

Angewandte Statistik

8. November 2001

Inhaltsverzeichnis

1	Information, Entscheidung, Statistik	1
1.1	Daten und Information	1
1.2	Entscheidung bei Unsicherheit	2
1.3	Kontinuierliche Entscheidungsprobleme	3
2	Beschreibende Statistik	4
2.1	Empirische Verteilungsfunktion	4
2.2	Diskrete Größen	4
2.3	Kontinuierliche Größen	4
2.4	Statistische Kenngrößen eindimensionaler Häufigkeitsverteilungen	4
2.4.1	Lageparameter	4
2.4.2	Streuungsparameter	7
2.5	Schiefe und Exzess eindimensionaler, empirischer Verteilungen	9
2.6	2-Dimensionale Häufigkeitsverteilungen	9
2.6.1	Kontingenzdaten	10
2.6.2	Empirischer Korrelationskoeffizient	10
2.7	Indexzahlen	11
2.7.1	Einfache Indexzahlen	11
2.7.2	Zeitreihen	11
2.7.3	Indexzahlen (zusammengesetzte Indizes)	12
3	Grafische Methoden zu Beurteilung empirischer Verteilungen	14
3.1	Normalverteilungsnetz	14
3.2	Wahrscheinlichkeitsnetze für logarithmische normalverteilte Merkmale	15
3.3	Zuverlässigkeitsanalyse und Ausfallraten	15

Inhaltsverzeichnis

3.4	Lebensdauer-Netz (Weibull-Verteilung)	17
3.5	Anpassung von Verteilungsmodellen	18
3.6	Ergänzung zu Schätzungen	20
4	Bayes'sche Anteilsschätzung	23
4.1	Gamma-Funktion	23
4.2	Gamma-Verteilung $\text{Gam}(\alpha, \beta)$ $\alpha > 0, \beta > 0$	23
4.3	Beta-Verteilung 1.Art $\text{Be}(\alpha, \beta)$ $\alpha > 0, \beta > 0$	24
4.4	Bayes'sches Theorem	25
4.5	A-posteriori Dichte für einfache Anteile	27
4.6	HPD-Intervalle für θ	28
4.7	Bayes'sche Anteilsschätzung bei mehrfachen Alternativen	31
4.8	Konjugierte Verteilungsfamilien	33
4.9	Suffizienz	34
5	Statistische Qualitätskontrolle	35
5.1	Einfache Stichprobenpläne	35
5.2	Festlegung einfacher Stichprobenpläne	38
5.3	Zweifache Stichprobenpläne	40
5.4	Sequentielle Stichprobenpläne	40
5.5	Kontrollkarten	41
6	Regressionsanalyse	46
6.1	Regression 1.Art	46
6.2	Regression 2.Art	50
6.3	Stochastische Regressionsanalyse	54
6.4	Regressionsgeraden 2.Art bei Normalverteilung	55
6.5	Tests für den Parameter von Regressionsgeraden bei Normalverteilung	60
6.6	Tests für Regressionskurven bei Normalverteilung	62
6.7	Test auf Regressionsgerade	62
6.8	Multiple lineare Regression	65
6.9	Multiple lineare Regression bei Normalverteilung	69
6.10	Bayes'sche Regressionsanalyse	72

Inhaltsverzeichnis

7	Varianzanalyse	75
7.1	Grundlagen der Varianzanalyse (ANOVA)	75
7.2	Einfache Varianzanalyse (1 Faktor)	76
7.3	Zweifache Varianzanalyse	78

1 Information, Entscheidung, Statistik

Für gute Organisation bzw. vernünftige Entscheidungen

Beispiel: Wasserbedarf eines Stadt: Zahlenangaben, Meßvorgänge

Statistik ist die zahlenmäßige Erfassung und Beschreibung von Phänomenen.

Beispiel:

- Volkszählung (Zensus)
- Stichprobenerhebung (Mikrozensus, Sozialstatistik)
- Prognosen (z.B. Bevölkerung)

1.1 Daten und Information

Aus Erhebungen bzw. empirischen Untersuchungen (Versuche) erhält man Rohdaten, oft auch nur Daten genannt. Beispiel:

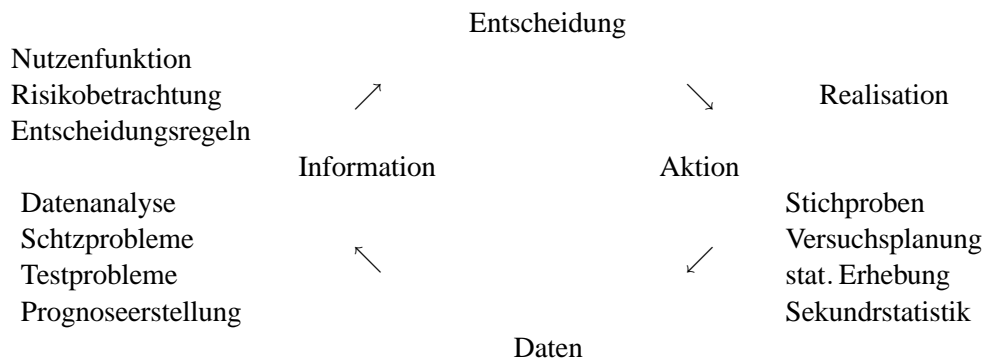
- Messung von Wartezeiten
- Messung von Lebensdauern

Informationsgewinnung: Daten oft in großer Zahl, daher unübersichtlich; schwierig zu vergleichen. Als Entscheidungsgrundlage nicht geeignet. Konzentration der Daten in Tabellen, Diagrammen und Ermittlung von Verteilungen.

Bemerkung: Kontrolle auf Vollständigkeit/Plausibilität

Raffung von Daten: z.B. Börsenindex

Information-Feedback-Cycle



1.2 Entscheidung bei Unsicherheit

Die stochastische Beschreibung von Entscheidungen, welche mit Nutzen bzw. Verlust verbunden sind.

Beispiel: Planung einer Straße, m mögliche Varianten, k verschiedene mögliche Klassen von Fahrzeughäufigkeiten.

$\theta_1, \dots, \theta_k$	k verschiedene Zustände (Fahrzeughäufigkeiten)
d_1, \dots, d_m	mögliche Entscheidungen (Fahrbahnbreiten)
$\tilde{\theta}$	stochastische Grö., die den Zustand beschreibt
$p(\theta_1), \dots, p(\theta_k)$	Wahrscheinlichkeiten der $\theta_1, \dots, \theta_k \hat{=} \text{Wahrscheinlichkeitsverteilung von } \tilde{\theta}$
$U(\theta_i, d_j) \geq 0$	Nutzenfunktion; Nutzen der Entscheidung d_j , falls System im Zustand θ_i ist.
$\bar{U}(d_j) = \mathbb{E}U(\tilde{\theta}, d_j) = \sum_{i=1}^k U(\theta_i, d_j) \cdot p(\theta_i)$	zu erwartender Nutzen der Entscheidung d_j .

Die optimale Entscheidung d^* ist jenes d_j , für das $\bar{U}(d_j)$ maximal ist (Bayes'sche Entscheidung).

Bemerkung: Man kann diese Probleme auch für kontinuierliche Zustände und Entscheidungsmöglichkeiten beschreiben (vgl. Abschnitt 1.3).

Beispiel: Ausnahme/Ablehnung einer Warensendung

N Losumfang

1 Information, Entscheidung, Statistik

θ	Anteil der schlechten Stecke im Los
d_0	Entscheidung Annahme
d_1	Entscheidung Ablehnung

Entscheidung ist abhängig von einer Stichprobe: x_1, \dots, x_n mit $n \ll N$, $d = \delta(x_1, \dots, x_n)$,
 $\delta \hat{=} \text{Entscheidungsregel}$

Aus der Stichprobe erhält man Aussagen über θ (vgl. 4).

K	=	Verlust fr ein schlechtes Stck
G		Gewinn fr ein gutes verkaufte Stck

Nutzen:

$$\begin{aligned}U(\theta, d) & \quad \text{falls Anteil } \theta \\U(\theta, d_0) & = (1 - \theta)NG - \theta NK \\U(\theta, d_1) & = 0\end{aligned}$$

Entscheidungskriterium: zu erwartender Nutzen $\mathbb{E}[U(\tilde{\theta}, d)]$
Erwartungsbildung mittels der Verteilung von $\tilde{\theta}$: in diesem Fall: a-posteriori Verteilung
 $\pi(\theta|x_1, \dots, x_n)$ (vgl. Abschnitt 4)

$$\tilde{\theta} \sim U_{0,1} \quad \Rightarrow \quad \tilde{\theta}|_{x_1, \dots, x_n} \hat{=} \text{Be}(\cdot, \cdot)$$

Entscheidung: $d_0 \hat{=} \bar{U}(d_0) > 0$

Frage: $\bar{U}(d_j) = ?$

1.3 Kontinuierliche Entscheidungsprobleme

Die Menge der möglichen Werte für θ ist kontinuierlich: $\theta \in \Theta$. Daher ist die Verteilung von $\tilde{\theta}$ ebenfalls kontinuierlich mit Dichtefunktion $g(\theta)$. Der erwartende Nutzen ist:

$$\bar{U}(d) = \mathbb{E}[U(\tilde{\theta}, d)] = \int_{\Theta} U(\theta, d) \cdot g(\theta) d\theta$$

Bemerkung: Häufig ist $g(\theta)$ eine a-posteriori-Dichte $\pi(\theta|D)$.

Die optimale Entscheidung d^* ist dann jene, mit größtem, erwartenden Nutzen, d.h.:

$$\bar{U}(d^*) = \max_{d \in \mathcal{D}} \bar{U}(d) \quad \mathcal{D} \hat{=} \text{Menge aller mglichen Entscheidungen}$$

2 Beschreibende Statistik

Erhält man aus Erhebungen oder Versuchen eine große Anzahl von Daten, so sollen diese zu einer überschaubaren Information zusammengefasst werden: Konzentration in Kennzahlen, Tabellen, Diagrammen und empirischen Verteilungen.

2.1 Empirische Verteilungsfunktion

2.2 Diskrete Größen

2.3 Kontinuierliche Größen

2.4 Statistische Kenngrößen eindimensionaler Häufigkeitsverteilungen

Für Beobachtungen x_1, \dots, x_n eines 1-dimensionalen Merkmales sucht man irgendwie charakterisierende Größen der Verteilung, z.B. für die Mitte (Lageparameter) oder für das Streuverhalten der Verteilung (Streuparameter).

2.4.1 Lageparameter

Mittlerer Wert

$$x_1, \dots, x_n \rightarrow \bar{x}_n := \frac{\sum_{i=1}^n x_i}{n} \quad \text{Stichprobenmittel}$$

Sind die Daten gruppiert mit Klassenmitten z_j und Häufigkeiten H_j so gilt:

$$\bar{z} = \frac{1}{n} \sum_{j=1}^k z_j \cdot H_j \quad j = 1..k \quad \text{empirischer Mittelwert}$$

Bei der Einteilung in Klassen kommt es zu einem Verlust von Information.

Geometrisches Mittel

$$\bar{x}_g := \sqrt[n]{x_1 \cdot x_2 \cdot \dots \cdot x_n} = \sqrt[n]{\prod_{i=1}^n x_i}$$

Beispiel: Wachstumsfaktoren bei Gemeindegrößen

Harmonisches Mittel

z.B. für Durchschnittsgeschwindigkeiten

$x_1, \dots, x_n \quad x_i > 0 \quad \text{oder} \quad x_i < 0 \quad \forall i$

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

Beispiel: Teilstrecken L_i mit x_i Geschwindigkeit durchfahren, Zeiten t_i

Gesamtstrecke $L = \sum_{i=1}^n L_i$ und Gesamtzeit $t = \sum_{i=1}^n t_i$

Durchschnittsgeschwindigkeit:

$$\begin{aligned} x &= \frac{l}{t} = \frac{\sum l_i}{\sum t_i} = \frac{\sum l_i}{\sum \frac{l_i}{x_i}} = \\ &= \frac{1}{\sum_{i=1}^n \frac{l_i}{x_i \cdot \sum l_i}} = \frac{1}{\sum \frac{w_i}{t_i}} \quad \text{mit } w_i = \frac{l_i}{\sum l_i} \end{aligned}$$

Sind die Daten in k Klassen gruppiert mit Klassenmitten z_j und Häufigkeiten H_j so gilt:

$$\bar{z}_{gh} := \frac{n}{\sum_{j=1}^k \left(\frac{1}{z_j} \cdot H_j \right)}$$

Empirischer Median

Für Einzeldaten x_1, \dots, x_n wird bei ungerader Anzahl, d.h. $n = 2 \cdot k + 1$ der mittlere Wert $x_{(k+1)}$ als Median definiert. Bei gerader Anzahl $n = 2 \cdot k$ ist er definiert als $\frac{x_{(k)} + x_{(k+1)}}{2}$.

Für gruppiert Daten ermittelt man das Summenpolygon. Der empirische Median $x_{0.5}$ ist definiert als das 0.5-Fraktile des Summenpolygons (vgl. Abbildung 2.1)

2 Beschreibende Statistik

Abbildung 2.1: Empirischer Median

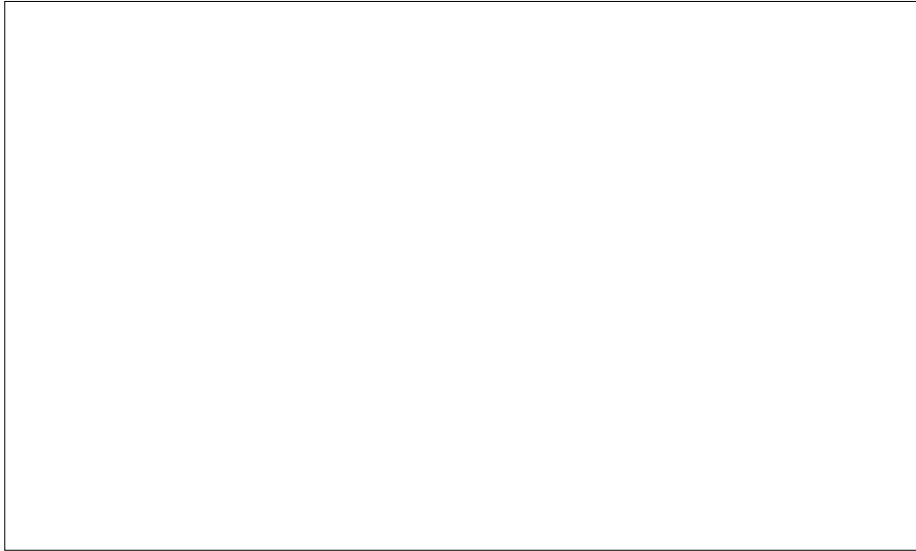


Abbildung 2.2: Empirischer Modus für diskrete, empirische Verteilungen

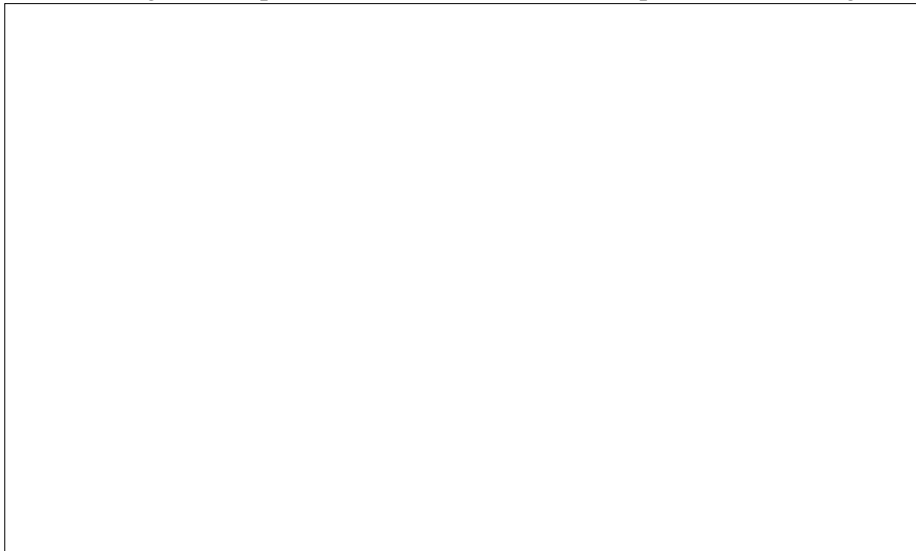
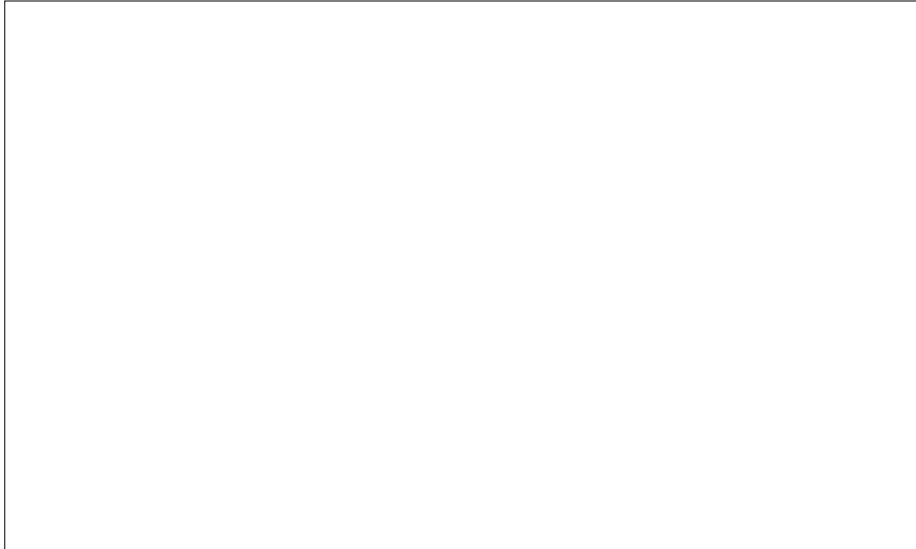


Abbildung 2.3: Empirischer Modus für kontinuierliche Größen



Empirischer Modus (Modalwert)

Für diskrete, empirische Verteilungen ist der empirische Modus als der häufigste Wert definiert, falls dieser existiert (vgl. Abbildung 2.2).

Für kontinuierliche Größen vgl. Abbildung 2.3.

2.4.2 Streuungsparameter

Spannweite

$$x_{(n)} - x_{(1)}$$

Quartilabstand

$$x_{0.75} - x_{0.25}$$

Bemerkung: Im Intervall $[x_{0.25}, x_{0.75}]$ liegen 50% aller Daten

Mittlere absolute Abweichung

$$MAD := \frac{1}{n} \sum_{i=1}^n |x_i - x_{0.5}|$$

2 Beschreibende Statistik

Bemerkung: Berechnet man die Summe $\sum_{i=1}^n |x_i - a|$ für ein beliebiges $a \in \mathbb{R}$, so gilt:

$$\sum_{i=1}^n |x_i - x_{0.5}| \leq \sum_{i=1}^n |x_i - a| \quad \forall a \in \mathbb{R}$$

Bei gruppierten Daten mit k Gruppen und den Gruppenmitten z_j sowie den Häufigkeiten H_j gilt:

$$MAD := \frac{1}{n} \sum_{j=1}^k |z_j - z_{0.5}| \cdot H_j$$

Empirische Varianz

Bei vollständigen Daten gilt:

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

Bei unvollständigen Daten gilt:

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

Bei Stichproben erfolgt die Division durch $n - 1$ zwecks Unverzerrtheit. Es gilt:

$$\sum_{i=1}^n (x_i - \bar{x})^2 \leq \sum_{i=1}^n (x_i - a)^2 \quad \forall a \in \mathbb{R}$$

Für gruppierte Daten mit Gruppenmitten z_j und Häufigkeiten H_j , $j = 1(1)k$ gilt:

$$S_g^2 = \frac{1}{n} \sum_{j=1}^k (z_j - \bar{z})^2 \cdot H_j = \sum_{j=1}^k (z_j - \bar{z})^2 \cdot h_j$$

Die empirische Streuung bzw. empirische Standardabweichung ist definiert als:

$$S := +\sqrt{S^2}$$

Empirischer Variationskoeffizient (dimensionslos)

$$VK = \frac{S}{\bar{x}}$$

2.5 Schiefe und Exzess eindimensionaler, empirischer Verteilungen

Die Abweichung einer „Verteilung“ von der symmetrischen Form wird manchmal durch die sog. Schiefe beschrieben. Für gruppierte Daten gilt:

$$\gamma_1 = \frac{1}{n \cdot S^3} \sum_{j=1}^k (z_j - \bar{z})^3 \cdot H_j$$

Der Wert von γ_1 kann positiv oder negativ werden:

rechtschief	$\gamma_1 > 0$
linksschief	$\gamma_1 < 0$
symmetrisch	$\gamma_1 = 0$

Es Maß für die „Spitzigkeit“ (Stauchung bzw. Überhöhung) der Verteilungsform ist der sog. Exzess:

$$\gamma_2 = \left[\frac{1}{n \cdot S^4} \sum_{j=1}^k (z_j - \bar{z})^4 \cdot H_j \right] - 3$$

Auch γ_2 kann positive oder negative Werte annehmen.

