_i}{\sum_{i=1}^k h_i a_i} - 1, \quad \text{wobei} \quad \tilde{u}_i = \sum_{j=1}^i h_j/n, \quad \tilde{v}_i = \sum_{j=1}^i f_j a_j / \sum_{j=1}^k f_j a_j.$$

Normierter Gini-Koeffizient (Lorenz-Münzner-Koeffizient):

$$G^* = \frac{n}{n-1} G, \quad G^* \in [0; 1]$$

3.2 Konzentrationsrate CR_g **Konzentrationsrate CR_g :**

Für vorgegebenes g und $x_1 \leq \dots \leq x_n$ bildet man

$$CR_g = \sum_{i=n-g+1}^n p_i, \quad \text{wobei} \quad p_i = \frac{x_i}{\sum_{j=1}^n x_j}$$

den Merkmalsanteil der i -ten Einheit bezeichnet.

3.3 Herfindahl-Index

Herfindahl-Index:

$$H = \sum_{i=1}^n p_i^2, \quad \text{wobei} \quad p_i = \frac{x_i}{\sum_{j=1}^n x_j}.$$

4 Elementare Wahrscheinlichkeitsrechnung

4.1 Rechenregeln für Wahrscheinlichkeiten

Axiome von Kolmogorow:

- $P(A) \geq 0$ für jedes Ereignis A
- $P(\Omega) = 1$
- $P(A_1 \cup A_2 \cup A_3 \cup \dots) = P(A_1) + P(A_2) + P(A_3) + \dots$
für endlich oder abzählbar unendlich viele paarweise disjunkte Ereignisse, d.h. Ereignisse mit $A_i \cap A_j = \emptyset$ für alle $i \neq j$

Folgerungen:

- $P(A) \leq 1$
- $P(\emptyset) = 0$
- Aus $A \subset B$ folgt $P(A) \leq P(B)$
- $P(\bar{A}) = 1 - P(A)$

Allgemeiner Additionssatz für $n = 2$:

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

Allgemeiner Additionssatz für $n = 3$:

$$P(A \cup B \cup C) = P(A) + P(B) + P(C) - P(A \cap B) - P(A \cap C) - P(B \cap C) + P(A \cap B \cap C)$$

Allgemeiner Additionssatz (Ein- und Ausschlußformel für abhängige Ereignisse):

$$\begin{aligned} P(E_1 \cup E_2 \cup \dots \cup E_n) &= \sum_{i=1}^n P(E_i) - \sum_{i_1 < i_2} P(E_{i_1} \cap E_{i_2}) + \dots \\ &+ (-1)^{r+1} \sum_{i_1 < i_2 < \dots < i_r} P(E_{i_1} \cap E_{i_2} \cap \dots \cap E_{i_r}) + \dots \\ &+ (-1)^{n+1} P(E_1 \cap E_2 \cap \dots \cap E_n) \end{aligned}$$

(Die Summation $\sum_{i_1 < i_2 < \dots < i_r} P(E_{i_1} \cap E_{i_2} \cap \dots \cap E_{i_r})$ läuft über alle $\binom{n}{r}$ möglichen Teilmengen von $\{1, 2, \dots, n\}$.)

4.2 Kombinatorik

Anzahl möglicher Stichproben vom Umfang N aus Grundgesamtheit vom Umfang N :

- Permutation ohne Wiederholung: $P(N) = N!$
- Permutation mit Wiederholung:

$$P^W(N | g_1, \dots, g_r) = \frac{N!}{g_1! \cdot \dots \cdot g_r!}$$

mit r Gruppen mit jeweils gleichen Elementen. Es muss gelten $g_1 + \dots + g_r = N$.

Anzahl möglicher Stichproben vom Umfang n aus Grundgesamtheit vom Umfang N :

- Kombination ohne Wiederholung: Modell ohne Zurücklegen ohne Berücksichtigung der Reihenfolge $K(N, n) = \binom{N}{n}$.
- Kombination mit Wiederholung: Modell mit Zurücklegen ohne Berücksichtigung der Reihenfolge $K^W(N, n) = \binom{N+n-1}{n}$.
- Variation ohne Wiederholung: Modell ohne Zurücklegen mit Berücksichtigung der Reihenfolge $V(N, n) = \frac{N!}{(N-n)!}$.
- Variation mit Wiederholung: Modell mit Zurücklegen mit Berücksichtigung der Reihenfolge $V^W(N, n) = N^n$.

4.3 Bedingte Wahrscheinlichkeiten

Definition:

$$P(A|B) = \frac{P(A \cap B)}{P(B)}, \quad \text{falls } P(B) > 0$$

Folgerungen:

- $P(A \cap B) = P(A|B)P(B)$, falls $P(B) > 0$
- $P(B \cap A) = P(B|A)P(A)$, falls $P(A) > 0$
- $P(A_1 \cap \dots \cap A_m) = P(A_1)P(A_2|A_1) \dots P(A_m|A_1 \cap \dots \cap A_{m-1})$, falls $P(A_1 \cap \dots \cap A_{m-1}) > 0$.

4.4 Satz von der totalen Wahrscheinlichkeit

Sei A_1, \dots, A_k eine disjunkte Zerlegung von Ω , dann gilt für jedes Ereignis $B \subset \Omega$

$$P(B) = \sum_{i=1}^k P(B|A_i)P(A_i).$$

4.5 Formel von Bayes

Sei A_1, \dots, A_k eine disjunkte Zerlegung von Ω , wobei für mindestens ein i , $i = 1, \dots, k$, $P(A_i) > 0$ und $P(B|A_i) > 0$ erfüllt ist, dann gilt

$$P(A_j|B) = \frac{P(B|A_j)P(A_j)}{\sum_{i=1}^k P(B|A_i)P(A_i)} = \frac{P(B|A_j)P(A_j)}{P(B)} \quad \text{für jedes } j = 1, \dots, k.$$

4.6 Unabhängigkeit zweier Ereignisse

Zwei Ereignisse A und B heißen (stochastisch) unabhängig, wenn gilt

$$P(A \cap B) = P(A) \cdot P(B)$$

bzw. $P(A|B) = P(A)$, falls $P(B) > 0$ bzw. $P(B|A) = P(B)$, falls $P(A) > 0$.

5 Eindimensionale Zufallsvariablen und ihre Verteilungen

5.1 Definition von diskreten und stetigen Zufallsvariablen und Dichten

Diskrete Zufallsvariable:

Eine ZV X heißt *diskret*, falls der Wertebereich von X nur endlich oder abzählbar unendlich viele Werte $x_1, x_2, \dots, x_k, \dots$ annehmen kann.

Wahrscheinlichkeitsfunktion (Dichte):

$$f(x) = \begin{cases} P(X = x_i) = p_i & \text{für } x = x_i \in \{x_1, x_2, \dots, x_k, \dots\} \\ 0 & \text{sonst} \end{cases}$$

Stetige Zufallsvariable:

Eine ZV X heißt *stetig*, wenn es eine Funktion $f(x) \geq 0$ gibt, so dass für jedes Intervall $[a, b]$ gilt:

$$P(a \leq X \leq b) = \int_a^b f(x) dx$$

- Die Funktion $f(x)$ heißt **Dichte** von X
- Es gilt: $P(X = x) = 0$ für alle $x \in \mathbb{R}$.

5.2 Die Verteilungsfunktion einer Zufallsvariablen

Verteilungsfunktion einer Zufallsvariablen:

$$F(x) = P(X \leq x), \quad x \in \mathbb{R}$$

Rechenregeln:

- $P(X = x) = F(x) - P(X < x)$
- $P(X > x) = 1 - F(x)$
- $P(a < X \leq b) = F(b) - F(a)$, falls $a < b$
- $P(a \leq X \leq b) = F(b) - F(a) + P(X = a)$
speziell: $P(a \leq X \leq b) = F(b) - F(a)$ für stetige Verteilungen
- $P(a < X < b) = F(b) - F(a) - P(X = b)$
- $P(a \leq X < b) = F(b) - F(a) + P(X = a) - P(X = b)$

5.3 Zusammenhänge zwischen Dichten und Verteilungsfunktionen

Im diskreten Fall:

$$F(x) = \sum_{x_i \leq x} f(x_i) = \sum_{x_i \leq x} P(X = x_i)$$

$$P(a \leq X \leq b) = \sum_{a \leq x_i \leq b} f(x_i) = \sum_{a \leq x_i \leq b} P(X = x_i)$$

$$f(x_i) = P(X = x_i) = F(x_i) - F(x_{i-1})$$

Im stetigen Fall:

$$F(x) = \int_{-\infty}^x f(t) dt$$

$$P(a \leq X \leq b) = \int_a^b f(t) dt$$

$$f(x) = F'(x) = \frac{dF(x)}{dx}, \quad \text{falls } F(x) \text{ an der Stelle } x \text{ differenzierbar ist.}$$

5.4 Modus, Median und Quantile

Modus: x_{mod} : Jeder Wert x , an dem $f(x)$ maximal ist.**Quantile:**

- Jeder Wert x_p mit $0 < p < 1$, für den

$$P(X \geq x_p) \geq 1 - p \quad \text{und} \quad P(X \leq x_p) \geq p$$

gilt, heißt p -Quantil einer diskreten Verteilung.

- Jeder Wert x_p mit $F(x_p) = p$ heißt p -Quantil einer stetigen Verteilung.

Median:Jedes 50%-Quantil heißt Median ($p = 0.5$).

5.5 Erwartungswert, Varianz und Standardabweichung

Erwartungswert:

$$\mu = E(X) = \begin{cases} \sum_x x f(x), & \text{falls } X \text{ diskret} \\ \int_{-\infty}^{+\infty} x f(x) dx, & \text{falls } X \text{ stetig} \end{cases}$$

Varianz:

$$\sigma^2 = \text{Var}(X) = \begin{cases} \sum_x (x - E(X))^2 f(x), & \text{falls } X \text{ diskret} \\ \int_{-\infty}^{+\infty} (x - E(X))^2 f(x) dx, & \text{falls } X \text{ stetig} \end{cases}$$

Standardabweichung:

$$\sigma = +\sqrt{\sigma^2} = +\sqrt{\text{Var}(X)}$$

5.6 Rechenregeln und Eigenschaften von Erwartungswerten

Transformation:

- Die Zufallsvariable $Y = g(X)$ besitzt den Erwartungswert

$$E(Y) = \begin{cases} \sum_x g(x) f(x), & \text{falls } X \text{ diskret} \\ \int_{-\infty}^{+\infty} g(x) f(x) dx, & \text{falls } X \text{ stetig} \end{cases}$$

- Spezialfall: lineare Transformation

$$E(aX + b) = aE(X) + b \quad \text{für alle } a, b \in \mathbb{R}$$

$$\text{Var}(aX + b) = a^2 \text{Var}(X)$$

Verschiebungssatz:

$$\text{Var}(X) = E(X^2) - [E(X)]^2$$

5.7 Spezielle diskrete Verteilungen

Verteilung	Wahrscheinlichkeitsfunktion	E(X)	Var(X)
$X \sim Ber(\pi)$ Bernoulliverteilung	$f(x) = \begin{cases} \pi & \text{für } x = 1, \\ 1 - \pi & \text{für } x = 0, \\ 0, & \text{sonst} \end{cases}$	π	$\pi(1 - \pi)$
$X \sim G(\pi)$ Geometrische Verteilung	$f(x) = (1 - \pi)^{x-1}\pi$ für $x = 1, 2, \dots$	$\frac{1}{\pi}$	$\frac{1 - \pi}{\pi^2}$
$X \sim B(n, \pi)$ Binomialverteilung	$f(x) = \begin{cases} \binom{n}{x}\pi^x(1 - \pi)^{n-x} & \text{für } x = 0, 1, \dots, n \\ 0, & \text{sonst} \end{cases}$	$n\pi$	$n\pi(1 - \pi)$
$X \sim NB(r, \pi)$ negative Binomialverteilung	$f(x) = \begin{cases} \binom{x-1}{r-1}\pi^r(1 - \pi)^{x-r} & \text{für } x = r, r + 1, \dots \\ 0, & \text{sonst} \end{cases}$	$\frac{r}{\pi}$	$\frac{r(1 - \pi)}{\pi^2}$
$X \sim H(n, M, N)$ Hypergeometrische Verteilung	$f(x) = \frac{\binom{M}{x}\binom{N-M}{n-x}}{\binom{N}{n}}$ für $x \in \tau$	$n\frac{M}{N}$	s.u.
$X \sim Po(\lambda)$ Poissonverteilung	$f(x) = \begin{cases} \frac{\lambda^x}{x!}e^{-\lambda} & \text{für } x = 0, 1, 2, \dots \\ 0, & \text{sonst} \end{cases}$ ($\lambda > 0$)	λ	λ

