



ulm university universität
uulm

Stochastik III

Vorlesungsskript

Prof. Dr. Evgeny Spodarev

ULM
2015

Inhaltsverzeichnis

1	Tests statistischer Hypothesen	3
1.1	Allgemeine Philosophie des Testens	3
1.2	Nichtrandomisierte Tests	12
1.2.1	Parametrische Signifikanztests	12
1.3	Randomisierte Tests	17
1.3.1	Grundlagen	17
1.3.2	Neyman-Pearson-Tests bei einfachen Hypothesen	19
1.3.3	Einseitige Neyman-Pearson-Tests	24
1.3.4	Unverfälschte zweiseitige Tests	30
1.4	Anpassungstests	36
1.4.1	χ^2 -Anpassungstest	36
1.4.2	χ^2 -Anpassungstest von Pearson-Fisher	42
1.4.3	Anpassungstest von Shapiro	48
1.5	Weitere, nicht parametrische Tests	50
1.5.1	Binomialtest	50
1.5.2	Iterationstests auf Zufälligkeit	52
2	Lineare Regression	55
2.1	Multivariate Normalverteilung	55
2.1.1	Eigenschaften der multivariaten Normalverteilung	59
2.1.2	Lineare und quadratische Formen von normalverteilten Zufallsvariablen	60
2.2	Multivariate lineare Regressionsmodelle mit vollem Rang	68
2.2.1	Methode der kleinsten Quadrate	68
2.2.2	Schätzer der Varianz σ^2	73
2.2.3	Maximum-Likelihood-Schätzer für β und σ^2	75
2.2.4	Tests für Regressionsparameter	78
2.2.5	Konfidenzbereiche	81
2.3	Multivariate lineare Regression mit $\text{Rang}(X) < m$	84
2.3.1	Verallgemeinerte Inverse	85
2.3.2	MKQ-Schätzer für β	86
2.3.3	Erwartungstreu schätzbare Funktionen	89
2.3.4	Normalverteilte Störgrößen	92
2.3.5	Hypothesentests	96
2.3.6	Konfidenzbereiche	98
2.3.7	Einführung in die Varianzanalyse	100

3	Verallgemeinerte lineare Modelle	103
3.1	Exponentialfamilie von Verteilungen	103
3.2	Linkfunktion	107
3.3	Maximum-Likelihood-Schätzung von β	109
3.4	Asymptotische Tests für β	116
3.5	Kriterien zur Modellwahl bzw. Modellanpassung	122
4	Hauptkomponentenanalyse	125
4.1	Einführung	125
4.2	Hauptkomponentenanalyse auf Modellebene	126
4.3	Hauptkomponentenanalyse auf Datenebene	134
4.4	Asymptotische Verteilung von HK bei normalverteilten Stichproben	138
4.5	Ausreißerererkennung	140
4.6	Hauptkomponentenanalyse und Regression	142
4.7	Numerische Berechnung der Hauptkomponenten	148
	Literatur	151

Vorwort

Dieses Skript entstand aus dem Zyklus der Vorlesungen über Statistik, die ich in den Jahren 2006-2010 an der Universität Ulm gehalten habe. Dabei handelt es sich um die aufbauende Vorlesung Stochastik III, die auf der Vorlesung *Stochastik I* basiert.

Ich möchte gerne meinen Kollegen aus dem Institut für Stochastik, Herrn Prof. Volker Schmidt und Herrn Dipl.-Math. Malte Spiess, für ihre Unterstützung und anregenden Diskussionen während der Entstehung des Skriptes danken. Herr Marco Baur hat eine hervorragende Arbeit beim Tippen des Skriptes und bei der Erstellung zahlreicher Abbildungen, die den Text begleiten, geleistet. Dafür gilt ihm mein herzlicher Dank.

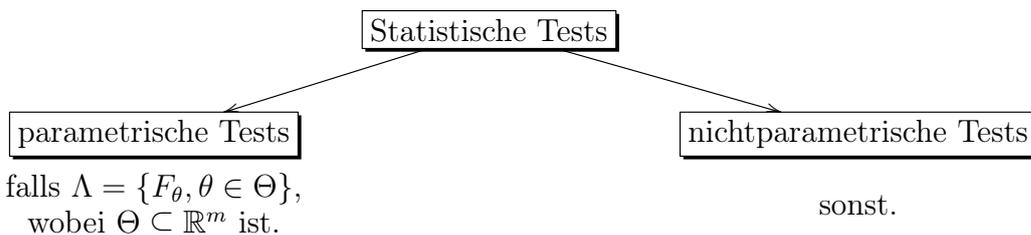
Ulm, den 20.10.2011
Evgeny Spodarev

1 Tests statistischer Hypothesen

In der Vorlesung Stochastik I haben wir schon Beispiele von statistischen Tests kennengelernt, wie etwa den Kolmogorow-Smirnow-Test (vergleiche Bemerkung 3.3.38, 3), Skript Stochastik I). Jetzt sollen statistische Signifikanztests formal eingeführt und ihre Eigenschaften untersucht werden.

1.1 Allgemeine Philosophie des Testens

Es sei eine Zufallsstichprobe (X_1, \dots, X_n) von unabhängigen, identisch verteilten Zufallsvariablen X_i gegeben, mit Verteilungsfunktion $F \in \Lambda$, wobei Λ eine Klasse von Verteilungsfunktionen ist. Es sei (x_1, \dots, x_n) eine konkrete Stichprobe, die als Realisierung von (X_1, \dots, X_n) interpretiert wird. In der Theorie des statistischen Testens werden Hypothesen über die Beschaffenheit der (unbekannten) Verteilungsfunktion F gestellt und geprüft. Dabei unterscheidet man



Bei parametrischen Tests prüft man, ob der Parameter θ bestimmte Werte annimmt (zum Beispiel $\theta = 0$). Bekannte Beispiele von nichtparametrischen Tests sind Anpassungstests, bei denen man prüft, ob die Verteilungsfunktion F gleich einer vorgegebenen Funktion F_0 ist.

Formalisieren wir zunächst den Begriff *Hypothese*. Die Menge Λ von zulässigen Verteilungsfunktionen F wird in zwei disjunkte Teilmengen Λ_0 und Λ_1 zerlegt, $\Lambda_0 \cup \Lambda_1 = \Lambda$. Die Aussage

„Man testet die *Haupthypothese* $H_0 : F \in \Lambda_0$ gegen die *Alternative* $H_1 : F \in \Lambda_1$,“

bedeutet, daß man an Hand der konkreten Stichprobe (x_1, \dots, x_n) versucht, eine Entscheidung zu fällen, ob die Verteilungsfunktion der Zufallsvariable X_i zu Λ_0 oder zu Λ_1 gehört. Dies passiert auf Grund einer statistischen *Entscheidungsregel*

$$\varphi : \mathbb{R}^n \rightarrow [0, 1],$$

die eine Statistik mit folgender Interpretation ist:

Der Stichprobenraum \mathbb{R}^n wird in drei disjunkte Bereiche K_0, K_{01} und K_1 unterteilt, sodaß $\mathbb{R}^n = K_0 \cup K_{01} \cup K_1$, wobei

$$\begin{aligned} K_0 &= \varphi^{-1}(\{0\}) &= \{x \in \mathbb{R}^n : \varphi(x) = 0\}, \\ K_1 &= \varphi^{-1}(\{1\}) &= \{x \in \mathbb{R}^n : \varphi(x) = 1\}, \\ K_{01} &= \varphi^{-1}((0, 1)) &= \{x \in \mathbb{R}^n : 0 < \varphi(x) < 1\}. \end{aligned}$$

Dementsprechend wird $H_0 : F \in \Lambda_0$

- verworfen, falls $\varphi(x) = 1$, also $x \in K_1$,
- nicht verworfen, falls $\varphi(x) = 0$, also $x \in K_0$;
- falls $\varphi(x) \in (0, 1)$, also $x \in K_{01}$, wird $\varphi(x)$ als Bernoulli-Wahrscheinlichkeit interpretiert, und es wird eine Zufallsvariable $Y \sim \text{Bernoulli}(\varphi(x))$ generiert, für die gilt:

$$Y = \begin{cases} 1 & \implies H_0 \text{ wird verworfen} \\ 0 & \implies H_0 \text{ wird nicht verworfen} \end{cases}$$

Falls $K_{01} \neq \emptyset$, wird eine solche Entscheidungsregel *randomisiert* genannt. Bei $K_{01} = \emptyset$, also $\mathbb{R}^n = K_0 \cup K_1$ spricht man dagegen von *nicht-randomisierten* Tests. Dabei heißt K_0 bzw. K_1 *Annahmehereich* bzw. *Ablehnungsbereich (kritischer Bereich)* von H_0 . K_{01} heißt *Randomisierungsbereich*.

Bemerkung 1.1.1. 1. Man sagt absichtlich „ H_0 wird nicht verworfen“, statt „ H_0 wird akzeptiert“, weil die schließende Statistik generell keine positiven, sondern nur negative Entscheidungen treffen kann. Dies ist generell ein philosophisches Problem der Falsifizierbarkeit von Hypothesen oder wissenschaftlichen Theorien, von denen aber keiner behaupten kann, daß sie der Wahrheit entsprechen (vergleiche die *wissenschaftliche Erkenntnistheorie von Karl Popper (1902-1994)*).

2. Die randomisierten Tests sind hauptsächlich von theoretischem Interesse (vergleiche Abschnitt 2.3). In der Praxis werden meistens nichtrandomisierte Regeln verwendet, bei denen man aus der Stichprobe (x_1, \dots, x_n) allein die Entscheidung über H_0 treffen kann. Hier gilt $\varphi(x) = \mathbb{1}_{K_1}, x = (x_1, \dots, x_n) \in \mathbb{R}^n$.

In diesem und in folgendem Abschnitt betrachten wir ausschließlich nichtrandomisierte Tests, um in Abschnitt 2.3 zu der allgemeinen Situation zurückzukehren.

Definition 1.1.1. Man sagt, daß die nicht-randomisierte Testregel $\varphi : \mathbb{R}^n \rightarrow \{0, 1\}$ einen (*nichtrandomisierten*) *statistischen Test zum Signifikanzniveau α* angibt, falls für $F \in \Lambda_0$ gilt

$$\mathbb{P}_F(\varphi(X_1, \dots, X_n) = 1) = P(H_0 \text{ verwerfen} \mid H_0 \text{ richtig}) \leq \alpha.$$

Definition 1.1.2. 1. Wenn man H_0 verwirft, obwohl H_0 richtig ist, begeht man den sogenannten *Fehler 1. Art*. Die Wahrscheinlichkeit

$$\alpha_n(F) = \mathbb{P}_F(\varphi(x_1, \dots, x_n) = 1), \quad F \in \Lambda_0$$

heißt die *Wahrscheinlichkeit des Fehlers 1. Art* und soll unter dem Niveau α bleiben.

2. Den *Fehler 2. Art* begeht man, wenn man die falsche Hypothese H_0 nicht verwirft. Dabei ist

$$\beta_n(F) = \mathbb{P}_F(\varphi(x_1, \dots, x_n) = 0), \quad F \in \Lambda_1$$

die *Wahrscheinlichkeit des Fehlers 2. Art*.

Eine Zusammenfassung aller Möglichkeiten wird in folgender Tabelle festgehalten:

	H_0 richtig	H_0 falsch
H_0 verwerfen	Fehler 1. Art, Wahrscheinlichkeit $\alpha_n(F) \leq \alpha$	richtige Entscheidung
H_0 nicht verwerfen	richtige Entscheidung	Fehler 2. Art mit Wahrscheinlichkeit $\beta_n(F)$

Dabei sollen α_n und β_n möglichst klein sein, was gegenläufige Tendenzen darstellt, weil beim Kleinwerden von α die Wahrscheinlichkeit des Fehlers 2. Art notwendigerweise wächst.

Definition 1.1.3. 1. Die Funktion

$$G_n(F) = \mathbb{P}_F(\varphi(X_1, \dots, X_n) = 1), \quad F \in \Lambda$$

heißt *Gütefunktion* eines Tests φ .

2. Die Einschränkung von G_n auf Λ_1 heißt *Stärke*, *Schärfe* oder *Macht* (englisch *power*) des Tests φ .

Es gilt

$$\begin{cases} G_n(F) = \alpha_n(F) \leq \alpha, & F \in \Lambda_0 \\ G_n(F) = 1 - \beta_n(F), & F \in \Lambda_1 \end{cases}$$

Beispiel 1.1.1. Parametrische Tests. Wie sieht ein parametrischer Test aus? Der Parameterraum Θ wird als $\Theta_0 \cup \Theta_1$ dargestellt, wobei $\Theta_0 \cap \Theta_1 = \emptyset$. Es gilt $\Lambda_0 = \{F_\theta : \theta \in \Theta_0\}$, $\Lambda_1 = \{F_\theta : \theta \in \Theta_1\}$. P_F wird zu P_θ , α_n , G_n und β_n werden statt auf Λ auf Θ definiert.

Welche Hypothesen H_0 und H_1 kommen oft bei parametrischen Tests vor? Zur Einfachheit betrachten wir den Spezialfall $\Theta = \mathbb{R}$.

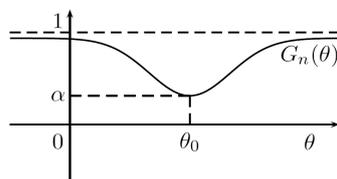
1. $H_0 : \theta = \theta_0$ vs. $H_1 : \theta \neq \theta_0$

2. $H_0 : \theta \geq \theta_0$ vs. $H_1 : \theta < \theta_0$
3. $H_0 : \theta \leq \theta_0$ vs. $H_1 : \theta > \theta_0$
4. $H_0 : \theta \in [a, b]$ vs. $H_1 : \theta \notin [a, b]$

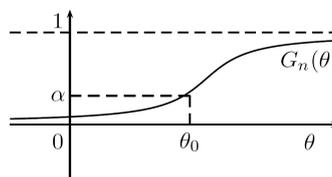
Im Fall (1) heißt der parametrische Test *zweiseitig*, in den Fällen (2) und (3) *einseitig* (*rechts-* bzw. *linksseitig*). In Fall (4) spricht man von der *Intervallhypothese* H_0 .

Bei einem zweiseitigen bzw. einseitigen Test kann die Gütefunktion wie in Abbildung 1.1 (a) bzw. 1.1 (b) aussehen,

Abbildung 1.1: Gütefunktion



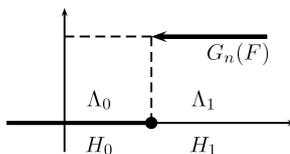
(a) eines zweiseitigen Tests



(b) eines einseitigen Tests

Bei einem allgemeinen (nicht notwendigerweise parametrischen) Modell kann man die ideale Gütefunktion wie in Abbildung 1.2 schematisch darstellen.