überhöht	$\gamma_2 > 0$
gestaucht	$\gamma_2 < 0$
Ähnlich einer Normalverteilung	$\gamma_2 = 0$

Bemerkung: Pearson-Diagramm: $\beta_2 = \gamma_2 + 3$, $\beta_1 = \gamma_1$ (vgl. Abschnitt 3)

2.6 2-Dimensionale Häufigkeitsverteilungen

Zur Untersuchung von Zusammenhängen zwischen 1-dimensionalen Merkmalen. Je nach der Art der Größen X und Y werden verschiedene Zusammenhangsmaße verwendet.

2.6.1 Kontingenzdaten

X und Y sind qualitativer Art, d.h. diskret aber nicht der Größe nach zu ordnen.

X l Ausprägungen

Y m Ausprägungen

Daten: Paare (x_k, y_k) , $k = 1(1)n$

x	y	b_1	b_2	b_j	b_m	$n_{i\bullet} = \sum_{j=1}^m n_{ij}$
a_1		n_{11}	n_{12}	n_{1j}	n_{1m}	$n_{1\bullet}$
a_2		n_{21}	n_{22}	n_{2j}	n_{2m}	$n_{2\bullet}$
a_i		n_{i1}	n_{i2}	n_{ij}	n_{im}	$n_{i\bullet}$
a_l		n_{l1}	n_{l2}	n_{lj}	n_{lm}	$n_{l\bullet}$
	$n_{\bullet j} = \sum_{i=1}^l n_{ij}$	$n_{\bullet 1}$	$n_{\bullet 2}$	$n_{\bullet j}$	$n_{\bullet m}$	n

n_{ij} =Anzahl jener Paare mit Ausprägung $x_k = a_i$, $y_k = b_j$

Bemerkung: Für Abhängigkeitsuntersuchungen.

Beispiel: Vier-Felder-Tafel

x	y	b_1	b_2	
a_1		n_{11}	n_{12}	$n_{1\bullet}$
a_2		n_{21}	n_{22}	$n_{2\bullet}$
		$n_{\bullet 1}$	$n_{\bullet 2}$	n

Der sog. Yule'sche Assoziationskoeffizient ist definiert als:

$$Q = \frac{n_{11} \cdot n_{22} - n_{12} \cdot n_{21}}{n_{11} \cdot n_{22} + n_{12} \cdot n_{21}} \quad -1 \leq Q \leq 1$$

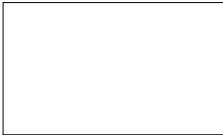

2.6.2 Empirischer Korrelationskoeffizient

X und Y sind metrische Merkmale.

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

Bemerkung:

- $-1 \leq r \leq 1$
- r ändert sich nicht, wenn man anstelle der Wertepaare (x_i, y_i) Lineartransformationen (u_i, v_i) mit $u_i = \frac{x_i - x_0}{c}$ und $v_i = \frac{y_i - y_0}{c}$ verwendet.
- Falls $|r| = 1$, so liegen alle Paare (x_i, y_i) auf einer Geraden in der (x, y) -Ebene.

$r = +1$		Positive lineare Abhängigkeit. Falls x steigt, so steigt auch y .
$r = -1$		Negative lineare Abhängigkeit. Falls x steigt, so sinkt y .
$r = 0$		Keine lineare Abhängigkeit. Bedeutet NICHT keine Abhängigkeit.

2.7 Indexzahlen

2.7.1 Einfache Indexzahlen

Statistische Kennzahlen, die nebengeordnete Größen auf eine von ihnen oder auf einen Durchschnitt beziehen.

- Maßzahlen des sachlichen Vergleiches.
- Maßzahlen des örtlichen Vergleiches.
- Maßzahlen des zeitlichen Vergleiches.

2.7.2 Zeitreihen

Zeitlich geordnete Folge statistische Beobachtungen $(A_t, t = 0, 1, 2, \dots)$. Darauf ermittelt man eine Folge von Maßzahlen.

Zeitreihe der Maßzahlen mit Basisjahr 0:

$$I_{0t} := \frac{A_t}{A_0}; \quad t = 0, 1, 2, 3, \dots$$

Abbildung 2.4: Meßzahlen gewonnen aus einer Zeitreihe



2.7.3 Indexpzahlen (zusammengesetzte Indizes)

Hat man mehrere sachlich zusammenhängende Reihen von Daten, so möchte man deren Verlauf oft jeweils durch eine Zahl beschreiben. Dies erfolgt durch sog. Indexpzahlen.

Wertindex: Aussage über die relative Änderung des Wertes eines Bündels von Waren oder Dienstleistungen.

$1, \dots, n$	Warennummern der n Waren des Warenkorbes
p_{t1}, \dots, p_{tn}	Preise der Waren zu Berichtsperiode t
q_{t1}, \dots, q_{tn}	Mengen der Waren zu Berichtsperioden t

Mit der Basisperiode $t = 0$ ergibt sich der Wert des Warenkorbes in der Periode t als:

$$\sum_{i=1}^n p_{ti} \cdot q_{ti}$$

Der Wertindex ist definiert als:

$$I_{0t}^W := \frac{\sum_{i=1}^n p_{ti} \cdot q_{ti}}{\sum_{i=1}^n p_{0i} \cdot q_{0i}}$$

Der Preisindex nach Paasche ist definiert als:

2 Beschreibende Statistik

$$pI_{0t}^P := \frac{\sum_{i=1}^n p_{ti} \cdot q_{ti}}{\sum_{i=1}^n p_{0i} \cdot q_{ti}}$$

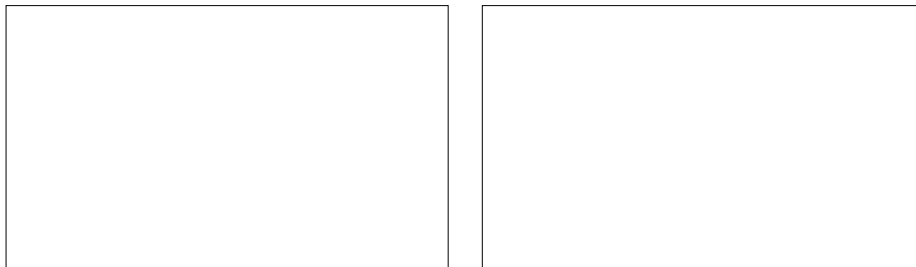
Der Preisindex nach Laspeyres (ÖSTAT) ist definiert als:

$$L I_{0t}^P := \frac{\sum_{i=1}^n p_{ti} \cdot q_{0i}}{\sum_{i=1}^n p_{0i} \cdot q_{0i}}$$

3 Grafische Methoden zu Beurteilung empirischer Verteilungen

Zur Beurteilung der Verteilungsart von empirisch gegebenen Merkmalen (stoch. Größen) konstruiert man sog. Wahrscheinlichkeitspapiere (oder Netze).

Dabei transformiert man die Darstellungsebene der Verteilungsfunktion so, dass die Bilder der entsprechenden Verteilungsfunktion in Geraden übergeführt werden.



$$(x, y) \rightarrow (t, z)$$

$$t = \Phi(x) \quad z = \psi(y)$$

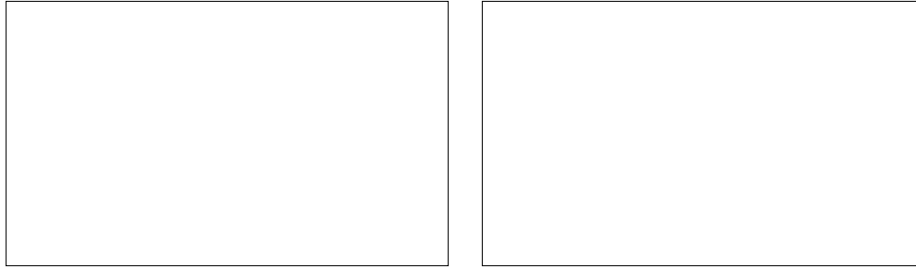
Das bekannteste solche Papier ist jenes für die Normalverteilung, genannt:

3.1 Normalverteilungsnetz

Ist $X \sim N(\mu, \sigma^2)$ mit Verteilungsfunktion $F(x) = W\{X \leq x\} = \Phi\left(\frac{x-\mu}{\sigma}\right)$, $\forall x \in \mathbb{R}$, so wird folgende Transformation durchgeführt:

$$\begin{aligned} t &= x \\ z &= \Phi^{-1}(y) \\ u_p &= p - \text{Fraktile der } N(0, 1) \end{aligned}$$

3 Grafische Methoden zu Beurteilung empirischer Verteilungen



Es gilt: Ist x_p das p -Fraktile der $N(\mu, \sigma^2)$ so gilt folgender Zusammenhang:

$$u_p = \frac{x_p - \mu}{\sigma} = \frac{1}{\sigma}x_p - \frac{\mu}{\sigma}$$

Beweis:

$$p = W\{X \leq x_p\} = W\left\{\frac{x - \mu}{\sigma} \leq \frac{x_p - \mu}{\sigma}\right\}$$

3.2 Wahrscheinlichkeitsnetze für logarithmische normalverteilte Merkmale

Das Bild der Verteilungsfunktion von $\ln X$ ist im Normalverteilungsnetz eine Gerade. Teilt man die Abszisse logarithmisch, so kann man direkt die Werte von X eintragen. Die Ordinate ist wie beim Wahrscheinlichkeitspapier für Normalverteilungen geteilt (Fraktile der $N(0, 1)$).

Trägt man die Punkte der empirischen Verteilungsfunktion in das logarithmische Wahrscheinlichkeitsnetz ein, so deutet ein annähernd geradliniger Verlauf auf eine $LN(\mu, \sigma^2)$ hin.

3.3 Zuverlässigkeitsanalyse und Ausfallraten

Lebensdauern X mit kontinuierlichen Verteilungen und Dichtefunktion $f(\cdot)$, $f(x) = 0 \quad \forall x < 0$. Gesucht ist die altersabhängige Ausfallswahrscheinlichkeit

$$W\{x \leq X \leq x + \Delta x \mid X > x\} = ?$$

Bemerkung: Bei Lebensdauern wird meist nicht die Verteilungsfunktion sondern die Zuverlässigkeitsfunktion

$$R(x) = W\{X > x\} = 1 - F(x) \quad \forall x \geq 0$$

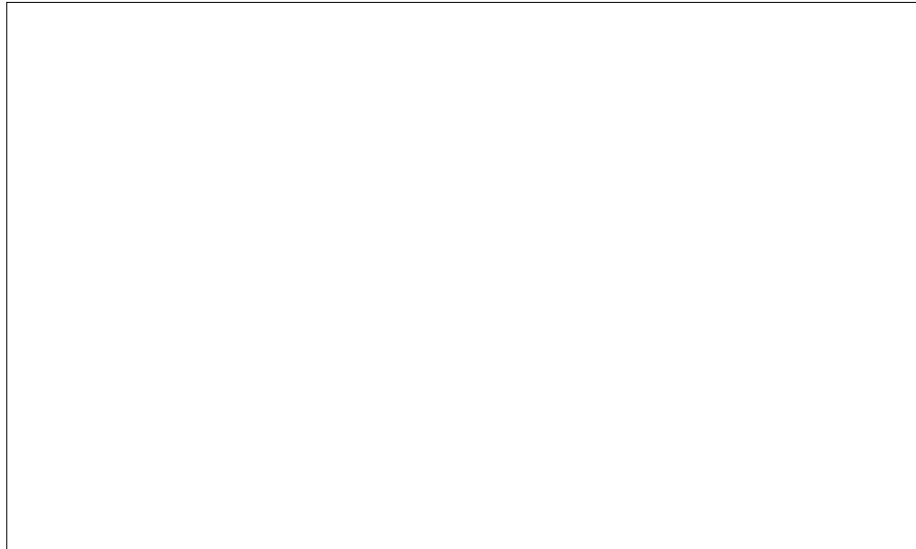
3 Grafische Methoden zu Beurteilung empirischer Verteilungen

betrachtet. Damit erhält man für obige altersabhängige Ausfallswahrscheinlichkeit:

$$\begin{aligned} W\{x < X \leq x + \Delta x \mid X > x\} &= \frac{W\{x < X \leq x + \Delta x\}}{R(x)} \\ &= \frac{\int_x^{x+\Delta x} f(t) dt}{R(x)} \\ &= \frac{f(\xi) \cdot \Delta x}{R(x)} \end{aligned}$$

Für kleine Δx ist der letzte Ausdruck ungefähr $\frac{f(x)}{R(x)} \cdot \Delta x$, d.h. die altersabhängige Ausfallswahrscheinlichkeit ist nahezu proportional zu $\frac{f(x)}{R(x)}$. Diese altersabhängige Größe $\lambda(x) := \frac{f(x)}{R(x)}$ heißt Ausfallsrate.

Abbildung 3.1: Ausfallsrate



Beispiel: Für die Exponentialverteilung gilt:

$$\begin{aligned} f(x) &= \frac{1}{\tau} \cdot e^{-\frac{x}{\tau}} \mathbf{I}_{(0, \infty)}(x) \\ R(x) &= e^{-\frac{x}{\tau}} \quad \forall x > 0 \\ \lambda(x) &= \frac{1}{\tau} \quad \text{konstant!} \end{aligned}$$

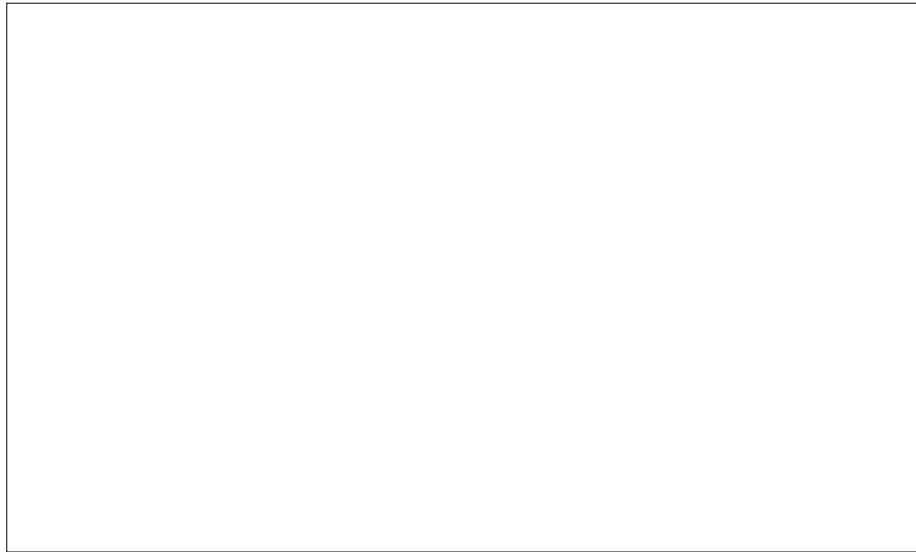
Für Ausfallsraten der Form $\lambda(x) = c \cdot x^\alpha$ verwendet man die Weibull-Verteilung.

3.4 Lebensdauer-Netz (Weibull-Verteilung)

$$X \sim \text{Wei}(\tau, \beta)$$

$$F(t) = W\{X \leq t\} = 1 - R(t) = 1 - e^{-\left(\frac{t}{\tau}\right)^\beta} \quad \forall x > 0$$

Abbildung 3.2: Weibull-Verteilung



Die Transformation des (t, z) -Koordinatensystems in das (x, y) -Koordinatensystem geschieht folgendermaßen:

$$\begin{aligned} x &= \ln t \\ y &= \ln \left[\ln \left(\frac{1}{1-z} \right) \right] \end{aligned}$$

Es gilt: Die Bilder von Verteilungsfunktionen von Weibull-verteilten stochastischen Größen sind im (x, y) -Koordinatensystem Geraden, denn:

$$\begin{aligned} y &= \ln \left[\ln \left(\frac{1}{1-F(t)} \right) \right] = \\ &= \beta \cdot (\ln t - \ln \tau) = \\ &= \beta \cdot x - \beta \cdot \ln \tau \end{aligned}$$

Abbildung 3.3: Transformierte Weibull-Verteilung



3.5 Anpassung von Verteilungsmodellen

Für Modellrechnungen und Prognosen benötigt man „die“ Wahrscheinlichkeitsverteilung der entsprechenden stochastischen Größe.

Bemerkung: Dazu dienen Tests, Parameterschätzungen (Konfidenzbereiche) und Prognoseverteilungen. Auch das sogenannte Pearson-Diagramm ist eine Entscheidungshilfe.

Zur Kontrolle des sog. Pearson-Diagramms benötigt man analog zur Schiefe und Exzess empirischer Verteilungen für Wahrscheinlichkeitsverteilungen stochastische Größen.

Definition: Ist X eine stochastische Größe mit existierendem $\mathbb{E}(X^r)$, so schreibt man:

$$\begin{aligned}\mu'_r &= \mathbb{E}(X^r) \\ \mu_r &= \mathbb{E}[(X - \mathbb{E}X)^r] \\ \mu'_1 &= \mu = \mathbb{E}(X) \\ \mu_2 &= \text{Var}X\end{aligned}$$

Vergleiche Abschnitt 2 (Schiefe, Spitzigkeit).

Schiefe: Es gilt:

$$\begin{aligned}\mu_3 &:= \mathbb{E}(X - \mu)^3 = \mu'_3 - 3 \cdot \mu'_2 \mu + 2 \cdot \mu^3 \\ \sqrt{\beta_1} &:= \frac{\mu_3}{(\mu_2)^{\frac{3}{2}}} \dots \text{Maß für die Schiefe relativ zu Spannweite}\end{aligned}$$

3 Grafische Methoden zu Beurteilung empirischer Verteilungen

Beispiel: Exponentialverteilung $\sqrt{\beta_1} = 2$

Für unimodale (1 Maximum) Verteilungen ist μ_4 ein Maß für die Spitzigkeit:

$$\beta_2 := \frac{\mu_4}{\mu_2} \dots \text{Maß für die relative Spitzigkeit (Wlbung)}$$

Exzess: $\gamma_2 = \beta_2 - 3$

Beispiel: Normalverteilung, Gleichverteilung (vgl. Abbildung 3.4)

Abbildung 3.4: β_2 ist ein Maß für die relative Spitzigkeit



Trägt man die Werte (β_1, β_2) in einem Diagramm auf, so erhält man Charakteristika für Verteilungstypen: sog. Pearson-Diagramm.

Definition: Ist X_1, \dots, X_n eine Stichprobe einer Verteilung W , so heißt $M'_r := \frac{1}{n} \sum_{i=1}^n X_i^r$ das r -te Stichprobenmoment bezüglich 0 und $M_r := \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^r$ das r -te Stichprobenmoment bezüglich \bar{X}_n .

Satz 3.1 Ist X_1, \dots, X_n eine Stichprobe von X mit $\exists \mu'_r$ und $\exists \mu'_{2r}$, so folgt:

- $\mathbb{E}M'_r = \mu'_r$ (unverzerrte Schätzfunktion)
- $\text{Var}M'_r = \frac{1}{n} \cdot \left[\mathbb{E}(X^{2r}) \cdot \mathbb{E}(X^r)^2 \right] = \frac{1}{n} \cdot \left[\mu'_{2r} - (\mu'_r)^2 \right]$

Korollar 3.1 Ist X_1, \dots, X_n eine Stichprobe von X und $\exists \text{Var}X = \sigma^2$ so folgt:

$$\mathbb{E}\bar{X}_n = \mathbb{E}X \quad \text{und} \quad \text{Var}\bar{X}_n = \frac{\sigma^2}{n}$$

Beweis: Sonderfall $r=1$

Anwendung des Pearsondiagrammes

Aus der Stichprobe x_1, \dots, x_n schätzt man die entsprechenden Momente. Daraus berechnet man Näherungswerte b_1 und b_2 für β_1 und β_2 und trägt das Paar (b_1, b_2) in das Pearsondiagramm ein. Aus der Lage des Punktes erhält man einen Anhaltspunkt für den Verteilungstyp.