Varianz der hypergeometrischen Verteilung: $Var(X) = n\frac{M}{N}\left(1 - \frac{M}{N}\right)\frac{N-n}{N-1}$

5.8 Spezielle stetige Verteilungen

<p>Stetige Gleichverteilung: $X \sim U[a; b]$</p> <p>Dichte und Verteilungsfunktion:</p> $f(x) = \begin{cases} \frac{1}{b-a}, & \text{für } a \leq x \leq b, \\ 0, & \text{sonst.} \end{cases} \quad F(x) = \begin{cases} 0, & x < a, \\ \frac{x-a}{b-a}, & a \leq x \leq b, \\ 1, & x > b. \end{cases}$ <p>Erwartungswert und Varianz:</p> $E(X) = \frac{a+b}{2}, \quad Var(X) = \frac{(b-a)^2}{12}.$
--

<p>Exponentialverteilung: $X \sim Exp(\lambda)$</p> <p>Dichte und Verteilungsfunktion:</p> $f(x) = \begin{cases} \lambda e^{-\lambda x} & \text{für } x \geq 0 \\ 0 & \text{sonst} \end{cases} \quad F(x) = \begin{cases} 1 - e^{-\lambda x}, & x \geq 0, \\ 0, & \text{sonst.} \end{cases} \quad (\lambda > 0)$ <p>Erwartungswert und Varianz:</p> $E(X) = \frac{1}{\lambda}, \quad Var(X) = \frac{1}{\lambda^2}.$
<p>Normalverteilung: $X \sim N(\mu, \sigma^2)$</p> <p>Dichte:</p> $f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(\frac{-(x - \mu)^2}{2\sigma^2}\right) \quad \text{für } x \in \mathbb{R}$ <p>Erwartungswert und Varianz:</p> $E(X) = \mu, \quad Var(X) = \sigma^2.$
<p>Weibullverteilung: $X \sim W(\lambda, \alpha)$</p> <p>Dichte und Verteilungsfunktion:</p> $f(x) = \begin{cases} \lambda\alpha(\lambda x)^{\alpha-1} \exp(-(\lambda x)^\alpha), & \text{für } x \geq 0, \\ 0, & \text{sonst.} \end{cases} \quad F(x) = \begin{cases} 1 - e^{-(\lambda x)^\alpha}, & x \geq 0, \\ 0, & \text{sonst.} \end{cases} \quad (\lambda, \alpha > 0)$ <p>Erwartungswert und Varianz:</p> $E(X) = \frac{1}{\lambda}\Gamma\left(\frac{\alpha+1}{\alpha}\right), \quad Var(X) = \frac{1}{\lambda^2}\left(\Gamma\left(\frac{\alpha+2}{\alpha}\right) - \left[\Gamma\left(\frac{\alpha+1}{\alpha}\right)\right]^2\right)$ <p>mit</p> $\Gamma(x) = \int_0^\infty e^{-t}t^{x-1}dt \quad (x > 0).$

5.9 Die Chi-Quadrat-, Student- und Fisher-Verteilung

Chi-Quadrat- (χ^2) -Verteilung:

$$Z = \sum_{i=1}^n X_i^2 \sim \chi_n^2 \quad \text{d.h. } \chi^2\text{-verteilt mit } n \text{ Freiheitsgraden}$$

falls X_1, \dots, X_n unabhängige, standardnormalverteilte Zufallsvariablen sind

Student- (t) -Verteilung:

$$T = \frac{X}{\sqrt{Z/n}} \sim t_n \quad \text{d.h. } t\text{-verteilt mit } n \text{ Freiheitsgraden}$$

falls X standardnormalverteilt, $Z \sim \chi_n^2$ -verteilt und X und Z unabhängig sind

Fisher- (F) -Verteilung:

$$Z = \frac{X/m}{Y/n} \sim F_{m,n} \quad \text{d.h. } F\text{-verteilt mit } m \text{ und } n \text{ Freiheitsgraden}$$

falls $X \sim \chi_m^2$ - und $Y \sim \chi_n^2$ -verteilt und unabhängig sind.

6 Zweidimensionale Zufallsvariablen und ihre Verteilungen

6.1 Definition zweidimensionaler Zufallsvariablen

Zweidimensionale diskrete Zufallsvariable:

Seien X und Y zwei diskrete ZV, wobei X die Werte x_1, x_2, \dots und Y die Werte y_1, y_2, \dots annehmen kann, so ist (X, Y) eine *zweidimensionale diskrete* Zufallsvariable mit Werten (x_i, y_j) , $i = 1, 2, \dots$, $j = 1, 2, \dots$

Gemeinsame Wahrscheinlichkeitsfunktion:

$$f(x, y) = \begin{cases} P(X = x, Y = y) & \text{für } (x, y) \in \{(x_1, y_1), (x_1, y_2), \dots\} \\ 0 & \text{sonst} \end{cases}$$

Randverteilungen:

$$f_X(x) = P(X = x) = \sum_j f(x, y_j) \quad \text{und} \quad f_Y(y) = P(Y = y) = \sum_i f(x_i, y)$$

Zweidimensionale stetige Zufallsvariable:

Die Zufallsvariablen X und Y sind *gemeinsam stetig verteilt*, wenn es eine **zweidimensionale Dichtefunktion** $f(x, y) \geq 0$ gibt, so daß gilt

$$P(a \leq X \leq b, c \leq Y \leq d) = \int_a^b \int_c^d f(x, y) dy dx.$$

Randdichten:

$$f_X(x) = \int_{-\infty}^{\infty} f(x, y) dy \quad \text{und} \quad f_Y(y) = \int_{-\infty}^{\infty} f(x, y) dx$$

Bedingte Wahrscheinlichkeitsfunktionen/Dichten:

$$f_X(x|y) = \frac{f(x, y)}{f_Y(y)} \quad \text{und} \quad f_Y(y|x) = \frac{f(x, y)}{f_X(x)}$$

Gemeinsame Verteilungsfunktion:

$$F(x, y) = P(X \leq x, Y \leq y) = \begin{cases} \sum_{x_i \leq x} \sum_{y_j \leq y} f(x_i, y_j) & \text{(diskret)} \\ \int_{-\infty}^x \int_{-\infty}^y f(u, v) dv du & \text{(stetig)} \end{cases}$$

6.2 Unabhängigkeit, Kovarianz und Korrelation

Unabhängigkeit von zwei Zufallsvariablen:

Zwei Zufallsvariablen X und Y heißen unabhängig, wenn für alle x und y gilt

$$f(x, y) = f_X(x)f_Y(y)$$

Kovarianz:

$$\begin{aligned} \text{Cov}(X, Y) &= E\{[X - E(X)][Y - E(Y)]\} \\ &= \begin{cases} \sum_i \sum_j f(x_i, y_j)(x_i - E(X))(y_j - E(Y)) & \text{(diskret)} \\ \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} f(x, y)(x - E(X))(y - E(Y)) dx dy & \text{(stetig)} \end{cases} \end{aligned}$$

Verschiebungssatz:

$$\text{Cov}(X, Y) = E(X \cdot Y) - E(X) \cdot E(Y)$$

$$\text{mit } E(X \cdot Y) = \begin{cases} \sum_i \sum_j f(x_i, y_j) x_i y_j & \text{(diskret)} \\ \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy f(x, y) dy dx & \text{(stetig)} \end{cases}$$

Lineare Transformation:

Für die Zufallsvariablen $\tilde{X} = a_X X + b_X$ und $\tilde{Y} = a_Y Y + b_Y$ gilt

$$\text{Cov}(\tilde{X}, \tilde{Y}) = a_X a_Y \text{Cov}(X, Y)$$

Korrelationskoeffizient:

$$\rho = \rho(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)} \sqrt{\text{Var}(Y)}} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

Unkorreliertheit:

Die Zufallsvariablen X und Y heißen unkorreliert, wenn gilt

$$\rho(X, Y) = 0 \quad \text{bzw.} \quad \text{Cov}(X, Y) = 0$$

Varianz der Summe zweier Zufallsvariablen:

$$\text{Var}(X + Y) = \text{Var}(X) + \text{Var}(Y) + 2\text{Cov}(X, Y)$$

Erwartungswert und Varianz von Linearkombinationen:

Die gewichtete Summe

$$X = a_1 X_1 + \dots + a_n X_n$$

der Zufallsvariablen X_1, \dots, X_n besitzt den Erwartungswert

$$E(X) = a_1 E(X_1) + \dots + a_n E(X_n)$$

und die Varianz

$$\text{Var}(X) = \sum_{i=1}^n a_i^2 \text{Var}(X_i) + 2 \sum_{i < j} a_i a_j \text{Cov}(X_i, X_j)$$

7 Ergänzungen zu Zufallsvariablen**7.1 Grenzwertsätze****Das Gesetz der großen Zahlen**

Sei X_1, \dots, X_n, \dots eine Folge von Zufallsvariablen mit Erwartungswert μ und Varianz σ^2 . Dann gilt für alle $c > 0$:

$$\lim_{n \rightarrow \infty} P(|\bar{X}_n - \mu| \leq c) = 1.$$

Man sagt: \bar{X}_n konvergiert nach Wahrscheinlichkeit gegen μ .

Satz von Glivenko–Cantelli

Sei X eine Zufallsvariable mit Verteilungsfunktion $F(x)$ und $F_n(x)$ die empirische Verteilungsfunktion bei einer Stichprobe vom Umfang n . Dann gilt für jedes $c > 0$:

$$\lim_{n \rightarrow \infty} P(\sup_x |F_n(x) - F(x)| \leq c) = 1, \quad x \in \mathbb{R}.$$

Zentraler Grenzwertsatz

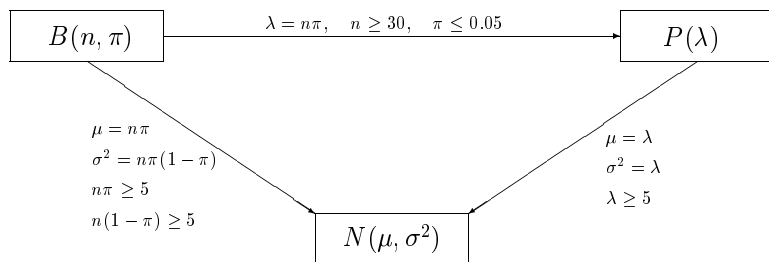
Seien X_1, \dots, X_n unabhängige, identisch verteilte Zufallsvariablen mit Erwartungswert μ und Varianz $\sigma^2 > 0$. Dann konvergiert die Verteilungsfunktion $F_n(z)$ der standardisierten Summe

$$Z_n = \frac{X_1 + \dots + X_n - n\mu}{\sigma \sqrt{n}} = \frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{X_i - \mu}{\sigma}$$

für $n \rightarrow \infty$ an jeder Stelle $z \in \mathbb{R}$ gegen die Verteilungsfunktion $\Phi(z)$ der Standardnormalverteilung, d.h.