Abbildung 1.2: Schematische Darstellung der idealen Gütefunktion



- Man sieht aus Definition 1.1.2, dem Fehler 1. und 2. Art und der Ablehnungsregel, daß die Hypothesen H_0 und H_1 nicht symmetrisch behandelt werden, denn nur die

Wahrscheinlichkeit des Fehlers 1. Art wird kontrolliert. Dies ist der Grund dafür, daß Statistiker die eigentlich interessierende Hypothese nicht als H_0 , sondern als H_1 formulieren, damit, wenn man sich für H_1 entscheidet, man mit Sicherheit sagen kann, daß die Wahrscheinlichkeit der Fehlentscheidung unter dem Niveau α liegt.

- Wie wird ein statistischer, nicht randomisierter Test praktisch konstruiert? Die Konstruktion der Ablehnungsregel φ ähnelt sich sehr der von Konfidenzintervallen:
 1. Finde eine Teststatistik $T : \mathbb{R}^n \rightarrow \mathbb{R}$, die unter H_0 eine (möglicherweise asymptotisch für $n \rightarrow \infty$) bestimmte Prüfverteilung hat.
 2. Definiere $B_0 = [t_{\alpha_1}, t_{1-\alpha_2}]$, wobei t_{α_1} und $t_{1-\alpha_2}$ Quantile der Prüfverteilung von T sind, $\alpha_1 + \alpha_2 = \alpha \in [0, 1]$.
 3. Falls $T(X_1, \dots, X_n) \in \mathbb{R} \setminus B_0 = B_1$, setze $\varphi(X_1, \dots, X_n) = 1$. H_0 wird verworfen. Ansonsten setze $\varphi(X_1, \dots, X_n) = 0$.
- Falls die Verteilung von T nur asymptotisch bestimmt werden kann, so heißt φ *asymptotischer Test*.
- Sehr oft aber ist auch die asymptotische Verteilung von T nicht bekannt. Dann verwendet man sogenannte *Monte-Carlo Tests*, in denen dann Quantile t_α näherungsweise aus sehr vielen Monte-Carlo-Simulationen von T (unter H_0) bestimmt werden: Falls t^i , $i = 1, \dots, m$ die Werte von T in m unabhängigen Simulationsvorgängen sind, das heißt $t^i = T(x_1^i, \dots, x_n^i)$, x_j^i sind unabhängige Realisierungen von $X_j \sim F \in \Lambda_0$, $j = 1, \dots, n$, $i = 1, \dots, m$ dann bildet man ihre Ordnungsstatistiken $t^{(1)}, \dots, t^{(m)}$ und setzt $t_\alpha \approx t^{(\lfloor \alpha \cdot m \rfloor)}$, $\alpha \in [0, 1]$, wobei $t^{(0)} = -\infty$.

Bemerkung 1.1.2. Man sieht deutlich, daß aus einem beliebigen Konfidenzintervall

$$I_\theta = \left[I_1^\theta(X_1, \dots, X_n), I_2^\theta(X_1, \dots, X_n) \right]$$

zum Niveau $1 - \alpha$ für einen Parameter $\theta \in \mathbb{R}$ ein Test für θ konstruierbar ist. Die Hypothese $H_0 : \theta = \theta_0$ vs. $H_1 : \theta \neq \theta_0$ wird mit folgender Entscheidungsregel getestet:

$$\varphi(X_1, \dots, X_n) = 1, \text{ falls } \theta_0 \notin \left[I_1^{\theta_0}(X_1, \dots, X_n), I_2^{\theta_0}(X_1, \dots, X_n) \right].$$

Das Signifikanzniveau des Tests ist α .

Beispiel 1.1.2. *Normalverteilung, Test des Erwartungswertes bei bekannter Varianz.* Es seien

$$X_1, \dots, X_n \sim N(\mu, \sigma^2)$$

mit bekannter Varianz σ^2 . Ein Konfidenzintervall für μ ist

$$I^\mu = \left[I_1^\mu(X_1, \dots, X_n), I_2^\mu(X_1, \dots, X_n) \right] = \left[\bar{X}_n - \frac{z_{1-\alpha/2} \cdot \sigma}{\sqrt{n}}, \bar{X}_n + \frac{z_{1-\alpha/2} \cdot \sigma}{\sqrt{n}} \right]$$

(vergleiche Stochastik I, 4.2.1) H_0 wird verworfen, falls $|\mu_0 - \bar{X}_n| > \frac{z_{1-\alpha/2} \cdot \sigma}{\sqrt{n}}$. In der Testsprache bedeutet es, dass

$$\varphi(x_1, \dots, x_n) = \mathbb{I}((x_1, \dots, x_n) \in K_1),$$

wobei

$$K_1 = \left\{ (x_1, \dots, x_n) \in \mathbb{R}^n : |\mu_0 - \bar{x}_n| > \frac{\sigma z_{1-\alpha/2}}{\sqrt{n}} \right\}$$

der Ablehnungsbereich ist. Für die Teststatistik $T(X_1, \dots, X_n)$ gilt:

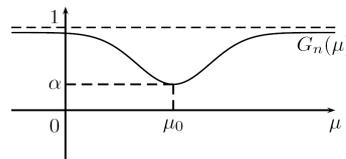
$$T(X_1, \dots, X_n) = \frac{\bar{X}_n - \mu_0}{\sigma} \sqrt{n} \sim N(0, 1) \mid \text{unter } H_0,$$

$$\alpha_n(\mu) = \alpha.$$

Berechnen wir nun die Gütefunktion (vergleiche Abbildung 1.3).

$$\begin{aligned} G_n(\mu) &= \mathbb{P}_\mu \left(|\mu_0 - \bar{X}_n| > \frac{z_{1-\alpha/2}}{\sqrt{n}} \right) = 1 - \mathbb{P}_\mu \left(|\bar{X}_n - \mu_0| \leq \frac{\sigma z_{1-\alpha/2}}{\sqrt{n}} \right) \\ &= 1 - \mathbb{P}_\mu \left(\left| \sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} + \frac{\mu - \mu_0}{\sigma} \sqrt{n} \right| \leq z_{1-\alpha/2} \right) \\ &= 1 - \mathbb{P}_\mu \left(-z_{1-\alpha/2} - \frac{\mu - \mu_0}{\sigma} \sqrt{n} \leq \sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} \leq z_{1-\alpha/2} - \frac{\mu - \mu_0}{\sigma} \sqrt{n} \right) \\ &= 1 - \Phi \left(z_{1-\alpha/2} - \frac{\mu - \mu_0}{\sigma} \sqrt{n} \right) + \Phi \left(-z_{1-\alpha/2} - \frac{\mu - \mu_0}{\sigma} \sqrt{n} \right) \\ &= \Phi \left(-z_{1-\alpha/2} + \frac{\mu - \mu_0}{\sigma} \sqrt{n} \right) + \Phi \left(-z_{1-\alpha/2} - \frac{\mu - \mu_0}{\sigma} \sqrt{n} \right). \end{aligned}$$

Abbildung 1.3: Gütefunktion für den zweiseitigen Test des Erwartungswertes einer Normalverteilung bei bekannter Varianz



Die „Ja-Nein“-Entscheidung des Testens wird oft als zu grob empfunden. Deswegen versucht man, ein feineres Maß der Verträglichkeit der Daten mit den Hypothesen H_0 und H_1 zu bestimmen. Dies ist der sogenannte p -Wert, der von den meisten Statistik-Softwarepaketen ausgegeben wird.

Definition 1.1.4. Es sei (x_1, \dots, x_n) die konkrete Stichprobe von Daten, die als Realisierung von (X_1, \dots, X_n) interpretiert wird und $T(X_1, \dots, X_n)$ die Teststatistik, mit deren Hilfe die Entscheidungsregel φ konstruiert wurde. Der p -Wert des statistischen Tests φ ist das kleinste Signifikanzniveau, zu dem der Wert $t = T(x_1, \dots, x_n)$ zur Verwerfung der Hypothese H_0 führt.

Im Beispiel eines einseitigen Tests mit $H_0 : \theta = \theta_0$ mit dem Ablehnungsbereich $B_1 = (t, \infty)$ sagt man grob, daß

$$p = \mathbb{P}(T(X_1, \dots, X_n) \geq t \mid H_0),$$

wobei die Anführungszeichen bedeuten, daß dies keine klassische, sondern eine bedingte Wahrscheinlichkeit ist, die später präzise angegeben wird.

Bei der Verwendung des p -Wertes verändert sich die Ablehnungsregel: die Hypothese $H_0 : \theta = \theta_0$ wird zum Signifikanzniveau α abgelehnt, falls $\alpha \geq p$. Früher hat man die Signifikanz der Testentscheidung (Ablehnung von H_0) an Hand folgender Tabelle festgesetzt:

p -Wert	Interpretation
$p \leq 0,001$	sehr stark signifikant
$0,001 < p \leq 0,01$	stark signifikant
$0,01 < p \leq 0,05$	schwach signifikant
$0,05 < p$	nicht signifikant

Da aber heute der p -Wert an sich verwendet werden kann, kann der Anwender der Tests bei vorgegebenem p -Wert selbst entscheiden, zu welchem Niveau er seine Tests durchführen will.

Bemerkung 1.1.3. 1. Das Signifikanzniveau darf nicht in Abhängigkeit von p festgelegt werden. Dies würde die allgemeine Testphilosophie zerstören!

2. Der p -Wert ist keine Wahrscheinlichkeit, sondern eine Zufallsvariable, denn er hängt von (X_1, \dots, X_n) ab. Der Ausdruck $p = \mathbb{P}(T(X_1, \dots, X_n) \geq t \mid H_0)$, der in Definition 1.1.4 für den p -Wert eines einseitigen Tests mit Teststatistik T gegeben wurde, soll demnach als *Überschreitungswahrscheinlichkeit* interpretiert werden, daß bei Wiederholung des Zufallsexperiments unter $H_0 : \theta = \theta_0$ der Wert $t = T(x_1, \dots, x_n)$ oder extremere Werte in Richtung der Hypothese H_1 betrachtet werden:

$$p = \mathbb{P}(T(X'_1, \dots, X'_n) \geq T(x_1, \dots, x_n) \mid H_0),$$

wobei $(X'_1, \dots, X'_n) \stackrel{d}{=} (X_1, \dots, X_n)$. Falls wir von einer konkreten Realisierung (x_1, \dots, x_n) zur Zufallsstichprobe (X_1, \dots, X_n) übergehen, erhalten wir

$$p = p(X_1, \dots, X_n) = \mathbb{P}(T(X'_1, \dots, X'_n) \geq T(X_1, \dots, X_n) \mid H_0)$$

3. Für andere Hypothesen H_1 wird der p -Wert auch eine andere Form haben. Zum Beispiel für

a) einen symmetrischen zweiseitigen Test ist

$$B_0 = [-t_{1-\alpha/2}, t_{1-\alpha/2}]$$

der Akzeptanzbereich für H_0 .

$$\Rightarrow p = P(|T(X'_1, \dots, X'_n)| \geq t | H_0), t = |T(X_1, \dots, X_n)|$$

b) einen rechtsseitigen Test mit $B_0 = [t_\alpha, \infty]$ gilt

$$p = P(T(X'_1, \dots, X'_n) \leq t | H_0), t = T(X_1, \dots, X_n)$$

c) Das Verhalten des p -Wertes kann folgendermaßen untersucht werden:

Lemma 1.1.1. Falls die Verteilungsfunktion F von X_i stetig und streng monoton steigend ist (die Verteilung von T ist absolut stetig mit zum Beispiel stetiger Dichte), dann ist $p \sim U[0, 1]$.

Beweis. Wir zeigen es am speziellen Beispiel des rechtsseitigen Tests.

$$\begin{aligned} \mathbb{P}(p \leq \alpha | H_0) &= \mathbb{P}(\overline{F}_T(T(X_1, \dots, X_n)) \leq \alpha | H_0) \\ &= \mathbb{P}(F_T(T(X_1, \dots, X_n)) \geq 1 - \alpha | H_0) \\ &= \mathbb{P}(U \geq 1 - \alpha) = 1 - (1 - \alpha) = \alpha, \quad \alpha \in [0, 1], \end{aligned}$$

da $F_T(T(X_1, \dots, X_n)) \stackrel{d}{=} U \sim U[0, 1]$ und F_T absolut stetig ist. \square

Übung 1.1.1. Zeigen Sie, daß für eine beliebige Zufallsvariable X mit absolut stetiger Verteilung und streng monoton steigender Verteilungsfunktion F_X gilt:

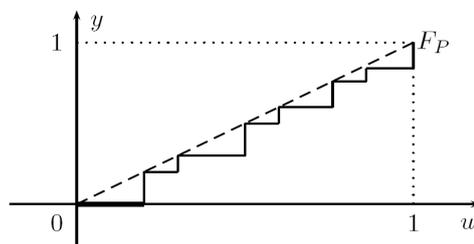
$$F_X(X) \sim U[0, 1]$$

Falls die Verteilung von T diskret ist, mit dem Wertebereich $\{t_1, \dots, t_n\}$, $t_i < t_j$ für $i < j$, so ist auch die Verteilung von p diskret, somit gilt nicht $p \sim U[0, 1]$. In diesem Fall ist $F_T(x)$ eine Treppenfunktion, die die Gerade $y = u$ in den Punkten $u = \sum_{i=1}^k \mathbb{P}(T(X_1, \dots, X_n) = t_i)$, $k = 1, \dots, n$ berührt (vgl. Abbildung 1.4).

Definition 1.1.5. 1. Falls die Macht $G_n(\cdot)$ eines Tests φ zum Niveau α die Ungleichung

$$G_n(F) \geq \alpha, \quad F \in \Lambda_1$$

erfüllt, dann heißt der Test *unverfälscht*.

Abbildung 1.4: Verteilung von p für diskrete T 

2. Es seien φ und φ^* zwei Tests zum Niveau α mit Gütefunktionen $G_n(\cdot)$ und $G_n^*(\cdot)$. Man sagt, daß der Test φ *besser* als φ^* ist, falls er eine größere Macht besitzt:

$$G_n(F) \geq G_n^*(F) \quad \forall F \in \Lambda_1$$

3. Der Test φ heißt konsistent, falls $G_n(F) \xrightarrow{n \rightarrow \infty} 1$ für alle $F \in \Lambda_1$.

Bemerkung 1.1.4. 1. Die einseitigen Tests haben oft eine größere Macht als ihre zweiseitigen Versionen.