3.6 Ergänzung zu Schätzungen

Für die Stichprobenvarianz

$$S_n^2 := \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2$$

gilt:

- $S_n^2 = \frac{n}{n-1} \cdot M_2$
- $S_n^2 = \frac{1}{2n \cdot (n-1)} \sum_{i=1}^n \sum_{j=1}^n (X_i - X_j)^2$
 $\sum_{i=1}^n (X_i - \bar{X}_n)^2 = \dots = \sum_{i=1}^n X_i^2 - \frac{1}{n} \cdot (\sum_{i=1}^n X_i)^2$
 $\frac{1}{2n} \sum_{i=1}^n \sum_{j=1}^n (X_i - X_j)^2 = -\| -$

Satz 3.2 Ist X_1, \dots, X_n eine Stichprobe von X mit $\exists \mathbb{E}X^4$ und $\exists \text{Var}X = \sigma^2$ so gilt:

- $\mathbb{E}S_n^2 = \sigma^2$
- $\text{Var}S_n^2 = \frac{1}{n} \left(\mu_4 - \frac{n-3}{n-1} \sigma^4 \right)$ für $n > 1$

Satz 3.3 Ist X eine 1-dimensionale stochastische Größe mit endlicher Varianz σ^2 und X_1, \dots, X_n eine Stichprobe von X , so ist

$$Z_n = \frac{\bar{X}_n - \mu}{\frac{\sigma}{\sqrt{n}}}$$

asymptotisch $N(0, 1)$ -verteilt.

Beweis: Zentraler Grenzwertungssatz

Bemerkung: \bar{X}_n ist näherungsweise normalverteilt mit Mittel μ und der Varianz $\left(\frac{\sigma^2}{n} \right)$.

Beispiel: $X \sim Ex_1, f(x) = e^{-x} \mathbf{I}_{(0, \infty)}(x)$

Anwendung Anteilsschätzung

Gesucht ist ein approximatives Konfidenzintervall für den Parameter θ einer Alternativverteilung. Für

$$\bar{X}_n = \frac{\sum_{i=1}^n X_i}{n}$$

gilt nach dem Satz von MOIVRE:

$$W \left\{ \frac{\bar{X}_n - \theta}{\sqrt{\theta(1-\theta)/n}} \leq x \right\} \rightarrow_{n \rightarrow \infty} \Phi(x)$$

Dies ist ein Sonderfall des zentralen Grenzwertungssatzes. Der Ausdruck

$$\frac{\bar{X}_n - \theta}{\sqrt{\theta(1-\theta)/n}}$$

ist approximativ $N(0, 1)$ -verteilt. Das bedeutet:

$$\bar{X}_n \sim N \left(\mathbb{E}\bar{X}_n = \theta, \text{Var}X = \frac{\theta(1-\theta)}{n} \right)$$

$$W \left\{ -u_{\frac{1+\gamma}{2}} \leq \frac{\bar{X}_n - \theta}{\sqrt{\theta(1-\theta)/n}} \leq +u_{\frac{1+\gamma}{2}} \right\} \doteq \gamma$$

mit $u_\alpha = \alpha$ -Fraktile der $N(0, 1)$. Man kann diese Doppelgleichung umschreiben als

$$(X_n - \theta)^2 \leq \frac{\theta(1-\theta)}{n} u_{\frac{1+\gamma}{2}}^2$$

$$(\bar{X}_n - \theta)^2 - \frac{\theta(1-\theta)}{n} \cdot u_{\frac{1+\gamma}{2}}^2 \leq 0$$

Nun sucht man die untere Nullstelle $\underline{\theta}(\bar{X}_n, \gamma)$ und die obere Nullstelle $\bar{\theta}(\bar{X}_n, \gamma)$:

$$\underline{\theta}(\bar{X}_n, \gamma) = \frac{2n\bar{X}_n + u^2 - u \cdot \sqrt{4n\bar{X}_n - 4n\bar{X}_n^2 + u^2}}{2(n + u^2)}$$

$$\bar{\theta}(\bar{X}_n, \gamma) = \frac{2n\bar{X}_n + u^2 + u \cdot \sqrt{4n\bar{X}_n - 4n\bar{X}_n^2 + u^2}}{2(n + u^2)}$$

Division durch $2n$ und Vernachlässigung aller Glieder der Form $\frac{c}{n}$ ergibt als Näherung das approximative Konfidenzintervall mit Überdeckungswahrscheinlichkeit γ für θ

$$\left[\bar{X}_n - u_{\frac{1+\gamma}{2}} \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}}, \bar{X}_n + u_{\frac{1+\gamma}{2}} \sqrt{\frac{\bar{X}_n(1-\bar{X}_n)}{n}} \right].$$

3 Grafische Methoden zu Beurteilung empirischer Verteilungen

Bemerkung: Dabei ist nur jener Teil des Intervalles zu berücksichtigen, der im Intervall $[0, 1]$ liegt. Die Form des Konfidenzintervalles ist naheliegend, da \bar{X}_n eine gute Schätzfunktion für den $\mathbb{E}X_n = \theta$ ist und

$$\frac{\bar{X}_n(1 - \bar{X}_n)}{n}$$

eine Schätzfunktion für

$$\text{Var}\bar{X}_n = \frac{\theta(1 - \theta)}{n}$$

ist.

Beispiel: Stichprobe von $X \sim A_\theta$ vom Umfang $n=130$ und $\sum_{i=1}^n x_i=28$ und $\gamma = 0.95$. Gesucht ist ein approximatives Konfidenzintervall für θ .

Bemerkung: Vergleich mit HPD-Intervall für $\pi(\cdot) \hat{=} U_{0,1} = \text{Be}(1, 1)$ (vgl. Abschnitt 4).

4 Bayes'sche Anteilsschätzung

In der Bayes'schen Statistik werden alle unbekanntes Größen durch stochastische Größen mit zugehöriger Wahrscheinlichkeitsverteilung beschrieben. Für ein stochastisches Modell $X \sim f(\cdot|\theta)$ mit Parameterraum Θ die Unsicherheit bezüglich des Parameters θ , der durch eine stochastische Größe $\hat{\theta}$ beschrieben wird, mittels einer sog. a-priori-Verteilung $\pi(\theta)$ von $\hat{\theta}$ ausgedrückt.

Beispiel: $\pi(\cdot) \hat{=} U_{0,1} \Rightarrow \pi(\theta) = I_{(0,1)}(\theta)$

Dies ist ein Sonderfall der Funktion

$$f(\theta) = \theta^{\alpha-1}(1-\theta)^{\beta-1}I_{(0,1)}(\theta) \quad \text{mit } \alpha = \beta = 1,$$

welches die Dichtefunktion der Beta-Verteilung $\text{Be}(\alpha, \beta)$ darstellt.

Frage: Kann man durch eine Stichprobe (Daten D) bessere Informationen über θ erhalten?

Bemerkung: Die Aktualisierung erfolgt über das Bayes'sche Theorem (vgl. Abschnitt 4.4).

4.1 Gamma-Funktion

$$\Gamma(x) \quad \forall x > 0$$

$$\Gamma(n+1) = n! \quad \Gamma(1) = 1 \quad \Gamma\left(\frac{1}{2}\right) = \sqrt{\pi}$$

4.2 Gamma-Verteilung $\text{Gam}(\alpha, \beta) \quad \alpha > 0, \beta > 0$

$$f(x|\alpha, \beta) = \frac{x^{\alpha-1} \cdot e^{-\frac{x}{\beta}}}{\Gamma(\alpha) \cdot \beta^\alpha} I_{(0,\infty)}(x)$$

Falls $X \sim \text{Gam}(\alpha, \beta)$ dann gilt:

$$\mathbb{E}X = \alpha \cdot \beta \quad \text{Var}X = \alpha \cdot \beta^2$$

Anwendung: Lebensdauern

Bemerkung: Die Gamma-Verteilung ist eine Verallgemeinerung der Exponential-Verteilung und der χ_n^2 -Verteilung.

Abbildung 4.1: Gamma-Verteilung



Satz (Additionstheorem) für Gamma-Verteilungen

Falls $X_i \sim \text{Gam}(\alpha_i, \beta)$, $i = 1, \dots, n$ so gilt:

$$\sum_{i=1}^n X_i \sim \text{Gam}\left(\sum_{i=1}^n \alpha_i, \beta\right)$$

Anwendung: Lebensdauern von erneuerbaren Systemen (Kalte Reserven).

Bemerkung: Die Verteilungsfunktion zur Gamma-Verteilung ist durch die sog. unvollständige Gamma-Funktion gegeben:

$$F_X(x|\alpha, \beta) = W\{X \leq x\} = \int_0^x \frac{t^{\alpha-1} e^{-t/\beta}}{\Gamma(\alpha)\beta^\alpha} dt$$

Diese Funktion ist tabelliert.

4.3 Beta-Verteilung 1.Art $\text{Be}(\alpha, \beta)$ $\alpha > 0, \beta > 0$

$$f(x|\alpha, \beta) \propto x^{\alpha-1} (1-x)^{\beta-1} I_{(0,1)}(x)$$

Damit eine Dichtefunktion entsteht, bedient man sich einer normierenden Konstante:

$$\int_0^1 x^{\alpha-1} (1-x)^{\beta-1} dx =: B(\alpha, \beta) \dots \text{Betafunktion}$$

4 Bayes'sche Anteilsschätzung

Für B gilt:

$$B(\alpha, \beta) = 2 \cdot \int_0^{\pi/2} \sin^{2\alpha-1} \varphi \cdot \cos^{2\beta-1} \varphi d\varphi$$

Es gilt:

$$B(x, y) = \frac{\Gamma(x)\Gamma(y)}{\Gamma(x+y)}$$

Bemerkung: Die Verteilungsfunktion der Be-Verteilung 1.Art ist durch die sog. unvollständige Beta-Funktion gegeben:

$$F_X(x|\alpha, \beta) = \int_0^x \frac{t^{\alpha-1}(1-t)^{\beta-1}}{B(\alpha, \beta)} dt \quad 0 < x < 1$$

Diese Funktion ist tabelliert.

Anwendung: Anteilsschätzung im Fall einfacher Alternativen.

Abbildung 4.2: Beta-Verteilung



4.4 Bayes'sches Theorem

Dient dazu, um von der a-priori Einschätzung zur a-posteriori Einschätzung zu gelangen.

4 Bayes'sche Anteilsschätzung

Vor Erhebung der Daten: a-priori Dichte $\pi(\theta)$, gemeinsame Dichte von $(X, \tilde{\theta})$:

$$g(x, \theta) = f(x|\theta) \cdot \pi(\theta)$$

Nach Erhebung *einer* Beobachtung x von X erhält man die durch $X = x$ bedingte Dichte von $\tilde{\theta}$:

$$\pi(\theta|x) := \frac{g(x, \theta)}{\int_{\Theta} g(x, \theta) d\theta} = \frac{f(x|\theta) \cdot \pi(\theta)}{\int_{\Theta} f(x|\theta) \cdot \pi(\theta) d\theta}$$

Für mehrere Beobachtungen x_1, \dots, x_n von X muß man die durch $X_1 = x_1$ bis $X_n = x_n$ bedingte Dicht von $\tilde{\theta}$ berechnen:

$$\pi(\theta|x_1, \dots, x_n) = \frac{g(x_1, \dots, x_n, \theta)}{\int_{\Theta} g(x_1, \dots, x_n, \theta) d\theta}$$

Für die gemeinsame Dichte $g(x_1, \dots, x_n, \theta)$ von $(X_1, \dots, X_n, \tilde{\theta})$ gilt:

$$g(x_1, \dots, x_n, \theta) = f_n(x_1, \dots, x_n|\theta) \cdot \pi(\theta)$$

und im Fall einer vollständigen Stichprobe (unabhängig wie X verteilte X_i) gilt für die Dichte f_n von (X_1, \dots, X_n)

$$f_n(x_1, \dots, x_n|\theta) = \prod_{i=1}^n f(x_i|\theta)$$

Damit folgt für die a-posteriori Dichte

$$\begin{aligned} \pi(\theta|x_1, \dots, x_n) &= \frac{[\prod_{i=1}^n f(x_i|\theta)] \cdot \pi(\theta)}{\int_{\Theta} \underbrace{\left[\prod_{i=1}^n f(x_i|\theta) \right]}_{l(\theta; x_1, \dots, x_n)} \cdot \pi(\theta) d\theta} \\ &= \frac{l(\theta; x_1, \dots, x_n) \cdot \pi(\theta)}{\underbrace{\int_{\Theta} l(\theta; x_1, \dots, x_n) \cdot \pi(\theta) d\theta}_{c \in \mathbf{R} \text{ (konstant)}}} \\ &= c \cdot l(\theta; x_1, \dots, x_n) \cdot \pi(\theta) \\ &\propto \pi(\theta) \cdot l(\theta; x_1, \dots, x_n) \end{aligned}$$

Bemerkung: Für nicht vollständige Daten D schreibt man

$$\pi(\theta|D) \propto \pi(\theta) \cdot l(\theta; D)$$

$$D = (x_1, \dots, x_n; w_1, \dots, w_n)$$

Dabei bezeichnet w_1, \dots, w_n sog. „Withdrawings“, z.B. zurückgezogene Flugzeugmotoren.

4.5 A-posteriori Dichte für einfache Anteile

Als a-priori Verteilung bietet sich die $\text{Be}(\alpha, \beta)$ an, insbesondere $\text{Be}(1, 1) = U_{0,1}$. Für die Punktwahrscheinlichkeiten der A_θ gilt

$$p(x|\theta) = \theta^x(1-\theta)^{1-x} \quad \text{für } x \in \{0, 1\}.$$

Für n Beobachtungen sind die gemeinsamen Punktwahrscheinlichkeiten einer vollständigen Stichprobe vom Umfang n gegeben durch

$$\begin{aligned} p_n(x_1, \dots, x_n|\theta) &= \prod_{i=1}^n p(x_i|\theta) \\ &= \prod_{i=1}^n \theta^{x_i}(1-\theta)^{1-x_i} \\ &= \theta^{\sum_{i=1}^n x_i} (1-\theta)^{n-\sum_{i=1}^n x_i} \\ &= l(\theta; x_1, \dots, x_n). \end{aligned}$$

Für a-priori Verteilungen $\pi(\theta) \hat{=} \text{Be}(\alpha, \beta)$ gilt

$$\pi(\theta) \propto \theta^{\alpha-1} \cdot (1-\theta)^{\beta-1} \mathbf{I}_{(0,1)}(\theta).$$

Aus dem Bayes'schem Theorem folgt:

$$\begin{aligned} \pi(\theta|x_1, \dots, x_n) &\propto \pi(\theta) \cdot l(\theta; x_1, \dots, x_n) \\ &\propto \theta^{\alpha-1} (1-\theta)^{\beta-1} \cdot \theta^{\sum_{i=1}^n x_i} \cdot (1-\theta)^{n-\sum_{i=1}^n x_i} \\ &\propto \theta^{(\alpha+\sum_{i=1}^n x_i)-1} \cdot (1-\theta)^{(\beta+n-\sum_{i=1}^n x_i)-1} \\ &\hat{=} \text{Be}\left(\alpha + \sum_{i=1}^n x_i, \beta + n - \sum_{i=1}^n x_i\right) \end{aligned}$$

Daraus berechnet man nun Punktschätzungen für den Anteil θ (falls ein solcher existiert).

Definition: Der a-posteriori Bayes-Schätzer $\hat{\theta}$ für θ ist der Erwartungswert von $\tilde{\theta}$ a-posteriori, d.h.

$$\hat{\theta} = \mathbb{E}_{\pi(\theta|\underline{x})} \tilde{\theta} = \int_0^1 \theta \cdot \pi(\theta|x_1, \dots, x_n) d\theta.$$

Es gilt: Für $\theta \sim \text{Be}(1 + \sum_{i=1}^n x_i, 1 + n - \sum_{i=1}^n x_i)$ ist

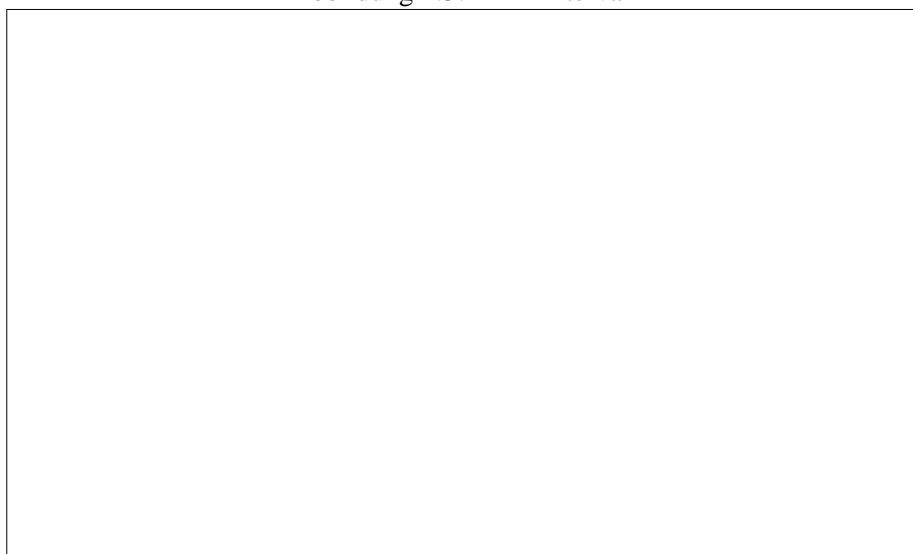
$$\mathbb{E}_{\pi(\theta|\underline{x})} \tilde{\theta} = \frac{\sum_{i=1}^n x_i + 1}{n + 2}.$$

4.6 HPD-Intervalle für θ

Gegeben eine a-posteriori Dichte $\pi(\theta|D)$, $\theta \in \Theta$ ist ein HPD-Intervall mit Sicherheit $1 - \alpha$ - bei Existenz - ein Intervall $\Theta^* \in \Theta$ mit:

- $W_{\pi(\theta|D)} \{ \tilde{\theta} \in \Theta^* \} = 1 - \alpha$
- $\pi(\theta|D) \geq c \quad \forall \theta \in \Theta^*$, wobei c die größtmögliche derartige Konstante ist

Abbildung 4.3: HPD-Intervall



Bemerkung: Analog für mehrdimensionale Parameter.

Bemerkung: Ein (einfacher) Alternativversuch kann auf zwei Arten ausgehen, aber die Wahrscheinlichkeit der zweiten Art ist durch jene der ersten Art bestimmt.

$$\theta = W \{ X = 1.\text{Art} \} \quad 1 - \theta = W \{ X = 2.\text{Art} \}$$

Auch darstellbar als Versuch, der auf zwei Arten ausgeht ist:

$$W \{ X = 1.\text{Art} \} = \theta_1 \quad W \{ X = 2.\text{Art} \} = \theta_2 \quad \text{mit } \theta_1 + \theta_2 = 1$$

Führt man jetzt diesen Versuch n -mal durch und zählt danach Y_1 der Ausgänge erster Art als auch Y_2 der Ausgänge zweiter Art, so erhält man eine 2-dimensionale stochastische Größe.

$$(Y_1, Y_2) \sim M_{n; \theta_1, \theta_2} \quad \text{Multinomialverteilung}$$

4 Bayes'sche Anteilsschätzung

Abbildung 4.4: HPD-Intervall für einen zweidimensionalen Parameter

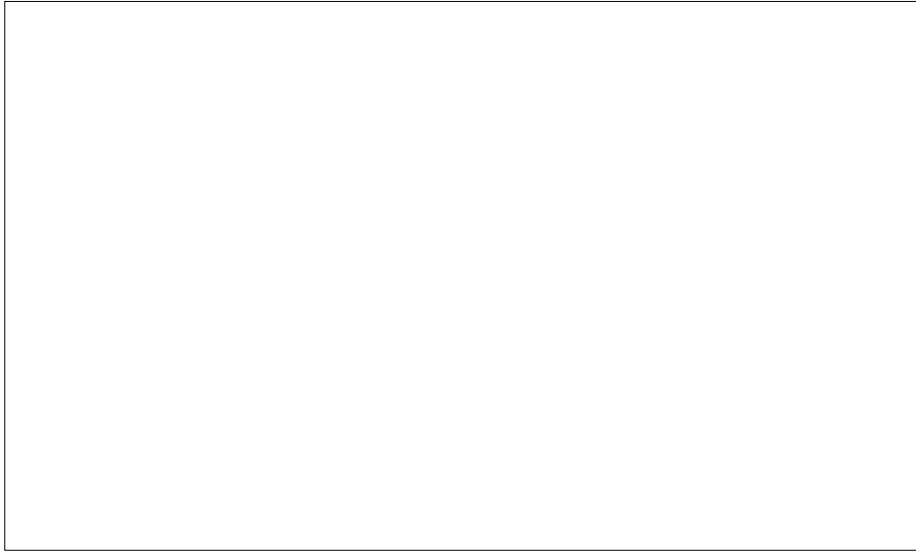
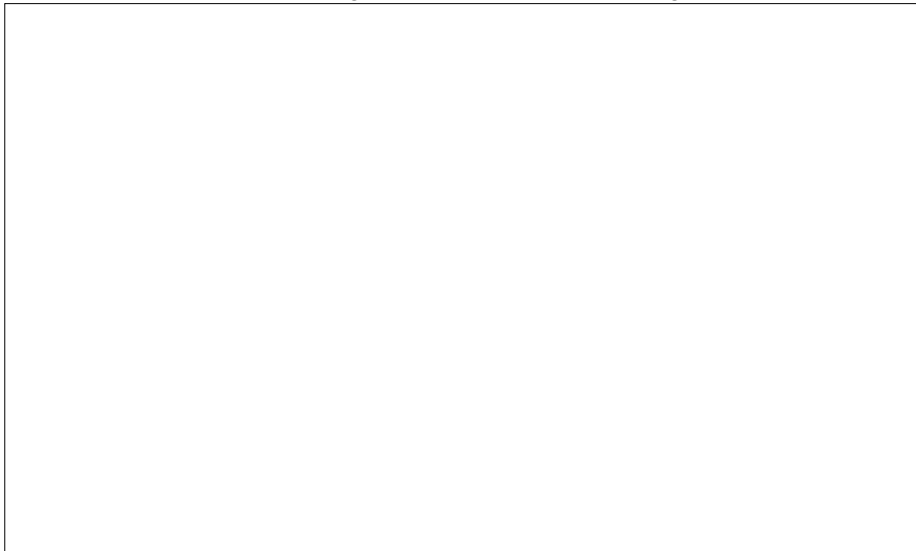


Abbildung 4.5: Multinomialverteilung



4 Bayes'sche Anteilsschätzung

Bemerkung: Der Merkmalraum von (Y_1, Y_2) liegt auf einer Geraden in \mathbf{R} (linearer Teilraum).