$$Z_n \stackrel{d}{\sim} N(0, 1).$$

7.2 Approximation von Verteilungen



7.3 Ungleichung von Tschebyschew

Für eine Zufallsvariable X mit Erwartungswert μ und Varianz σ^2 gelten für jedes $c > 0$ die folgenden Ungleichungen:

$$P(|X - \mu| \geq c) \leq \frac{\sigma^2}{c^2} \quad \text{und} \quad P(|X - \mu| < c) \geq 1 - \frac{\sigma^2}{c^2}.$$

8 Testen und Schätzen

Bezeichnungen:

- Merkmal X : metrisch oder dichotom (Bernoulliverteilt)
- Unbekannter Parameter der Verteilung von X : θ
- Stichprobenvariablen: X_1, X_2, \dots, X_n
- Annahme: X_1, \dots, X_n unabhängig und identisch verteilt wie X
- Stichprobenmittelwert: $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ (entspricht der relativen Häufigkeit, falls X dichotom)
- Stichprobenvarianz: $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$
- Realisationen: x_1, x_2, \dots, x_n
- Schätzfunktion/Schätzstatistik/Teststatistik: $T = g(X_1, \dots, X_n)$
- Schätzwert: $\hat{\theta} = g(x_1, \dots, x_n)$

Abkürzungen der Quantile:

z_p ... p -Quantil der Standardnormalverteilung (siehe S. 48)
 $t_{p,k}$... p -Quantil der t -Verteilung mit k Freiheitsgraden (siehe S. 49)
 $\chi_{p,k}^2$... p -Quantil der χ^2 -Verteilung mit k Freiheitsgraden (siehe S. 50)

8.1 Punktschätzung

Erwartungstreue:

Eine Schätzstatistik T heißt erwartungstreu für θ , wenn gilt

$$E_\theta(T) = \theta$$

Asymptotische Erwartungstreue:

Eine Schätzstatistik T heißt asymptotisch erwartungstreu für θ , wenn gilt

$$\lim_{n \rightarrow \infty} E_\theta(T) = \theta$$

Bias:

Eine nicht erwartungstreu Schätzstatistik heißt verzerrt. Die Stärke der Verzerrung wird durch den Bias angegeben

$$\text{Bias}_\theta(T) = E_\theta(T) - \theta$$

Erwartete mittlere quadratische Abweichung (MSE):

Die erwartete mittlere quadratische Abweichung (mean squared error) ist bestimmt durch

$$\text{MSE} = E_\theta([T - \theta]^2) = \text{Var}_\theta(T) + \text{Bias}_\theta(T)^2$$

MSE-Konsistenz (Konsistenz im quadratischen Mittel):

Eine Schätzstatistik heißt MSE-konsistent, wenn gilt

$$\lim_{n \rightarrow \infty} \text{MSE} = 0.$$

Eine Schätzstatistik T heißt schwach konsistent, wenn zu beliebigem $\epsilon > 0$ gilt

$$\lim_{n \rightarrow \infty} P(|T - \theta| < \epsilon) = 1 \quad \text{bzw.} \quad \lim_{n \rightarrow \infty} P(|T - \theta| \geq \epsilon) = 0.$$

MSE-Wirksamkeit (MSE-Effizienz):

Von zwei Schätzstatistiken T_1 und T_2 heißt T_1 MSE-wirksamer (MSE-effizient), wenn gilt

$$\text{MSE}(T_1) \leq \text{MSE}(T_2)$$

Momenten-Schätzung:

Die Parameter $\theta_1, \dots, \theta_k$ der theoretischen Verteilung von X werden als Funktionen der Momente

$$\theta_i = h_i(\mu_1, \dots, \mu_\ell), \quad \mu_j = E(X^j), \quad j = 1, \dots, \ell, \quad i = 1, \dots, k \quad (1)$$

angegeben. Die Momentenschätzer $\hat{\theta}_i, i = 1, \dots, k$ werden berechnet, indem in (1) die empirischen Momente $\hat{\mu}_j = \frac{1}{n} \sum_{i=1}^n x_i^j$ für die Momente μ_j eingesetzt werden.

Maximum-Likelihood-Schätzung:

Der Maximum-Likelihood-Schätzwert $\hat{\theta}$ ist die Lösung der Gleichung

$$L(\hat{\theta}) = \max_{\theta} L(\theta)$$

bzw.

$$f(x_1, \dots, x_n | \hat{\theta}) = \max_{\theta} f(x_1, \dots, x_n | \theta)$$

Zur praktischen Berechnung von $\hat{\theta}$ wird üblicherweise die *Log-Likelihood* $\ln L(\theta)$ gebildet und diese bezüglich θ maximiert.

Bayes-Schätzung:

- Bayes-Inferenz

Sei

- $f(x|\theta)$ die Wahrscheinlichkeitsfunktion bzw. Dichte von X , gegeben θ
- $L(\theta) = f(x_1, \dots, x_n|\theta)$ die gemeinsame Dichte bzw. Likelihoodfunktion für n unabhängige Wiederholungen
- $f(\theta)$ eine *a priori Dichte* für den unbekannt Parameter

Dann ist die *a posteriori Dichte* als Basis zur Bestimmung eines Bayes-Schätzers für θ gegeben durch

$$f(\theta|x_1, \dots, x_n) = \frac{f(x_1|\theta) \cdots f(x_n|\theta)f(\theta)}{\int f(x_1|\theta) \cdots f(x_n|\theta)f(\theta)d\theta} = \frac{L(\theta)f(\theta)}{\int L(\theta)f(\theta)d\theta}$$

- Bayes-Schätzer

Mögliche Bayes-Schätzer für θ basierend auf der *a posteriori Dichte* sind

- *a posteriori Erwartungswert:*

$$\hat{\theta} = E(\theta|x_1, \dots, x_n) = \int \theta f(\theta|x_1, \dots, x_n)d\theta$$

- *a posteriori Modus, maximum a posteriori Schätzer:*
Wähle für $\hat{\theta}$ denjenigen Parameterwert, für den die *a posteriori Dichte* maximal wird, d.h.

$$L(\hat{\theta})f(\hat{\theta}) = \max_{\theta} L(\theta)f(\theta)$$

bzw.

$$\ln L(\hat{\theta}) + \ln f(\hat{\theta}) = \max_{\theta} \{\ln L(\theta) + \ln f(\theta)\}$$

8.2 Intervallschätzung**(1 - α)-Konfidenzintervall:**

Die beiden Schätzstatistiken

$$G_u = g_u(X_1, \dots, X_n) \quad \text{und} \quad G_o = g_o(X_1, \dots, X_n)$$

bilden ein $(1 - \alpha)$ -Konfidenzintervall für θ , falls gilt

$$P(G_u \leq G_o) = 1 \quad \text{und} \quad P(G_u \leq \theta \leq G_o) = 1 - \alpha$$

Das Konfidenzintervall besitzt dann die Gestalt

$$\text{KI} = [g_u(x_1, \dots, x_n), g_o(x_1, \dots, x_n)]$$

Einseitige $(1 - \alpha)$ -Konfidenzintervalle:

Für $G_u = -\infty$ bzw. $G_o = \infty$ ergibt sich

$$P(\theta \leq G_o) = 1 - \alpha \quad \text{bzw.} \quad P(G_u \leq \theta) = 1 - \alpha$$

und man erhält die Konfidenzintervalle

$$\text{KI} = (-\infty, g_o(x_1, \dots, x_n)] \quad \text{bzw.} \quad \text{KI} = [g_u(x_1, \dots, x_n), \infty)$$

8.3 Spezielle Schätzprobleme

Verteilung	θ	$\hat{\theta}$	$(1 - \alpha)$ -Konfidenzintervall
$X \sim N(\mu, \sigma^2)$, σ^2 bekannt	μ	\bar{X}	$\left[\bar{X} - z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \right]$
$X \sim N(\mu, \sigma^2)$, σ^2 unbekannt	μ	\bar{X}	$\left[\bar{X} - t_{1-\frac{\alpha}{2}, n-1} \frac{S}{\sqrt{n}}, \bar{X} + t_{1-\frac{\alpha}{2}, n-1} \frac{S}{\sqrt{n}} \right]$
$X \sim N(\mu, \sigma^2)$, μ beliebig	σ^2	S^2	$\left[(n-1)S^2 \frac{1}{\chi_{1-\frac{\alpha}{2}, n-1}^2}, (n-1)S^2 \frac{1}{\chi_{\frac{\alpha}{2}, n-1}^2} \right]$

Verteilung	θ	$\hat{\theta}$	approximatives $(1 - \alpha)$ -Konfidenzintervall
X beliebig verteilt, $n \geq 30$, $E(X) = \mu$, $\text{Var}(X) = \sigma^2$ bekannt	μ	\bar{X}	$\left[\bar{X} - z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}}, \bar{X} + z_{1-\frac{\alpha}{2}} \frac{\sigma}{\sqrt{n}} \right]$
X beliebig verteilt, $n \geq 30$, $E(X) = \mu$, $\text{Var}(X) = \sigma^2$ unbekannt	μ	\bar{X}	$\left[\bar{X} - z_{1-\frac{\alpha}{2}} \frac{S}{\sqrt{n}}, \bar{X} + z_{1-\frac{\alpha}{2}} \frac{S}{\sqrt{n}} \right]$
X dichotom, $n \geq 30$ $\Leftrightarrow \sum_{i=1}^n X_i \sim B(n, \pi)$, $n \geq 30$	π	\bar{X}	$\left[\hat{\pi} - z_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{\pi}(1-\hat{\pi})}{n}}, \hat{\pi} + z_{1-\frac{\alpha}{2}} \sqrt{\frac{\hat{\pi}(1-\hat{\pi})}{n}} \right]$

8.4 Testen von Hypothesen

Statistisches Testproblem:

Nullhypothese H_0 und Alternative H_1 treffen Aussagen über θ

$$H_0 : \theta \in \Theta_0 \quad \text{vs.} \quad H_1 : \theta \in \Theta_1$$

Die Entscheidung für oder gegen H_0 wird anhand einer Prüfgröße (Teststatistik) getroffen.

Fehlentscheidung:

- Fehler 1.Art: H_0 wird verworfen, obwohl H_0 zutrifft
- Fehler 2.Art: H_0 wird beibehalten, obwohl H_1 zutrifft

Signifikanztest:

Falls gilt

$$P(H_0 \text{ verwerfen} | H_0 \text{ trifft zu}) \leq \alpha \quad (\text{d.h. } P(\text{Fehler 1.Art}) \leq \alpha),$$

dann heißt der Test Signifikanztest, oder Test zum *Signifikanzniveau* α

p-Wert:

Der p -Wert ist definiert als die Wahrscheinlichkeit, unter H_0 den beobachteten Prüfgrößenwert oder einen in Richtung der Alternative extremeren Wert zu erhalten.

Gütefunktion:

Für vorgegebenes Signifikanzniveau α und festen Stichprobenumfang n gibt die Gütefunktion g die Wahrscheinlichkeit für einen statistischen Test an, die Nullhypothese über θ zu verwerfen, d.h.

$$g(\theta) = P(H_0 \text{ verwerfen} | \theta).$$

Gilt $\theta \in H_0$, so ist $g(\theta) \leq \alpha$. Falls $\theta \in H_1$, so ist $1 - g(\theta)$ die Wahrscheinlichkeit β für den Fehler 2.Art.