Beispiel 1.1.3. Betrachten wir zum Beispiel den Gauß-Test des Erwartungswertes der Normalverteilung bei bekannter Varianz. Beim zweiseitigen Test

$$H_0 : \mu = \mu_0 \text{ vs. } H_1 : \mu \neq \mu_0.$$

erhalten wir die Gütefunktion

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

Beim einseitigen Test φ^* der Hypothesen

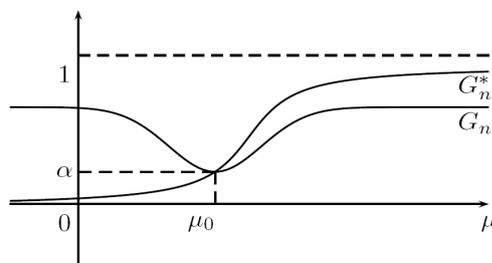
$$H_0^* : \mu \leq \mu_0 \text{ vs. } H_1^* : \mu > \mu_0$$

ist seine Gütefunktion gleich

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

Beide Tests sind offensichtlich konsistent, denn $G_n(\mu) \xrightarrow{n \rightarrow \infty} 1$, $G_n^*(\mu) \xrightarrow{n \rightarrow \infty} 1$. Dabei ist φ^* besser als φ . Beide Tests sind unverfälscht (vergleiche Abbildung 1.5).

Abbildung 1.5: Gütefunktionen eines ein- bzw. zweiseitigen Tests der Erwartungswertes einer Normalverteilung



2. Beim Testen einer Intervallhypothese $H_0 : \theta \in [a, b]$ vs. $H_1 : \theta \notin [a, b]$ zum Niveau α kann man wie folgt vorgehen: Teste

a) $H_0^a : \theta \geq a$ vs. $H_1^a : \theta < a$ zum Niveau $\alpha/2$.

b) $H_0^b : \theta \leq b$ vs. $H_1^b : \theta > b$ zum Niveau $\alpha/2$.

H_0 wird nicht abgelehnt, falls H_0^a und H_0^b nicht abgelehnt werden. Die Wahrscheinlichkeit des Fehlers 1. Art ist hier $\leq \alpha$. Die Macht dieses Tests ist im Allgemeinen schlecht.

3. Je mehr Parameter für den Aufbau der Teststatistik T geschätzt werden müssen, desto kleiner wird in der Regel die Macht.

1.2 Nichtrandomisierte Tests

1.2.1 Parametrische Signifikanztests

In diesem Abschnitt geben wir Beispiele einiger Tests, die meistens aus den entsprechenden Konfidenzintervallen für die Parameter von Verteilungen entstehen. Deshalb werden wir sie nur kurz behandeln.

1. Tests für die Parameter der Normalverteilung $N(\mu, \sigma^2)$

a) Test von μ bei unbekannter Varianz

- Hypothesen: $H_0 : \mu = \mu_0$ vs. $H_1 : \mu \neq \mu_0$.

- Teststatistik:

$$T(X_1, \dots, X_n) = \frac{\bar{X}_n - \mu_0}{S_n} \sim t_{n-1} \quad | H_0$$

- Entscheidungsregel:

$$\varphi(X_1, \dots, X_n) = 1, \text{ falls } |T(X_1, \dots, X_n)| > t_{n-1, 1-\alpha/2}.$$

b) Test von σ^2 bei bekanntem μ

- Hypothesen: $H_0 : \sigma^2 = \sigma_0^2$ vs. $H_1 : \sigma^2 \neq \sigma_0^2$.
- Teststatistik:

$$T(X_1, \dots, X_n) = \frac{n\tilde{S}_n^2}{\sigma_0^2} \sim \chi_n^2 \quad | H_0$$

$$\text{mit } \tilde{S}_n^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2.$$

- Entscheidungsregel:

$$\varphi(X_1, \dots, X_n) = 1, \text{ falls } T(X_1, \dots, X_n) \notin \left[\chi_{n, \alpha/2}^2, \chi_{n, 1-\alpha/2}^2 \right].$$

- Gütefunktion:

$$\begin{aligned} G_n(\sigma^2) &= 1 - \mathbb{P}_{\sigma^2} \left(\chi_{n, \alpha/2}^2 \leq \frac{n\tilde{S}_n^2}{\sigma_0^2} \leq \chi_{n, 1-\alpha/2}^2 \right) \\ &= 1 - \mathbb{P}_{\sigma^2} \left(\frac{\chi_{n, \alpha/2}^2 \sigma_0^2}{\sigma^2} \leq \frac{n\tilde{S}_n^2}{\sigma^2} \leq \frac{\chi_{n, 1-\alpha/2}^2 \sigma_0^2}{\sigma^2} \right) \\ &= 1 - F_{\chi_n^2} \left(\chi_{n, 1-\alpha/2}^2 \frac{\sigma_0^2}{\sigma^2} \right) + F_{\chi_n^2} \left(\chi_{n, \alpha/2}^2 \frac{\sigma_0^2}{\sigma^2} \right) \end{aligned}$$

c) Test von σ^2 bei unbekanntem μ

- Hypothesen: $H_0 : \sigma^2 = \sigma_0^2$ vs. $H_1 : \sigma^2 \neq \sigma_0^2$.
- Teststatistik:

$$T(X_1, \dots, X_n) = \frac{(n-1)S_n^2}{\sigma_0^2} \sim \chi_{n-1}^2 \quad | H_0,$$

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

- Entscheidungsregel:

$$\varphi(X_1, \dots, X_n) = 1, \text{ falls } T(X_1, \dots, X_n) \notin \left[\chi_{n-1, \alpha/2}^2, \chi_{n-1, 1-\alpha/2}^2 \right].$$

Übung 1.2.1. (i) Finden Sie $G_n(\cdot)$ für die einseitige Version der obigen Tests.

(ii) Zeigen Sie, daß diese einseitigen Tests unverfälscht sind, die zweiseitigen aber nicht.

2. Asymptotische Tests

Bei asymptotischen Tests ist die Verteilung der Teststatistik nur näherungsweise (für große n) bekannt. Ebenso asymptotisch wird das Konfidenzniveau α erreicht. Ihre Konstruktion basiert meistens auf Verwendung der Grenzwertsätze.

Die allgemeine Vorgehensweise wird im sogenannten *Wald-Test* (genannt nach dem Statistiker Abraham Wald (1902-1980)) fixiert:

- Sei (X_1, \dots, X_n) eine Zufallsstichprobe, X_i seien unabhängig und identisch verteilt für $i = 1, \dots, n$, mit $X_i \sim F_\theta$, $\theta \in \Theta \subseteq \mathbb{R}$.
- Wir testen $H_0 : \theta = \theta_0$ vs. $H_1 : \theta \neq \theta_0$. Es sei $\hat{\theta}_n = \hat{\theta}(X_1, \dots, X_n)$ ein erwartungstreuer, asymptotisch normalverteilter Schätzer für θ .

$$\frac{\hat{\theta}_n - \theta_0}{\hat{\sigma}_n} \xrightarrow[n \rightarrow \infty]{d} Y \sim N(0, 1) \quad | H_0,$$

wobei $\hat{\sigma}_n^2$ ein konsistenter Schätzer für die Varianz von $\hat{\theta}_n$ sei.

Die Teststatistik ist

$$T(X_1, \dots, X_n) = \frac{\hat{\theta}_n(X_1, \dots, X_n) - \theta_0}{\hat{\sigma}_n}.$$

- Die Entscheidungsregel lautet: H_0 wird abgelehnt, wenn $|T(X_1, \dots, X_n)| > z_{1-\alpha/2}$, wobei $z_{1-\alpha/2} = \Phi^{-1}(1 - \alpha/2)$. Diese Entscheidungsregel soll nur bei großen n verwendet werden. Die Wahrscheinlichkeit des Fehlers 1. Art ist asymptotisch gleich α , denn $\mathbb{P}(|T(X_1, \dots, X_n)| > z_{1-\alpha/2} | H_0) \xrightarrow[n \rightarrow \infty]{} \alpha$ wegen der asymptotischen Normalverteilung von T .

Die Gütefunktion des Tests ist asymptotisch gleich

$$\lim_{n \rightarrow \infty} G_n(\theta) = 1 - \Phi\left(z_{1-\alpha/2} + \frac{\theta_0 - \theta}{\sigma}\right) + \Phi\left(-z_{1-\alpha/2} + \frac{\theta_0 - \theta}{\sigma}\right),$$

wobei $\hat{\sigma}_n^2 \xrightarrow[n \rightarrow \infty]{P} \sigma^2$.

Spezialfälle des Wald-Tests sind asymptotische Tests der Erwartungswerte bei einer Poisson- oder Bernoulliverteilten Stichprobe.

Beispiel 1.2.1. a) Bernoulliverteilung

Es seien $X_i \sim \text{Bernoulli}(p)$, $p \in [0, 1]$ unabhängige, identisch verteilte Zufallsvariablen.

- Hypothesen: $H_0 : p = p_0$ vs. $H_1 : p \neq p_0$.
- Teststatistik:

$$T(X_1, \dots, X_n) = \begin{cases} \sqrt{n} \frac{\bar{X}_n - p_0}{\sqrt{\bar{X}_n(1-\bar{X}_n)}}, & \text{falls } \bar{X}_n \neq 0, 1, \\ 0, & \text{sonst.} \end{cases}$$

Unter H_0 gilt: $T(X_1, \dots, X_n) \xrightarrow[n \rightarrow \infty]{d} Y \sim N(0, 1)$.

b) Poissonverteilung

Es seien $X_i \sim \text{Poisson}(\lambda)$, $\lambda > 0$ unabhängige, identisch verteilte Zufallsvariablen.

- Hypothesen: $H_0 : \lambda = \lambda_0$ vs. $H_1 : \lambda \neq \lambda_0$
- Teststatistik:

$$T(X_1, \dots, X_n) = \begin{cases} \sqrt{n} \frac{\bar{X}_n - \lambda_0}{\sqrt{\bar{X}_n}}, & \text{falls } \bar{X}_n > 0, \\ 0, & \text{sonst.} \end{cases}$$

Unter H_0 gilt: $T(X_1, \dots, X_n) \xrightarrow[n \rightarrow \infty]{d} Y \sim N(0, 1)$

3. Zwei-Stichproben-Probleme

Gegeben seien zwei Zufallsstichproben

$$Y_1 = (X_{11}, \dots, X_{1n_1}), Y_2 = (X_{21}, \dots, X_{2n_2}), n = \max\{n_1, n_2\}.$$

X_{ij} seien unabhängig für $j = 1, \dots, n_i$, $X_{ij} \sim F_{\theta_i}$, $i = 1, 2$.

a) Test der Gleichheit zweier Erwartungswerte bei normalverteilten Stichproben

- bei bekannten Varianzen

Es seien $X_{ij} \sim N(\mu_i, \sigma_i^2)$, $i = 1, 2$, $j = 1, \dots, n$. Dabei seien σ_1^2, σ_2^2 bekannt, X_{ij} seien unabhängig voneinander für alle i, j .

Die Hypothesen sind $H_0 : \mu_1 = \mu_2$ vs. $H_1 : \mu_1 \neq \mu_2$. Wir betrachten die Teststatistik:

$$T(Y_1, Y_2) = \frac{\bar{X}_{1n_1} - \bar{X}_{2n_2}}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

Unter H_0 gilt: $T(Y_1, Y_2) \sim N(0, 1)$. Als Entscheidungsregel gilt: H_0 wird abgelehnt, falls $|T(Y_1, Y_2)| > z_{1-\alpha/2}$.

- **bei unbekanntem (jedoch gleichen) Varianzen**

Es seien $X_{ij} \sim N(\mu_i, \sigma_i^2)$, $i = 1, 2$, $j = 1, \dots, n$. Dabei seien σ_1^2, σ_2^2 unbekannt, $\sigma_1^2 = \sigma_2^2$ und X_{ij} seien unabhängig voneinander für alle i, j .

Die Hypothesen sind: $H_0 : \mu_1 = \mu_2$ vs. $H_1 : \mu_1 \neq \mu_2$. Wir betrachten die Teststatistik

$$T(Y_1, Y_2) = \frac{\bar{X}_{1n_1} - \bar{X}_{2n_2}}{S_{n_1 n_2}} \sqrt{\frac{n_1 n_2}{n_1 + n_2}},$$

wobei

$$S_{n_1 n_2}^2 = \frac{1}{n_1 + n_2 - 2} \cdot \left(\sum_{j=1}^{n_1} (X_{1j} - \bar{X}_{1n_1})^2 + \sum_{j=1}^{n_2} (X_{2j} - \bar{X}_{2n_2})^2 \right).$$

Man kann zeigen, daß unter H_0 gilt: $T(Y_1, Y_2) \sim t_{n_1+n_2-2}$. Die Entscheidungsregel lautet: H_0 ablehnen, falls $|T(Y_1, Y_2)| > t_{n_1+n_2-2, 1-\alpha/2}$.

b) **Test der Gleichheit von Erwartungswerten bei verbundenen Stichproben**

Es seien $Y_1 = (X_{11}, \dots, X_{1n})$ und $Y_2 = (X_{21}, \dots, X_{2n})$, $n_1 = n_2 = n$,

$$Z_j = X_{1j} - X_{2j} \sim N(\mu_1 - \mu_2, \sigma^2), \quad j = 1, \dots, n$$

unabhängig und identisch verteilt mit $\mu_i = \mathbb{E} X_{ij}$, $i = 1, 2$. Die Hypothesen sind: $H_0 : \mu_1 = \mu_2$ vs. $H_1 : \mu_1 \neq \mu_2$ bei unbekannter Varianz σ^2 . Als Teststatistik verwenden wir

$$T(Z_1, \dots, Z_n) = \sqrt{n} \frac{\bar{Z}_n}{S_n},$$

wobei

$$S_n^2 = \frac{1}{n-1} \sum_{j=1}^n (Z_j - \bar{Z}_n)^2.$$

Unter H_0 gilt dann: $T(Z_1, \dots, Z_n) \sim t_{n-1}$. Die Entscheidungsregel lautet: H_0 wird abgelehnt, falls $|T(z_1, \dots, z_n)| > t_{n-1, 1-\alpha/2}$.

c) **Test der Gleichheit von Varianzen bei unabhängigen Gaußschen Stichproben**

Es seien $Y_1 = (X_{11}, \dots, X_{1n_1})$ und $Y_2 = (X_{21}, \dots, X_{2n_2})$ unabhängig und identisch verteilt mit $X_{ij} \sim N(\mu_i, \sigma_i^2)$, wobei μ_i und σ_i^2 beide unbekannt sind. Die Hypothesen sind: $H_0 : \sigma_1^2 = \sigma_2^2$ vs. $H_1 : \sigma_1^2 \neq \sigma_2^2$. Als Teststatistik verwenden wir

$$T(Y_1, Y_2) = \frac{S_{2n_2}^2}{S_{1n_1}^2},$$

wobei

$$S_{in_i}^2 = \frac{1}{n_i - 1} \sum_{j=1}^n (X_{ij} - \bar{X}_{in_i})^2, \quad i = 1, 2.$$

Unter H_0 gilt: $T(Y_1, Y_2) \sim F_{n_2-1, n_1-1}$. Die Entscheidungsregel lautet: H_0 wird abgelehnt, falls $T(Y_1, Y_2) \notin [F_{n_2-1, n_1-1, \alpha/2}, F_{n_2-1, n_1-1, 1-\alpha/2}]$.

d) Asymptotische Zwei-Stichproben-Tests

- bei Bernoulli-verteilten Stichproben

Es gilt $X_{ij} \sim \text{Bernoulli}(p_i)$, $j = 1, \dots, n_i$, $i = 1, 2$. Die Hypothesen sind $H_0 : p_1 = p_2$ vs. $H_1 : p_1 \neq p_2$. Als Teststatistik verwenden wir

$$T(Y_1, Y_2) = \frac{(\bar{X}_{1n_1} - \bar{X}_{2n_2})(1 - \mathbb{I}(\bar{X}_{1n_1} = 0, \bar{X}_{2n_2} = 0))}{\sqrt{\frac{\bar{X}_{1n_1}(1-\bar{X}_{1n_1})}{n_1} + \frac{\bar{X}_{2n_2}(1-\bar{X}_{2n_2})}{n_2}}}$$

Unter H_0 gilt: $T(Y_1, Y_2) \xrightarrow[n_1, n_2 \rightarrow \infty]{d} Y \sim N(0, 1)$. Die Entscheidungsregel lautet: H_0 wird verworfen, falls $|T(Y_1, Y_2)| > z_{1-\alpha/2}$. Dies ist ein Test zum asymptotischen Signifikanzniveau α .