Es gilt:

$$\mathbf{W} \left\{ \begin{pmatrix} Y_1 \\ Y_2 \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \end{pmatrix} \right\} = \frac{n!}{y_1! \cdot y_2!} \cdot \theta_1^{y_1} \cdot \theta_2^{y_2}$$

mit $y_1 + y_2 = n$ und $\theta_1 + \theta_2 = 1$.

Bemerkung: Für den allgemeinen Fall „mehrfacher“ Alternativen (K mögliche Ausgänge eines Einzelversuches) ist die Beschreibung der a-priori Verteilung in etwas anderer Form nützlich.

Dazu betrachtet man (als 2-dimensionale Form der $\text{Be}(\alpha_1, \alpha_2)$) die sog. Dirichlet-Verteilung (2-dimensional) mit Dichte

$$\pi(\theta_1, \theta_2 | \alpha_1, \alpha_2) = \frac{\Gamma(\alpha_1 + \alpha_2)}{\Gamma(\alpha_1) \cdot \Gamma(\alpha_2)} \cdot \theta_1^{\alpha_1 - 1} \cdot \theta_2^{\alpha_2 - 1}$$

für $\theta_i \geq 0$ und $\theta_1 + \theta_2 = 1$ und $\alpha_i = 0$.

Abbildung 4.6: Dirichlet-Verteilung



Bemerkung: Im Fall einfacher Alternativen betrachtet man meist das 1-dimensionale Problem Y_1 bzw. θ_1 .

Im Fall von drei möglichen Versuchsausgängen eines Einzelversuches mit den entsprechenden Wahrscheinlichkeiten $\theta_1, \theta_2, \theta_3$ mit $\theta_1 + \theta_2 + \theta_3 = 1$ ist der Merkmalraum von $(\tilde{\theta}_1, \tilde{\theta}_2, \tilde{\theta}_3)$ eine Ebene des \mathbf{R}^3 .

Abbildung 4.7: Drei mögliche Versuchsausgänge eines Einzelversuches



4.7 Bayes'sche Anteilsschätzung bei mehrfachen Alternativen

Ein Einzelversuch kann auf K Arten ausgehen mit $\mathbf{W} = \theta_j$, $k = 2$, $\sum_{j=1}^k \theta_j = 1$. Bei n -facher Durchführung eines solchen Versuches erhält man X_1, \dots, X_n . Danach stellt man fest, wie oft die k verschiedenen Arten beobachtet wurden:

$$\begin{pmatrix} Y_1 \\ \vdots \\ Y_k \end{pmatrix} \quad Y_j = \text{Anzahl der Ausgänge auf die } j\text{-te Art}$$

$$\begin{pmatrix} Y_1 \\ \vdots \\ Y_k \end{pmatrix} \sim \mathbf{M}_{n; \theta_1, \dots, \theta_k} \quad \text{Multinomialverteilung}$$

$$\mathbb{P} \left\{ \begin{pmatrix} Y_1 \\ \vdots \\ Y_k \end{pmatrix} = \begin{pmatrix} y_1 \\ \vdots \\ y_k \end{pmatrix} \right\} = \frac{n!}{y_1! \cdot \dots \cdot y_k!} \cdot \theta_1^{y_1} \cdot \dots \cdot \theta_k^{y_k} \quad \text{für } \begin{cases} \sum_{j=1}^k y_j = n \\ y_j \in \mathbf{N}_0 \end{cases}$$

Bemerkung: Der Merkmalraum von diesem stochastischen Vektor $(Y_1, \dots, Y_k)^T$ ist ein Teilraum einer $(k - 1)$ -dimensionalen Hyperebene des \mathbf{R}^k .

4 Bayes'sche Anteilsschätzung

Abbildung 4.8: Spezialfall k=3 (2-dimensionale diskrete Verteilung)



Spezialfall k=3: Wegen $y_3 = n - y_1 - y_2$ kann dies als zweidimensionales Problem behandelt werden.

Als a-priori Verteilung für $\tilde{\theta} = (\tilde{\theta}_1, \dots, \tilde{\theta}_k)$ verwendet man die k -dimensionale Dirichlet-Verteilung $\text{Dir}(\alpha_1, \dots, \alpha_k)$.

$$\pi(\theta_1, \dots, \theta_k | \alpha_1, \dots, \alpha_k) = \frac{\Gamma\left(\sum_{j=1}^k \alpha_j\right)}{\prod_{j=1}^k \Gamma(\alpha_j)} \cdot \prod_{j=1}^k \theta_j^{\alpha_j - 1}$$

$$\begin{cases} \alpha_j \geq 0 \\ \theta_j \geq 0 \\ \sum_{j=1}^k \theta_j = 1 \end{cases}$$

Der Merkmalraum von $(\tilde{\theta}_1, \dots, \tilde{\theta}_k)$ ist eine Teilmenge einer $(k-1)$ -dimensionalen Hyperebene des \mathbf{R}^k . Die Parameter $\alpha_1, \dots, \alpha_k$ der a-priori Verteilung heißen Hyperparameter.

Die k -dimensionale uniforme Verteilung auf der $(k-1)$ -dimensionalen Hyperebene $\sum_{j=1}^k \theta_j = 1$ von $\Theta \subseteq \mathbf{R}^k$ enthält man mit $\alpha_1 = \alpha_2 = \dots = \alpha_k = 1$.

Spezialfall k=3:

$$\pi(\theta_1, \theta_2, \theta_3 | \alpha_1, \alpha_2, \alpha_3) > 0 \quad \text{nur fuer} \quad \theta_1 + \theta_2 + \theta_3 = 1$$

Für die Aktualisierung der Verteilung $\pi(\theta_1, \dots, \theta_k | \alpha_1, \dots, \alpha_k)$ nach Beobachtung einer Stichprobe x_1, \dots, x_n und den entsprechenden Werten y_1, \dots, y_k verwendet man das Bayes'sche Theorem:

$$\pi(\theta_1, \dots, \theta_k | \alpha_1, \dots, \alpha_k) \propto \prod_{j=1}^k \theta_j^{\alpha_j - 1}$$

4 Bayes'sche Anteilsschätzung

$$\begin{aligned}l(\theta_1, \dots, \theta_k; y_1, \dots, y_k) &\propto \prod_{j=1}^k \theta_j^{y_j} \\ \pi(\theta_1, \dots, \theta_k | D) &\propto \prod_{j=1}^k \theta_j^{\alpha_j - 1} \cdot \prod_{j=1}^k \theta_j^{y_j} \\ &\propto \prod_{j=1}^k \theta_j^{\alpha_j + y_j - 1} \\ \pi(\theta_1, \dots, \theta_k | D) &\hat{=} \text{Dir}(\alpha_1 + y_1, \dots, \alpha_k + y_k)\end{aligned}$$

Spezialfall k=3:

$$\pi(\theta_1, \theta_2, \theta_3 | D) \hat{=} \text{Dir}(\alpha_1 + y_1, \alpha_2 + y_2, \alpha_3 + y_3)$$

Daraus ermittelt man eine Verteilung für $(\tilde{\theta}_1, \tilde{\theta}_2)$ durch $\theta_3 = 1 - \theta_1 - \theta_2$.

Beispiel: Meinungsumfrage Hainburg

- a-priori von $\theta_1 = \text{pro}$, $\theta_2 = \text{contra}$, $\theta_3 = \text{unentschlossen}$
- a-posteriori Dichte
- HPD-Bereiche
- $W\{\text{pro} > \text{contra}\}$

4.8 Konjugierte Verteilungsfamilien

Anteilsschätzung:

- a-priori $\text{Be}(\alpha, \beta)$
- a-posteriori $\text{Be}(\alpha + \sum_{i=1}^n x_i, \beta + n - \sum_{i=1}^n x_i)$

Analog bei mehrfachen Alternativen:

- a-priori $\text{Dir}(\alpha_1, \dots, \alpha_k)$
- a-posteriori $\text{Dir}(\alpha_1 + y_1, \dots, \alpha_k + y_k)$

Allgemeiner für die stochastischen Modelle und zugehörigen a-priori Familien.

Definition: Ist $X \sim W_\theta$, $\theta \in \Theta$ ein stochastisches Modell und \mathcal{P} eine Familie von a-priori Verteilungen für $\tilde{\theta}$, sodaß für beliebige Stichproben von X die a-posteriori Verteilung von $\tilde{\theta}$ wieder ein Element von \mathcal{P} ist, so spricht man von konjugierten Verteilungsfamilien.

4 Bayes'sche Anteilsschätzung

Beispiel: $X \sim \text{Ex}_\tau, \tau > 0; \mathcal{P} = \{\text{Gam}(\alpha, \beta) : \alpha > 0, \beta > 0\}$

Beispiel: $X \sim P_\mu, \mu > 0; \mathcal{P} = \{\text{Gam}(\alpha, \beta) : \alpha > 0, \beta > 0\}$

Bemerkung: Man spricht auch von „einer zu einem stochastischen Modell konjugierten a-priori Familie“.

4.9 Suffizienz

Eine aus einer Stichprobe X_1, \dots, X_n berechnete Statistik $s(X_1, \dots, X_n)$ heißt **suffizient** für einen Parameter θ der Verteilung von X , falls zur Berechnung der a-posteriori Verteilung $\pi(\theta|x_1, \dots, x_n)$ die Kenntnis von $s(x_1, \dots, x_n)$ hinreicht.

Beispiel: $\sum_{i=1}^n X_i$ für die Alternativ-Verteilung

Beispiel: $(Y_1, \dots, Y_k), M_{n;\theta_1, \dots, \theta_k}$

$$\begin{aligned} W \left\{ \left(\begin{array}{c} Y_1 \\ \vdots \\ Y_k \end{array} \right) = \left(\begin{array}{c} y_1 \\ \vdots \\ y_k \end{array} \right) \middle| n, \theta_1, \dots, \theta_k \right\} &= p(y_1, \dots, y_k | \theta_1, \dots, \theta_k) = \\ &= \frac{n!}{y_1! \cdot \dots \cdot y_k!} \cdot \theta_1^{y_1} \cdot \dots \cdot \theta_k^{y_k} \cdot \mathbf{I}_{\{\sum_{j=1}^k y_j = n\} \cap \{y_j \in \mathbf{N}_0\}}(y_1, \dots, y_k) \end{aligned}$$

Dazu konjugiert die Dirichlet-Verteilung

$$\begin{aligned} f(x_1, \dots, x_k | \alpha_1, \dots, \alpha_k) &= \frac{\Gamma(\alpha_1 + \dots + \alpha_k)}{\Gamma(\alpha_1) \cdot \dots \cdot \Gamma(\alpha_k)} \cdot x_1^{\alpha_1-1} \cdot \dots \cdot x_k^{\alpha_k-1} \cdot \mathbf{I}_{D^*}(x_1, \dots, x_k) \\ D^* &= \left\{ (x_1, \dots, x_k) : x_i \geq 0 \wedge \sum_{i=1}^k x_i = 1 \right\} \end{aligned}$$

$$\pi(\theta_1, \dots, \theta_k | y_1, \dots, y_k) = \text{Dir}(\alpha_1 + y_1, \dots, \alpha_k + y_k)$$

In diesem Beispiel ist (Y_1, \dots, Y_k) eine **suffiziente Statistik** für $\underline{\theta} = (\theta_1, \dots, \theta_k)$.

Bemerkung: Auch (Y_1, \dots, Y_{k-1}) ist **suffizient**.

5 Statistische Qualitätskontrolle

Entscheidung, ob eine Warensendung (Los) genügend gute Stücke enthält.

N	Losumfang
n	Stichprobenumfang

Entscheidungsfindung:

- Risiko für Lieferanten und Käufer zumutbar
- Laufende Überwachung eines Produktionsvorganges

5.1 Einfache Stichprobenpläne

θ	Anteil der schlechten Stücke (bzw. A =Anzahl der schlechten Stücke)
x	Anzahl der schlechten Stücke in der Stichprobe
X	stochastische Größe, die x vor der Erhebung der Stichprobe beschreibt
(n, c)	c =Annahmekennzahl (wieviel schlechte Stücke maximal)

Entscheidung: Nehme Los an, falls $x \leq c$.

Bemerkung: c so zu wählen, dass gute Lose mit großer Wahrscheinlichkeit angenommen werden und schlechte Lose mit sehr kleiner Wahrscheinlichkeit angenommen werden. Dazu dient die sog. „Operationscharakteristik“ OC; eine Funktion von θ (auch Kennkurve).

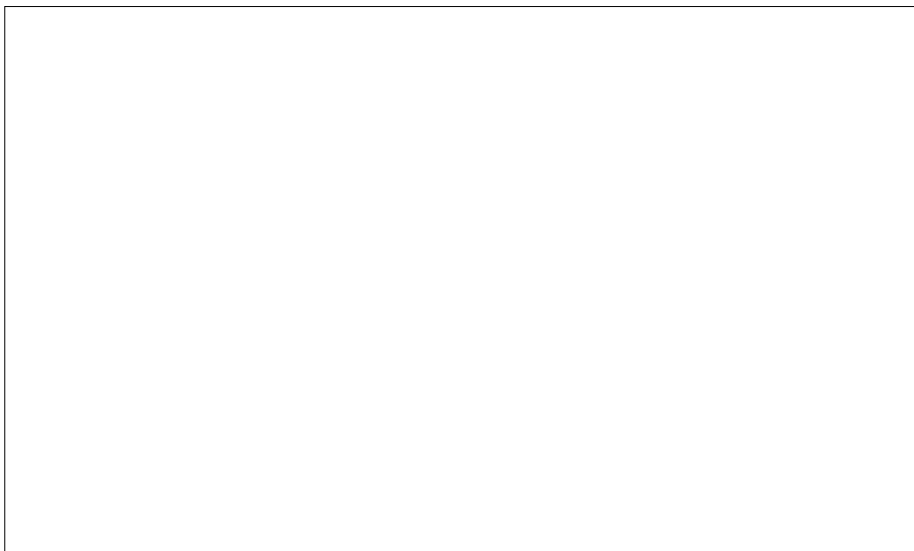
$$\begin{aligned} W(\theta) &= W \{ \text{Los angenommen} \mid \text{Schlechtanteil} = \theta \} \\ &= W \{ X \leq c \mid \text{Schlechtanteil} = \theta \} \end{aligned}$$

Abbildung 5.1: Operationscharakteristik



Eigenschaften der Operationscharakteristik

In konkreten Produktionen ist oft ein bestimmter (kleiner) Schlechtanteil p zulässig. Lose mit Schlechtanteil $\theta \leq p$ sollten angenommen werden, solche mit $\theta > p$ sollten abgelehnt werden. Eine ideale OC wäre daher:



Dies ist mit konkreten Stichproben im allgemeinen nicht erreichbar. Die OC kann folgenderma-

5 Statistische Qualitätskontrolle

ßen berechnet werden:

$$W(\theta) = W\{X \leq c | \theta\} = \sum_{a=0}^c W\{X = a | \theta\}$$

Bei klassischen Stichproben gilt:

$$X \sim H_{N,A,n} \quad \text{mit } \theta = \frac{A}{N}$$

und

$$W\{X = a\} = \frac{\binom{A}{a} \binom{N-A}{n-a}}{\binom{N}{n}} \quad \text{für } a_1 \leq a \leq a_2$$

$$\text{mit } a_1 = \max(0, n - (N - A))$$

$$a_2 = \min(n, A)$$

Wenn der Losumfang N groß ist ($n < \frac{N}{10}$) kann für praktische Zwecke die $H_{N,A,n}$ durch die $B_{n, \frac{A}{N}}$ approximiert werden, d.h.

$$W\{X = a\} \approx \binom{n}{a} \cdot \theta^a \cdot (1 - \theta)^{n-a} \quad \text{für } a = 0(1)n.$$

Wenn n relativ groß ist und obiges gilt, sowie $n \cdot \theta \leq 5$ so kann man die Binomialverteilung $B_{n,\theta}$ durch die Poisson-Verteilung $P_{n \cdot \theta}$ approximieren:

$$W\{X = a\} \approx \frac{(n \cdot \theta)^a \cdot e^{-n\theta}}{a!}$$

Bemerkung: Die OC ist abhängig von n , c und θ .

Beispiel: Große Warensendung von Schrauben, einfacher Stichprobenplan $n = 80$, $c = 2$.

1. Falls der Schlechtanteil $p=0.01$ ist, untersuche man ob dieser Stichprobenplan Lose mit guter Qualität, d.h. $\theta \leq p$, mit hoher Wahrscheinlichkeit annimmt.

Hier ist die Poisson-Approximation zulässig, da $n < \frac{N}{10}$ und $n \cdot \theta = 0.8 < 5$.

$$\begin{aligned} W(\theta) &= W\{X \leq 2\} \approx \sum_{a=0}^2 \frac{(80 \cdot \theta)^a}{a!} \cdot e^{-80 \cdot \theta} \\ &= e^{-80 \cdot \theta} \left[1 + 80 \cdot \theta + \frac{(80 \cdot \theta)^2}{2!} \right] \end{aligned}$$

Für $\theta = p = 0.01$ ergibt sich $W(0.01) = 0.953$. Solche Lose werden also mit hoher Wahrscheinlichkeit angenommen.

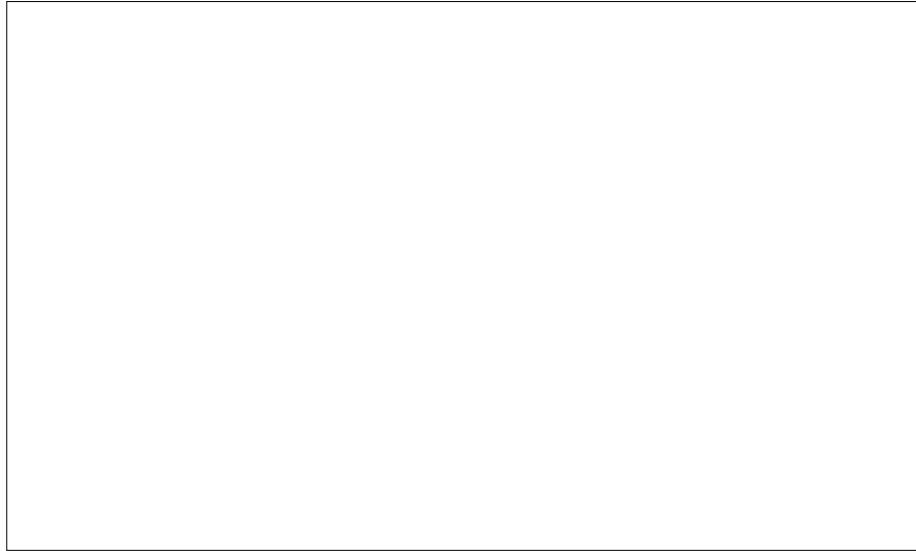
5 Statistische Qualitätskontrolle

2. Wie groß ist die Wahrscheinlichkeit der Annahme des Loses, falls $\theta = 0.05$ ist ?

Für $\theta = 0.05$ ergibt sich $W(0.05) = 0.238$, d.h. ein solcher Los wird nur mit einer Wahrscheinlichkeit von 23.8% angenommen.