8.5 Spezielle Testprobleme

8.5.1 Einstichproben-Testprobleme

Formulierung der Hypothesen:

- Zweiseitiges Testproblem:

$$(a) \quad H_0 : \theta = \theta_0 \quad \text{vs.} \quad H_1 : \theta \neq \theta_0$$
- Einseitige Testprobleme:

$$(b) \quad H_0 : \theta \geq \theta_0 \quad \text{vs.} \quad H_1 : \theta < \theta_0$$

$$(c) \quad H_0 : \theta \leq \theta_0 \quad \text{vs.} \quad H_1 : \theta > \theta_0$$

Verteilung	θ	Teststatistik	Ablehnbereiche
$X \sim N(\mu, \sigma^2)$, σ^2 bekannt	μ	$Z = \frac{\bar{X} - \mu_0}{\sigma} \sqrt{n}$	(a) $ Z > z_{1-\frac{\alpha}{2}}$ (b) $Z < -z_{1-\alpha}$ (c) $Z > z_{1-\alpha}$
$X \sim N(\mu, \sigma^2)$, σ^2 unbekannt	μ	$T = \frac{\bar{X} - \mu_0}{S} \sqrt{n}$	(a) $ T > t_{1-\frac{\alpha}{2}, n-1}$ (b) $T < -t_{1-\alpha, n-1}$ (c) $T > t_{1-\alpha, n-1}$
X beliebig verteilt, $n \geq 30$, $E(X) = \mu$, $\text{Var}(X) = \sigma^2$ bekannt	μ	$Z = \frac{\bar{X} - \mu_0}{\sigma} \sqrt{n}$	(a) $ Z > z_{1-\frac{\alpha}{2}}$ (b) $Z < -z_{1-\alpha}$ (c) $Z > z_{1-\alpha}$
X beliebig verteilt, $n \geq 30$, $E(X) = \mu$, $\text{Var}(X) = \sigma^2$ unbekannt	μ	$T = \frac{\bar{X} - \mu_0}{S} \sqrt{n}$	(a) $ T > z_{1-\frac{\alpha}{2}}$ (b) $T < -z_{1-\alpha}$ (c) $T > z_{1-\alpha}$
X dichotom $\Leftrightarrow \sum_{i=1}^n X_i \sim B(n, \pi)$	π	$Z = \frac{\hat{\pi} - \pi_0}{\sqrt{\pi_0(1-\pi_0)}} \sqrt{n}$	(a) $ Z > z_{1-\frac{\alpha}{2}}$ (b) $Z < -z_{1-\alpha}$ (c) $Z > z_{1-\alpha}$

Gütefunktion für Spezialfall Einstichproben-Gauß-Test

Gegeben seien iid Zufallsvariablen X_1, \dots, X_n mit $X_i \sim N(\mu, \sigma^2)$, σ^2 bekannt, bzw. mit beliebiger stetiger Verteilung und $E(X_i) = \mu$, $\text{Var}(X_i) = \sigma^2$, $n \geq 30$.

Für die Gütefunktion $g(\mu)$ ergibt sich dann im Fall des Testproblems

(a) $H_0 : \mu = \mu_0$ vs. $H_1 : \mu \neq \mu_0$

$$g(\mu) = \Phi\left(-z_{1-\alpha/2} + \frac{\mu - \mu_0}{\sigma} \sqrt{n}\right) + \Phi\left(-z_{1-\alpha/2} - \frac{\mu - \mu_0}{\sigma} \sqrt{n}\right)$$

(b) $H_0 : \mu \geq \mu_0$ vs. $H_1 : \mu < \mu_0$

$$g(\mu) = \Phi\left(z_\alpha - \frac{\mu - \mu_0}{\sigma} \sqrt{n}\right)$$

(c) $H_0 : \mu \leq \mu_0$ vs. $H_1 : \mu > \mu_0$

$$g(\mu) = 1 - \Phi\left(z_{1-\alpha} - \frac{\mu - \mu_0}{\sigma} \sqrt{n}\right)$$

wobei Φ die Verteilungsfunktion der $N(0,1)$ -Verteilung ist.

8.5.2 Zweistichproben-Mittelwertsvergleiche

Bezeichnungen:

- Metrische Merkmale X und Y
- Unbekannte Parameter: $E(X) = \mu_X$ und $E(Y) = \mu_Y$
- Stichprobenvariablen: X_1, X_2, \dots, X_n und Y_1, Y_2, \dots, Y_m
- Annahmen: X_1, \dots, X_n unabhängig und identisch verteilt wie X
 Y_1, \dots, Y_m unabhängig und identisch verteilt wie Y
 $X_1, \dots, X_n, Y_1, \dots, Y_m$ unabhängig

Formulierung der Hypothesen:

- Zweiseitiges Testproblem:

(a) $H_0 : \mu_X - \mu_Y = \delta_0$ vs. $H_1 : \mu_X - \mu_Y \neq \delta_0$

- Einseitige Testprobleme:

(b) $H_0 : \mu_X - \mu_Y \geq \delta_0$ vs. $H_1 : \mu_X - \mu_Y < \delta_0$

(c) $H_0 : \mu_X - \mu_Y \leq \delta_0$ vs. $H_1 : \mu_X - \mu_Y > \delta_0$

Verteilung	Teststatistik	Ablehnbereiche
$X \sim N(\mu_X, \sigma_X^2)$, $Y \sim N(\mu_Y, \sigma_Y^2)$ σ_X^2, σ_Y^2 bekannt	$Z = \frac{\bar{X} - \bar{Y} - \delta_0}{\sqrt{\frac{\sigma_X^2}{n} + \frac{\sigma_Y^2}{m}}}$	(a) $ Z > z_{1-\frac{\alpha}{2}}$ (b) $Z < -z_{1-\alpha}$ (c) $Z > z_{1-\alpha}$
$X \sim N(\mu_X, \sigma_X^2)$, $Y \sim N(\mu_Y, \sigma_Y^2)$ $\sigma_X^2 = \sigma_Y^2$ unbekannt	$T = \frac{\bar{X} - \bar{Y} - \delta_0}{\sqrt{\left(\frac{1}{n} + \frac{1}{m}\right) \frac{(n-1)S_X^2 + (m-1)S_Y^2}{n+m-2}}}$	(a) $ T > t_{1-\frac{\alpha}{2}, n+m-2}$ (b) $T < -t_{1-\alpha, n+m-2}$ (c) $T > t_{1-\alpha, n+m-2}$
X, Y beliebig verteilt σ_X^2, σ_Y^2 unbekannt, $n, m > 30$	$T = \frac{\bar{X} - \bar{Y} - \delta_0}{\sqrt{\frac{S_X^2}{n} + \frac{S_Y^2}{m}}}$	(a) $ T > z_{1-\frac{\alpha}{2}}$ (b) $T < -z_{1-\alpha}$ (c) $T > z_{1-\alpha}$

8.5.3 Weitere Testprobleme

Vorzeichen-Test:

- Annahmen: X_1, \dots, X_n unabhängige Wiederholungen, X besitzt stetige Verteilungsfunktion

- Hypothesen:

$$\begin{aligned} (a) \quad & H_0 : x_{med} = \delta_0 \quad \text{vs.} \quad H_1 : x_{med} \neq \delta_0 \\ (b) \quad & H_0 : x_{med} \geq \delta_0 \quad \text{vs.} \quad H_1 : x_{med} < \delta_0 \\ (c) \quad & H_0 : x_{med} \leq \delta_0 \quad \text{vs.} \quad H_1 : x_{med} > \delta_0 \end{aligned}$$

- Teststatistik:

$A =$ Anzahl der Stichprobenvariablen mit einem Wert kleiner als δ_0

- Ablehnungsbereiche:

$$\begin{aligned} (a) \quad & A \leq b_{\alpha/2} \quad \text{oder} \quad n - A \leq b_{\alpha/2} \\ (b) \quad & A > o_{1-\alpha} \\ (c) \quad & A \leq b_{\alpha} \end{aligned}$$

Die kritischen Schranken $b_{\alpha/2}$, $o_{1-\alpha}$ und b_{α} sind bestimmt durch

$$\begin{aligned} (a) \quad & B(b_{\alpha/2}) \leq \alpha/2 < B(b_{\alpha/2} + 1) \\ (b) \quad & B(o_{1-\alpha}) < 1 - \alpha \leq B(o_{1-\alpha} + 1) \\ (c) \quad & B(b_{\alpha}) \leq \alpha < B(b_{\alpha} + 1) \end{aligned}$$

wobei B die Verteilungsfunktion der $B(n, 0.5)$ -Verteilung bezeichnet.

- Bemerkung: Für $n \geq 25$ ist die Teststatistik unter H_0 approximativ $N(0.5n, 0.25n)$ -verteilt.

Wilcoxon-Vorzeichen-Rang-Test:

- Annahmen: X_1, \dots, X_n unabhängig und identisch verteilt wie X
 X metrisch skaliert mit stetiger und symmetrischer Verteilungsfunktion.

- Hypothesen:

$$\begin{aligned} (a) \quad & H_0 : x_{med} = \delta_0 \quad \text{vs.} \quad H_1 : x_{med} \neq \delta_0 \\ (b) \quad & H_0 : x_{med} \geq \delta_0 \quad \text{vs.} \quad H_1 : x_{med} < \delta_0 \\ (c) \quad & H_0 : x_{med} \leq \delta_0 \quad \text{vs.} \quad H_1 : x_{med} > \delta_0 \end{aligned}$$

- Teststatistik:

$$W^+ = \sum_{i=1}^n \text{rg}(|D_i|) Z_i \quad \text{mit} \quad D_i = X_i - \delta_0, \quad Z_i = \begin{cases} 1, & D_i > 0 \\ 0, & D_i < 0 \end{cases}$$

- Ablehnungsbereiche:

$$\begin{aligned} (a) \quad & W^+ < w_{\alpha/2, n}^+ \quad \text{oder} \quad W^+ > w_{1-\alpha/2, n}^+ \\ (b) \quad & W^+ < w_{\alpha, n}^+ \\ (c) \quad & W^+ > w_{1-\alpha, n}^+ \end{aligned}$$

wobei $w_{\alpha, n}^+$ das tabellierte α -Quantil der Verteilung von W^+ ist.

- Bemerkung: Für $n > 20$ ist die Teststatistik unter H_0 approximativ $N(\frac{n(n+1)}{4}, \frac{n(n+1)(2n+1)}{24})$ -verteilt.

Wilcoxon-Rangsummen-Test:

- Annahmen: X_1, \dots, X_n unabhängig und identisch verteilt wie X
 Y_1, \dots, Y_m unabhängig und identisch verteilt wie Y
 X_1, \dots, X_n und Y_1, \dots, Y_m unabhängig
 X und Y besitzen stetige Verteilungsfunktionen F bzw. G

- Hypothesen:
 - (a) $H_0 : x_{med} = y_{med}$ vs. $H_1 : x_{med} \neq y_{med}$
 - (b) $H_0 : x_{med} \geq y_{med}$ vs. $H_1 : x_{med} < y_{med}$
 - (c) $H_0 : x_{med} \leq y_{med}$ vs. $H_1 : x_{med} > y_{med}$

- Teststatistik:

$$T_W = \sum_{i=1}^n \text{rg}(X_i) = \sum_{i=1}^{n+m} iV_i$$

mit
$$V_i = \begin{cases} 1, & i\text{-te Beobachtung der gepoolten Stichprobe ist } X\text{-Variable} \\ 0, & \text{sonst} \end{cases}$$

- Ablehnungsbereiche:
 - (a) $T_W < w_{\alpha/2;n,m}$ oder $T_W > w_{1-\alpha/2;n,m}$
 - (b) $T_W < w_{\alpha;n,m}$
 - (c) $T_W > w_{1-\alpha;n,m}$

wobei $w_{\tilde{\alpha}}$ das tabellierte $\tilde{\alpha}$ -Quantil der Verteilung von T_W ist.

- Bemerkung:
 Für m oder $n > 25$ ist die Teststatistik unter H_0 approximativ $N(\frac{n(n+m+1)}{2}, \frac{nm(n+m+1)}{12})$ -verteilt.