- bei Poisson-verteilten Stichproben

Es seien X_{ij} unabhängig, $X_{ij} \sim \text{Poisson}(\lambda_i)$, $i = 1, 2$. Die Hypothesen sind: $H_0 : \lambda_1 = \lambda_2$ vs. $H_1 : \lambda_1 \neq \lambda_2$. Als Teststatistik verwenden wir:

$$T(Y_1, Y_2) = \frac{\bar{X}_{1n_1} - \bar{X}_{2n_2}}{\sqrt{\frac{\bar{X}_{1n_1}}{n_1} + \frac{\bar{X}_{2n_2}}{n_2}}}$$

Die Entscheidungsregel lautet: H_0 ablehnen, falls $|T(Y_1, Y_2)| > z_{1-\alpha/2}$. Dies ist ein Test zum asymptotischen Niveau α .

Bemerkung 1.2.1. Asymptotische Tests dürfen nur für große Stichprobenumfänge verwendet werden. Bei ihrer Verwendung für kleine Stichproben kann das asymptotische Signifikanzniveau nicht garantiert werden.

1.3 Randomisierte Tests

In diesem Abschnitt werden wir klassische Ergebnisse von Neyman-Pearson über die besten Tests präsentieren. Dabei werden randomisierte Tests eine wichtige Rolle spielen.

1.3.1 Grundlagen

Gegeben sei eine Zufallsstichprobe (X_1, \dots, X_n) von unabhängigen und identisch verteilten Zufallsvariablen X_i mit konkreter Ausprägung (x_1, \dots, x_n) . Sei unser Stichprobenraum (B, \mathcal{B}) entweder $(\mathbb{R}^n, \mathcal{B}_{\mathbb{R}^n})$ oder $(\mathbb{N}_0^n, \mathcal{B}_{\mathbb{N}_0^n})$, je nachdem, ob die Stichprobenvariablen X_i , $i = 1, \dots, n$ absolut stetig oder diskret verteilt sind.

Hier wird zur Einfachheit im Falle einer diskret verteilten Zufallsvariable X_i ihr diskreter Wertebereich mit $\mathbb{N}_0 = \mathbb{N} \cup \{0\}$ gleichgesetzt. Der Wertebereich sei mit einem Maß μ versehen, wobei

$$\mu = \begin{cases} \text{Lebesgue-Ma\ss auf } \mathbb{R}, & \text{falls } X_i \text{ als stetig verteilt} \\ \text{Z\ahhlma\ss auf } \mathbb{N}_0, & \text{falls } X_i \text{ diskret verteilt.} \end{cases}$$

Dementsprechend gilt

$$\int g(x)\mu(dx) = \begin{cases} \int_{\mathbb{R}} g(x)dx, & \text{im absolut stetigen Fall,} \\ \sum_{x \in \mathbb{N}_0} g(x), & \text{im diskreten Fall.} \end{cases}$$

Es sei zus\atztlich $X_i \sim F_\theta$, $\theta \in \Theta \subseteq \mathbb{R}^m$, $i = 1, \dots, n$ (parametrisches Modell). F\ur $\Theta = \Theta_0 \cup \Theta_1$, $\Theta_0 \cap \Theta_1 = \emptyset$ formulieren wir die Hypothesen $H_0 : \theta \in \Theta_0$ vs. $H_1 : \theta \in \Theta_1$, die mit Hilfe eines randomisierten Tests

$$\varphi(x) = \begin{cases} 1, & x \in K_1, \\ \gamma \in (0, 1), & x \in K_{01} \\ 0, & x \in K_0 \end{cases} \quad x = (x_1, \dots, x_n),$$

getestet werden.

Im Falle $x \in K_{01}$ wird mit Hilfe einer Zufallsvariable $Y \sim \text{Bernoulli}(\varphi(x))$ entschieden, ob H_0 verworfen wird ($Y = 1$) oder nicht ($Y = 0$).

Definition 1.3.1. 1. Die *G\utefunktion* eines randomisierten Tests φ sei

$$G_n(\theta) = G_n(\varphi, \theta) = \mathbb{E}_\theta \varphi(X_1, \dots, X_n), \theta \in \Theta.$$

2. Der Test φ hat das *Signifikanzniveau* $\alpha \in [0, 1]$, falls $G_n(\varphi, \theta) \leq \alpha$, $\forall \theta \in \Theta_0$ ist. Die Zahl

$$\sup_{\theta \in \Theta_0} G_n(\varphi, \theta)$$

wird *Umfang* des Tests φ genannt. Offensichtlich ist der Umfang eines Niveau- α -Tests kleiner gleich α .

3. Sei $\Psi(\alpha)$ die Menge aller Tests zum Niveau α . Der Test $\varphi_1 \in \Psi(\alpha)$ ist (*gleichm\ai\ssig*) *besser* als Test $\varphi_2 \in \Psi(\alpha)$, falls $G_n(\varphi_1, \theta) \geq G_n(\varphi_2, \theta)$, $\theta \in \Theta_1$, also falls φ_1 eine gr\o\ssere Macht besitzt.
4. Ein Test $\varphi^* \in \Psi(\alpha)$ ist (*gleichm\ai\ssig*) *bester Test* in $\Psi(\alpha)$, falls

$$G_n(\varphi^*, \theta) \geq G_n(\varphi, \theta), \text{ f\ur alle Tests } \varphi \in \Psi(\alpha), \theta \in \Theta_1.$$

Bemerkung 1.3.1. 1. Definition 1.3.1 1) ist eine offensichtliche Verallgemeinerung der Definition 1.1.3 der Gütefunktion eines nicht-randomisierten Tests φ . Nämlich, für $\varphi(x) = \mathbb{I}(x \in K_1)$ gilt:

$$\begin{aligned} G_n(\varphi, \theta) &= \mathbb{E}_\theta \varphi(X_1, \dots, X_n) \\ &= \mathbb{P}_\theta((X_1, \dots, X_n) \in K_1) \\ &= \mathbb{P}_\theta(H_0 \text{ ablehnen}), \theta \in \Theta. \end{aligned}$$

2. Ein bester Test φ^* in $\Psi(\alpha)$ existiert nicht immer, sondern nur unter gewissen Voraussetzungen an $\mathbb{P}_\theta, \Theta_0, \Theta_1$ und $\Psi(\alpha)$.

1.3.2 Neyman-Pearson-Tests bei einfachen Hypothesen

In diesem Abschnitt betrachten wir einfache Hypothesen

$$H_0 : \theta = \theta_0 \quad \text{vs.} \quad H_1 : \theta = \theta_1 \tag{1.3.1}$$

wobei $\theta_0, \theta_1 \in \Theta, \theta_1 \neq \theta_0$.

Dementsprechend sind $\Theta_0 = \{\theta_0\}, \Theta_1 = \{\theta_1\}$. Wir setzen voraus, daß F_{θ_i} eine Dichte $g_i(x)$ bezüglich μ besitzt, $i = 0, 1$. Führen wir einige abkürzende Bezeichnungen $\mathbb{P}_0 = \mathbb{P}_{\theta_0}, \mathbb{P}_1 = \mathbb{P}_{\theta_1}, \mathbb{E}_0 = \mathbb{E}_{\theta_0}, \mathbb{E}_1 = \mathbb{E}_{\theta_1}$ ein. Sei $f_i(x) = \prod_{j=1}^n g_i(x_j), x = (x_1, \dots, x_n), i = 0, 1$ die Dichte der Stichprobe unter H_0 bzw. H_1 .

Definition 1.3.2. Ein *Neyman-Pearson-Test (NP-Test)* der einfachen Hypothesen in (1.3.1) ist gegeben durch die Regel

$$\varphi(x) = \varphi_K(x) = \begin{cases} 1, & \text{falls } f_1(x) > K f_0(x), \\ \gamma, & \text{falls } f_1(x) = K f_0(x), \\ 0, & \text{falls } f_1(x) < K f_0(x) \end{cases} \tag{1.3.2}$$

für Konstanten $K > 0$ und $\gamma \in [0, 1]$.

Bemerkung 1.3.2. 1. Manchmal werden $K = K(x)$ und $\gamma = \gamma(x)$ als Funktionen von x und nicht als Konstanten betrachtet.

2. Der *Ablehnungsbereich* des Neyman-Pearson-Tests φ_K ist

$$K_1 = \{x \in B : f_1(x) > K f_0(x)\}.$$

3. Der *Umfang* des Neyman-Pearson-Tests φ_K ist

$$\begin{aligned} \mathbb{E}_0 \varphi_K(X_1, \dots, X_n) &= \mathbb{P}_0(f_1(X_1, \dots, X_n) > K f_0(X_1, \dots, X_n)) \\ &\quad + \gamma \mathbb{P}_0(f_1(X_1, \dots, X_n) = K f_0(X_1, \dots, X_n)) \end{aligned}$$

4. Die Definition 1.3.2 kann man äquivalent folgendermaßen geben: Wir definieren eine Teststatistik

$$T(x) = \begin{cases} \frac{f_1(x)}{f_0(x)}, & x \in B : f_0(x) > 0, \\ \infty, & x \in B : f_0(x) = 0. \end{cases}$$

Dann wird der neue Test

$$\tilde{\varphi}_K(x) = \begin{cases} 1, & \text{falls } T(x) > K, \\ \gamma, & \text{falls } T(x) = K, \\ 0, & \text{falls } T(x) < K \end{cases}$$

eingeführt, der für P_0 - und P_1 - fast alle $x \in B$ äquivalent zu φ_K ist. In der Tat gilt $\varphi_K(x) = \tilde{\varphi}_K(x) \forall x \in B \setminus C$, wobei $C = \{x \in B : f_0(x) = f_1(x) = 0\}$ das \mathbb{P}_0 - bzw. \mathbb{P}_1 -Maß Null besitzt.

In der neuen Formulierung ist der Umfang von φ bzw. $\tilde{\varphi}_K$ gleich

$$\mathbb{E}_0 \tilde{\varphi}_K = \mathbb{P}_0(T(X_1, \dots, X_n) > K) + \gamma \cdot \mathbb{P}_0(T(X_1, \dots, X_n) = K).$$

Satz 1.3.1. Optimalitätssatz

Es sei φ_K ein Neyman-Pearson-Test für ein $K > 0$ und $\gamma \in [0, 1]$. Dann ist φ_K der beste Test zum Niveau $\alpha = \mathbb{E}_0 \varphi_K$ seines Umfangs.

Beweis. Sei $\varphi \in \Psi(\alpha)$, also $\mathbb{E}_0(\varphi(X_1, \dots, X_n)) \leq \alpha$. Um zu zeigen, daß φ_K besser als φ ist, genügt es bei einfachen Hypothesen H_0 und H_1 zu zeigen, daß $\mathbb{E}_1 \varphi_K(X_1, \dots, X_n) \geq \mathbb{E}_1 \varphi(X_1, \dots, X_n)$. Wir führen dazu die folgenden Mengen ein:

$$M^+ = \{x \in B : \varphi_K(x) > \varphi(x)\}$$

$$M^- = \{x \in B : \varphi_K(x) < \varphi(x)\}$$

$$M^= = \{x \in B : \varphi_K(x) = \varphi(x)\}$$

Es gilt offensichtlich $x \in M^+ \Rightarrow \varphi_K(x) > 0 \Rightarrow f_1(x) \geq K f_0(x)$,

$$x \in M^- \Rightarrow \varphi_K(x) < 1 \Rightarrow f_1(x) \leq K f_0(x) \text{ und } B = M^+ \cup M^- \cup M^=.$$

Als Folgerung erhalten wir

$$\begin{aligned} \mathbb{E}_1(\varphi_K(X_1, \dots, X_n) - \varphi(X_1, \dots, X_n)) &= \int_B (\varphi_K(x) - \varphi(x)) f_1(x) \mu(dx) \\ &= \left(\int_{M^+} + \int_{M^-} + \int_{M^=} \right) (\varphi_K(x) - \varphi(x)) f_1(x) \mu(dx) \\ &\geq \int_{M^+} (\varphi_K(x) - \varphi(x)) K f_0(x) \mu(dx) \\ &\quad + \int_{M^-} (\varphi_K(x) - \varphi(x)) K f_0(x) \mu(dx) \\ &= \int_B (\varphi_K(x) - \varphi(x)) K f_0(x) \mu(dx) \\ &= K [\mathbb{E}_0 \varphi_K(X_1, \dots, X_n) - \mathbb{E}_0 \varphi(X_1, \dots, X_n)] \\ &\geq K(\alpha - \alpha) = 0, \end{aligned}$$

weil beide Tests das Niveau α haben. Damit ist die Behauptung bewiesen. \square

Bemerkung 1.3.3. 1. Da im Beweis γ nicht vorkommt, wird derselbe Beweis im Falle von $\gamma(x) \neq \text{const}$ gelten.

2. Aus dem Beweis folgt die Gültigkeit der Ungleichung

$$\int_B (\varphi_K(x) - \varphi(x)) (f_1(x) - K f_0(x)) \mu(dx) \geq 0$$

im Falle des konstanten K , bzw.

$$\mathbb{E}_1 (\varphi_K(X_1, \dots, X_n) - \varphi(X_1, \dots, X_n)) \geq \int_B (\varphi_K(x) - \varphi(x)) K f_0(x) \mu(dx)$$

im allgemeinen Fall.