Übung: Im obigen Beispiel berechne man $W(\theta)$ für $\theta = 0(0.01)0.10$ und zeichne die OC. Für größerer Werte von θ sollten die Binomial-Wahrscheinlichkeiten berechnet werden ($80 \cdot 0.1 = 8 > 5$).

Abbildung 5.2: Operationscharakteristik für $W(\theta)$ mit $\theta = 0(0.01)0.10$



5.2 Festlegung einfacher Stichprobenpläne

Man gibt zwei Punkte der OC vor und zwar für p_1 (gute Lose, $\theta \leq p_1$)

$$WS(p_1) = 1 - \alpha \quad (5.1)$$

und p_2 (schlechte Lose, $\theta \geq p_2$)

$$WS(p_2) = \beta \quad (5.2)$$

wobei α und β kleine Zahlen sind.

- Produzentenrisiko α : der Produzent möchte das Lose von hoher Qualität ($\theta \leq p_1$) mit hoher Wahrscheinlichkeit $WS(1 - \alpha)$ angenommen werden.

Bemerkung: Beziehung zur Irrtumswahrscheinlichkeit α bei statistischen Tests (Fehler 1. Art)

5 Statistische Qualitätskontrolle

- Konsumentenrisiko β : der Käufer möchte das Lose von schlechter Qualität ($\theta \geq p_2$) mit kleiner Wahrscheinlichkeit $WS(\beta)$ angenommen werden.

Bemerkung: Fehler 2.Art, Schärfe $1 - \beta$ bei statistischen Tests

Die Werte für α und β sind vom Losumfang N abhängig; je größer das N ist, desto wichtiger ist es keine falsche Entscheidung zu treffen. Da n und $c \in \mathbf{N}_0$, können i.a. die Gleichungen 5.1 und 5.2 nicht exakt erfüllt sein. Man wählt für n und c solche Werte, dass die OC möglichst nahe an den Punkten $(p_1, 1 - \alpha)$ und $(p_2, 1 - \beta)$ liegt und ferner gilt:

$$W(p_1) \geq 1 - \alpha \quad \text{und} \quad W(p_2) \leq \beta.$$

Zur Bestimmung von n und c ist folgendes hilfreich:

$$W(\theta) = W\{Y > 2 \cdot n \cdot \theta\} \quad \text{wobei} \quad Y \sim \chi_{2(c+1)}^2$$

Beweis: Durch sukzessive partielle Integration sieht man die Gültigkeit von

$$\int_{2n\theta}^{\infty} \frac{y^c e^{-y/2}}{c! \cdot 2^{c+1}} dy = \sum_{k=0}^c \frac{(n\theta)^k}{k!} \cdot e^{-n\theta} = W(\theta)$$

$$\frac{y^c e^{-y/2}}{c! \cdot 2^{c+1}} = \frac{y^{\frac{2(c+1)}{2}-1} \cdot e^{-y/2}}{\underbrace{\Gamma(c+1) \cdot 2^{\frac{2(c+1)}{2}}}_{c!}}$$

Dichte von $\chi_{2(c+1)}^2$

Zur Bestimmung von n und c hat man zwei Gleichungen:

$$1 - \alpha \leq W(p_1) = W\left\{Y \geq 2np_1 \mid Y \sim \chi_{2(c+1)}^2\right\} \quad (5.3)$$

$$\beta \geq W(p_2) = W\left\{Y \geq 2np_2 \mid Y \sim \chi_{2(c+1)}^2\right\} \quad (5.4)$$

In einer Tabelle der χ^2 -Verteilung sucht man Fraktile:

$$\chi_{\bullet, \alpha}^2 \quad \text{und} \quad \chi_{\bullet, 1-\beta}^2$$

sodaß gilt:

$$2np_1 = \chi_{2(c+1); \alpha}^2$$

$$2np_2 = \chi_{2(c+1); 1-\beta}^2$$

Dies ist dann der Fall, wenn

$$\frac{\chi_{\bullet, 1-\beta}^2}{\chi_{\bullet, \alpha}^2} = \frac{p_2}{p_1}$$

Man sucht jene Anzahl von Freiheitsgraden, sodaß der Quotient der Fraktilen möglichst nahe bei $\frac{p_2}{p_1}$ liegt (es kommen nur geraden Anzahlen in Betracht).

Zur Bestimmung von n müssen die Gleichungen 5.3 und 5.4 erfüllt sein. Daraus erhält man für n zwei Ungleichungen. Die kleinste ganze Zahl, die beide Ungleichungen erfüllt, ist der gesuchte Stichprobenumfang.

Bemerkung: Für die praktische Anwendung gibt es Tabellen und sog. „Normogramme“.

5.3 Zweifache Stichprobenpläne

Hier wird eine Stichprobe vom Umfang n_1 gezogen (x_1 =Anzahl der schlechten Stücke in dieser Stichprobe). Falls:

$$\begin{array}{ll} x_1 \leq c_1 & \text{Annahme des Loses} \\ x_1 \geq d_1 & \text{Ablehnung} \\ c_1 < x_1 < d_1 & \text{Zweite Stichprobe mit Umfang } n_2 \end{array}$$

Dabei bezeichnet n_2 die Anzahl der schlechten Stücke in der zweiten Stichprobe. Falls:

$$\begin{array}{ll} x_1 + x_2 \leq c_2 & \text{Annahme des Loses} \\ \text{andererseits} & \text{Ablehnung} \end{array}$$

Man erhält einen zweifachen Stichprobenplan $(n_1, c_1, d_1, n_2, c_2)$.

Bemerkung: Die Berechnung mehrfacher Stichprobenpläne ist kompliziert. Es gibt aber Tabellen. Der Vorteil ist, dass der zu erwartende Stichprobenumfang bei praktisch gleicher OC kleiner ist.

5.4 Sequentielle Stichprobenpläne

Man entnimmt hintereinander einzelne Elemente des Loses. Dabei bezeichnet x_n die Anzahl der schlechten Stücke, nachdem man insgesamt n Stücke geprüft hat.

Man benötigt einen sog. Prüfplan (a, b, c)

$$\begin{array}{ll} x_n \leq c \cdot n - a & \text{Annahme} \\ x_n \geq c \cdot n + b & \text{Ablehnung} \\ c \cdot n - a < x_n < c \cdot n + b & \text{Weiterprüfen} \end{array}$$

Die Berechnung von a, b, c erfolgt aus den vorgegebenen Werte p_1, p_2, α, β der OC mit Hilfe des sog. approximativen Sequentialtest von A. Wald:

Abbildung 5.3: Sequentieller Stichprobenplan



$$A = \ln \frac{1-\alpha}{\beta} \quad B = \ln \frac{1-\beta}{\alpha} \quad P = \ln \frac{p_2}{p_1} \quad Q = \ln \frac{1-p_1}{1-p_2}$$

$$a = \frac{A}{P+Q} \quad b = \frac{B}{P+Q} \quad c = \frac{Q}{P+Q}$$

5.5 Kontrollkarten

Zur laufenden Überwachung eines Produktionsprozesses. Dazu wird eine Größe x gemessen, die in Abhängigkeit der Zeit variiert. Dabei nimmt man meist an, dass diese eine Realisierung einer stochastischen Größe X ist, die (approximativ) normalverteilt ist.

Dabei sind gewisse Schwankungen von X durchaus zulässig, solange sie in einem gewissen Bereich (Intervall) bleiben.

Beispiel:

1. Druckfestigkeit von Beton
2. Feuchtigkeit von Bauholz
3. Durchmesser von Schrauben

Qualitätsänderungen entsprechen Änderungen von zumindest einem Parameter. Hier kann sich eine kleine Streuung (Varianz) insofern günstig auswirken, da man in diesem Fall mit dem Erwartungswert μ_0 der Produktion näher an die Qualitätsgrenze herangehen kann.

Es gibt verschiedene Arten von Kontrollkarten, je nachdem was aufgetragen wird:

- Einzelmeßwerte: x -Karten
- Mittelwerte: \bar{x} -Karten
- Streuungen: s -Karten

Damit können Schwankungen im Produktionsprozeß erkannt werden.

Bezeichnungen:

ML	Mittellinie	
T_o	obere Toleranzgrenze	
T_u	untere Toleranzgrenze	
OKG	obere Kontrollgrenze	} Eingriff notwendig
UKG	untere Kontrollgrenze	
OWG	obere Warngrenze	
UWG	untere Warngrenze	

Ur- oder Einzelwertkarte

In regelmäßigen Abständen werden kleine Stichproben vom Umfang n (etwa 5) gezogen und die Messwerte in der Kontrollkarte eingetragen. ML , KG_n und WG_n werden aus einem Verlauf von k (20-30) Stichproben aus je m Einzelwerten bestimmt.

x_{j1}, \dots, x_{jm}	Einzelwerte der j – Stichprobe	
\bar{x}_j	Stichprobenmittel	} der j – ten Stichprobe
s_j	Stichprobenstreuung	
R_j	Spannweite $x_{max} - x_{min}$	
$\bar{\bar{x}} = \frac{\bar{x}_1 + \dots + \bar{x}_k}{k}$	Gesamtmittel	
$\bar{R} = \frac{R_1 + \dots + R_k}{k}$	Mittelwert der Spannweiten	
$\bar{s} = \frac{s_1 + \dots + s_k}{k}$	Mittelwert der Stichprobenstreuungen	

Die verschiedenen Grenzen ergeben sich aus den Bedingungen für die Messgrößen zum Zeitpunkt t :

$$W\{G_u \leq X_1(t), \dots, X_n(t) \leq G_o\} \geq 1 - \alpha$$

Bei einseitigen Fragestellungen (z.B. Mindestgewicht) gilt:

$$W\{G_u \leq X_1(t), \dots, X_n(t)\} \geq 1 - \alpha$$

Dabei bezeichnet α die Verwerfungswahrscheinlichkeit, falls der Prozess unter Kontrolle ist.

Bemerkung: Für $\alpha = 0.01$ Kontrollgrenzen, für $\alpha = 0.05$ Warngrenzen.

Bestimmung der Grenzen:

$$\begin{aligned}
 W\{G_u \leq X_1(t), \dots, X_n(t) \leq G_o\} &= W\left(\bigcap_{i=1}^n \{G_u \leq X_i(t) \leq G_o\}\right) \\
 &= \prod_{i=1}^n W\{G_u \leq X_i(t) \leq G_o\} \\
 &= (W\{G_u \leq X(t) \leq G_o\})^n = 1 - \alpha \\
 W\{G_u \leq X \leq G_o\} &= \sqrt[n]{1 - \alpha} \\
 G_u &= \frac{1 - \sqrt[n]{1 - \alpha}}{2} - \text{Fraktilen der Vtlg. von } X \\
 G_o &= 1 - \frac{1 - \sqrt[n]{1 - \alpha}}{2} - \text{Fraktilen der Vtlg. von } X
 \end{aligned}$$

Diese Fraktilwerte können wegen der Normalverteilungsannahme leicht berechnet werden. Praktisch werden die Grenzen mittels Tabellen ermittelt.

\bar{x} -Karte

Wenn der Prozess unter Kontrolle ist, gilt

$$\bar{X}_n \sim N\left(\mu_0, \frac{\sigma_0^2}{n}\right), \sigma^2 = \frac{\sigma_0^2}{n}.$$

Damit kann man Warn- und Kontrollgrenzen ermitteln.

Die KG_n werden entweder mittels der Sicherheitswahrscheinlichkeit α (meist 0.01) oder mittels der 3σ -Regel festgelegt

$$\bar{\bar{x}} \pm \frac{3\sigma}{\sqrt{n}}.$$

Für zweiseitige Kontrollkarten gilt:

$$\begin{aligned}
 OKG &:= \inf_y \{W\{\bar{X}_n > y\}\} \leq \frac{\alpha}{2} \\
 UKG &:= \sup_y \{W\{\bar{X}_n < y\}\} \leq \frac{\alpha}{2}
 \end{aligned}$$

Bei Normalverteilungen verwendet man oft

$$OKG - ML = ML - UKG = 3\sqrt{\text{Var}\bar{X}_n}.$$

Dabei wird $\text{Var}\bar{X}_n$ geschätzt bzw. wird eine Schätzung \bar{s} für $\sqrt{\text{Var}\bar{X}_n}$ benutzt.

5 Statistische Qualitätskontrolle

ML	WG	KG	KG
	0.05	$\alpha = 0.01$	3σ
$\bar{\bar{x}}$	$\bar{\bar{x}} \pm \frac{1.96}{\sqrt{n}}\sigma$	$\bar{\bar{x}} \pm \frac{2.58}{\sqrt{n}}\sigma$	$\bar{\bar{x}} \pm \frac{3}{\sqrt{n}}\sigma$

Begründung:

$$\begin{aligned}
 \bar{X}_n \sim N(\mu, \sigma^2) &\rightarrow \mathbb{W}\{|\bar{X}_n - \mu| \leq G\} = 1 - \alpha \\
 \mathbb{W}\{|\bar{X}_n - \mu| \leq G\} &= \mathbb{W}\{\mu - G \leq \bar{X}_n \leq \mu + G\} \\
 &= \mathbb{W}\left\{\frac{\mu - G - \mu}{\sigma} \leq \underbrace{\frac{\bar{X}_n - \mu}{\sigma}}_{\sim N(0,1)} \leq \frac{\mu + G - \mu}{\sigma}\right\} \\
 &= \Phi\left(\frac{G}{\sigma}\right) - \Phi\left(-\frac{G}{\sigma}\right) \\
 &= 2\Phi\left(\frac{G}{\sigma}\right) - 1 \\
 &= 1 - \alpha \\
 \Phi\left(\frac{G}{\sigma}\right) &= 1 - \frac{\alpha}{2} \\
 \frac{G}{\sigma} &= z_{1-\frac{\alpha}{2}} \quad \text{Fraktile der } N(0, 1)
 \end{aligned}$$

Für $\alpha = 0.05$ gilt $z_{1-\frac{\alpha}{2}} = z_{0.975} = 1.96$.

Für $\alpha = 0.01$ gilt $z_{1-\frac{\alpha}{2}} = z_{0.995} = 2.58$.

s-Karte

Zur Überwachung der Streuung (Varianz) von X gilt bei Normalverteilung

$$\frac{(n-1)S_n^2}{\sigma^2} \sim \chi_{n-1}^2$$

$$\begin{aligned}
 \mathbb{W}\left\{\chi_{n-1; \frac{\alpha}{2}}^2 \leq \frac{(n-1)S_n^2}{\sigma^2} \leq \chi_{n-1; 1-\frac{\alpha}{2}}^2\right\} &= 1 - \alpha \\
 \mathbb{W}\left\{\frac{\sigma^2}{n-1}\chi_{n-1; \frac{\alpha}{2}}^2 \leq S_n^2 \leq \frac{\sigma^2}{n-1}\chi_{n-1; 1-\frac{\alpha}{2}}^2\right\} &= 1 - \alpha \\
 \mathbb{W}\left\{\sigma\sqrt{\frac{\chi_{n-1; \frac{\alpha}{2}}^2}{n-1}} \leq S_n \leq \sigma\sqrt{\frac{\chi_{n-1; 1-\frac{\alpha}{2}}^2}{n-1}}\right\} &= 1 - \alpha
 \end{aligned}$$

5 Statistische Qualitätskontrolle

Für die Grenzen ergibt sich somit

$$\begin{aligned} \text{obere} & \quad \sigma \sqrt{\frac{\chi_{n-1; 1-\frac{\alpha}{2}}^2}{n-1}} \\ \text{untere} & \quad \sigma \sqrt{\frac{\chi_{n-1; \frac{\alpha}{2}}^2}{n-1}} \end{aligned}$$

Da σ mittels \bar{s} geschätzt wird, erhält man folgende Grenzen:

- Warngrenzen $\alpha = 0.05$

$$\frac{\bar{s}}{c_n} \sqrt{\frac{\chi_{n-1; 0.025}^2}{n-1}} \quad \frac{\bar{s}}{c_n} \sqrt{\frac{\chi_{n-1; 0.975}^2}{n-1}}$$

- Kontrollgrenzen

$$\frac{\bar{s}}{c_n} \sqrt{\frac{\chi_{n-1; 0.005}^2}{n-1}} \quad \frac{\bar{s}}{c_n} \sqrt{\frac{\chi_{n-1; 0.995}^2}{n-1}}$$

Bemerkung: Es werden verschiedene Typen von QR-Karten (QS-Karten) verwendet (\tilde{x} -Karten, R -Karten, Summenkarten CUSUM).

6 Regressionsanalyse

Anwendung zur Beschreibung (kausaler) aber nicht deterministischer Zusammenhänge.

- Zwischen stochastischen Größen
- Zwischen Einstellgrößen (unabhängige Größen) und stochastischen Größen

Beispiel: Ist (X, Y) ein 2-dimensionaler Vektor der Körpergröße und -gewicht beschreibt.

Beispiel: Versuch, dessen Ausgang von einer deterministischen einstellbaren Größe abhängt.

Einstellgröße x : Abhängige stochastische Größe Y_x .

6.1 Regression 1.Art

$$\begin{aligned}\mathbb{E}[Y|X = x] &= y(x) \\ y(x) &= \sum_{j=0}^k \theta_j \cdot x^j \quad \theta_j \text{ Parameter}\end{aligned}$$

Dies ist eine sog. „polynomische Regression der Ordnung k “.

Sonderfall $k=1$:

$$\mathbb{E}[Y|X = x] = \theta_0 + \theta_1 \cdot x \dots \text{Regressionsgeraden}$$

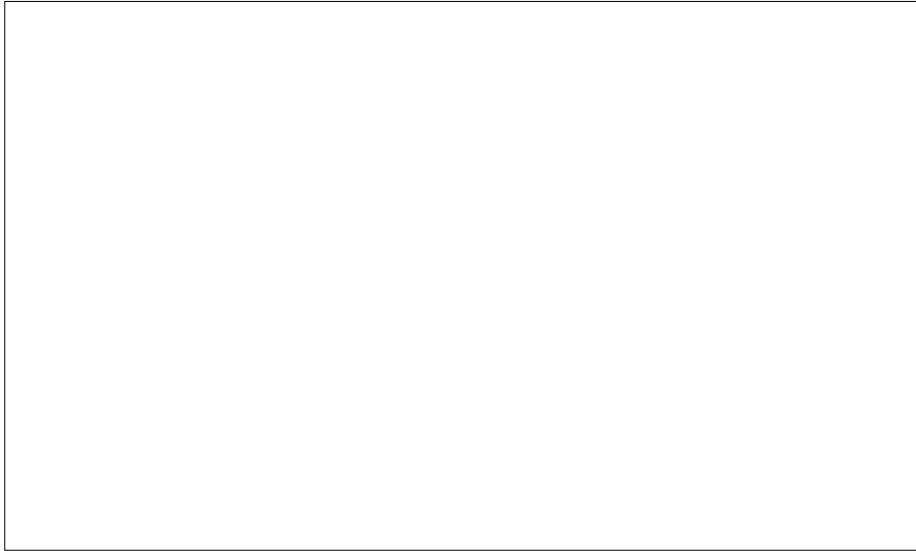
Satz 6.1

$$\begin{aligned}(X, Y) &\sim N(\mu_1, \mu_2, \sigma_1^2, \sigma_2^2) \\ \mathbb{E}[Y|X = x] &= \mu_2 + \rho \frac{\sigma_2}{\sigma_1} (x - \mu_1) \\ \mathbb{E}[X|Y = y] &= \mu_1 + \rho \frac{\sigma_1}{\sigma_2} (y - \mu_2)\end{aligned}$$

d.h. bei 2-dimensionalen Normalverteilungen hat man immer Regressionsgeraden.

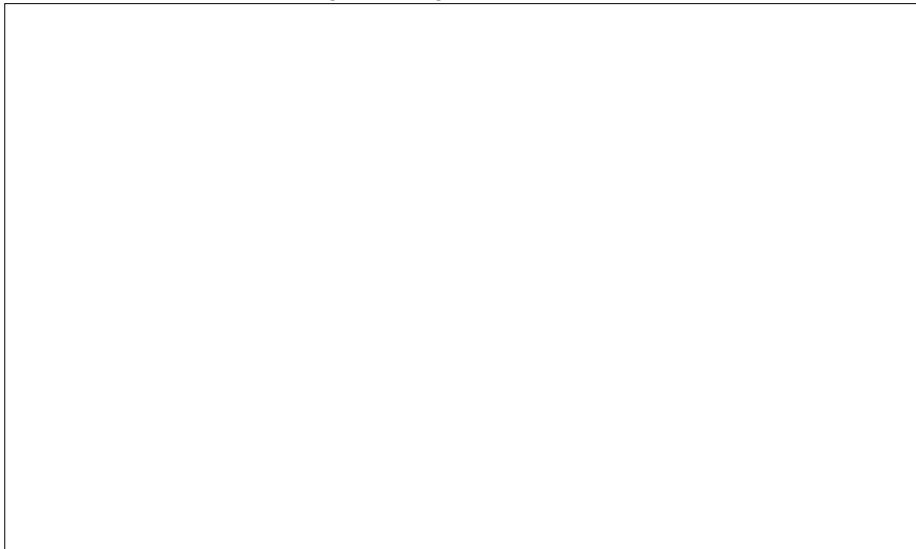
6 Regressionsanalyse

Abbildung 6.1: Regressionsfunktion 2.Art



$$\begin{aligned}\psi(x) &= \mathbb{E}Y_x \\ Y_x &= \psi(x) + \underbrace{U_x}_{\mathbb{E}U_x=0}\end{aligned}$$

Abbildung 6.2: Regressionsfunktion 1.Art



Bemerkung:

- $\tan(\phi) = \frac{\sigma_1\sigma_2}{\sigma_1^2 + \sigma_2^2} \cdot \frac{1-\rho^2}{\rho}$



- Falls X unabhängig von Y ist, sind die Regressionslinien parallel zu den Regressionsachsen.