χ^2 -Anpassungstest:

- Annahmen: X_1, \dots, X_n unabhängig und identisch verteilt wie X
 Einteilung der Daten in k disjunkte Klassen

- Hypothesen:

$$H_0 : P(X = i) = \pi_i \quad \text{für } i = 1, \dots, k$$

$$H_1 : P(X = i) \neq \pi_i \quad \text{für mindestens ein } i$$

- Teststatistik:

$$\chi^2 = \sum_{i=1}^k \frac{(h_i - n\pi_i)^2}{n\pi_i}$$

- Ablehnungsbereich:

$$\chi^2 > \chi_{1-\alpha, k-1}^2$$
 -Anzahl der unter H_0 zu schätzenden Parameter

- Faustregel:
 Approximative Verteilung der Teststatistik gilt, wenn $n\pi_i \geq 1$ für alle Klassen und $n\pi_i \geq 5$ für mindestens 80% der Klassen erfüllt ist.

χ^2 -Homogenitätstest:

- Annahmen: Unabhängige Stichproben aus k Populationen mit den Stichprobenumfängen n_1, \dots, n_k

- Hypothesen:

$$H_0 : P(X_1 = j) = \dots = P(X_k = j) \quad \text{für } j = 1, \dots, m$$

$$H_1 : P(X_{i_1} = j) \neq P(X_{i_2} = j) \quad \text{für mindestens ein Tupel } (i_1, i_2, j)$$

	X					X			
	1	...	m			1	...	m	
1	h_{11}	...	h_{1m}	n_1	$\xrightarrow{\text{unter } H_0}$	$\frac{n_1 h_{\bullet 1}}{n}$...	$\frac{n_1 h_{\bullet m}}{n}$	n_1
:	:		:	:		:		:	:
k	h_{k1}	...	h_{km}	n_k		$\frac{n_k h_{\bullet 1}}{n}$...	$\frac{n_k h_{\bullet m}}{n}$	n_k
	$h_{\bullet 1}$...	$h_{\bullet m}$	n		$h_{\bullet 1}$...	$h_{\bullet m}$	n

- Teststatistik:

$$\chi^2 = \sum_{i=1}^k \sum_{j=1}^m \frac{(h_{ij} - \hat{h}_{ij})^2}{\hat{h}_{ij}} \quad \text{mit } \hat{h}_{ij} = \frac{n_i h_{\bullet j}}{n}$$

- Ablehnungsbereich:

$$\chi^2 > \chi_{1-\alpha, (k-1) \cdot (m-1)}^2$$

χ^2 -Unabhängigkeitstest:

- Annahmen: Unabhängige Stichprobenvariablen (X_i, Y_i) , $i = 1, \dots, n$
- Hypothesen:

$$H_0 : P(X = i, Y = j) = P(X = i) \cdot P(Y = j) \quad \text{für alle } i, j$$

$$H_1 : P(X = i, Y = j) \neq P(X = i) \cdot P(Y = j) \quad \text{für mindestens ein Paar } (i, j)$$

		Y			
		1	...	m	
1	h_{11}	...	h_{1m}	$h_{1\bullet}$	unter H_0 →
X	⋮	⋮	⋮	⋮	
k	h_{k1}	...	h_{km}	$h_{k\bullet}$	
	$h_{\bullet 1}$...	$h_{\bullet m}$	n	

		Y			
		1	...	m	
1	$\frac{h_{1\bullet}h_{\bullet 1}}{n}$...	$\frac{h_{1\bullet}h_{\bullet m}}{n}$	$h_{1\bullet}$	
X	⋮	⋮	⋮	⋮	
k	$\frac{h_{k\bullet}h_{\bullet 1}}{n}$...	$\frac{h_{k\bullet}h_{\bullet m}}{n}$	$h_{k\bullet}$	
	$h_{\bullet 1}$...	$h_{\bullet m}$	n	

- Teststatistik:

$$\chi^2 = \sum_{i=1}^k \sum_{j=1}^m \frac{(h_{ij} - \tilde{h}_{ij})^2}{\tilde{h}_{ij}} \quad \text{mit} \quad \tilde{h}_{ij} = \frac{h_{i\bullet}h_{\bullet j}}{n}$$

- Ablehnungsbereich:

$$\chi^2 > \chi^2_{1-\alpha, (k-1) \cdot (m-1)}$$

Korrelationstest:

- Annahmen: Unabhängige gemeinsam normalverteilte Stichprobenvariablen (X_i, Y_i) , $i = 1, \dots, n$
- Hypothesen:

- (a) $H_0 : \rho_{XY} = \rho_0$ gegen $H_1 : \rho_{XY} \neq \rho_0$
- (b) $H_0 : \rho_{XY} \geq \rho_0$ gegen $H_1 : \rho_{XY} < \rho_0$
- (c) $H_0 : \rho_{XY} \leq \rho_0$ gegen $H_1 : \rho_{XY} > \rho_0$

- Teststatistik für $\rho_0 = 0$

$$T = \frac{r_{XY}}{\sqrt{1-r_{XY}^2}} \sqrt{n-2} \stackrel{H_0}{\sim} t_{n-2},$$

bzw. für $\rho_0 = \rho_{XY}$

$$Z = \frac{1}{2} \left(\ln \frac{1+r_{XY}}{1-r_{XY}} - \ln \frac{1+\rho_0}{1-\rho_0} \right) \sqrt{n-3} \stackrel{H_0, n > 25}{\sim} N(0, 1).$$

- Ablehnungsbereiche:

- (a) $|T| > t_{1-\alpha/2, n-2}$ bzw. $|Z| > z_{1-\alpha/2}$
- (b) $T < -t_{1-\alpha, n-2}$ bzw. $Z < -z_{1-\alpha}$
- (c) $T > -t_{1-\alpha, n-2}$ bzw. $Z > z_{1-\alpha}$.

9 Regressionsanalyse

9.1 Lineare Einfachregression

Lineare Einfachregression:

Y abhängige (zu erklärende) Variable, Zielgröße, Regressand
 X unabhängige (erklärende) Variable, Einflussgröße, Regressor

Regressionsansatz:

$$Y = f(X) + \epsilon = \alpha + \beta X + \epsilon$$

Bezeichnungen:

- geschätzte Regressionsgerade: $\hat{Y} = \hat{\alpha} + \hat{\beta}x$
- Regressionskoeffizienten: $\hat{\alpha}$ und $\hat{\beta}$
- Residuen: $\hat{\epsilon}_i = Y_i - \hat{Y}_i = Y_i - (\hat{\alpha} + \hat{\beta}x_i)$

Normalverteilungsannahme:

$$\epsilon_i \sim N(0, \sigma^2) \Leftrightarrow Y_i \sim N(\alpha + \beta x_i, \sigma^2), \quad i = 1, \dots, n.$$

Kleinste-Quadrate-Schätzer:

$$\hat{\alpha} = \bar{Y} - \hat{\beta}\bar{x}, \quad \hat{\beta} = \frac{\sum_{i=1}^n (x_i - \bar{x})(Y_i - \bar{Y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{\sum_{i=1}^n x_i Y_i - n\bar{x}\bar{Y}}{\sum_{i=1}^n x_i^2 - n\bar{x}^2}$$

Schätzer für die Varianz σ^2 :

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n \hat{\epsilon}_i^2 = \frac{1}{n-2} \sum_{i=1}^n (Y_i - (\hat{\alpha} + \hat{\beta}x_i))^2$$

Verteilung der geschätzten Regressionskoeffizienten:

$$\hat{\alpha} \sim N(\alpha, \sigma_{\hat{\alpha}}^2) \quad \text{mit} \quad \text{Var}(\hat{\alpha}) = \sigma_{\hat{\alpha}}^2 = \sigma^2 \frac{\sum x_i^2}{n \sum (x_i - \bar{x})^2} = \sigma^2 \frac{\sum x_i^2}{n(\sum x_i^2 - n\bar{x}^2)}$$

$$\hat{\beta} \sim N(\beta, \sigma_{\hat{\beta}}^2) \quad \text{mit} \quad \text{Var}(\hat{\beta}) = \sigma_{\hat{\beta}}^2 = \frac{\sigma^2}{\sum (x_i - \bar{x})^2} = \frac{\sigma^2}{\sum x_i^2 - n\bar{x}^2}$$

Verteilung der standardisierten Schätzfunktionen:

$$\frac{\hat{\alpha} - \alpha}{\hat{\sigma}_{\hat{\alpha}}} \sim t_{n-2} \quad \text{mit} \quad \hat{\sigma}_{\hat{\alpha}} = \hat{\sigma} \frac{\sqrt{\sum x_i^2}}{\sqrt{n \sum (x_i - \bar{x})^2}} = \hat{\sigma} \frac{\sqrt{\sum x_i^2}}{\sqrt{n(\sum x_i^2 - n\bar{x}^2)}}$$

$$\frac{\hat{\beta} - \beta}{\hat{\sigma}_{\hat{\beta}}} \sim t_{n-2} \quad \text{mit} \quad \hat{\sigma}_{\hat{\beta}} = \frac{\hat{\sigma}}{\sqrt{\sum (x_i - \bar{x})^2}} = \frac{\hat{\sigma}}{\sqrt{\sum x_i^2 - n\bar{x}^2}}$$

$(1 - \alpha)$ -Konfidenzintervalle für α und β :

$$\text{für } \alpha: \quad \left[\hat{\alpha} - \hat{\sigma}_{\hat{\alpha}} t_{1-\frac{\alpha}{2}, n-2}, \hat{\alpha} + \hat{\sigma}_{\hat{\alpha}} t_{1-\frac{\alpha}{2}, n-2} \right]$$

$$\text{für } \beta: \quad \left[\hat{\beta} - \hat{\sigma}_{\hat{\beta}} t_{1-\frac{\alpha}{2}, n-2}, \hat{\beta} + \hat{\sigma}_{\hat{\beta}} t_{1-\frac{\alpha}{2}, n-2} \right]$$

Teststatistiken:

$$T_{\alpha_0} = \frac{\hat{\alpha} - \alpha_0}{\hat{\sigma}_{\hat{\alpha}}} \quad \text{und} \quad T_{\beta_0} = \frac{\hat{\beta} - \beta_0}{\hat{\sigma}_{\hat{\beta}}}$$

Hypothesen und Ablehnbereiche:

Hypothesen		Ablehnbereich
$H_0: \alpha = \alpha_0$	vs. $H_1: \alpha \neq \alpha_0$	$ T_{\alpha_0} > t_{1-\frac{\alpha}{2}, n-2}$
$H_0: \beta = \beta_0$	vs. $H_1: \beta \neq \beta_0$	$ T_{\beta_0} > t_{1-\frac{\alpha}{2}, n-2}$
$H_0: \alpha \geq \alpha_0$	vs. $H_1: \alpha < \alpha_0$	$T_{\alpha_0} < -t_{1-\alpha, n-2}$
$H_0: \beta \geq \beta_0$	vs. $H_1: \beta < \beta_0$	$T_{\beta_0} < -t_{1-\alpha, n-2}$
$H_0: \alpha \leq \alpha_0$	vs. $H_1: \alpha > \alpha_0$	$T_{\alpha_0} > t_{1-\alpha, n-2}$
$H_0: \beta \leq \beta_0$	vs. $H_1: \beta > \beta_0$	$T_{\beta_0} > t_{1-\alpha, n-2}$

Prognose:

$$\hat{Y}_0 = \hat{\alpha} + \hat{\beta}x_0$$

Konfidenzintervall für Y_0 :

$$\left[\hat{Y}_0 - t_{1-\frac{\alpha}{2}, n-2} \hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum x_i^2 - n\bar{x}^2}}, \hat{Y}_0 + t_{1-\frac{\alpha}{2}, n-2} \hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum x_i^2 - n\bar{x}^2}} \right]$$

Quadratsummenzerlegung:

$$\underbrace{\sum_{i=1}^n (Y_i - \bar{Y})^2}_{SQT} = \underbrace{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}_{SQE} + \underbrace{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}_{SQR}$$

SQT: Gesamtabweichungsquadratsumme in Y -Richtung

SQE: Durch die Regression erklärter Teil von *SQT*

SQR: Trotz der Regression unerklärt bleibender Teil von *SQT*

Bestimmtheitsmaß:

$$R^2 = \frac{SQE}{SQT} = 1 - \frac{SQR}{SQT}, \quad \text{Berechnung: } R^2 = \frac{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}{\sum_{i=1}^n (Y_i - \bar{Y})^2} = \frac{\sum_{i=1}^n \hat{Y}_i^2 - n\bar{Y}^2}{\sum_{i=1}^n Y_i^2 - n\bar{Y}^2}$$