Satz 1.3.2. (Fundamentallemma von Neyman-Pearson)

1. Zu einem beliebigen $\alpha \in (0, 1)$ gibt es einen Neyman-Pearson-Test φ_K mit Umfang α , der dann nach Satz 1.3.1 der beste Niveau- α -Test ist.
2. Ist φ ebenfalls bester Test zum Niveau α , so gilt $\varphi(x) = \varphi_K(x)$ für μ -fast alle $x \in K_0 \cup K_1 = \{x \in B : f_1(x) \neq K f_0(x)\}$ und φ_K aus Teil 1).

Beweis. 1. Für $\varphi_K(x)$ gilt

$$\varphi_K(x) = \begin{cases} 1, & \text{falls } x \in K_1 = \{x : f_1(x) > K \cdot f_0(x)\}, \\ \gamma, & \text{falls } x \in K_{01} = \{x : f_1(x) = K \cdot f_0(x)\}, \\ 0, & \text{falls } x \in K_0 = \{x : f_1(x) < K \cdot f_0(x)\}. \end{cases}$$

Der Umfang von φ_K ist

$$\mathbb{P}_0(T(X_1, \dots, X_n) > K) + \gamma \mathbb{P}_0(T(X_1, \dots, X_n) = K) = \alpha, \quad (1.3.3)$$

wobei

$$T(x_1, \dots, x_n) = \begin{cases} \frac{f_1(x_1, \dots, x_n)}{f_0(x_1, \dots, x_n)}, & \text{falls } f_0(x_1, \dots, x_n) > 0, \\ \infty, & \text{sonst.} \end{cases}$$

Nun suchen wir ein $K > 0$ und ein $\gamma \in [0, 1]$, sodaß Gleichung (1.3.3) stimmt. Es sei $\tilde{F}_0(x) = \mathbb{P}_0(T(X_1, \dots, X_n) \leq x)$, $x \in \mathbb{R}$ die Verteilungsfunktion von T . Da $T \geq 0$ ist, gilt $\tilde{F}_0(x) = 0$, falls $x < 0$. Außerdem ist $\mathbb{P}_0(T(X_1, \dots, X_n) < \infty) = 1$, das heißt $\tilde{F}_0^{-1}(\alpha) \in [0, \infty)$, $\alpha \in (0, 1)$. Die Gleichung (1.3.3) kann dann folgendermaßen umgeschrieben werden:

$$1 - \tilde{F}_0(K) + \gamma (\tilde{F}_0(K) - \tilde{F}_0(K-)) = \alpha, \quad (1.3.4)$$

wobei $\tilde{F}_0(K-) = \lim_{x \rightarrow K-0} \tilde{F}_0(x)$.

Sei $K = \tilde{F}_0^{-1}(1 - \alpha)$, dann gilt:

- a) Falls K ein Stetigkeitspunkt von \tilde{F}_0 ist, ist Gleichung (1.3.4) erfüllt für alle $\gamma \in [0, 1]$, zum Beispiel $\gamma = 0$.
- b) Falls K kein Stetigkeitspunkt von \tilde{F}_0 ist, dann ist $\tilde{F}_0(K) - \tilde{F}_0(K-) > 0$, woraus folgt

$$\gamma = \frac{\alpha - 1 + \tilde{F}_0(K)}{\tilde{F}_0(K) - \tilde{F}_0(K-)}$$

\Rightarrow es gibt einen Neyman-Pearson-Test zum Niveau α .

2. Wir definieren $M^\neq = \{x \in B : \varphi(x) \neq \varphi_K(x)\}$. Es muss gezeigt werden, daß

$$\mu\left((K_0 \cup K_1) \cap M^\neq\right) = 0.$$

Dazu betrachten wir

$$\begin{aligned} \mathbb{E}_1 \varphi(X_1, \dots, X_n) - \mathbb{E}_1 \varphi_K(X_1, \dots, X_n) &= 0 && (\varphi \text{ und } \varphi_K \text{ sind beste Tests}) \\ \mathbb{E}_0 \varphi(X_1, \dots, X_n) - \mathbb{E}_0 \varphi_K(X_1, \dots, X_n) &\leq 0 && (\varphi \text{ und } \varphi_K \text{ sind } \alpha\text{-Tests} \\ &&& \text{mit Umfang von } \varphi_K = \alpha) \end{aligned}$$

$$\Rightarrow \int_B (\varphi - \varphi_K) \cdot (f_1 - K \cdot f_0) \mu(dx) \geq 0.$$

In Bemerkung 1.3.3 wurde bewiesen, daß

$$\begin{aligned} \int_B (\varphi - \varphi_K)(f_1 - K \cdot f_0) d\mu &\leq 0 \\ \Rightarrow \int_B (\varphi - \varphi_K)(f_1 - K \cdot f_0) d\mu &= 0 = \int_{M^\neq \cap (K_0 \cup K_1)} (\varphi - \varphi_K)(f_1 - K f_0) d\mu. \end{aligned}$$

Es gilt $\mu(M^\neq \cap (K_0 \cup K_1)) = 0$, falls der Integrand $(\varphi_K - \varphi)(f_1 - K f_0) > 0$ auf M^\neq ist. Wir zeigen, daß

$$(\varphi_K - \varphi)(f_1 - K f_0) > 0 \text{ für } x \in M^\neq \quad (1.3.5)$$

ist. Es gilt

$$\begin{aligned} f_1 - K f_0 > 0 &\Rightarrow \varphi_K - \varphi > 0, \\ f_1 - K f_0 < 0 &\Rightarrow \varphi_K - \varphi < 0, \end{aligned}$$

weil

$$\begin{aligned} f_1(x) > K f_0(x) &\Rightarrow \varphi_K(x) = 1 \\ &\text{und mit } \varphi(x) < 1 \Rightarrow \varphi_K(x) - \varphi(x) > 0 \text{ auf } M^\neq. \\ f_1(x) < K f_0(x) &\Rightarrow \varphi_K(x) = 0 \\ &\text{und mit } \varphi(x) > 0 \Rightarrow \varphi_K(x) - \varphi(x) < 0 \text{ auf } M^\neq. \end{aligned}$$

Daraus folgt die Gültigkeit der Ungleichung (1.3.5) und somit

$$\mu \left((K_0 \cup K_1) \cap M^{\neq} \right) = 0.$$

□

Bemerkung 1.3.4. Falls φ und φ_K beste α -Neyman-Pearson-Tests sind, dann sind sie P_0 - bzw. P_1 - fast sicher gleich.

Beispiel 1.3.1 (Neyman-Pearson-Test für den Parameter der Poissonverteilung). Es sei (X_1, \dots, X_n) eine Zufallsstichprobe mit $X_i \sim \text{Poisson}(\lambda)$, $\lambda > 0$, wobei X_i unabhängig und identisch verteilt sind für $i = 1, \dots, n$. Wir testen die Hypothesen $H_0 : \lambda = \lambda_0$ vs. $H_1 : \lambda = \lambda_1$. Dabei ist

$$g_i(x) = e^{-\lambda_i} \frac{\lambda_i^x}{x!}, \quad x \in \mathbb{N}_0, \quad i = 0, 1,$$

$$f_i(x) = f_i(x_1, \dots, x_n) = \prod_{j=1}^n g_i(x_j) = \prod_{j=1}^n e^{-\lambda_i} \frac{\lambda_i^{x_j}}{x_j!} = e^{-n\lambda_i} \cdot \frac{\lambda_i^{\sum_{j=1}^n x_j}}{(x_1! \cdot \dots \cdot x_n!)}$$

für $i = 0, 1$. Die Neyman-Pearson-Teststatistik ist

$$T(x_1, \dots, x_n) = \begin{cases} \frac{f_1(x)}{f_0(x)} = e^{-n(\lambda_1 - \lambda_0)} \cdot (\lambda_1/\lambda_0)^{\sum_{j=1}^n x_j}, & \text{falls } x_1, \dots, x_n \in \mathbb{N}_0, \\ \infty, & \text{sonst.} \end{cases}$$

Die Neyman-Pearson-Entscheidungsregel lautet

$$\varphi_K(x_1, \dots, x_n) = \begin{cases} 1, & \text{falls } T(x_1, \dots, x_n) > K, \\ \gamma, & \text{falls } T(x_1, \dots, x_n) = K, \\ 0, & \text{falls } T(x_1, \dots, x_n) < K. \end{cases}$$

Wir wählen $K > 0$, $\gamma \in [0, 1]$, sodaß φ_K den Umfang α hat. Dazu lösen wir

$$\alpha = \mathbb{P}_0(T(X_1, \dots, X_n) > K) + \gamma \mathbb{P}_0(T(X_1, \dots, X_n) = K)$$

bezüglich γ und K auf.

$$\begin{aligned} \mathbb{P}_0(T(X_1, \dots, X_n) > K) &= \mathbb{P}_0(\log T(X_1, \dots, X_n) > \log K) \\ &= \mathbb{P}_0 \left(-n(\lambda_1 - \lambda_0) + \sum_{j=1}^n X_j \cdot \log \left(\frac{\lambda_1}{\lambda_0} \right) > \log K \right) = \mathbb{P}_0 \left(\sum_{j=1}^n X_j > A \right) \end{aligned}$$

$$\text{wobei } A := \left\lceil \frac{\log K + n \cdot (\lambda_1 - \lambda_0)}{\log \frac{\lambda_1}{\lambda_0}} \right\rceil,$$

falls zum Beispiel $\lambda_1 > \lambda_0$. Im Falle $\lambda_1 < \lambda_0$ ändert sich das $>$ auf $<$ in der Wahrscheinlichkeit.

Wegen der Faltungsstabilität der Poissonverteilung ist unter H_0

$$\sum_{j=1}^n X_j \sim \text{Poisson}(n\lambda_0),$$

also wählen wir K als minimale, nichtnegative Zahl, für die gilt: $\mathbb{P}_0\left(\sum_{j=1}^n X_j > A\right) \leq \alpha$, und setzen

$$\gamma = \frac{\alpha - \mathbb{P}_0(\sum_{j=1}^n X_j > A)}{\mathbb{P}_0(\sum_{j=1}^n X_j = A)},$$

wobei

$$\begin{aligned} \mathbb{P}_0\left(\sum_{j=1}^n X_j > A\right) &= 1 - \sum_{j=0}^A e^{-\lambda_0 n} \frac{(\lambda_0 n)^j}{j!}, \\ \mathbb{P}_0\left(\sum_{j=1}^n X_j = A\right) &= e^{-\lambda_0 n} \frac{(\lambda_0 n)^A}{A!}. \end{aligned}$$

Somit haben wir die Parameter K und γ gefunden und damit einen Neyman-Pearson-Test φ_K konstruiert.

1.3.3 Einseitige Neyman-Pearson-Tests

Bisher betrachteten wir Neyman-Pearson-Tests für einfache Hypothesen der Form $H_i : \theta = \theta_i$, $i = 0, 1$. In diesem Abschnitt wollen wir einseitige Neyman-Pearson-Tests einführen, für Hypothesen der Form $H_0 : \theta \leq \theta_0$ vs. $H_1 : \theta > \theta_0$.

Zunächst konstruieren wir einen Test für diese Hypothesen: Sei (X_1, \dots, X_n) eine Zufallsstichprobe, X_i seien unabhängig und identisch verteilt mit

$$X_i \sim F_\theta \in \Lambda = \{F_\theta : \theta \in \Theta\},$$

wobei $\Theta \subset \mathbb{R}$ offen ist und Λ eindeutig parametrisiert, das heißt

$$\theta \neq \theta' \Rightarrow F_\theta \neq F_{\theta'}.$$

Ferner besitze F_θ eine Dichte g_θ bezüglich des Lebesgue-Maßes (bzw. Zählmaßes) auf \mathbb{R} (bzw. \mathbb{N}_0). Dann ist

$$f_\theta(x) = \prod_{j=1}^n g_\theta(x_j), \quad x = (x_1, \dots, x_n)$$

eine Dichte von (X_1, \dots, X_n) bezüglich μ auf B .

Definition 1.3.3. Eine Verteilung auf B mit Dichte f_θ gehört zur Klasse von Verteilungen mit monotonen Dichtekoeffizienten in T , falls es für alle $\theta < \theta'$ eine Funktion $h : \mathbb{R} \times \Theta^2 \rightarrow \mathbb{R} \cup \infty$, die monoton wachsend in $t \in \mathbb{R}$ ist und eine Statistik $T : B \rightarrow \mathbb{R}$ gibt, mit der Eigenschaft

$$\frac{f_{\theta'}(x)}{f_\theta(x)} = h(T(x), \theta, \theta'),$$

wobei

$$h(T(x), \theta, \theta') = \infty \quad \text{für alle } x \in B : f_\theta(x) = 0, f_{\theta'}(x) > 0.$$

Der Fall $f_\theta(x) = f_{\theta'}(x) = 0$ tritt mit \mathbb{P}_θ - bzw. $\mathbb{P}_{\theta'}$ -Wahrscheinlichkeit 0 auf.

Definition 1.3.4. Es sei Q_θ eine Verteilung auf (B, \mathcal{B}) mit der Dichte f_θ bzgl. μ . Q_θ gehört zur *einparametrischen Exponentialklasse* ($\theta \in \Theta \subset \mathbb{R}$ offen), falls die Dichte folgende Form hat:

$$f_\theta(x) = \exp \{c(\theta) \cdot T(x) + a(\theta)\} \cdot l(x), \quad x = (x_1, \dots, x_n) \in B,$$

wobei $c(\theta)$ eine monoton steigende Funktion ist, und $\text{Var}_\theta T(X_1, \dots, X_n) > 0$, $\theta \in \Theta$.

Lemma 1.3.1. Verteilungen aus der einparametrischen Exponentialfamilie besitzen einen monotonen Dichtekoeffizienten.

Beweis. Es sei Q_θ aus der einparametrischen Exponentialfamilie mit der Dichte

$$f_\theta(x) = \exp \{c(\theta) \cdot T(x) + a(\theta)\} \cdot l(x).$$

Für $\theta < \theta'$ ist dann

$$\frac{f_{\theta'}(x)}{f_\theta(x)} = \exp \{(c(\theta') - c(\theta)) \cdot T(x) + a(\theta') - a(\theta)\}$$

monoton bezüglich T , weil $c(\theta') - c(\theta) > 0$ wegen der Monotonie von $c(\theta)$. Also besitzt f_θ einen monotonen Dichtekoeffizienten. \square

Beispiel 1.3.2. 1. *Normalverteilte StichprobenvARIABLEN*

Es seien $X_i \sim N(\mu, \sigma_0^2)$, $i = 1, \dots, n$, unabhängige, identisch verteilte Zufallsvariablen, mit unbekanntem Parameter μ und bekannter Varianz σ_0^2 (Hier wird μ für die Bezeichnung des Erwartungswertes von X_i und nicht des Maßes auf \mathbb{R}^n verwendet.