Aus einer Stichprobe (X_i, Y_i) , $i = 1(1)n$ kann die Regressionsfunktion $X \rightarrow \mathbb{E}[Y|X = x]$ geschätzt werden (vgl. Abbildung 6.3).

Definition: Bedingte Mittelwerte $g(x_i)$

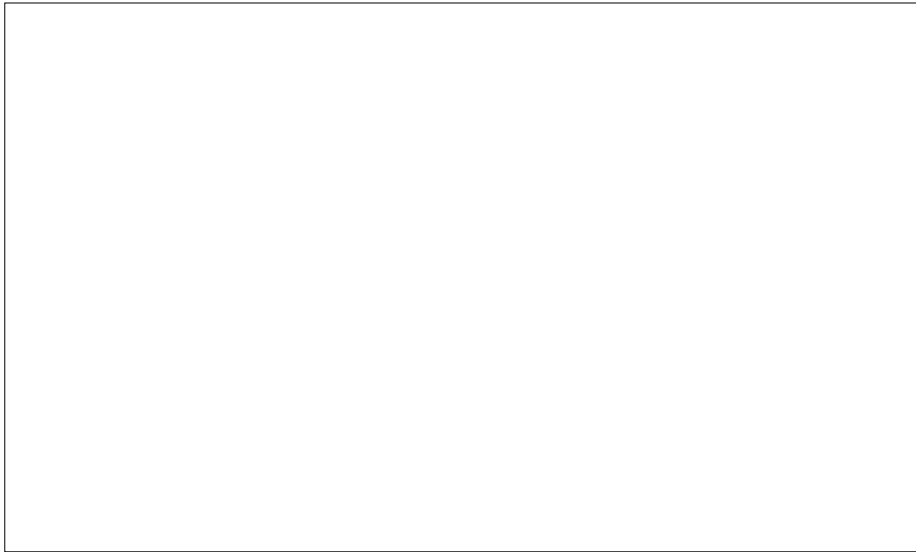
$$g(x_i) = \bar{y}_j = \frac{1}{n_j} \cdot \sum_{k=1}^{n_j} y_{jk} = \frac{y_{j\bullet}}{n_j} \quad j = 1(1)n$$

Abbildung 6.3 zeigt eine empirische Regressionsfunktion 1.Art von Y bezüglich X (auf den Werten x_1, \dots, x_m) definiert). Durch die Punkte $(x_j, g(x_j))$, $j = 1(1)n$ gehende Kurven heißen *empirische Regressionskurven 1.Art*.

Bemerkung: Für größeren Datenumfang n bildet man Klasseneinteilungen für die x -Werte und trägt die bedingten y -Mittelwerte über den Klassenmitten auf. So erhält man oft u.U. glattere Regressionskurven.

6 Regressionsanalyse

Abbildung 6.3:



Konkrete Stichprobe (x_i, y_i) , $i = 1(1)n$, m Werte von x_i verschieden.

$$\begin{array}{l} (x_1, y_{11}), \dots, (x_1, y_{1n_1}) \quad x_1, \dots, x_m \text{ verschieden} \\ (x_2, y_{21}), \dots, (x_2, y_{2n_2}) \\ (x_m, y_{m1}), \dots, (x_m, y_{mn_m}) \quad \sum_{j=1}^n n_j = n \end{array}$$

Satz 6.2 Die empirische Regressionsfunktion von y bezüglich x ist unter allen (x_1, \dots, x_m) definierten Funktionen $f(\cdot)$ jene, für welche die Summe der Quadrate aller y -Differenzen minimal ist, d.h.

$$\sum_{j=1}^m \sum_{k=1}^{n_j} [y_{jk} - g(x_j)]^2 \equiv \sum_{j=1}^m \sum_{k=1}^{n_j} [y_{jk} - f(x_j)]^2$$

Beweis: $g(x_j)$ ist das arithmetische Mittel der y_{j1}, \dots, y_{jn_j} ; für festes j gilt wegen der bekannten Eigenschaften der Abstandsquadratsumme:

$$\sum_{k=1}^{n_j} [y_{jk} - g(x_j)]^2 \leq \sum_{k=1}^{n_j} [y_{jk} - c_j]^2 \quad \forall c_j \neq g(x_j)$$

Da alle m Summanden nicht negativ sind, folgt die Behauptung.

Bemerkung: Empirische Regressionsfunktion 1.Art von X bezüglich Y werden analog berechnet.

6.2 Regression 2.Art

Abhängigkeit einer Größe Y von einer nicht stochastischen Größe (Einstellgröße) X .

Ausgleichsrechnung

Beobachtete Paare (x_i, y_i) , $i = 1(1)n$ liegen annähernd auf einer Geraden

$$y = \alpha + \beta x.$$

Die „Parameter“ α und β können mit der Methode der kleinstenm Abstandsquadratsumme bestimmt werden:

$$QS = \sum_{i=1}^n (y_i - \alpha - \beta x_i)^2 \rightarrow \text{Min} \quad (6.1)$$

Jene Werte $\hat{\alpha}$ und $\hat{\beta}$ für die QS in Gleichung 6.1 minimal ist, sind die Parameter der Ausgleichsgeraden $y = \hat{\alpha} + \hat{\beta}x$.

Die Werte $\hat{\alpha}$ und $\hat{\beta}$ erhält man durch die notwendigen Bedingungen

$$\frac{\partial QS}{\partial \alpha} = 0 \quad \frac{\partial QS}{\partial \beta} = 0$$

6 Regressionsanalyse

Satz 6.3 Als Lösung ergibt sich

$$\hat{\alpha} = \frac{(\sum_{i=1}^n x_i^2) \cdot (\sum_{i=1}^n y_i)^2 - \sum_{i=1}^n x_i \cdot \sum_{i=1}^n x_i \cdot y_i}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}$$

$$\hat{\beta} = \frac{n (\sum_{i=1}^n x_i y_i) - (\sum_{i=1}^n x_i) (\sum_{i=1}^n y_i)}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}$$

Mit der Bezeichnung

$$s_{x,y} = \frac{1}{n-1} \left[\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y} \right] = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

gilt:

$$\hat{\beta} = \frac{s_{x,y}}{s_x^2} \quad \hat{\alpha} = \bar{y} - \hat{\beta} \bar{x}.$$

Die Ausgleichsgerade ergibt sich zu:

$$y - \bar{y} = \hat{\beta} (x - \bar{x}) = \frac{s_{x,y}}{s_x^2} (x - \bar{x})$$

Bemerkung: Die Ausgleichsgerade (Regressionsgerade) geht durch den Punkt (\bar{x}, \bar{y}) und hat die Steigung $\frac{s_{x,y}}{s_x^2}$.

Bemerkung: P.S. de Laplace hat für das Regressionsmodell $y = \theta x_i + \epsilon_i$ die Absolutbeträge herangezogen:

$$\hat{\theta} = \arg_{\theta} \min \sum_{i=1}^n |y_i - \theta x_i|$$

Allgemeine Ausgleichskurven

$$y = \psi(x; \alpha_1, \dots, \alpha_l)$$

Mittels der Methode der kleinsten Abstandsquadratsumme aus Daten $x_i, y_i, i = 1(1)n$

$$QS = \sum_{i=1}^n [y_i - \psi(x_i; \alpha_1, \dots, \alpha_l)]^2 \rightarrow \text{Min}$$

ermittelt man die „Parameter“ $\hat{\alpha}_1, \dots, \hat{\alpha}_l$ der Ausgleichskurve

$$y = \psi(x; \hat{\alpha}_1, \dots, \hat{\alpha}_l).$$

Notwendige Bedingung:

$$\frac{\partial QS}{\partial \alpha_k} = -2 \sum_{i=1}^n [y_i - \psi(x_i; \alpha_1, \dots, \alpha_l)] \frac{\partial \psi}{\partial \alpha_k} = 0 \quad k = 1(1)l$$

Sonderfall: $y = a + bx + cx^2$

Multiple Ausgleichsrechnung

Für funktionale Abhängigkeiten von k reellen Variablen

$$y = \psi(x_1, \dots, x_k; \theta) \quad \theta = \begin{pmatrix} \theta_1 \\ \vdots \\ \theta_m \end{pmatrix} \in \mathbf{R}^m$$

und Beobachtungen $(x_{i1}, \dots, x_{ik}; y_i); i = 1(1)n$ sind Ausgleichsflächen gesucht.

Sonderfall k=1:**Spezialfall: Lineare Funktion**

$$y = \theta_0 + \sum_{j=1}^k \theta_j x_j \quad \theta = \begin{pmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_k \end{pmatrix} \quad m = k + 1$$

Notation:

$$\begin{aligned} \underline{x} &= (x_1, \dots, x_k) \\ \dot{\underline{x}} &= (1, x_1, \dots, x_k) \\ y &= \dot{\underline{x}} \cdot \theta \end{aligned}$$

Spezialfall k=2:



Mit Hilfe der Methode der kleinsten Abstandsquadratsumme für n Beobachtungen

$$QS = \sum_{i=1}^n \left[y_i - \theta_0 - \sum_{j=1}^k \theta_j x_{ij} \right]^2 \rightarrow \text{Min} \quad (6.2)$$

können über die sog. Gauß'schen Normalgleichungen die Lösungen $\hat{\theta}_0, \hat{\theta}_1, \dots, \hat{\theta}_k$ gefunden werden. Schreibt man für die Matrix der Einstellgrößen

$$\mathcal{X} = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & & & & \vdots \\ \vdots & & & & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nk} \end{pmatrix}$$

$$\theta = \begin{pmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_k \end{pmatrix}$$

$$\underline{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix}$$

so erhält 6.2 die Form

$$QS = (\underline{y} - \mathcal{X}\theta)^T (\underline{y} - \mathcal{X}\theta).$$

Satz 6.4 Als Lösungsvektor θ ergibt sich für $\text{Rang } \mathcal{X} = k + 1$ und $n > k + 1$

$$\hat{\theta} = (\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T \underline{Y}.$$

6.3 Stochastische Regressionsanalyse

Für Streuungsuntersuchungen betrachtet man y als stochastische Größe Y_x , die von der Einstellgröße x abhängt:

$$Y_x = \psi(x) + U_x$$

Die Werte y_i in den Datenpaaren (x_i, y_i) werden als Realisierungen von der stochastischen Größen Y_i , die den Y_{x_i} entsprechen, betrachtet. Meist wird identisches Streuverhalten vorausgesetzt, d.h.

$$\text{Var}Y_x \equiv \sigma^2.$$

1-dimensionaler Fall $x_i \in \mathbb{R}$

Satz 6.5 (Gauß-Markoff) Sind Y_1, \dots, Y_n unkorreliert und gilt $\mathbb{E}Y_i = \alpha + \beta \cdot x_i$, sowie $\text{Var}Y_i \equiv \sigma^2$, $\forall i = 1(1)n$ und sind die x_1, \dots, x_n bekannt und nicht alle gleich und sind α, β und σ^2 unbekannt, so gilt:

Die mittels der Ausgleichsparameter $\hat{\alpha}$ und $\hat{\beta}$ konstruierten Schätzfunktionen A und B für α und β

$$A = \frac{(\sum_{i=1}^n x_i^2) (\sum_{i=1}^n Y_i) - (\sum_{i=1}^n x_i) (\sum_{i=1}^n x_i Y_i)}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}$$

↑
SG!
↓

$$B = \frac{n (\sum_{i=1}^n x_i Y_i) - (\sum_{i=1}^n x_i) (\sum_{i=1}^n Y_i)}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}$$

sind die effizienten linearen Schätzfunktionen (linear in Y_i) für α bzw. β und es gilt:

$$\begin{aligned} \text{Var}A &= \frac{\sum_{i=1}^n x_i^2}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \sigma^2 \\ \text{Var}B &= \frac{n}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \sigma^2 \\ \text{Cor}(A, B) &= -\frac{\sum_{i=1}^n x_i}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \sigma^2 \end{aligned}$$

Eine unverzerrte Schätzfunktion für σ^2 ist

$$S^2 = \frac{\sum_{i=1}^n (Y_i - A - Bx_i)^2}{n - 2}$$

Bemerkung:

1. $A = \bar{Y} - B\bar{x}$

6 Regressionsanalyse

2. $n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2 = n \sum_{i=1}^n (x_i - \bar{x}_n)^2$
3. Matrixschreibweise

$$\underline{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} \quad \mathcal{X} = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \quad \theta = \begin{pmatrix} \alpha \\ \beta \end{pmatrix}$$

$$\begin{pmatrix} A \\ B \end{pmatrix} \doteq \hat{\theta} - \begin{pmatrix} \hat{\alpha} \\ \hat{\beta} \end{pmatrix} = (\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T \underline{Y}$$

4. Der Wert $\mathbb{E}Y_x$ der Regressionsfunktion an der Stelle x wird durch $\widehat{\mathbb{E}_Q Y_x} = A + B \cdot x$ unverzerrt geschätzt.
Für diese Schätzfunktion gilt: $A + B \cdot x = \bar{Y} + B(x - \bar{x})$

6.4 Regressionsgeraden 2.Art bei Normalverteilung

Sind in Satz 6.5 die stochastischen Größen Y_i normalverteilt, d.h.

$$Y_i \sim N(\alpha + \beta x_i, \sigma^2),$$

so kann man einfache Konfidenzbereiche für α , β , σ^2 angeben und stochastische Tests konstruieren.

Satz 6.6 Gilt zu den Voraussetzungen von Satz 6.5 noch $Y_i \sim N(\alpha + \beta x_i, \sigma^2)$ so folgt:

1. $\hat{\alpha}$ und $\hat{\beta}$ sind auch die plausiblen Schätzwerte für α und β .
2. $A \sim N(\alpha, \text{Var}A)$ und $B \sim N(\beta, \text{Var}B)$
3. $\frac{\sum_{i=1}^n (Y_i - A - Bx_i)^2}{\sigma^2} = \underbrace{\frac{(n-2)S^2}{\sigma^2}}_{\text{unabh. von } A \text{ und } B} \sim \chi_{n-2}^2$

Beweis:

1. Lösungen der Plausibilitätsgleichungen

$$\frac{\partial}{\partial \alpha} \ln l(\alpha, \beta, \sigma^2; y_1, \dots, y_n, x_1, \dots, x_n) = 0$$

2. Y_1, \dots, Y_n sind unabhängig normalverteilt, daher gilt:

a) $\underline{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}$ ist normalverteilt

6 Regressionsanalyse

- b) $\begin{pmatrix} A \\ B \end{pmatrix} = (\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T \cdot \underline{Y}$ ist normalverteilt
 c) A und B sind normalverteilt

In Verbindung mit Satz 6.5 ergibt sich die Behauptung 2.

Konstruktion von Konfidenzintervallen für $\alpha, \beta, \sigma^2, \sigma$:

Satz 6.7 Unter den Voraussetzungen von Satz 6.6 erhält man folgende Konfidenzintervalle mit Überdeckungswahrscheinlichkeit $1 - \gamma$:

1. Für α

$$\left[A - t_{n-2; 1-\frac{\gamma}{2}} \sqrt{s^2 \frac{\sum_{i=1}^n x_i^2}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}}; A + t \dots \right]$$
2. Für β

$$\left[B - t_{n-2; 1-\frac{\gamma}{2}} \sqrt{s^2 \frac{n}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2}}; B + t \dots \right]$$
3. Für σ^2

$$\left[\frac{(n-2) s^2}{\chi_{n-2; 1-\frac{\gamma}{2}}^2}; \frac{(n-2) s^2}{\chi_{n-2; \frac{\gamma}{2}}^2} \right]$$
4. Für σ

$$\left[\sqrt{\frac{(n-2) s^2}{\chi_{n-2; 1-\frac{\gamma}{2}}^2}}; \sqrt{\frac{(n-2) s^2}{\chi_{n-2; \frac{\gamma}{2}}^2}} \right]$$

Beweis:

1. $\frac{A-\alpha}{\sqrt{\text{Var}A}} \sim N(0, 1)$ und Satz 6.6 Punkt 3:

$$\frac{(n-2) S^2}{\sigma^2} \sim \chi_{n-2}^2$$

Wenn $X \sim N(0, 1)$ und $Y \sim \chi_n^2$ dann sind die beiden Größen unabhängig und somit gilt:

$$\frac{X}{\sqrt{\frac{Y}{n}}} \sim t_n$$

$$\frac{(A-\alpha)\sigma}{\sqrt{\text{Var}A}\sqrt{S^2}} \sim t_{n-2}$$

2. Analog zu 1

6 Regressionsanalyse

3. $\frac{(n-2)S^2}{\sigma^2}$ ist eine Pivot-Größe

$$\mathbb{W} \left\{ \chi_{n-2; \frac{\gamma}{2}}^2 \leq \frac{(n-2)S^2}{\sigma^2} \leq \chi_{n-2; 1-\frac{\gamma}{2}}^2 \right\} = 1 - \gamma$$

4.

$$\mathbb{W} \left\{ \sqrt{\chi_{n-2; \frac{\gamma}{2}}^2} \leq \sqrt{\frac{(n-2)S^2}{\sigma^2}} \leq \sqrt{\chi_{n-2; 1-\frac{\gamma}{2}}^2} \right\} = 1 - \gamma$$

Bemerkung:

1. $\frac{n}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} = \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2}$
2. Mann kann auch 2-dimensionale Konfidenzbereiche für den Parametervektor (α, β) konstruieren.

Unter den Voraussetzungen von Satz 6.6 gilt

$$A + Bx - \alpha - \beta \sim N \left(0, \frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right),$$

da

$$\begin{aligned} \text{Var}(A + Bx - \alpha - \beta) &= \text{Var}(A + Bx) \\ \text{Var}A + \text{Var}(Bx) + 2x\text{Cov}(A, B) &= \frac{nx^2 - 2x \sum_{i=1}^n x_i + \sum_{i=1}^n x_i^2}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \sigma^2 \\ &= \sigma^2 \left(\frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) \end{aligned}$$

Gesucht ist nun ein Konfidenzintervall für $\mathbb{E}Y_x$ zu gegebenen x (vgl. Abbildung 6.4)

Satz 6.8 Unter den Voraussetzungen von Satz 6.6 gilt:

$$\left. \begin{array}{c} \left(\begin{array}{c} Y_1 - A - Bx_1 \\ \vdots \\ Y_n - A - Bx_n \end{array} \right) \\ \left(\begin{array}{c} A \\ B \end{array} \right) \end{array} \right\} \text{unabhaengig normalverteilt}$$

Folgerung:

$$\begin{aligned} Y_x - A - Bx &\sim N \left(0, 1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) \text{ da } \mathbb{E}(Y_x - A - Bx) = 0 \\ \text{Var}(Y_x - A - Bx) &= \text{Var}Y_x + \text{Var}(A + Bx) \\ &= \sigma^2 + \sigma^2 \left(\frac{1}{n} + \frac{(x - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right) \end{aligned}$$

Abbildung 6.4: Konfidenzintervall für $\mathbb{E}Y_x$



Satz 6.9 Mit der Bezeichnung

$$\hat{Y}_x := \bar{Y} + B(x - \bar{x}) = A + Bx$$

$$D^2 := \sum_{i=1}^n [(Y_i - \bar{Y}) - B(x_i - \bar{x})]^2 = (n - 2) S^2$$

gilt:

$$\begin{aligned} T &= \frac{\sqrt{n-2} \bar{Y} + B(x - \bar{x}) - \alpha - \beta x}{\sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{(n-1)S_x^2}} \cdot D} \\ &= \frac{\sqrt{n-2} \hat{Y}_x - \alpha - \beta x}{\sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{(n-1)S_x^2}} \cdot D} \\ &\sim t_{n-2} \end{aligned}$$

Daraus konstruiert man ein Konfidenzintervall für $\alpha + \beta x$ an der Stelle x mit Überdeckungswahrscheinlichkeit $1 - \gamma$:

$$\left[\hat{Y}_x - t_{n-2; 1-\frac{\gamma}{2}} \cdot \frac{D}{\sqrt{n-2}} \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{(n-1)S_x^2}}; \hat{Y}_x + \dots \right]$$

Ist $y = \hat{\alpha} + \hat{\beta}x$ die empirische Regressionsgerade, so erhält man für jedes x das empirische Konfidenzintervall für $\mathbb{E}Y_x$:

$$\left[\hat{\alpha} + \hat{\beta}x - t_{n-2; 1-\frac{\gamma}{2}} \cdot \frac{d}{\sqrt{n-2}} \sqrt{\frac{1}{n} + \frac{(x - \bar{x})^2}{(n-1)s_x^2}}; \hat{\alpha} + \hat{\beta}x + \dots \right]$$

6 Regressionsanalyse

Lässt man x variieren, dann erhält man einen sog. Konfidenzgürtel (vgl. Abbildung 6.5)

Abbildung 6.5: Konfidenzgürtel



Prognoseintervalle für einen zu x gehörigen Wert der stochastischen Größe Y_x .