9.2 Multiple lineare Regression in Summennotation**Regressionsansatz:**

$$Y_i = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip} + \epsilon_i, \quad i = 1, \dots, n.$$

Normalverteilungsannahme:

$$\epsilon_i \sim N(0, \sigma^2), \quad \Leftrightarrow \quad Y_i \sim N(\beta_0 + \beta_1 x_{i1} + \dots + \beta_p x_{ip}, \sigma^2), \quad i = 1, \dots, n.$$

Gefittete Werte:

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_{i1} + \hat{\beta}_2 x_{i2} + \dots + \hat{\beta}_p x_{ip}$$

Residuen:

$$\hat{\epsilon}_i = Y_i - \hat{Y}_i, \quad i = 1, \dots, n.$$

Schätzer für die Varianz σ^2 :

$$\hat{\sigma}^2 = \frac{1}{n-p-1} \sum_{i=1}^n \hat{\epsilon}_i^2 = \frac{1}{n-p-1} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$$

Verteilung der standardisierten Schätzfunktionen:

$$\frac{\hat{\beta}_j - \beta_j}{\hat{\sigma}_j} \sim t_{n-p-1}, \quad j = 0, \dots, p$$

(1 - α)-Konfidenzintervalle für β_j :

$$\left[\hat{\beta}_j - \hat{\sigma}_j t_{1-\frac{\alpha}{2}, n-p-1}, \hat{\beta}_j + \hat{\sigma}_j t_{1-\frac{\alpha}{2}, n-p-1} \right]$$

Teststatistiken:

$$T_j = \frac{\hat{\beta}_j - \beta_{0j}}{\hat{\sigma}_j}, \quad j = 0, \dots, p$$

Hypothesen und Ablehnbereiche:

Hypothesen		Ablehnbereich
$H_0: \beta_j = \beta_{0j}$	vs. $H_1: \beta_j \neq \beta_{0j}$	$ T_j > t_{1-\frac{\alpha}{2}, n-p-1}$
$H_0: \beta_j \geq \beta_{0j}$	vs. $H_1: \beta_j < \beta_{0j}$	$T_j < -t_{1-\alpha, n-p-1}$
$H_0: \beta_j \leq \beta_{0j}$	vs. $H_1: \beta_j > \beta_{0j}$	$T_j > t_{1-\alpha, n-p-1}$

Overall-F-Test:

• Hypothesen:

$$H_0: \beta_1 = \dots = \beta_p = 0$$

$$H_1: \beta_j \neq 0 \quad \text{für mindestens ein } j$$

• Teststatistik:

$$F = \frac{R^2}{1-R^2} \frac{n-p-1}{p} = \frac{SQE}{SQR} \frac{n-p-1}{p}$$

• Ablehnungsbereich:

$$F > F_{1-\alpha, p, n-p-1}$$

9.3 Multiple lineare Regression in Matrixnotation

• Modell in Matrixnotation:

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \mathbf{E}(\boldsymbol{\epsilon}) = \mathbf{0}, \quad \text{Var}(\boldsymbol{\epsilon}) = \mathbf{E}(\boldsymbol{\epsilon}\boldsymbol{\epsilon}') = \sigma^2 \mathbf{I}_n$$

mit

$$\mathbf{y} = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1p} \\ 1 & x_{21} & \cdots & x_{2p} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \cdots & x_{np} \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}, \quad \boldsymbol{\epsilon} = \begin{pmatrix} \epsilon_1 \\ \epsilon_2 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

und \mathbf{I}_n der n -dimensionalen Einheitsmatrix.

KQ-Schätzer für β :

- Aus dem KQ-Ansatz

$$(\mathbf{y} - \mathbf{X}\beta)'(\mathbf{y} - \mathbf{X}\beta) \rightarrow \min_{\beta}$$

ergibt sich durch Nullsetzen der ersten Ableitung nach β und, falls $\mathbf{X}'\mathbf{X}$ invertierbar ist, anschließendem Lösen des resultierenden Gleichungssystems der KQ-Schätzer

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

- Varianz der Schätzer $\hat{\beta}_j$:

Mit den Diagonalelementen v_j aus $(\mathbf{X}'\mathbf{X})^{-1}$ erhält man für bekanntes σ^2 als Varianz von $\hat{\beta}_j$

$$\sigma_j^2 = \text{Var}(\hat{\beta}_j) = \sigma^2 v_j$$

bzw. für unbekanntes σ^2 die geschätzte Varianz von $\hat{\beta}_j$ gemäß

$$\hat{\sigma}_j^2 = \hat{\sigma}^2 v_j.$$

Zusammenfassende Darstellung in Vektornotation

$$\text{Var}(\hat{\beta}) = \sigma^2(\mathbf{X}'\mathbf{X})^{-1} \quad \text{bzw.} \quad \widehat{\text{Var}}(\hat{\beta}) = \hat{\sigma}^2(\mathbf{X}'\mathbf{X})^{-1}$$

KQ-Schätzer für σ^2 :

$$\hat{\sigma}^2 = \frac{\hat{\epsilon}'\hat{\epsilon}}{n-p-1} = \frac{\mathbf{y}'\mathbf{y} - \hat{\beta}'\mathbf{X}'\mathbf{y}}{n-p-1}$$

Prognose:

$$\hat{y}_0 = \mathbf{x}'_0 \hat{\beta}$$

Konfidenzintervall für \hat{y}_0 :

$$\left[\hat{y}_0 - t_{1-\frac{\alpha}{2}, n-p-1} \hat{\sigma} \sqrt{\mathbf{x}'_0 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_0 + 1}, \hat{y}_0 + t_{1-\frac{\alpha}{2}, n-p-1} \hat{\sigma} \sqrt{\mathbf{x}'_0 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_0 + 1} \right]$$

Bestimmtheitsmaß und korrigiertes Bestimmtheitsmaß:

Bestimmtheitsmaß:

$$R^2 = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{\hat{\beta}'\mathbf{X}'\mathbf{y} - n\bar{y}^2}{\mathbf{y}'\mathbf{y} - n\bar{y}^2}$$

korrigiertes Bestimmtheitsmaß:

$$\bar{R}^2 = 1 - \frac{n-1}{n-p-1} (1 - R^2)$$

Hypothesentests:

Sei $R_1 := (r_0, \dots, r_p)$. Man betrachte folgende Testprobleme

- (a) $H_0 : R_1\beta = r$ gegen $H_1 : R_1\beta \neq r$
- (b) $H_0 : R_1\beta \geq r$ gegen $H_1 : R_1\beta < r$
- (c) $H_0 : R_1\beta \leq r$ gegen $H_1 : R_1\beta > r$.

Teststatistik:

$$T = \frac{R_1\hat{\beta} - r}{\hat{\sigma} \sqrt{R_1(\mathbf{X}'\mathbf{X})^{-1}R_1'}} \stackrel{H_0}{\sim} t_{n-p-1}$$

Ablehnungsbereiche:

- (a), $|T| > t_{1-\frac{\alpha}{2}, n-p-1}$
- (b), $T < -t_{1-\alpha, n-p-1}$
- (c), $T > t_{1-\alpha, n-p-1}$

Kodierung kategorialer Einflussgrößen:

Sei $M \in \{1, \dots, m\}$ eine mehrkategoriale erklärende Variable mit m Kategorien.

Dummy-Kodierung:

$$x_i^M := \begin{cases} 1, & M = i, \\ 0, & \text{sonst.} \end{cases} \quad i = 1, \dots, m-1$$

Effekt-Kodierung:

$$x_i^M := \begin{cases} 1, & M = i, \\ -1, & M = m, \\ 0, & \text{sonst.} \end{cases} \quad i = 1, \dots, m-1$$

mit m als Referenzkategorie.

SPSS-Output einer multiplen Regression:

		Coefficients ^a				
Model		Unstandardized Coefficients		t	Sig.	
		B	Std. Error			
1	(Constant)	$\hat{\beta}_0$	$\hat{\sigma}_0$	T_0	$P(T \geq T_0)$	
	X_1	$\hat{\beta}_1$	$\hat{\sigma}_1$	T_1	$P(T \geq T_1)$	
	X_2	$\hat{\beta}_2$	$\hat{\sigma}_2$	T_2	$P(T \geq T_2)$	
	\vdots	\vdots	\vdots	\vdots	\vdots	
	X_p	$\hat{\beta}_p$	$\hat{\sigma}_p$	T_p	$P(T \geq T_p)$	

a Dependent Variable: Y

Testprobleme:

- t ist der Wert der t-Statistik für

$$H_0 : \beta_j = 0 \quad \text{vs.} \quad H_1 : \beta_j \neq 0,$$

$$\text{d.h. } T_j = \frac{\hat{\beta}_j}{\hat{\sigma}_j}, \quad j = 0, \dots, p$$

- Sig. ist der zugehörige p-Wert

10 Varianzanalyse

10.1 Einfaktorielle Varianzanalyse

Modell 1:

$$Y_{ij} = \mu_i + \epsilon_{ij}, \quad \epsilon_{ij} \sim N(0, \sigma^2), \text{ unabhängig, } i = 1, \dots, I, j = 1, \dots, n_i.$$

Modell 2:

$$Y_{ij} = \mu + \alpha_i + \epsilon_{ij}, \quad \sum_{i=1}^I n_i \alpha_i = 0, \quad \epsilon_{ij} \sim N(0, \sigma^2), \text{ unabhängig, } i = 1, \dots, I, j = 1, \dots, n_i.$$

Schätzer für Modell 2:

$$\hat{\mu} = \frac{1}{n} \sum_{i=1}^I \sum_{j=1}^{n_i} Y_{ij} = \bar{Y}_{\bullet\bullet} \quad \text{und} \quad \hat{\alpha}_i = \bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet}, \quad \text{mit} \quad \bar{Y}_{i\bullet} = \frac{1}{n_i} \sum_{j=1}^{n_i} Y_{ij}$$

Die Prüfgröße für das Testproblem

$$H_0 : \alpha_1 = \dots = \alpha_I = 0 \quad \text{gegen} \quad H_1 : \text{mind. zwei } \alpha_i \neq 0$$

ist gegeben als

$$F = \frac{MQE}{MQR} = \frac{\sum_{i=1}^I n_i (\bar{Y}_{i\bullet} - \bar{Y}_{\bullet\bullet})^2 / (I - 1)}{\sum_{i=1}^I \sum_{j=1}^{n_i} (Y_{ij} - \bar{Y}_{i\bullet})^2 / (n - I)},$$

Ablehnungsbereich:

$$C = \{F : F > F_{1-\alpha; I-1, n-I}\}.$$

10.2 Zweifaktorielle Varianzanalyse

Modelle:

Modell 1:

$$Y_{ijk} = \mu_{ij} + \epsilon_{ijk}, \quad \epsilon_{ijk} \sim N(0, \sigma^2), \text{ unabhängig, } i = 1, \dots, I, j = 1, \dots, J, k = 1, \dots, K.$$

Modell 2:

$$Y_{ijk} = \mu + \alpha_i + \beta_j + (\alpha\beta)_{ij} + \epsilon_{ijk}, \quad \sum_{i=1}^I \alpha_i = 0, \quad \sum_{j=1}^J \beta_j = 0, \quad \sum_{i=1}^I (\alpha\beta)_{ij} = 0, \quad \sum_{j=1}^J (\alpha\beta)_{ij} = 0,$$

$$\epsilon_{ijk} \sim N(0, \sigma^2), \text{ unabhängig, } i = 1, \dots, I, j = 1, \dots, J, k = 1, \dots, K.$$

Schätzer:

Die Schätzer der Parameter in Modell 2 sind gegeben als

$$\begin{aligned}\hat{\mu} &= \frac{1}{IJK} \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K Y_{ijk} = \bar{Y}_{\dots}, \\ \hat{\alpha}_i &= \bar{Y}_{i\bullet\bullet} - \bar{Y}_{\dots} \quad \text{mit} \quad \bar{Y}_{i\bullet\bullet} = \frac{1}{JK} \sum_{j=1}^J \sum_{k=1}^K Y_{ijk} \\ \hat{\beta}_j &= \bar{Y}_{\bullet j \bullet} - \bar{Y}_{\dots} \quad \text{mit} \quad \bar{Y}_{\bullet j \bullet} = \frac{1}{IK} \sum_{i=1}^I \sum_{k=1}^K Y_{ijk} \\ (\widehat{\alpha\beta})_{ij} &= \bar{Y}_{ij\bullet} - \bar{Y}_{i\bullet\bullet} - \bar{Y}_{\bullet j \bullet} + \bar{Y}_{\dots} \quad \text{mit} \quad \bar{Y}_{ij\bullet} = \frac{1}{K} \sum_{k=1}^K Y_{ijk}.\end{aligned}$$

Prüfgrößen:

Vorliegen von Wechselwirkungen:

$$F_{A \times B} = \frac{MQ(A \times B)}{MQR} = \frac{K \sum_{i=1}^I \sum_{j=1}^J (\bar{Y}_{ij\bullet} - \bar{Y}_{i\bullet\bullet} - \bar{Y}_{\bullet j \bullet} + \bar{Y}_{\dots})^2 / ((I-1)(J-1))}{\sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K (Y_{ijk} - \bar{Y}_{ij\bullet})^2 / (IJ(K-1))}$$

mit Ablehnungsbereich

$$F_{A \times B} > F_{1-\alpha; (I-1)(J-1), IJ(K-1)}.$$

Vorliegen von Haupteffekten bedingt durch Faktor A

$$F_A = \frac{MQA}{MQR} = \frac{KJ \sum_{i=1}^I (\bar{Y}_{i\bullet\bullet} - \bar{Y}_{\dots})^2 / (I-1)}{\sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K (Y_{ijk} - \bar{Y}_{ij\bullet})^2 / (IJ(K-1))}$$

mit Ablehnungsbereich

$$F_A > F_{1-\alpha; I-1, IJ(K-1)}.$$

Vorliegen von Haupteffekten bedingt durch Faktor B

$$F_B = \frac{MQB}{MQR} = \frac{KI \sum_{j=1}^J (\bar{Y}_{\bullet j \bullet} - \bar{Y}_{\dots})^2 / (J-1)}{\sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K (Y_{ijk} - \bar{Y}_{ij\bullet})^2 / (IJ(K-1))}$$

mit Ablehnungsbereich

$$F_B > F_{1-\alpha; J-1, IJ(K-1)}.$$

11 Zeitreihenanalyse**Einfacher gleitender Durchschnitt:**

Gleitender Durchschnitt ungerader Ordnung $p = 2q + 1$:

$$\hat{g}_t = \frac{1}{2q+1} \sum_{i=t-q}^{t+q} y_i, \quad t = q+1, \dots, n-q$$

Gleitender Durchschnitt gerader Ordnung $p = 2q$:

$$\hat{g}_t = \frac{1}{2q} \left[\frac{1}{2} y_{t-q} + \frac{1}{2} y_{t+q} + \sum_{i=t-(q-1)}^{t+(q-1)} y_i \right], \quad t = q+1, \dots, n-q$$

Globale Trendmodelle:

$g_t = \beta_0 + \beta_1 t$	linearer Trend
$g_t = \beta_0 + \beta_1 t + \beta_2 t^2$	quadratischer Trend
$g_t = \beta_0 + \beta_1 t + \dots + \beta_q t^q$	polynomialer Trend
$g_t = \beta_0 \beta_1^t$	Exponentialtrend
$g_t = \beta_0 \exp\{\beta_1 t\}$	exponentielles Wachstum
$g_t = \frac{\beta_0}{\beta_1 + \exp\{-\beta_2 t\}}$	logistische Sättigungskurve

Schätzung von Trendfunktionen:

lineare Trendfunktion $\hat{g}_t = \hat{\beta}_0 + \hat{\beta}_1 t$ mit

$$\hat{\beta}_0 = \frac{\sum_{t=1}^n g_t \sum_{t=1}^n t^2 - \sum_{t=1}^n t \sum_{t=1}^n g_t \cdot t}{n \sum_{t=1}^n t^2 - (\sum_{t=1}^n t)^2} \quad \hat{\beta}_1 = \frac{n \sum_{t=1}^n g_t \cdot t - \sum_{t=1}^n g_t \sum_{t=1}^n t}{n \sum_{t=1}^n t^2 - (\sum_{t=1}^n t)^2}$$

Exponentialtrend $\hat{g}_t = \hat{\beta}_0 \hat{\beta}_1^t \iff \ln \hat{g}_t = \ln \hat{\beta}_0 + t \cdot \ln \hat{\beta}_1$ mit

$$\ln \hat{\beta}_0 = \frac{\sum_{t=1}^n \ln g_t \sum_{t=1}^n t^2 - \sum_{t=1}^n t \sum_{t=1}^n \ln g_t \cdot t}{n \sum_{t=1}^n t^2 - (\sum_{t=1}^n t)^2} \quad \ln \hat{\beta}_1 = \frac{n \sum_{t=1}^n \ln g_t \cdot t - \sum_{t=1}^n \ln g_t \sum_{t=1}^n t}{n \sum_{t=1}^n t^2 - (\sum_{t=1}^n t)^2}$$

12.2 Student's t -Verteilung

Tabelliert sind die Quantile für n Freiheitsgrade.

Für das Quantil $t_{1-\alpha,n}$ gilt $F(t_{1-\alpha,n}) = 1 - \alpha$.

Links vom Quantil $t_{1-\alpha,n}$ liegt die Wahrscheinlichkeitsmasse $1 - \alpha$.

Ablesebeispiel: $t_{0,99,20} = 2.528$

Die Quantile für $0 < 1 - \alpha < 0.5$ erhält man aus $t_{\alpha,n} = -t_{1-\alpha,n}$

Approximation für $n > 30$:

$$t_{\alpha,n} \approx z_{\alpha} \quad (z_{\alpha} \text{ ist das } (\alpha)\text{-Quantil der Standardnormalverteilung})$$

n	0.6	0.8	0.9	0.95	0.975	0.99	0.995	0.999	0.9995
1	0.3249	1.3764	3.0777	6.3138	12.706	31.821	63.657	318.31	636.62
2	0.2887	1.0607	1.8856	2.9200	4.3027	6.9646	9.9248	22.327	31.599
3	0.2767	0.9785	1.6377	2.3534	3.1824	4.5407	5.8409	10.215	12.924
4	0.2707	0.9410	1.5332	2.1318	2.7764	3.7469	4.6041	7.1732	8.6103
5	0.2672	0.9195	1.4759	2.0150	2.5706	3.3649	4.0321	5.8934	6.8688
6	0.2648	0.9057	1.4398	1.9432	2.4469	3.1427	3.7074	5.2076	5.9588
7	0.2632	0.8960	1.4149	1.8946	2.3646	2.9980	3.4995	4.7853	5.4079
8	0.2619	0.8889	1.3968	1.8595	2.3060	2.8965	3.3554	4.5008	5.0413
9	0.2610	0.8834	1.3830	1.8331	2.2622	2.8214	3.2498	4.2968	4.7809
10	0.2602	0.8791	1.3722	1.8125	2.2281	2.7638	3.1693	4.1437	4.5869
11	0.2596	0.8755	1.3634	1.7959	2.2010	2.7181	3.1058	4.0247	4.4370
12	0.2590	0.8726	1.3562	1.7823	2.1788	2.6810	3.0545	3.9296	4.3178
13	0.2586	0.8702	1.3502	1.7709	2.1604	2.6503	3.0123	3.8520	4.2208
14	0.2582	0.8681	1.3450	1.7613	2.1448	2.6245	2.9768	3.7874	4.1405
15	0.2579	0.8662	1.3406	1.7531	2.1314	2.6025	2.9467	3.7328	4.0728
16	0.2576	0.8647	1.3368	1.7459	2.1199	2.5835	2.9208	3.6862	4.0150
17	0.2573	0.8633	1.3334	1.7396	2.1098	2.5669	2.8982	3.6458	3.9651
18	0.2571	0.8620	1.3304	1.7341	2.1009	2.5524	2.8784	3.6105	3.9216
19	0.2569	0.8610	1.3277	1.7291	2.0930	2.5395	2.8609	3.5794	3.8834
20	0.2567	0.8600	1.3253	1.7247	2.0860	2.5280	2.8453	3.5518	3.8495
21	0.2566	0.8591	1.3232	1.7207	2.0796	2.5176	2.8314	3.5272	3.8193
22	0.2564	0.8583	1.3212	1.7171	2.0739	2.5083	2.8188	3.5050	3.7921
23	0.2563	0.8575	1.3195	1.7139	2.0687	2.4999	2.8073	3.4850	3.7676
24	0.2562	0.8569	1.3178	1.7109	2.0639	2.4922	2.7969	3.4668	3.7454
25	0.2561	0.8562	1.3163	1.7081	2.0595	2.4851	2.7874	3.4502	3.7251
26	0.2560	0.8557	1.3150	1.7056	2.0555	2.4786	2.7787	3.4350	3.7066
27	0.2559	0.8551	1.3137	1.7033	2.0518	2.4727	2.7707	3.4210	3.6896
28	0.2558	0.8546	1.3125	1.7011	2.0484	2.4671	2.7633	3.4082	3.6739
29	0.2557	0.8542	1.3114	1.6991	2.0452	2.4620	2.7564	3.3962	3.6594
30	0.2556	0.8538	1.3104	1.6973	2.0423	2.4573	2.7500	3.3852	3.6460
∞	0.2533	0.8416	1.2816	1.6449	1.9600	2.3263	2.5758	3.0903	3.2906

12.3 χ^2 -Verteilung

Tabelliert sind die Quantile für n Freiheitsgrade.

Für das Quantil $\chi^2_{1-\alpha,n}$ gilt $F(\chi^2_{1-\alpha,n}) = 1 - \alpha$.

Links vom Quantil $\chi^2_{1-\alpha,n}$ liegt die Wahrscheinlichkeitsmasse $1 - \alpha$.