(wie früher)). Die Dichte des Zufallsvektors $X = (X_1, \dots, X_n)^\top$ ist gleich

$$\begin{aligned}
 f_\mu(x) &= \prod_{i=1}^n g_\mu(x_i) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma_0^2}} e^{-\frac{(x_i-\mu)^2}{2\sigma_0^2}} \\
 &= \frac{1}{(2\pi\sigma_0^2)^{n/2}} \exp\left\{-\frac{1}{2\sigma_0^2} \sum_{i=1}^n (x_i - \mu)^2\right\} \\
 &= \frac{1}{(2\pi\sigma_0^2)^{n/2}} \exp\left\{-\frac{1}{2\sigma_0^2} \left(\sum_{i=1}^n x_i^2 - 2\mu \sum_{i=1}^n x_i + \mu^2 n\right)\right\} \\
 &= \exp\left(\underbrace{\frac{\mu}{\sigma_0^2} \cdot \sum_{i=1}^n x_i}_{c(\mu)} - \underbrace{\frac{\mu^2 n}{2\sigma_0^2}}_{a(\mu)}\right) \cdot \underbrace{\frac{1}{(2\pi\sigma_0^2)^{n/2}} \exp\left(-\frac{\sum_{i=1}^n x_i^2}{2\sigma_0^2}\right)}_{l(x)}.
 \end{aligned}$$

Also gehört $N(\mu, \sigma_0^2)$ zur einparametrischen Exponentialklasse mit $c(\mu) = \frac{\mu}{\sigma_0^2}$ und $T(x) = \sum_{i=1}^n x_i$.

2. Binomialverteilte Stichprobenvariablen

Es seien $X_i \sim \text{Bin}(k, p)$ unabhängig und identisch verteilt, $i = 1, \dots, n$. Der Parameter p sei unbekannt. Die Zähldichte des Zufallsvektors $X = (X_1, \dots, X_n)^\top$ ist

$$\begin{aligned}
 f_p(x) &= \mathbb{P}_p(X_i = x_i, i = 1, \dots, n) \\
 &= \prod_{i=1}^n \binom{k}{x_i} p^{x_i} (1-p)^{k-x_i} = p^{\sum_{i=1}^n x_i} \cdot \frac{(1-p)^{nk}}{(1-p)^{\sum_{i=1}^n x_i}} \cdot \prod_{i=1}^n \binom{k}{x_i} \\
 &= \exp\left\{\underbrace{\left(\sum_{i=1}^n x_i\right)}_{T(x)} \cdot \underbrace{\log\left(\frac{p}{1-p}\right)}_{c(p)} + \underbrace{nk \cdot \log(1-p)}_{a(p)}\right\} \cdot \underbrace{\prod_{i=1}^n \binom{k}{x_i}}_{l(x)},
 \end{aligned}$$

also gehört $\text{Bin}(n, p)$ zur einparametrischen Exponentialklasse mit

$$c(p) = \log\left(\frac{p}{1-p}\right)$$

und

$$T(x) = \sum_{i=1}^n x_i.$$

Lemma 1.3.2. Falls φ_K der Neyman-Pearson-Test der Hypothesen $H_0 : \theta = \theta_0$ vs. $H_1 : \theta = \theta_1$ ist, dann gilt:

$$\mu(\underbrace{\{x \in B : f_1(x) \neq K f_0(x)\}}_{K_0 \cup K_1}) > 0.$$

Beweis. Wegen $\theta_0 \neq \theta_1$ und der eindeutigen Parametrisierung gilt $f_0 \neq f_1$ auf einer Menge mit μ -Maß > 0 .

Nun sei $\mu(K_0 \cup K_1) = 0$. Daraus folgt, daß $f_1(x) = K \cdot f_0(x)$ μ -fast sicher. Das heißt

$$1 = \int_B f_1(x) dx = K \cdot \int_B f_0(x) dx,$$

woraus folgt, daß $K = 1$ und $f_1(x) = f_0(x)$ μ -fast sicher, was aber ein Widerspruch zur eindeutigen Parametrisierung ist. \square

Im Folgenden sei (X_1, \dots, X_n) eine Stichprobe von unabhängigen, identisch verteilten Zufallsvariablen mit $X_i \sim$ Dichte g_θ , $i = 1, \dots, n$ und

$$(X_1, \dots, X_n) \sim \text{Dichte } f_\theta(x) = \prod_{i=1}^n g_\theta(x_i)$$

aus der Klasse der Verteilungen mit monotonen Dichtekoeffizienten und einer Statistik $T(X_1, \dots, X_n)$.

Wir betrachten die Hypothesen $H_0 : \theta \leq \theta_0$ vs. $H_1 : \theta > \theta_0$ und den Neyman-Pearson-Test:

$$\varphi_{K^*}^*(x) = \begin{cases} 1, & \text{falls } T(x) > K^*, \\ \gamma^*, & \text{falls } T(x) = K^*, \\ 0, & \text{falls } T(x) < K^* \end{cases} \quad (1.3.6)$$

für $K^* \in \mathbb{R}$ und $\gamma^* \in [0, 1]$. Die Gütefunktion von $\varphi_{K^*}^*$ bei θ_0 ist

$$G_n(\theta_0) = \mathbb{E}_0 \varphi_{K^*}^* = \mathbb{P}_0(T(X_1, \dots, X_n) > K^*) + \gamma^* \cdot \mathbb{P}_0(T(X_1, \dots, X_n) = K^*)$$

Satz 1.3.3. 1. Falls $\alpha = \mathbb{E}_0 \varphi_{K^*}^* > 0$, dann ist der soeben definierte Neyman-Pearson-Test ein bester Test der einseitigen Hypothesen H_0 vs. H_1 zum Niveau α .

2. Zu jedem Konfidenzniveau $\alpha \in (0, 1)$ gibt es ein $K^* \in \mathbb{R}$ und $\gamma^* \in [0, 1]$, sodaß $\varphi_{K^*}^*$ ein bester Test zum Umfang α ist.

3. Die Gütefunktion $G_n(\theta)$ von $\varphi_{K^*}^*(\theta)$ ist monoton wachsend in θ . Falls $0 < G_n(\theta) < 1$, dann ist sie sogar streng monoton wachsend.

Beweis. 1. Wähle $\theta_1 > \theta_0$ und betrachte die einfachen Hypothesen $H'_0 : \theta = \theta_0$ und $H'_1 : \theta = \theta_1$. Sei

$$\varphi_K(x) = \begin{cases} 1, & f_1(x) > K f_0(x), \\ \gamma, & f_1(x) = K f_0(x), \\ 0, & f_1(x) < K f_0(x) \end{cases}$$

der Neyman-Pearson-Test für H'_0, H'_1 mit $K > 0$. Da f_θ den monotonen Dichtekoeffizienten mit Statistik T besitzt,

$$\frac{f_1(x)}{f_0(x)} = h(T(x), \theta_0, \theta_1),$$

existiert ein $K > 0$, so dass

$$\left\{ x : \frac{f_1(x)}{f_0(x)} > K \right\} \subset \left\{ T(x) > K^* \right\} \quad \text{mit } K = h(K^*, \theta_0, \theta_1).$$

φ_K ist ein bester Neyman-Pearson-Test zum Niveau $\alpha = \mathbb{E}_0 \varphi_K = \mathbb{E}_0 \varphi_{K^*}$. Aus $\alpha > 0$ folgt $K < \infty$, denn aus $K = \infty$ würde folgen

$$\begin{aligned} 0 < \alpha = \mathbb{E}_0 \varphi_K &\leq \mathbb{P}_0 \left(\frac{f_1(X_1, \dots, X_n)}{f_0(X_1, \dots, X_n)} \geq K^* \right) \leq \mathbb{P}_0 \left(\frac{f_1(X_1, \dots, X_n)}{f_0(X_1, \dots, X_n)} = \infty \right) \\ &= \mathbb{P}_0 (f_1(X_1, \dots, X_n) > 0, f_0(X_1, \dots, X_n) = 0) \\ &= \int_B \mathbb{I}(f_1(x) > 0, f_0(x) = 0) \cdot f_0(x) \mu(dx) = 0. \end{aligned}$$

Für den Test φ_{K^*} aus (1.3.6) gilt dann

$$\varphi_{K^*}^*(x) = \begin{cases} 1, & \text{falls } f_1(x)/f_0(x) > K, \\ \gamma^*(x), & \text{falls } f_1(x)/f_0(x) = K, \\ 0, & \text{falls } f_1(x)/f_0(x) < K, \end{cases}$$

wobei $\gamma^*(x) \in \{\gamma^*, 0, 1\}$. Daraus folgt, daß $\varphi_{K^*}^*$ ein bester Neyman-Pearson-Test ist für H'_0 vs. H'_1 (vergleiche Bemerkung 1.3.2, 1.) und Bemerkung 1.3.3) für beliebige $\theta_1 > \theta_0$. Deshalb ist $\varphi_{K^*}^*$ ein bester Neyman-Pearson-Test für $H''_0 : \theta = \theta_0$ vs. $H''_1 : \theta > \theta_0$ ist.

Die selbe Behauptung erhalten wir aus dem Teil 3. des Satzes für $H_0 : \theta \leq \theta_0$ vs. $H_1 : \theta > \theta_0$, weil dann $G_n(\theta) \leq G_n(\theta_0) = \alpha$ für alle $\theta < \theta_0$.

2. Siehe Beweis zu Satz 1.3.2, 1.).

3. Wir müssen zeigen, daß $G_n(\theta)$ monoton ist. Dazu wählen wir $\theta_1 < \theta_2$ und zeigen, daß $\alpha_1 = G_n(\theta_1) \leq G_n(\theta_2)$. Wir betrachten die neuen, einfachen Hypothesen $H''_0 : \theta = \theta_1$ vs. $H''_1 : \theta = \theta_2$. Der Test $\varphi_{K^*}^*$ kann genauso wie in 1. als Neyman-Pearson-Test dargestellt werden (für die Hypothesen H''_0 und H''_1), der ein bester Test zum Niveau α_1 ist. Betrachten wir einen weiteren konstanten Test $\varphi(x) = \alpha_1$. Dann ist $\alpha_1 = \mathbb{E}_{\theta_2} \varphi \leq \mathbb{E}_{\theta_2} \varphi_{K^*}^* = G_n(\theta_2)$. Daraus folgt, daß $G_n(\theta_1) \leq G_n(\theta_2)$.

Nun zeigen wir, daß für $G_n(\theta) \in (0, 1)$ gilt: $G_n(\theta_1) < G_n(\theta_2)$. Wir nehmen an, daß $\alpha_1 = G_n(\theta_1) = G_n(\theta_2)$ und $\theta_1 < \theta_2$ für $\alpha \in (0, 1)$. Es folgt, daß $\varphi(x) = \alpha_1$ auch ein bester Test für H''_0 und H''_1 ist. Aus Satz 1.3.2, 2.) folgt

$$\mu(\underbrace{\{x \in B : \varphi(x) \neq \varphi_{K^*}^*(x)\}}_{=\alpha_1}) = 0 \text{ auf } K_0 \cup K_1 = \{f_1(x) \neq K f_0(x)\},$$

was ein Widerspruch zur Bauart des Tests φ_{K^*} ist, der auf $K_0 \cup K_1$ nicht gleich $\alpha_1 \in (0, 1)$ sein kann. □

Bemerkung 1.3.5. 1. Der Satz 1.3.3 ist genauso auf Neyman-Pearson-Tests der einseitigen Hypothesen

$$H_0 : \theta \geq \theta_0 \text{ vs. } H_1 : \theta < \theta_0$$

anwendbar, mit dem entsprechenden Unterschied

$$\begin{aligned} \theta &\mapsto -\theta \\ T &\mapsto -T \end{aligned}$$

Somit existiert der beste α -Test auch in diesem Fall.

2. Man kann zeigen, daß die Gütefunktion $G_n(\varphi_{K^*}^*, \theta)$ des besten Neyman-Pearson-Tests auf $\Theta_0 = (-\infty, \theta_0)$ folgende Minimalitätseigenschaft besitzt:

$$G_n(\varphi_{K^*}^*, \theta) \leq G_n(\varphi, \theta) \quad \forall \varphi \in \Psi(\alpha), \theta \leq \theta_0$$

Beispiel 1.3.3. Wir betrachten eine normalverteilte Stichprobe (X_1, \dots, X_n) von unabhängigen und identisch verteilten Zufallsvariablen X_i , wobei $X_i \sim N(\mu, \sigma_0^2)$ und σ_0^2 sei bekannt. Es werden die Hypothesen

$$H_0 : \mu \leq \mu_0 \text{ vs. } H_1 : \mu > \mu_0,$$

getestet. Aus Beispiel 1.1.2 kennen wir die Testgröße

$$T(X_1, \dots, X_n) = \sqrt{n} \frac{\bar{X}_n - \mu_0}{\sigma_0},$$

wobei unter H_0 gilt: $T(X_1, \dots, X_n) \sim N(0, 1)$. H_0 wird verworfen, falls

$$T(X_1, \dots, X_n) > z_{1-\alpha}, \quad \text{wobei } \alpha \in (0, 1).$$

Wir zeigen jetzt, daß dieser Test der beste Neyman-Pearson-Test zum Niveau α ist. Aus Beispiel 1.3.2 ist bekannt, daß die Dichte f_n von (X_1, \dots, X_n) aus der einparametrischen Exponentialklasse ist, mit

$$\tilde{T}(X_1, \dots, X_n) = \sum_{i=1}^n X_i.$$

Dann gehört f_μ von (x_1, \dots, x_n) zur einparametrischen Exponentialklasse auch bezüglich der Statistik

$$T(X_1, \dots, X_n) = \sqrt{n} \frac{\bar{X}_n - \mu}{\sigma_0}$$

Es gilt nämlich

$$\begin{aligned} f_\mu(x) &= \exp\left(\underbrace{\frac{\mu}{\sigma_0^2}}_{\tilde{c}(\mu)} \cdot \underbrace{\sum_{i=1}^n x_i}_{\tilde{T}} - \underbrace{\frac{\mu^2 n}{2\sigma_0^2}}_{\tilde{a}(\mu)}\right) \cdot l(x) \\ &= \exp\left(\underbrace{\frac{\mu\sqrt{n}}{\sigma_0}}_{c(\mu)} \cdot \underbrace{\sqrt{n}\bar{x}_n - \mu}_{T} + \underbrace{\frac{\mu^2 n}{2\sigma_0^2}}_{a(\mu)}\right) \cdot l(x). \end{aligned}$$

Die Statistik T kann also in der Konstruktion des Neyman-Pearson-Tests (Gleichung (1.3.6)) verwendet werden:

$$\varphi_{K^*}(x) = \begin{cases} 1, & \text{falls } T(x) > z_{1-\alpha}, \\ 0, & \text{falls } T(x) = z_{1-\alpha}, \\ 0, & \text{falls } T(x) < z_{1-\alpha} \end{cases}$$

(mit $K^* = z_{1-\alpha}$ und $\gamma^* = 0$). Nach Satz 1.3.3 ist dieser Test der beste Neyman-Pearson-Test zum Niveau α für unsere Hypothesen:

$$\begin{aligned} G_n(\varphi_{K^*}, \mu_0) &= \mathbb{P}_0(T(X_1, \dots, X_n) > z_{1-\alpha}) + 0 \cdot \mathbb{P}_0(T(X_1, \dots, X_n) \leq z_{1-\alpha}) \\ &= 1 - \Phi(z_{1-\alpha}) = 1 - (1 - \alpha) = \alpha. \end{aligned}$$

1.3.4 Unverfälschte zweiseitige Tests

Es sei (X_1, \dots, X_n) eine Stichprobe von unabhängigen und identisch verteilten Zufallsvariablen mit der Dichte

$$f_\theta(x) = \prod_{i=1}^n g_\theta(x_i).$$

Es wird ein zweiseitiger Test der Hypothesen

$$H_0 : \theta = \theta_0 \text{ vs. } H_1 : \theta \neq \theta_0$$

betrachtet. Für alle $\alpha \in (0, 1)$ kann es jedoch keinen besten Neyman-Pearson-Test φ zum Niveau α für H_0 vs. H_1 geben. Denn, nehmen wir an, φ wäre der beste Test zum Niveau α für H_0 vs. H_1 , dann wäre φ der beste Test für die Hypothesen

1. $H'_0 : \theta = \theta_0$ vs. $H'_1 : \theta > \theta_0$
2. $H''_0 : \theta = \theta_0$ vs. $H''_1 : \theta < \theta_0$.