Satz 6.10 Unter den Voraussetzungen dieses Abschnittes gilt

$$T = \frac{Y_x - A - Bx}{S \sqrt{1 + \frac{1}{n} + \frac{(x-\bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}} \sim t_{n-2}$$

Beweis: Vergleiche Folgerung von Satz 6.8. Daraus erhält man durch Umformung von

$$\mathbb{W} \left\{ -t_{n-2; 1-\frac{\gamma}{2}} \leq T \leq t_{n-2; 1-\frac{\gamma}{2}} \right\} = 1 - \gamma$$

ein Konfidenzintervall für Y_x :

$$\left[\hat{\alpha} + \hat{\beta}x - t_{n-2; 1-\frac{\gamma}{2}} \cdot s \cdot \sqrt{1 + \frac{1}{n} + \frac{(x-\bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2}}; \hat{\alpha} + \hat{\beta}x + \dots \right]$$

Bemerkung: Es gibt auch (breitere) simultane Konfidenzbänder.

6.5 Tests für den Parameter von Regressionsgeraden bei Normalverteilung

Voraussetzungen: $\mathbb{E}Y_x = \alpha + \beta x$, $\forall x \in$ Versuchsbereich und n Beobachtungen (x_i, y_i) , $i = 1(1)n$ gilt:

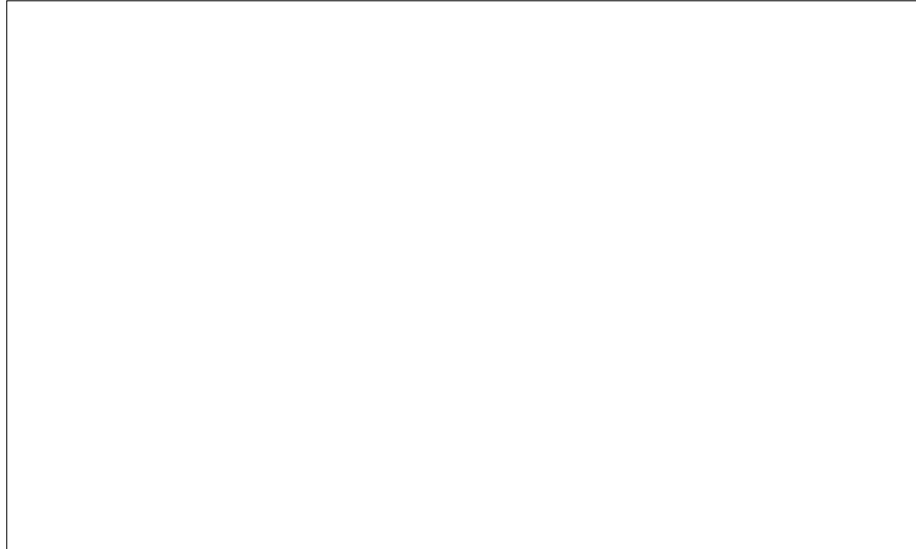
$$S_A^2 = \frac{\sum_{i=1}^n x_i^2}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \cdot S^2$$
$$S_B^2 = \frac{n}{n \sum_{i=1}^n x_i^2 - (\sum_{i=1}^n x_i)^2} \cdot S^2$$

Tests für verschieden Hypothesen sind in folgender Tabelle angegeben.

T1

6.6 Tests für Regressionskurven bei Normalverteilung

Voraussetzungen: $Y_x \sim N(\mu(x), \sigma^2)$, σ^2 ist unabhängig von x , (x_ν, y_ν) , $\nu = 1(1)n$, $x_i^* \rightarrow m$ verschieden Werte für x mit Vielfachheit n_i



$$\bar{y}_i = \frac{1}{n_i} \sum_{k=1}^{n_i} y_{ik}$$

6.7 Test auf Regressionsgerade

Null-Hypothese: $\mathcal{H}_0 : \mu(x) = \alpha + \beta x$, wobei α und β nicht bekannt sein müssen.

$$\begin{aligned} \tilde{y} &= A + Bx = \bar{y} + B(x - \bar{x}) \\ \tilde{y}_i &= \bar{y} + B(x_i^* - \bar{x}) \end{aligned}$$

T2

6 Regressionsanalyse

Zerlegung der Summe:

$$\begin{aligned}
 (y_{ik} - \tilde{y}_i)^2 &= [(y_{ik} - \bar{y}_i) + (\bar{y}_i - \tilde{y}_i)]^2 \\
 &= (y_{ik} - \bar{y}_i)^2 + (\bar{y}_i - \tilde{y}_i)^2 + 2(y_{ik} - \bar{y}_i)(\bar{y}_i - \tilde{y}_i) \\
 q &= \underbrace{\sum_{i=1}^m \sum_{k=1}^{n_i} (y_{ik} - \bar{y}_i)^2}_{q_2} + \underbrace{\sum_{i=1}^m \sum_{k=1}^{n_i} (\bar{y}_i - \tilde{y}_i)^2}_{q_1} + 2 \sum_{i=1}^m (\bar{y}_i - \tilde{y}_i) \underbrace{\sum_{k=1}^{n_i} (y_{ik} - \bar{y}_i)}_{=0}
 \end{aligned}$$

q_2 =Summe der quadrierten Abweichungen vom jeweiligen Mittel \bar{y}_i

q_1 =Summe der quadrierten Abweichungen der Schätzwerte von den Mittelwerten

Satz 6.11 von Cockran: Sind X_1, \dots, X_n unabhängig standard-normalverteilt und gilt

$$\sum_{k=1}^n X_k^2 = \sum_{i=1}^m Q_i(X_1, \dots, X_n)$$

mit Q_i sind nicht negativ symmetrische quadratische Formen mit $\text{Rang}(Q_i) \leq r_i$ und $\sum_{i=1}^m r_i = n$ dann gilt:

$$\left. \begin{aligned}
 Q_i(X_1, \dots, X_n) &\sim \chi_{r_i}^2 \\
 Q_1(X_1, \dots, X_n), \dots, Q_m(X_1, \dots, X_n) &\} \text{unabhaengig}
 \end{aligned} \right\}$$

Satz 6.12 Für die den Beobachtungen (x_ν, y_ν) entsprechenden stochastischen Größen Y_{ik}, \bar{Y}_i und \tilde{Y}_i gilt:

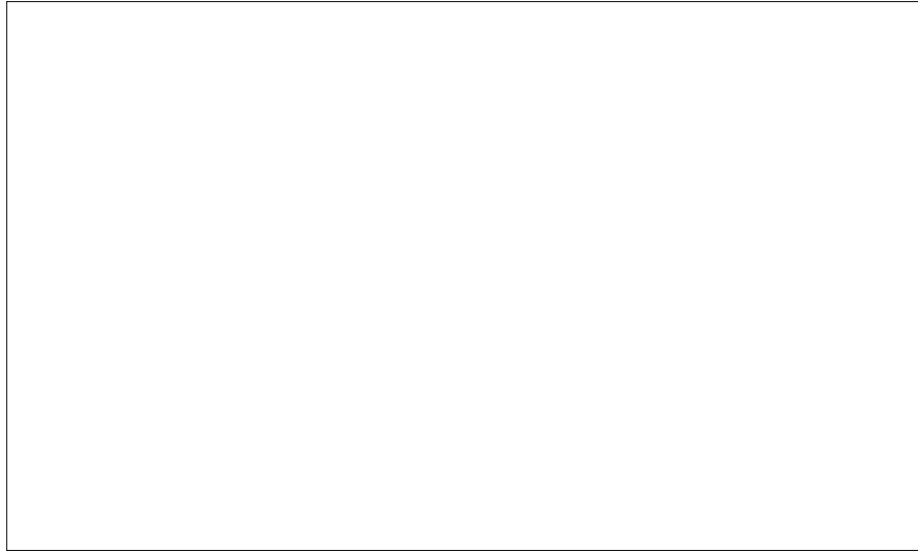
$$\begin{aligned}
 \frac{1}{\sigma^2} \sum_{i=1}^m n_i (\bar{Y}_i - \tilde{Y}_i)^2 &= \frac{Q_1}{\sigma^2} \sim \chi_{m-2}^2 \\
 \frac{1}{\sigma^2} \sum_{i=1}^m \sum_{k=1}^{n_i} (Y_{ik} - \bar{Y}_i)^2 &= \frac{Q_2}{\sigma^2} \sim \chi_{n-m}^2
 \end{aligned}$$

Diese beiden Größen sind unabhängig. Daraus folgt, falls \mathcal{H}_0 richtig ist:

$$Z = \frac{Q_1 / (m - 2)}{Q_2 / (n - m)} \sim F_{m-2; n-m}$$

Test: \mathcal{H}_0 wird abgelehnt, falls

$$\frac{q_1 / (m - 2)}{q_2 / (n - m)} > F_{m-2; n-m; 1-\gamma}$$



6.8 Multiple lineare Regression

k Einstellgrößen x_1, \dots, x_k , $k+1$ unbekannte Parameter $\theta_0, \theta_1, \dots, \theta_k$, Regressionsfunktion $\psi(x_1, \dots, x_k) = \theta_0 + \sum_{j=1}^k \theta_j \cdot x_j$

$$\begin{aligned} \underline{x} &= (x_1, \dots, x_k) \\ \dot{\underline{x}} &= (1, x_1, \dots, x_k) \\ \theta &= \begin{pmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_k \end{pmatrix} \end{aligned}$$

Y : abhängige Größe

$$\begin{aligned} Y_{\underline{x}} &= \dot{\underline{x}} \cdot \theta + U_{\underline{x}} \text{ mit } \mathbb{E}U_{\underline{x}} = 0 \\ \text{Var}Y_{\underline{x}} &\equiv \sigma^2 \end{aligned}$$

Mit den Bezeichnungen der multiplen Ausgleichsrechnung (vgl. Abschnitt 6.2) gilt:

n Beobachtungen $(x_{i1}, x_{i2}, \dots, x_{ik}; y_i)$ mit $i = 1(1)n$

$$\underline{y} = \mathcal{X}\theta + \underline{u} \quad \text{mit} \quad \mathcal{X} = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1k} \\ 1 & x_{21} & x_{22} & & x_{2k} \\ \vdots & & & & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{nk} \end{pmatrix}$$

$$\underline{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix}$$

Schätzungen aus dem Ausgleichsproblem:

Frage: Statistische Qualität der Schätzer? Schätzung der Varianz σ^2 ?

Das System der Gleichungen hat die Gestalt

$$\underline{Y} = \mathcal{X}\theta + \underline{U}$$

Vorraussetzungen:

1. Der Rang $\mathcal{X} = k + 1$, daher $n > k$
2. $\text{VCov}(U_1, \dots, U_n) = \sigma^2 \cdot I_n$ ($I_n=1$ -Matrix)
3. Der Parameterraum darf keine Hyperebene der \mathbf{R}^{k+1} sein.

Definition: Schätzfunktionen T_j für θ_j in einem multiplen linearen Regressionsmodell heißen effiziente lineare Schätzfunktionen, wenn gilt:

1. $T_j = \sum_{i=1}^n c_{ji} \cdot Y_i$, wobei die c_{ji} nur von den (bekannten) x_{ij} abhängen.
2. $\mathbb{E}_{\theta, \sigma^2} T_j = \theta_j$ (unverzerrt, erwartungstreu)
3. $\text{Var}_{\theta, \sigma^2} T_j$: minimal unter allen Schätzfunktionen, welche die ersten beiden Punkte erfüllen

Analog zu Satz 6.6:

Satz 6.13 (Gauß-Markoff) Unter den Voraussetzungen:

$$\begin{aligned} Y_1, \dots, Y_n & \text{ paarweise unkorreliert} \\ \text{Var} Y_i & \equiv \sigma^2 \quad \forall i = 1(1)n \quad (\text{Homoskedastizität}) \\ \mathcal{X} & = \begin{pmatrix} 1 & \underline{x}_1 \\ 1 & \underline{x}_2 \\ \vdots & \vdots \\ 1 & \underline{x}_n \end{pmatrix} \quad \text{bekannte Einstellgrößen } \underline{x}_i, i = 1(1)n \\ \theta_0, \theta_1, \dots, \theta_k & \text{ linear unabhängig (nicht in einer Hyperebene der } \mathbf{R}^{k+1}) \\ \mathbb{E} Y_i & = \underline{\dot{x}}_i \cdot \theta = \theta_0 + \sum_{j=1}^k x_{ij} \cdot \theta_j \\ \text{Rang } \mathcal{X} & = k + 1 \end{aligned}$$

gilt:

$$1. \underline{T} = \begin{pmatrix} T_0 \\ T_1 \\ \vdots \\ T_k \end{pmatrix} = (\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T \cdot \underline{Y}$$

sind die eindeutig bestimmten effizienten linearen Schätzfunktionen T_j für θ_j $j = 0(1)k$.

$$2. \text{VCov}(T_0, T_1, \dots, T_k) = (\mathcal{X}^T \mathcal{X})^{-1} \cdot \sigma^2$$

$$3. S^2 := \frac{(\underline{Y} - \mathcal{X}\underline{T})^T (\underline{Y} - \mathcal{X}\underline{T})}{n-k-1} \text{ ist eine unverzerrte Schätzfunktion für } \sigma^2.$$

Beweis:

- Extrema unter Nebenbedingungen. Unter den gegebenen Annahmen hat das Ausgleichsproblem $\exists!$ Lösung:

$$\hat{\theta}_0, \hat{\theta}_1, \dots, \hat{\theta}_k \leftrightarrow T_j, j = 0(1)k$$

- Kovarianzen berechnet man mit Hilfe von

$$\text{Cov}(Y_r, Y_s) = \sigma^2 \delta_{r,s} \dots$$

-

$$\begin{aligned} \mathbb{E}[(n-k-1)S^2] &= \mathbb{E}[(\underline{Y} - \mathcal{X}\underline{T})^T (\underline{Y} - \mathcal{X}\underline{T})] \\ &= \mathbb{E}\left[\sum_{i=1}^n \left(Y_i - T_0 - \sum_{j=1}^k x_{ij} T_j\right)^2\right] \\ &= \mathbb{E}\left[\sum_{i=1}^n \left\langle Y_i^2 - 2Y_i \left\{T_0 + \sum_{j=1}^k x_{ij} T_j\right\} + \left(T_0 + \sum_{j=1}^k x_{ij} T_j\right)^2 \right\rangle\right] \end{aligned}$$

Einsetzen von $\underline{T} = (\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T \underline{Y}, \dots$

Bemerkung: Die Varianzschätzung ist intuitiv; gemittelte quadrierte Abweichungen von der geschätzten Regressions-Hyperebene:

$$\begin{aligned} y &= \hat{\theta}_0 + \sum_{j=1}^k x_j \cdot \hat{\theta}_j, x_j \in \mathbf{R} \\ s^2 &= \frac{1}{n-k-1} \cdot \sum_{i=1}^n \left[y_i - \hat{\theta}_0 - \sum_{j=1}^k x_{ij} \cdot \hat{\theta}_j \right]^2 \end{aligned}$$

Prognose im linearen Regressionsmodell

$$Y_{\underline{x}} = \underline{\dot{x}}\theta + U_{\underline{x}}$$

Aus Daten geschätzte Parameter $\hat{\theta} = \begin{pmatrix} \hat{\theta}_0 \\ \hat{\theta}_1 \\ \vdots \\ \hat{\theta}_k \end{pmatrix}$

Prognose zum Einstell-Vektor $\underline{x} = (x_1, \dots, x_k)$:

$$\hat{Y}_{\underline{x}} := \underline{\dot{x}} \cdot \hat{\theta} = \hat{\theta}_0 + x_1 \cdot \hat{\theta}_1 + \dots + x_k \cdot \hat{\theta}_k$$

Für diese Prognose gilt:

1. $\mathbb{E}\hat{Y}_{\underline{x}} = \mathbb{E}Y_{\underline{x}}$ (unverzerrt)
2. $\text{Var}\hat{Y}_{\underline{x}} = \underline{\dot{x}} (\mathcal{X}^T \mathcal{X})^{-1} \cdot \underline{\dot{x}}^T \cdot \sigma^2$

Beweis:

1.

$$\begin{aligned} \mathbb{E}\hat{Y}_{\underline{x}} &= \mathbb{E}(\underline{\dot{x}}\hat{\theta}) \\ &= \mathbb{E}\left[\hat{\theta}_0 + \sum_{j=1}^k x_j \cdot \hat{\theta}_j\right] \\ &= \mathbb{E}\hat{\theta}_0 + \sum_{j=1}^k x_j \cdot \mathbb{E}\hat{\theta}_j \\ &= \theta_0 + \sum_{j=1}^k x_j \cdot \theta_j \\ &= \mathbb{E}Y_{\underline{x}} \end{aligned}$$

2.

$$\begin{aligned} \text{Var}\hat{Y}_{\underline{x}} &= \text{Var}(\underline{\dot{x}}\hat{\theta}) \\ &= \text{Var}\left(\underbrace{\underline{\dot{x}} (\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T}_{\text{Vektor } \underline{z}} \cdot \underline{Y}\right) \end{aligned}$$

6 Regressionsanalyse

$$\begin{aligned}
 &= \text{Var}(\underline{z} \cdot \underline{Y}) = \text{Var}\left(\sum_{i=1}^n z_i \cdot Y_i\right) \\
 &= \sum_{i=1}^n \text{Var}(z_i Y_i) = \sum_{i=1}^n z_i^2 \cdot \underbrace{\text{Var} Y_i}_{\sigma^2} \\
 &= \sigma^2 \cdot \underline{z} \cdot \underline{z}^T \\
 &= \sigma^2 \left(\dot{\underline{x}} (\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T\right) \cdot \left(\dot{\underline{x}} (\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T\right)^T \\
 &= \sigma^2 \cdot \dot{\underline{x}} (\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T \mathcal{X} \left((\mathcal{X}^T \mathcal{X})^{-1}\right)^T \dot{\underline{x}}^T \\
 &= \sigma^2 \cdot \underbrace{\left(\dot{\underline{x}} \cdot (\mathcal{X}^T \mathcal{X})^{-1} \cdot \dot{\underline{x}}^T\right)^T}_{\in \mathbb{R}}
 \end{aligned}$$

6.9 Multiple lineare Regression bei Normalverteilung

Falls die abhängigen Größen Y_1, \dots, Y_n im Satz von Gauß-Markoff normalverteilt sind, gilt:

Satz 6.14 Wird zu den Voraussetzungen von Gauß-Markoff vorausgesetzt, dass

$$\begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} \sim N(\mathcal{X}\theta, \sigma^2 I_n)$$

so folgt:

- Die T_j sind auch die plausiblen Schätzfunktionen für die Regressionsparameter θ_j , $j = 0(1)k$

$$2. \begin{pmatrix} T_0 \\ T_1 \\ \vdots \\ T_k \end{pmatrix} \sim N\left(\begin{pmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_k \end{pmatrix}, (\mathcal{X}^T \mathcal{X})^{-1} \cdot \sigma^2\right)$$

$$3. \left. \begin{array}{l} \frac{(n-k-1)S^2}{\sigma^2} \sim \chi_{n-k-1}^2 \\ \frac{(\underline{T}-\theta)^T \mathcal{X}^T \mathcal{X} (\underline{T}-\theta)}{\sigma^2} \sim \chi_{k+1}^2 \end{array} \right\} \text{unabhaengig}$$

Beweis:

- Plausibilitätsfunktion

$$\prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{(y_i - \theta_0 - \sum_{j=1}^k x_{ij}\theta_j)^2}{2\sigma^2}}$$

2. $\underline{T} = \underbrace{(\mathcal{X}^T \mathcal{X})^{-1} \mathcal{X}^T \underline{Y}}_{\text{konstant}}$ ist als Lineartransformation einer n -dimensionalen Normalverteilung auch normalverteilt.
3. Quadratsummen von normalverteilten stochastischen Größen (Division durch σ) ergeben χ^2 -Verteilungen.