Ablesebeispiel: $\chi^2_{0,95,10} = 18.307$

Approximation für $n > 30$:

$$\chi^2_{\alpha,n} \approx \frac{1}{2}(z_{\alpha} + \sqrt{2n-1})^2 \quad (z_{\alpha} \text{ ist das } \alpha\text{-Quantil der Standardnormalverteilung})$$

n	0.01	0.025	0.05	0.1	0.5	0.9	0.95	0.975	0.99
1	0.0002	0.0010	0.0039	0.0158	0.4549	2.7055	3.8415	5.0239	6.6349
2	0.0201	0.0506	0.1026	0.2107	1.3863	4.6052	5.9915	7.3778	9.2103
3	0.1148	0.2158	0.3518	0.5844	2.3660	6.2514	7.8147	9.3484	11.345
4	0.2971	0.4844	0.7107	1.0636	3.3567	7.7794	9.4877	11.143	13.277
5	0.5543	0.8312	1.1455	1.6103	4.3515	9.2364	11.070	12.833	15.086
6	0.8721	1.2373	1.6354	2.2041	5.3481	10.645	12.592	14.449	16.812
7	1.2390	1.6899	2.1674	2.8331	6.3458	12.017	14.067	16.013	18.475
8	1.6465	2.1797	2.7326	3.4895	7.3441	13.362	15.507	17.535	20.090
9	2.0879	2.7004	3.3251	4.1682	8.3428	14.684	16.919	19.023	21.666
10	2.5582	3.2470	3.9403	4.8652	9.3418	15.987	18.307	20.483	23.209
11	3.0535	3.8157	4.5748	5.5778	10.341	17.275	19.675	21.920	24.725
12	3.5706	4.4038	5.2260	6.3038	11.340	18.549	21.026	23.337	26.217
13	4.1069	5.0088	5.8919	7.0415	12.340	19.812	22.362	24.736	27.688
14	4.6604	5.6287	6.5706	7.7895	13.339	21.064	23.685	26.119	29.141
15	5.2293	6.2621	7.2609	8.5468	14.339	22.307	24.996	27.488	30.578
16	5.8122	6.9077	7.9616	9.3122	15.338	23.542	26.296	28.845	32.000
17	6.4078	7.5642	8.6718	10.085	16.338	24.769	27.587	30.191	33.409
18	7.0149	8.2307	9.3905	10.865	17.338	25.989	28.869	31.526	34.805
19	7.6327	8.9065	10.117	11.651	18.338	27.204	30.144	32.852	36.191
20	8.2604	9.5908	10.851	12.443	19.337	28.412	31.410	34.170	37.566
21	8.8972	10.283	11.591	13.240	20.337	29.615	32.671	35.479	38.932
22	9.5425	10.982	12.338	14.041	21.337	30.813	33.924	36.781	40.289
23	10.196	11.689	13.091	14.848	22.337	32.007	35.172	38.076	41.638
24	10.856	12.401	13.848	15.659	23.337	33.196	36.415	39.364	42.980
25	11.524	13.120	14.611	16.473	24.337	34.382	37.652	40.646	44.314
26	12.198	13.844	15.379	17.292	25.336	35.563	38.885	41.923	45.642
27	12.879	14.573	16.151	18.114	26.336	36.741	40.113	43.195	46.963
28	13.565	15.308	16.928	18.939	27.336	37.916	41.337	44.461	48.278
29	14.256	16.047	17.708	19.768	28.336	39.087	42.557	45.722	49.588
30	14.953	16.791	18.493	20.599	29.336	40.256	43.773	46.979	50.892

12.4 Poissonverteilung

Tabelliert sind die Werte der Verteilungsfunktion $F(x) = P(X \leq x) = \sum_{k=0}^x P(X = k)$.

Ablesebeispiel: $X \sim Po(4)$ $F(5) = P(X \leq 5) = 0.7851$

Approximation für $\lambda \geq 10$:

$$Po(\lambda) \overset{a}{\sim} N(\lambda, \lambda)$$

x	λ									
	1	2	3	4	5	6	7	8	9	10
0	0.3679	0.1353	0.0498	0.0183	0.0067	0.0025	0.0009	0.0003	0.0001	0.0000
1	0.7358	0.4060	0.1991	0.0916	0.0404	0.0174	0.0073	0.0030	0.0012	0.0005
2	0.9197	0.6767	0.4232	0.2381	0.1247	0.0620	0.0296	0.0138	0.0062	0.0028
3	0.9810	0.8571	0.6472	0.4335	0.2650	0.1512	0.0818	0.0424	0.0212	0.0103
4	0.9963	0.9473	0.8153	0.6288	0.4405	0.2851	0.1730	0.0996	0.0550	0.0293
5	0.9994	0.9834	0.9161	0.7851	0.6160	0.4457	0.3007	0.1912	0.1157	0.0671
6	0.9999	0.9955	0.9665	0.8893	0.7622	0.6063	0.4497	0.3134	0.2068	0.1301
7	1.0000	0.9989	0.9881	0.9489	0.8666	0.7440	0.5987	0.4530	0.3239	0.2202
8		0.9998	0.9962	0.9786	0.9319	0.8472	0.7291	0.5925	0.4557	0.3328
9		1.0000	0.9989	0.9919	0.9682	0.9161	0.8305	0.7166	0.5874	0.4579
10			0.9997	0.9972	0.9863	0.9574	0.9015	0.8159	0.7060	0.5830
11			0.9999	0.9991	0.9945	0.9799	0.9467	0.8881	0.8030	0.6968
12			1.0000	0.9997	0.9980	0.9912	0.9730	0.9362	0.8758	0.7916
13				0.9999	0.9993	0.9964	0.9872	0.9658	0.9261	0.8645
14				1.0000	0.9998	0.9986	0.9943	0.9827	0.9585	0.9165
15					0.9999	0.9995	0.9976	0.9918	0.9780	0.9513
16					1.0000	0.9998	0.9990	0.9963	0.9889	0.9730
17						0.9999	0.9996	0.9984	0.9947	0.9857
18						1.0000	0.9999	0.9993	0.9976	0.9928
19							1.0000	0.9997	0.9989	0.9965
20								0.9999	0.9996	0.9984
21								1.0000	0.9998	0.9993
22									0.9999	0.9997
23									1.0000	0.9999
24										1.0000

12.5 F-Verteilung

Tabelliert sind die rechtsseitigen Quantile für (n_1, n_2) Freiheitsgrade.

Für das Quantil $f_{1-\alpha; n_1, n_2}$ gilt $F(f_{1-\alpha; n_1, n_2}) = 1 - \alpha$. Links vom Quantil $f_{1-\alpha; n_1, n_2}$ liegt die Wahrscheinlichkeitsmasse $1 - \alpha$.

Ablesebeispiel: $f_{0.99; 15, 8} = 5.5151$

Linksseitige Quantile: $f_{\alpha; n_1, n_2} = \frac{1}{f_{1-\alpha; n_1, n_2}}$

n_1	α	n_2								
		1	2	3	4	5	6	7	8	9
1	0.9	39.863	8.5263	5.5383	4.5448	4.0604	3.7759	3.5894	3.4579	3.3603
	0.95	161.45	18.513	10.128	7.7086	6.6079	5.9874	5.5914	5.3177	5.1174
	0.975	647.79	38.506	17.443	12.218	10.007	8.8131	8.0727	7.5709	7.2093
	0.99	4052.2	98.502	34.116	21.198	16.258	13.745	12.246	11.259	10.561
2	0.9	49.500	9.0000	5.4624	4.3246	3.7797	3.4633	3.2574	3.1131	3.0065
	0.95	199.50	19.000	9.5521	6.9443	5.7861	5.1433	4.7374	4.4590	4.2565
	0.975	799.50	39.000	16.044	10.649	8.4336	7.2599	6.5415	6.0595	5.7147
	0.99	4999.5	99.000	30.817	18.000	13.274	10.925	9.5466	8.6491	8.0215
3	0.9	53.593	9.1618	5.3908	4.1909	3.6195	3.2888	3.0741	2.9238	2.8129
	0.95	215.71	19.164	9.2766	6.5914	5.4095	4.7571	4.3468	4.0662	3.8625
	0.975	864.16	39.165	15.439	9.9792	7.7636	6.5988	5.8898	5.4160	5.0781
	0.99	5403.4	99.166	29.457	16.694	12.060	9.7795	8.4513	7.5910	6.9919
4	0.9	55.833	9.2434	5.3426	4.1072	3.5202	3.1808	2.9605	2.8064	2.6927
	0.95	224.58	19.247	9.1172	6.3882	5.1922	4.5337	4.1203	3.8379	3.6331
	0.975	899.58	39.248	15.101	9.6045	7.3879	6.2272	5.5226	5.0526	4.7181
	0.99	5624.6	99.249	28.710	15.977	11.392	9.1483	7.8466	7.0061	6.4221
5	0.9	57.240	9.2926	5.3092	4.0506	3.4530	3.1075	2.8833	2.7264	2.6106
	0.95	230.16	19.296	9.0135	6.2561	5.0503	4.3874	3.9715	3.6875	3.4817
	0.975	921.85	39.298	14.885	9.3645	7.1464	5.9876	5.2852	4.8173	4.4844
	0.99	5763.6	99.299	28.237	15.522	10.967	8.7459	7.4604	6.6318	6.0569
6	0.9	58.204	9.3255	5.2847	4.0097	3.4045	3.0546	2.8274	2.6683	2.5509
	0.95	233.99	19.330	8.9406	6.1631	4.9503	4.2839	3.8660	3.5806	3.3738
	0.975	937.11	39.331	14.735	9.1973	6.9777	5.8198	5.1186	4.6517	4.3197
	0.99	5859.0	99.333	27.911	15.207	10.672	8.4661	7.1914	6.3707	5.8018
7	0.9	58.906	9.3491	5.2662	3.9790	3.3679	3.0145	2.7849	2.6241	2.5053
	0.95	236.77	19.353	8.8867	6.0942	4.8759	4.2067	3.7870	3.5005	3.2927
	0.975	948.22	39.355	14.624	9.0741	6.8531	5.6955	4.9949	4.5286	4.1970
	0.99	5928.4	99.356	27.672	14.976	10.456	8.2600	6.9928	6.1776	5.6129
8	0.9	59.439	9.3668	5.2517	3.9549	3.3393	2.9830	2.7516	2.5893	2.4694
	0.95	238.88	19.371	8.8452	6.0410	4.8183	4.1468	3.7257	3.4381	3.2296
	0.975	956.66	39.373	14.540	8.9796	6.7572	5.5996	4.8993	4.4333	4.1020
	0.99	5981.1	99.374	27.489	14.799	10.289	8.1017	6.8400	6.0289	5.4671
9	0.9	59.858	9.3805	5.2400	3.9357	3.3163	2.9577	2.7247	2.5612	2.4403
	0.95	240.54	19.385	8.8123	5.9988	4.7725	4.0990	3.6767	3.3881	3.1789
	0.975	963.28	39.387	14.473	8.9047	6.6811	5.5234	4.8232	4.3572	4.0260
	0.99	6022.5	99.388	27.345	14.659	10.158	7.9761	6.7188	5.9106	5.3511
10	0.9	60.195	9.3916	5.2304	3.9199	3.2974	2.9369	2.7025	2.5380	2.4163
	0.95	241.88	19.396	8.7855	5.9644	4.7351	4.0600	3.6365	3.3472	3.1373
	0.975	968.63	39.398	14.419	8.8439	6.6192	5.4613	4.7611	4.2951	3.9639
	0.99	6055.8	99.399	27.229	14.546	10.051	7.8741	6.6201	5.8143	5.2565
11	0.9	60.473	9.4006	5.2224	3.9067	3.2816	2.9195	2.6839	2.5186	2.3961
	0.95	242.98	19.405	8.7633	5.9358	4.7040	4.0274	3.6030	3.3130	3.1025
	0.975	973.03	39.407	14.374	8.7935	6.5678	5.4098	4.7095	4.2434	3.9121
	0.99	6083.3	99.408	27.133	14.452	9.9626	7.7896	6.5382	5.7343	5.1779
12	0.9	60.705	9.4081	5.2156	3.8955	3.2682	2.9047	2.6681	2.5020	2.3789
	0.95	243.91	19.413	8.7446	5.9117	4.6777	3.9999	3.5747	3.2839	3.0729

Table with columns n1, alpha, and n2 (30-110). It provides critical values for the Wilcoxon signed-rank test.

12.6 Wilcoxon-Vorzeichen-Rang-Test

Tabelliert sind die kritischen Werte w_{alpha,n}^+. Ablesebeispiel: w_{0,95,12}^+ = 59

Table with columns n\alpha and values for alpha = 0.01, 0.025, 0.05, 0.10, 0.90, 0.95, 0.975, 0.99. It provides critical values for the Wilcoxon signed-rank test.

12.7 Wilcoxon-Rangsummen-Test

Tabelliert sind die kritischen Werte w_alpha fuer alpha = 0.05 (1. Zeile) und alpha = 0.10 (2. Zeile).

Ablesebeispiel: Fuer n = 3 und m = 7 ist w_{0,10} = 11.

Es ist w_{1-alpha;n,m} = n(n + m + 1) - w_{alpha;n,m}

Large table with columns n\m and values for m from 2 to 20. It provides critical values for the Wilcoxon rank-sum test.