Dann ist nach Satz 1.3.3, 3. die Gütefunktion

1. $G_n(\varphi, \theta) < \alpha$ auf $\theta < \theta_0$, bzw.

2. $G_n(\varphi, \theta) > \alpha$ auf $\theta < \theta_0$,

was ein Widerspruch ist!

Darum werden wir die Klasse aller möglichen Tests auf unverfälschte Niveau- α -Tests (Definition 1.1.5) eingrenzen. Der Niveau- α -Test φ ist unverfälscht genau dann, wenn

$$G_n(\varphi, \theta) \leq \alpha \text{ für } \theta \in \Theta_0$$

$$G_n(\varphi, \theta) \geq \alpha \text{ für } \theta \in \Theta_1$$

Beispiel 1.3.4. 1. $\varphi(x) \equiv \alpha$ ist unverfälscht.

2. Der zweiseitige Gauß-Test ist unverfälscht, vergleiche Beispiel 1.1.2: $G_n(\varphi, \mu) \geq \alpha$ für alle $\mu \in \mathbb{R}$.

Im Folgenden seien X_i unabhängig und identisch verteilt. Die Dichte f_θ von (X_1, \dots, X_n) gehöre zur einparametrischen Exponentialklasse:

$$f_\theta(x) = \exp \{c(\theta) \cdot T(x) + a(\theta)\} \cdot l(x), \quad (1.3.7)$$

wobei $c(\theta)$ und $a(\theta)$ stetig differenzierbar auf Θ sein sollen, mit

$$c'(\theta) > 0 \quad \text{und} \quad \text{Var}_\theta T(X_1, \dots, X_n) > 0$$

für alle $\theta \in \Theta$. Sei $f_\Phi(x)$ stetig in (x, Θ) auf $B \times \Theta$.

Übungsaufgabe 1.3.1. Zeigen Sie, daß folgende Relation gilt:

$$a'(\theta) = -c'(\theta) \mathbb{E}_\theta T(X_1, \dots, X_n).$$

Lemma 1.3.3. Es sei φ ein unverfälschter Test zum Niveau α für

$$H_0 : \theta = \theta_0 \text{ vs. } H_1 : \theta \neq \theta_0.$$

Dann gilt:

1. $\alpha = \mathbb{E}_0 \varphi(X_1, \dots, X_n) = G_n(\varphi, \theta_0)$
2. $\mathbb{E}_0 [T(X_1, \dots, X_n) \varphi(X_1, \dots, X_n)] = \alpha \cdot \mathbb{E}_0 T(X_1, \dots, X_n)$

Beweis. 1. Die Gütefunktion von φ ist

$$G_n(\varphi, \theta) = \int_B \varphi(x) f_\theta(x) \mu(dx)$$

Da f_θ aus der einparametrischen Exponentialklasse ist, ist $G_n(\varphi, \theta)$ differenzierbar (unter dem Integral) bezüglich θ . Wegen der Unverfälschtheit von φ gilt

$$G_n(\varphi, \theta_0) \leq \alpha, \quad G_n(\varphi, \theta) \geq \alpha, \quad \theta \neq \theta_0$$

und daraus folgt $G_n(\varphi, \theta_0) = \alpha$ und θ_0 ist ein Minimumpunkt von G_n . Somit ist 1) bewiesen.

2. Da θ_0 der Minimumpunkt von G_n ist, gilt

$$\begin{aligned} 0 &= G'_n(\varphi, \theta_0) = \int_B \varphi(x)(c'(\theta_0)T(x) + a'(\theta_0))f_0(x)\mu(dx) \\ &= c'(\theta_0) \cdot \mathbb{E}_0 [\varphi(X_1, \dots, X_n)T(X_1, \dots, X_n)] + a'(\theta) \cdot G_n(\varphi, \theta_0) \\ &= c'(\theta_0) \cdot \mathbb{E}_0 [\varphi(X_1, \dots, X_n)T(X_1, \dots, X_n)] + \alpha a'(\theta_0) \\ &\stackrel{\text{(Übung 1.3.1)}}{=} c'(\theta_0) (\mathbb{E}_0 (\varphi \cdot T) - \alpha \mathbb{E}_0 T) \end{aligned}$$

Daraus folgt $\mathbb{E}_0 (\varphi T) = \alpha \mathbb{E}_0 T$ und damit ist das Lemma bewiesen. \square

Wir definieren jetzt die modifizierten Neyman-Pearson-Tests für einfache Hypothesen

$$H_0 : \theta = \theta_0 \text{ vs. } H_1' : \theta = \theta_1, \quad \theta_1 \neq \theta_0.$$

Für $\lambda, K \in \mathbb{R}$, $\gamma : B \rightarrow [0, 1]$ definieren wir

$$\varphi_{K,\lambda}(x) = \begin{cases} 1, & \text{falls } f_1(x) > (K + \lambda T(x))f_0(x), \\ \gamma(x), & \text{falls } f_1(x) = (K + \lambda T(x))f_0(x), \\ 0, & \text{falls } f_1(x) < (K + \lambda T(x))f_0(x), \end{cases} \quad (1.3.8)$$

wobei $T(x)$ die Statistik aus der Darstellung (1.3.7) ist.

Es sei $\tilde{\Psi}(\alpha)$ die Klasse aller Tests, die Aussagen 1) und 2) des Lemmas 1.3.3 erfüllen. Aus Lemma 1.3.3 folgt dann, daß die Menge der unverfälschten Tests zum Niveau α eine Teilmenge von $\tilde{\Psi}(\alpha)$ ist.

Satz 1.3.4. Der modifizierte Neyman-Pearson-Test $\varphi_{K,\lambda}$ ist der beste α -Test in $\tilde{\Psi}(\alpha)$ für Hypothesen H_0 vs. H_1' zum Niveau $\alpha = \mathbb{E}_0 \varphi_{K,\lambda}$, falls $\varphi_{K,\lambda} \in \tilde{\Psi}(\alpha)$.

Beweis. Es ist zu zeigen, daß $\mathbb{E}_1 \varphi_{K,\lambda} \geq \mathbb{E}_1 \varphi$ für alle $\varphi \in \tilde{\Psi}(\alpha)$, bzw. $\mathbb{E}_1 (\varphi_{K,\lambda} - \varphi) \geq 0$. Es gilt

$$\begin{aligned} \mathbb{E}_1 (\varphi_{K,\lambda} - \varphi) &= \int_B (\varphi_{K,\lambda}(x) - \varphi(x))f_1(x)\mu(dx) \\ &\stackrel{\text{(Bem. 1.3.3, 2.)}}{\geq} \int_B (\varphi_{K,\lambda}(x) - \varphi(x))(K + \lambda T(x))f_0(x)\mu(dx) \\ &= K \left(\underbrace{\mathbb{E}_0 \varphi_{K,\lambda}}_{=\alpha} - \underbrace{\mathbb{E}_0 \varphi}_{=\alpha} \right) + \lambda \left(\underbrace{\mathbb{E}_0 (\varphi_{K,\lambda} \cdot T)}_{\alpha \mathbb{E}_0 T} - \underbrace{\mathbb{E}_0 (\varphi \cdot T)}_{=\alpha \mathbb{E}_0 T} \right) \\ &= 0, \end{aligned}$$

weil $\varphi, \varphi_{K,\lambda} \in \tilde{\Psi}(\alpha)$. \square

Wir definieren folgende Entscheidungsregel, die später zum Testen der zweiseitigen Hypothesen

$$H_0 : \theta = \theta_0 \text{ vs. } H_1 : \theta \neq \theta_0$$

verwendet wird:

$$\varphi_c(x) = \begin{cases} 1, & \text{falls } T(x) \notin [c_1, c_2], \\ \gamma_1, & \text{falls } T(x) = c_1, \\ \gamma_2, & \text{falls } T(x) = c_2, \\ 0, & \text{falls } T(x) \in (c_1, c_2), \end{cases} \quad (1.3.9)$$

für $c_1 \leq c_2 \in \mathbb{R}$, $\gamma_1, \gamma_2 \in [0, 1]$ und die Statistik $T(x)$, $x = (x_1, \dots, x_n) \in B$, die in der Dichte (1.3.7) vorkommt. Zeigen wir, daß φ_c sich als modifizierter Neyman-Pearson-Test schreiben lässt.

Für die Dichte

$$f_\theta(x) = \exp\{c(\theta)T(x) + a(\theta)\} \cdot l(x)$$

wird (wie immer) vorausgesetzt, daß $l(x) > 0$, $c'(\theta) > 0$ und $a'(\theta)$ existiert für $\theta \in \Theta$.

Lemma 1.3.4. Es sei (X_1, \dots, X_n) eine Stichprobe von unabhängigen, identisch verteilten Zufallsvariablen mit gemeinsamer Dichte $f_\theta(x)$, $x \in B$, die zur einparametrischen Exponentialfamilie gehört. Sei $T(x)$ die dazugehörige Statistik, die im Exponenten der Dichte f_θ vorkommt. Für beliebige reelle Zahlen $c_1 \leq c_2$, $\gamma_1, \gamma_2 \in [0, 1]$ und Parameterwerte $\theta_0, \theta_1 \in \Theta$: $\theta_0 \neq \theta_1$ läßt sich der Test φ_c aus (1.3.9) als modifizierter Neyman-Pearson-Test $\varphi_{K,\lambda}$ aus (1.3.8) mit gegebenen $K, \lambda \in \mathbb{R}$, $\gamma(x) \in [0, 1]$ schreiben.

Beweis. Falls wir die Bezeichnung

$$f_{\theta_i}(x) = f_i(x), \quad i = 0, 1$$

verwenden, dann gilt

$$\frac{f_1(x)}{f_0(x)} = \exp \left\{ \underbrace{(c(\theta_1) - c(\theta_0))}_{c} T(x) + \underbrace{a(\theta_1) - a(\theta_0)}_a \right\},$$

und somit

$$\{x \in B : f_1(x) > (K + \lambda T(x)) f_0(x)\} = \{x \in B : \exp(cT(x) + a) > K + \lambda T(x)\}.$$

Finden wir solche K und λ aus \mathbb{R} , für die die Gerade $K + \lambda t$, $t \in \mathbb{R}$ die konvexe Kurve $\exp\{ct + a\}$ genau an den Stellen c_1 und c_2 schneidet (falls $c_1 \neq c_2$) bzw. an der Stelle $t = c_1$ berührt (falls $c_1 = c_2$). Dies ist immer möglich, siehe Abbildung 1.1.

Ferner setzen wir $\gamma(x) = \gamma_i$ für $\{x \in B : T(x) = c_i\}$. Insgesamt gilt dann

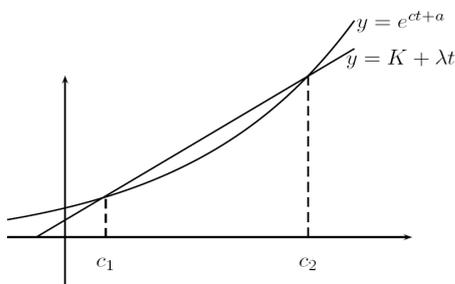
$$\{x : \exp(cT(x) + a) > K + \lambda T(x)\} = \{x : T(x) \notin [c_1, c_2]\}$$

und

$$\{x : \exp(cT(x) + a) < K + \lambda T(x)\} = \{x : T(x) \in (c_1, c_2)\}.$$

Damit ist das Lemma bewiesen. □

Abbildung 1.1:



Bemerkung 1.3.6. 1. Die Umkehrung des Lemmas stimmt nicht, denn bei vorgegebenen Kurven $y = K + \lambda t$ und $y = \exp\{ct + a\}$ muss es die Schnittpunkte c_1 und c_2 nicht unbedingt geben. So kann die Gerade vollständig unter der Kurve $y = \exp\{ct + a\}$ liegen.

2. Der Test φ_c macht von den Werten θ_0 und θ_1 nicht explizit Gebrauch. Dies unterscheidet ihn vom Test $\varphi_{K,\lambda}$, für den die Dichten f_0 und f_1 gebraucht werden.

Jetzt sind wir bereit, den Hauptsatz über zweiseitige Tests zum Prüfen der Hypothesen

$$H_0 : \theta = \theta_0 \text{ vs. } H_1 : \theta \neq \theta_0$$

zu formulieren und zu beweisen.