Satz 6.15 Unter den Voraussetzungen des vorausstehenden Satzes gilt:

$$\left. \begin{array}{l} \underline{Y} - \mathcal{X}\underline{T} = \begin{pmatrix} Y_1 - T_0 - \sum_{j=1}^k x_{1j}T_j \\ \vdots \\ Y_n - T_0 - \sum_{j=1}^k x_{nj}T_j \end{pmatrix} \\ \hat{\theta} = \underline{T} = \begin{pmatrix} T_0 \\ T_1 \\ \vdots \\ T_k \end{pmatrix} \end{array} \right\} \text{unabhaengig normalverteilt}$$

Satz 6.16 Unter den Voraussetzungen dieses Abschnittes und der Bezeichnung $s^{jj}=j$ -tes Hauptdiagonalelement der Matrix $(\mathcal{X}^T \mathcal{X})^{-1}$ gilt:

1. Konfidenzintervalle für θ_j mit Überdeckungswahrscheinlichkeit $1 - \delta$:

$$\left[T_j - t_{n-k-1; 1-\frac{\delta}{2}} \cdot \sqrt{s^{jj}} \cdot S, T_j + \dots \right]$$

2. Konfidenzintervall für σ^2 mit Überdeckungswahrscheinlichkeit $1 - \delta$:

$$\left[\frac{(n-k-1) S^2}{\chi_{n-k-1; 1-\frac{\delta}{2}}^2}, \frac{(n-k-1) S^2}{\chi_{n-k-1; \frac{\delta}{2}}^2} \right]$$

Beweis:

1. Laut Satz 6.14 und Satz 6.15 gilt:

$$\left. \begin{array}{l} \frac{T_j - \theta_j}{\sqrt{s^{jj}} \cdot \sigma} \sim N(0, 1) \\ \frac{(n-k-1)S^2}{\sigma^2} \sim \chi_{n-k-1}^2 \end{array} \right\} \text{unabhaengig}$$

$$\frac{T_j - \theta_j}{\sqrt{s^{jj}} \cdot S} \sim t_{n-k-1}$$

2. Laut Satz 6.14 gilt:

$$\frac{(n-k-1) S^2}{\sigma^2} \sim \chi_{n-k-1}^2$$

Rest analog zu ersterem.

Test auf Linearität einer Regressionsfunktion

Für normalverteilte abhängige Größen $Y_i \sim N(\mu(\underline{x}_i), \sigma^2)$ kann man auf Linearität testen, wenn mindestens ein Einstellvektor $\underline{x}_i = (x_{i1}, \dots, x_{ik})$, $i = 1(1)n$ vorliegt, zu dem mindestens zwei Beobachtungen $Y_{i,1}$ und $Y_{i,2}$ vorliegen.

Satz 6.17 Gilt für die abhängigen Größen $Y_i \sim N(\underline{\dot{x}}_i \cdot \theta, \sigma^2)$ mit $\underline{\dot{x}} = (1, x_{i1}, \dots, x_{ik})$

$$\theta = \begin{pmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_k \end{pmatrix} \quad \mathcal{X} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1k} \\ 1 & x_{21} & \cdots & x_{2k} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \cdots & x_{nk} \end{pmatrix} \quad \text{Rang } \mathcal{X} = k + 1$$

und bildet man analog zu Abschnitt 6.6 die zu gleichen Einstellvektoren \underline{x}_j gehörenden Mittelwerte der zugehörigen Y -Werte

$$\begin{array}{ll} \underline{x}_a, a = 1(1)l & \text{verschiedene } \underline{x}_j \text{ - Werte} \\ n_a, \sum_{a=1}^l n_a = n & \text{mindestens } 1 \ n_a \geq 2 \end{array}$$

$$\begin{aligned} \bar{Y}_a &= \frac{1}{n_a} \sum_{\underline{x}_i = \underline{x}_a} Y_i \\ Z &= \frac{(n-l) \sum_{a=1}^l n_a \cdot \left(\bar{Y}_a - T_0 - \sum_{j=1}^k x_{ij} T_j \right)^2}{(l-k-1) \sum_{a=1}^l \sum_{\{\underline{x}_i = \underline{x}_a\}} (Y_i - \bar{Y}_a)^2} \sim F_{l-k-1, n-l} \end{aligned}$$

Test: Unter den Voraussetzungen von Satz 6.17 ist ein Test für die Hypothese \mathcal{H}_0

$$\mathcal{H}_0 : \mathbb{E}Y_{\underline{x}} = \theta_0 + \sum_{j=1}^k x_{ij} \cdot \theta_j = \underline{\dot{x}}\theta$$

mit Überdeckungswahrscheinlichkeit α eines Fehlers 1.Art durch folgenden Verwerfungsraum V für die Teststatistik Z gegeben:

$$V = [F_{k-k-1, n-l; 1-\alpha}, \infty]$$

Begründung:



Bemerkung: Es ist sinnvoll für große Werte von Z zu verwerfen, da der Erwartungswert des Nenners unabhängig von der besonderen Art der Regressionsfunktion gleich σ^2 ist. Der Erwartungswert des Zählers ist umso größer, je stärker die tatsächlichen $\mathbb{E}\mu_i$ von den hypothetischen Erwartungswerten abweichen.

Bemerkung: Der Satz 6.12 ergibt sich als Sonderfall.

6.10 Bayes'sche Regressionsanalyse

Auch die Parameter werden durch stochastische Größen beschrieben.

		a – priori
α	$\tilde{\alpha}$	$\pi(\alpha)$
β	$\tilde{\beta}$	$\pi(\beta)$
σ^2	$\tilde{\sigma}^2$	$\pi(\sigma^2)$
θ	$\tilde{\theta}$	$\pi(\theta)$

Regressionsgeraden

$$\mathbb{E}Y_{\underline{x}} = \alpha + \beta x, \quad \mathbf{Var}Y_x = \sigma^2$$

A-priori Verteilung $\pi(\alpha, \beta, \sigma^2)$

habe die Dichte $f(y|x, \alpha, \beta, \sigma^2)$

Daten $D = (x_i, y_i), i = 1(1)n$

6 Regressionsanalyse

Plausibilitätsfunktion $l(\alpha, \beta, \sigma^2; D)$

$$l(\alpha, \beta, \sigma^2; x_1, y_1, \dots, x_n, y_n) = \prod_{i=1}^m f(y_i | x_i, \alpha, \beta, \sigma^2)$$

Mittels des Bayes'schen Theorems erhält man die a-posteriori Verteilung von $(\tilde{\alpha}, \tilde{\beta}, \tilde{\sigma}^2)$

$$\pi(\alpha, \beta, \sigma^2 | D) \propto \pi(\alpha, \beta, \sigma^2) \cdot l(\alpha, \beta, \sigma^2; D)$$

Damit kann man HPD-Bereiche

- für einzelne Parameter
- gemeinsame HPD-Bereiche

konstruieren.

Außerdem: Prädiktivverteilungen (Prognoseverteilungen) $Y_x | D$ mit Prognosedichte $f_x(\bullet | D)$

$$f_x(y | D) = \int \int \int f(y | x, \alpha, \beta, \sigma^2) \cdot \pi(\alpha, \beta, \sigma^2 | D) \, d\alpha d\beta d\sigma^2$$

Bemerkung: Mittels $f_x(\bullet | D)$ kann man Prognoseintervalle für Y_x ermitteln.

Spezialfall: $Y_x \sim N(\alpha + \beta x, \sigma^2)$

$$\begin{aligned} f(y | x, \alpha, \beta, \sigma^2) &= \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{y_i - \alpha - \beta x_i}{\sigma}\right)^2} \\ \text{Daten } D &= (x_1, y_1, \dots, x_n, y_n) \\ l(\alpha, \beta, \sigma^2; D) &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{y_i - \alpha - \beta x_i}{\sigma}\right)^2} \\ &= \frac{1}{(\sqrt{2\pi}\sigma)^n} e^{-\frac{1}{2}\sum_{i=1}^n \left(\frac{y_i - \alpha - \beta x_i}{\sigma}\right)^2} \end{aligned}$$

Die Anwendung des Bayes'schen Theorems liefert die a-posteriori Verteilung

$$\pi(\alpha, \beta, \sigma^2 | D) \propto \pi(\alpha, \beta, \sigma^2) \cdot \frac{1}{(\sigma^2)^{n/2}} e^{-\frac{1}{2}\sum_{i=1}^n \left(\frac{y_i - \alpha - \beta x_i}{\sigma}\right)^2}$$

Daraus berechnet man HPD-Bereiche und Prognosedichten sowie Prognoseintervalle.

Bemerkung: Häufig $\pi(\alpha, \beta, \sigma^2) = p(\alpha, \beta) \cdot q(\sigma^2)$

Multiple lineare Regression

$$\mathbb{E}Y_{\underline{x}} = \underline{\dot{x}}\theta \quad \text{mit} \quad \theta = \begin{pmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_k \end{pmatrix}$$

$$Y_{\underline{x}} \sim f(\bullet | \underline{x}, \theta, \sigma^2) \quad \text{Var}Y_{\underline{x}} = \sigma^2$$

$$\underline{x} = (x_1, \dots, x_k) \quad \underline{\dot{x}} = (1, x_1, \dots, x_k)$$

$$\text{Daten } D = (x_i, y_i), i = 1(1)n \quad \text{mit } \underline{x}_i = (x_{i1}, \dots, x_{ik})$$

Plausibilitätsfunktion $l(\theta, \sigma^2; D)$

Für vollständige Daten gilt:

$$l(\theta, \sigma^2; \underline{x}_1, \underline{y}_1, \dots, \underline{x}_n, \underline{y}_n) = \prod_{i=1}^n f(y_i | \underline{x}_i, \theta, \sigma^2)$$

Aufgrund des Bayes'schen Theorems ergibt sich

$$\pi(\theta, \sigma^2 | D) \propto \pi(\theta, \sigma^2) \cdot l(\theta, \sigma^2; D)$$

Daraus:

- HPD-Bereiche für Parameter
- Prognosedichten
- Prognoseintervalle

Sonderfall: $Y_{\underline{x}} \sim N(\underline{\dot{x}}\theta, \sigma^2)$

$$f(y | \underline{x}, \theta, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{(y - \underline{\dot{x}}\theta)^2}{2\sigma^2}}$$

7 Varianzanalyse

Es soll untersucht werden, ob ein oder mehrere Größen (genannt Faktoren) Einfluß auf ein beobachtetes Merkmal haben. Die Varianzanalyse umfaßt Testverfahren, die auf einer Analyse des Streuverhaltens empirisch gegebener stochastischer Größen beruhen.

Annahme: Alle beobachtbaren stochastischen Größen werden als normalverteilt mit identischen Varianzen vorausgesetzt.

7.1 Grundlagen der Varianzanalyse (ANOVA)

Mit jeder symmetrischen, quadratischen n -reihigen Matrix $A = (a_{ij})$ ist durch

$$\underline{x} = (x_1, \dots, x_n)^T \rightarrow Q(\underline{x}) = \sum_{i=1}^n \sum_{j=1}^n a_{ij} \cdot x_i \cdot x_j = \underline{x}^T \cdot A \cdot \underline{x}$$

eine symmetrische quadratische Form Q in den n Variablen x_1, \dots, x_n festgelegt.

Der Rang von Q ist definiert als der Rang der Matrix A .

Eine symmetrische Matrix B - bzw. die zugeordnete lineare Transformation - heißt orthogonal, wenn

$$B^{-1} = B^T$$

Satz 7.1 (Hauptachsentransformation) Jede symmetrische quadratische Form $Q(\underline{x})$ von Rang r kann durch eine orthogonale Transformation in die Form

$$Q^*(\underline{y}) = \sum_{i=1}^r c_i y_i^2 \quad \text{mit } c_1, \dots, c_r \neq 0$$

übergeführt werden. Ist Q nichtnegativ, d.h. $Q : \mathbf{R}^n \rightarrow [0, \infty)$, so gilt $c_1, \dots, c_r > 0$.

Satz 7.2 (Satz von Cockran) Sind X_1, \dots, X_n unabhängig verteilt nach $N(0, 1)$ und

$$\sum_{k=1}^n X_k^2 = \sum_{i=1}^m Q_i(X_1, \dots, X_n)$$

wobei die $Q_i(\dots)$ nichtnegative symmetrische quadratische Formen mit $\text{Rang}(Q_i) \leq r_i$ und $\sum_{i=1}^m r_i = n$, so folgt:

1. $Q_i(X_1, \dots, X_n) \sim \chi_{r_i}^2$
2. $Q_1(X_1, \dots, X_n), \dots, Q_m(X_1, \dots, X_n)$ sind unabhängig.

7.2 Einfache Varianzanalyse (1 Faktor)

Es soll getestet werden, ob m stochastische Größen X_1, \dots, X_m identische Erwartungswerte haben.

Nullhypothese: $\mathcal{H}_0 : \mu_1 = \mu_2 = \dots = \mu_m$

Zur stochastischen Größe X_i konkrete Stichprobe x_{i1}, \dots, x_{in_i} mit $n_1, \dots, n_m, \sum_{i=1}^m n_i = n$

Gruppe	Stichprobenwerte	Summe
1. Gruppe	$x_{11}, x_{12}, \dots, x_{1n_1}$	$x_{1\bullet}$
2. Gruppe	$x_{21}, x_{22}, \dots, x_{2n_2}$	$x_{2\bullet}$
\vdots		
i . Gruppe	$x_{i1}, x_{i1}, \dots, x_{in_i}$	$x_{i\bullet}$
\vdots		
m . Gruppe	$x_{m1}, x_{m2}, \dots, x_{mn_m}$	$x_{m\bullet}$
		$x_{\bullet\bullet}$

Gruppenwerte $x_{i\bullet} = \sum_{k=1}^{n_i} x_{ik} \quad i = 1(1)m$

Gruppenmittel $\bar{x}_i = \frac{x_{i\bullet}}{n_i}$

Gesamtmittel $\bar{x} = \frac{1}{n} \sum_{i=1}^m n_i \cdot \bar{x}_i = \frac{x_{\bullet\bullet}}{n}$

$s^2 = \frac{1}{n-1} \sum_{i=1}^m \sum_{k=1}^{n_i} (x_{ik} - \bar{x})^2$

Quadratsumme $q = (n-1) \cdot s^2$

$$\begin{aligned}
 q &= \sum_{i=1}^m \sum_{k=1}^{n_i} (x_{ik} - \bar{x})^2 \\
 &= \sum_{i=1}^m \sum_{k=1}^{n_i} [(x_{ik} - \bar{x}_i) + (\bar{x}_i - \bar{x})]^2 \\
 &= \sum_{i=1}^m \sum_{k=1}^{n_i} (x_{ik} - \bar{x}_i)^2 + \sum_{i=1}^m \sum_{k=1}^{n_i} (\bar{x}_i - \bar{x})^2 + 2 \underbrace{\sum_{i=1}^m \sum_{k=1}^{n_i} (x_{ik} - \bar{x}_i) (\bar{x}_i - \bar{x})}_0
 \end{aligned}$$

$$\sum_{i=1}^m \sum_{k=1}^{n_i} (x_{ik} - \bar{x}) (\bar{x}_i - \bar{x}) = \sum_{i=1}^m (\bar{x}_i - \bar{x}) \cdot \sum_{k=1}^{n_i} (x_{ik} - \bar{x}_i)$$

7 Varianzanalyse

$$= \underbrace{\sum_{i=1}^m (\bar{x}_i - \bar{x}) \left[\underbrace{\left(\sum_{k=1}^{n_i} x_{ik} \right)}_{n_i \cdot \bar{x}_i} - n_i \cdot \bar{x}_i \right]}_0$$

Wegen

$$\sum_{k=1}^{n_i} (\bar{x}_i - \bar{x})^2 = n_i \cdot (\bar{x}_i - \bar{x})^2$$

gilt:

$$\begin{aligned} q &= \sum_{i=1}^m \sum_{k=1}^{n_i} (x_{ik} - \bar{x})^2 \\ &= \underbrace{\sum_{i=1}^m n_i (\bar{x}_i - \bar{x})^2}_{q_1} + \underbrace{\sum_{i=1}^m \sum_{k=1}^{n_i} (x_{ik} - \bar{x}_i)^2}_{q_2} \end{aligned}$$

q_1 = Summe der Abweichungsquadrate zwischen den Gruppen

q_2 = Summe der Abweichungsquadrate innerhalb der Gruppen

Falls die Hypothese $\mathcal{H}_0 : \mu_1 = \mu_2 = \dots = \mu_m = \mu$ richtig ist, gilt $X_{ik} \sim N(\mu, \sigma^2)$ und $\underline{X} = (X_{11}, \dots, X_{1n_1}; \dots; X_{m1}, \dots, X_{mn_m})$ ist eine Stichprobe von X mit Umfang $n = \sum_{i=1}^m n_i$.

Bemerkung: Eine unerzerzte Schätzfunktion für σ^2 :

$$\frac{1}{n-1} Q = \frac{1}{n-1} \sum_{i=1}^m \sum_{k=1}^{n_i} (X_{ik} - \bar{X}_n)^2$$

$$\frac{Q}{\sigma^2} \sim \chi_{n-1}^2$$

$\frac{Q}{\sigma^2}$ hat den Rang $n-1$ und wegen $Q = Q_1 + Q_2$ folgt: $\frac{Q}{\sigma^2} = \frac{Q_1}{\sigma^2} + \frac{Q_2}{\sigma^2}$ und $\text{Rang} \frac{Q_1}{\sigma^2} \leq m-1$ und $\text{Rang} \frac{Q_2}{\sigma^2} \leq n-m$.

Nach dem Satz von Cockran folgt, dass $\frac{Q_1}{\sigma^2} \sim \chi_{m-1}^2$ und $\frac{Q_2}{\sigma^2} \sim \chi_{n-m}^2$ unabhängig sind.

Daraus folgt:

$$T = \frac{Q_1 / (m-1) \sigma^2}{Q_2 / (n-m) \sigma^2} = \frac{Q_1 / (m-1)}{Q_2 / (n-m)} \sim F_{m-1; n-m}$$

Verwerfungsraum für \mathcal{H}_0

$$V = \left\{ (x_{11}, \dots, x_{mn_m}) : \frac{q_1 / (m-1)}{q_2 / (n-m)} \geq F_{m-1, n-m; 1-\alpha} \right\}$$

mit einer Wahrscheinlichkeit α eines Fehlers 1. Art.

7.3 Zweifache Varianzanalyse

Es soll gleichzeitig der Einfluß zweier Faktoren („Einflußgrößen“) auf ein beobachtetes Merkmal untersucht werden.

Beispiel: Betonqualität abhängig von Kornmischung des Schotters und der Wasserbeigabe.

- Kornverteilung: Faktor A (Zeilenfaktor)
mögliche Werte $\alpha_1, \dots, \alpha_a$
- Wasserbeigabe: Faktor B (Spaltenfaktor)
mögliche Werte β_1, \dots, β_b

Zu jeder Variante des Zeilenfaktors A und des Spaltenfaktors B werden gleich viele c Beobachtungen gemacht.

x_{ijk} : k -ter beobachteter Wert zur i -ten Variante von A und j -ten Variante von B.

Beispiel: Jeder von a Geodäten führt mit jedem von b Meßgeräten genau c Messungen am selben Objekt durch. Frage:

1. Gibt es bezüglich der Geodäten (Zeilenfaktor) bzw. zwischen den Geräten (Spaltenfaktor) Unterschiede?
2. Gibt es Geodäten, denen ein bestimmtes Gerät mehr liegt als ein anderes, d.h. ob es eine Wechselwirkung zwischen Geodäten und Messgeräten gibt.

Modell:

$$X_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + \epsilon_{ijk} \quad i = 1(1)a, j = 1(1)b, k = 1(1)c$$

Dabei bezeichnet γ_{ij} die Wechselwirkung.

Stochastisches Modell:

$$X_{ijk} = \mu + \alpha_i + \beta_j + \gamma_{ij} + U_{ijk} \quad U_{ijk} \sim N(0, \sigma^2)$$

Die Parameter genügen folgenden Bedingungen:

$$\begin{aligned} \mu \in \mathbf{R} \quad & \sum_{i=1}^a \alpha_i = 0 \quad \sum_{j=1}^b \beta_j = 0 \\ & \sum_{i=1}^a \gamma_{ij} = 0 \quad \forall j = 1(1)b \\ & \sum_{j=1}^b \gamma_{ij} = 0 \quad \forall i = 1(1)a \end{aligned}$$

Kopien 6-7 bis 6-11