Satz 1.3.5. (Hauptsatz über zweiseitige Tests)

Unter den Voraussetzungen des Lemmas 1.3.4 sei φ_c ein Test aus (1.3.9), für den $\varphi_c \in \tilde{\Psi}(\alpha)$ gilt. Dann ist φ_c bester unverfälschter Test zum Niveau α (und dadurch bester Test in $\tilde{\Psi}(\alpha)$) der Hypothesen

$$H_0 : \theta = \theta_0 \text{ vs. } H_1 : \theta \neq \theta_0.$$

Beweis. Wählen wir ein beliebiges $\theta_1 \in \Theta$, $\theta_1 \neq \theta_0$. Nach Lemma 1.3.4 ist φ_c ein modifizierter Neyman-Pearson-Test $\varphi_{K,\lambda}$ für eine spezielle Wahl von K und $\lambda \in \mathbb{R}$. $\varphi_{K,\lambda}$ ist aber nach Satz 1.3.4 bester Test in $\tilde{\Psi}(\alpha)$ für $H_0 : \theta = \theta_0$ vs. $H_1' : \theta = \theta_1$. Da φ_c nicht von θ_1 abhängt, ist es bester Test in $\tilde{\Psi}(\alpha)$ für $H_1 : \theta \neq \theta_0$. Da unverfälschte Niveau- α -Tests in $\tilde{\Psi}(\alpha)$ liegen, müssen wir nur zeigen, daß φ_c unverfälscht ist. Da φ_c der beste Test ist, ist er nicht schlechter als der konstante unverfälschte Test $\varphi = \alpha$, das heißt

$$G_n(\varphi_c, \theta) \geq G_n(\varphi, \theta) = \alpha, \quad \theta \neq \theta_0.$$

Somit ist auch φ_c unverfälscht. Der Beweis ist beendet. \square

Bemerkung 1.3.7. Wir haben gezeigt, daß φ_c der beste Test seines Umfangs ist. Es wäre jedoch noch zu zeigen, daß für beliebiges $\alpha \in (0, 1)$ Konstanten $c_1, c_2, \gamma_1, \gamma_2$ gefunden werden, für die $\mathbb{E}_0 \varphi_c = \alpha$ gilt. Da der Beweis schwierig ist, wird er hier ausgelassen. Im folgenden Beispiel jedoch wird es klar, wie die Parameter $c_1, c_2, \gamma_1, \gamma_2$ zu wählen sind.

Beispiel 1.3.5 (Zweiseitiger Gauß-Test). Im Beispiel 1.1.2 haben wir folgenden Test des Erwartungswertes einer normalverteilten Stichprobe (X_1, \dots, X_n) mit unabhängigen und identisch verteilten X_i und $X_i \sim N(\mu, \sigma_0^2)$ bei bekannten Varianzen σ_0^2 betrachtet. Getestet werden die Hypothesen

$$H_0 : \mu = \mu_0 \text{ vs. } H_1 : \mu \neq \mu_0.$$

Der Test $\varphi(x)$ lautet

$$\varphi(x) = \mathbb{I}(x \in \mathbb{R}^n : |T(x)| > z_{1-\alpha/2}),$$

wobei

$$T(x) = \sqrt{n} \frac{\bar{x}_n - \mu_0}{\sigma_0}.$$

Zeigen wir, daß φ der beste Test zum Niveau α in $\tilde{\Psi}(\alpha)$ (und somit bester unverfälschter Test) ist. Nach Satz 1.3.5 müssen wir lediglich prüfen, daß φ als φ_c mit (1.3.9) dargestellt werden kann, weil die n -dimensionale Normalverteilung mit Dichte f_μ (siehe Beispiel 1.3.3) zu der einparametrischen Exponentialfamilie mit Statistik

$$T(x) = \sqrt{n} \frac{\bar{x}_n - \mu}{\sigma_0}$$

gehört. Setzen wir $c_1 = z_{1-\alpha/2}$, $c_2 = -z_{1-\alpha/2}$, $\gamma_1 = \gamma_2 = 0$. Damit ist

$$\varphi(x) = \varphi_c(x) = \begin{cases} 1, & \text{falls } |T(x)| > z_{1-\alpha/2}, \\ 0, & \text{falls } |T(x)| \leq z_{1-\alpha/2}. \end{cases}$$

und die Behauptung ist bewiesen, weil aus der in Beispiel 1.1.2 ermittelten Gütefunktion $G_n(\varphi, \theta)$ von φ ersichtlich ist, daß φ ein unverfälschter Test zum Niveau α ist (und somit $\varphi \in \tilde{\Psi}(\alpha)$).

Bemerkung 1.3.8. Bisher haben wir immer vorausgesetzt, daß nur *ein* Parameter der Verteilung der Stichprobe (X_1, \dots, X_n) unbekannt ist, um die Theorie des Abschnittes 1.3 über die besten (Neyman-Pearson-) Tests im Fall der einparametrischen Exponentialfamilie aufstellen zu können. Um jedoch den Fall weiterer unbekannter Parameter betrachten zu können (wie im Beispiel der zweiseitigen Tests des Erwartungswertes der normalverteilten Stichprobe bei unbekannter Varianz (der sog. t -Test, vergleiche Abschnitt 1.2.1, 1 (a)), bedarf es einer tiefergehenden Theorie, die aus Zeitgründen in dieser Vorlesung nicht behandelt wird. Der interessierte Leser findet das Material dann im Buch [14].

1.4 Anpassungstests

Sei eine Stichprobe von unabhängigen, identisch verteilten Zufallsvariablen (X_1, \dots, X_n) gegeben mit $X_i \sim F$ (Verteilungsfunktion) für $i = 1, \dots, n$. Bei den Anpassungstests wird die Hypothese

$$H_0 : F = F_0 \text{ vs. } H_1 : F \neq F_0$$

überprüft, wobei F_0 eine vorgegebene Verteilungsfunktion ist.

Einen Test aus dieser Klasse haben wir bereits in der Vorlesung Stochastik I kennengelernt: den Kolmogorow-Smirnov-Test (vergleiche Bemerkung 3.3.8. 3), Vorlesungsskript Stochastik I).

Jetzt werden weitere nichtparametrische Anpassungstests eingeführt. Der erste ist der χ^2 -Anpassungstest von K. Pearson.

1.4.1 χ^2 -Anpassungstest

Der Test von Kolmogorow-Smirnov basiert auf dem Abstand

$$D_n = \sup_{x \in \mathbb{R}} | \hat{F}_n(x) - F_0(x) |$$

zwischen der empirischen Verteilungsfunktion der Stichprobe (X_1, \dots, X_n) und der Verteilungsfunktion F_0 . In der Praxis jedoch erscheint dieser Test zu feinfühlig, denn er ist zu sensibel gegenüber Unregelmäßigkeiten in den Stichproben und verwirft H_0 zu oft. Einen Ausweg aus dieser Situation stellt die Vergrößerung der Haupthypothese H_0 dar, auf welcher der folgende χ^2 -Anpassungstest beruht.

Man zerlegt den Wertebereich der Stichprobenvariablen X_i in r Klassen $(a_j, b_j]$, $j = 1, \dots, r$ mit der Eigenschaft

$$-\infty \leq a_1 < b_1 = a_2 < b_2 = \dots = a_r < b_r \leq \infty.$$

Anstelle von $X_i, i = 1, \dots, n$ betrachten wir die sogenannten *Klassenstärken* $Z_j, j = 1, \dots, r$, wobei

$$Z_j = \#\{i : a_j < X_i \leq b_j, 1 \leq i \leq n\}.$$

Lemma 1.4.1. Der Zufallsvektor $Z = (Z_1, \dots, Z_r)^\top$ ist *multinomialverteilt* mit Parametervektor

$$p = (p_1, \dots, p_{r-1})^\top \in [0, 1]^{r-1},$$

wobei

$$p_j = \mathbb{P}(a_j < X_1 \leq b_j) = F(b_j) - F(a_j), \quad j = 1, \dots, r-1, \quad p_r = 1 - \sum_{j=1}^{r-1} p_j.$$

Schreibweise:

$$Z \sim M_{r-1}(n, p)$$

Beweis. Es ist zu zeigen, daß für alle Zahlen $k_1, \dots, k_r \in \mathbb{N}_0$ mit $k_1 + \dots + k_r = n$ gilt:

$$\mathbb{P}(Z_i = k_i, i = 1, \dots, r) = \frac{n!}{k_1! \cdot \dots \cdot k_r!} p_1^{k_1} \cdot \dots \cdot p_r^{k_r}. \quad (1.4.1)$$

Da X_i unabhängig und identisch verteilt sind, gilt

$$\mathbb{P}(X_j \in (a_{i_j}, b_{i_j}], j = 1, \dots, n) = \prod_{j=1}^n \mathbb{P}(a_{i_j} < X_1 \leq b_{i_j}) = p_1^{k_1} \cdot \dots \cdot p_r^{k_r},$$

falls die Folge von Intervallen $(a_{i_j}, b_{i_j}]_{j=1, \dots, n}$ das Intervall $(a_i, b_i]$ k_i Mal enthält, $i = 1, \dots, r$. Die Formel (1.4.1) ergibt sich aus dem Satz der totalen Wahrscheinlichkeit als Summe über die Permutationen von Folgen $(a_{i_j}, b_{i_j}]_{j=1, \dots, n}$ dieser Art. \square

Im Sinne des Lemmas 1.4.1 werden neue Hypothesen über die Beschaffenheit von F geprüft.

$$H_0 : p = p_0 \text{ vs. } H_1 : p \neq p_0,$$

wobei $p = (p_1, \dots, p_{r-1})^\top$ der Parametervektor der Multinomialverteilung von Z ist, und $p_0 = (p_{01}, \dots, p_{0,r-1})^\top \in (0, 1)^{r-1}$ mit $\sum_{i=1}^{r-1} p_{0i} < 1$. In diesem Fall ist

$$\Lambda_0 = \{F \in \Lambda : F(b_j) - F(a_j) = p_{0j}, \quad j = 1, \dots, r-1\},$$

$\Lambda_1 = \Lambda \setminus \Lambda_0$, wobei Λ die Menge aller Verteilungsfunktionen ist. Um H_0 vs. H_1 zu testen, führen wir die *Pearson-Teststatistik*

$$T_n(x) = \sum_{j=1}^r \frac{(z_j - np_{0j})^2}{np_{0j}}$$

ein, wobei $x = (x_1, \dots, x_n)$ eine konkrete Stichprobe der Daten ist und $z_j, j = 1, \dots, r$ ihre Klassenstärken sind.

Unter H_0 gilt

$$\mathbb{E} Z_j = np_{0j}, \quad j = 1, \dots, r,$$

somit soll H_0 abgelehnt werden, falls $T_n(X)$ ungewöhnlich große Werte annimmt.

Im nächsten Satz zeigen wir, daß $T(X_1, \dots, X_n)$ asymptotisch (für $n \rightarrow \infty$) χ_{r-1}^2 -verteilt ist, was zu folgendem Anpassungstest (χ^2 -Anpassungstest) führt:

$$H_0 \text{ wird verworfen, falls } T_n(x_1, \dots, x_n) > \chi_{r-1, 1-\alpha}^2.$$

Dieser Test ist nach seinem Entdecker *Karl Pearson* (1857-1936) benannt worden.

Satz 1.4.1. Unter H_0 gilt

$$\lim_{n \rightarrow \infty} \mathbb{P}_{p_0}(T_n(X_1, \dots, X_n) > \chi_{r-1, 1-\alpha}^2) = \alpha, \quad \alpha \in (0, 1),$$

das heißt, der χ^2 -Pearson-Test ist ein asymptotischer Test zum Niveau α .

Beweis. Führen wir die Bezeichnung $Z_{nj} = Z_j(X_1, \dots, X_n)$ der Klassenstärken ein, die aus der Stichprobe (X_1, \dots, X_n) entstehen. Nach Lemma 1.4.1 ist

$$Z_n = (Z_{n1}, \dots, Z_{nr}) \sim M_{r-1}(n, p_0) \text{ unter } H_0.$$

Insbesondere soll $\mathbb{E} Z_{nj} = np_{0j}$ und

$$\text{Cov}(Z_{ni}, Z_{nj}) = \begin{cases} np_{0j}(1 - p_{0j}), & i = j, \\ -np_{0i}p_{0j}, & i \neq j \end{cases}$$

für alle $i, j = 1, \dots, r$ gelten. Da

$$Z_{nj} = \sum_{i=1}^n \mathbb{I}(a_j < X_i \leq b_j), \quad j = 1, \dots, r,$$

ist $Z_n = (Z_{n1}, \dots, Z_{n,r-1})$ eine Summe von n unabhängigen und identisch verteilten Zufallsvektoren $Y_i \in \mathbb{R}^{r-1}$ mit Koordinaten $Y_{ij} = \mathbb{I}(a_j < X_i \leq b_j)$, $j = 1, \dots, r-1$. Daher gilt nach dem multivariaten Grenzwertsatz (der in Lemma 1.4.2 bewiesen wird), daß

$$Z'_n = \frac{Z_n - \mathbb{E} Z_n}{\sqrt{n}} = \frac{\sum_{i=1}^n Y_i - n\mathbb{E} Y_1}{\sqrt{n}} \xrightarrow[n \rightarrow \infty]{d} Y \sim N(0, K),$$

mit $N(0, K)$ eine $(r-1)$ -dimensionale multivariate Normalverteilung (vergleiche Vorlesungsskript WR, Beispiel 3.4.1. 3.) mit Erwartungswertvektor Null und Kovarianzmatrix $K = (\sigma_{ij}^2)$, wobei

$$\sigma_{ij}^2 = \begin{cases} -p_{0i}p_{0j}, & i \neq j, \\ p_{0i}(1 - p_{0j}), & i = j \end{cases}$$

für $i, j = 1, \dots, r-1$ ist. Diese Matrix K ist invertierbar mit $K^{-1} = A = (a_{ij})$,

$$a_{ij} = \begin{cases} \frac{1}{p_{0r}}, & i \neq j, \\ \frac{1}{p_{0i}} + \frac{1}{p_{0r}}, & i = j. \end{cases}$$

Außerdem ist K (als Kovarianzmatrix) symmetrisch und positiv definit. Aus der linearen Algebra ist bekannt, daß es eine invertierbare $(r-1) \times (r-1)$ -Matrix $A^{1/2}$ gibt, mit der Eigenschaft $A = A^{1/2}(A^{1/2})^\top$. Daraus folgt,

$$K = A^{-1} = ((A^{1/2})^\top)^{-1} \cdot (A^{1/2})^{-1}.$$

Wenn wir $(A^{1/2})^\top$ auf Z'_n anwenden, so bekommen wir

$$(A^{1/2})^\top \cdot Z'_n \xrightarrow[n \rightarrow \infty]{d} (A^{1/2})^\top \cdot Y,$$

wobei

$$(A^{1/2})^\top \cdot Y \sim N\left(0, (A^{1/2})^\top \cdot K \cdot A^{1/2}\right) = N(0, \mathcal{I}_{r-1})$$

nach der Eigenschaft der multivariaten Normalverteilung, die im Kapitel 2, Satz 2.1.3 behandelt wird. Des Weiteren wurde hier der Stetigkeitssatz aus der Wahrscheinlichkeitsrechnung benutzt, daß

$$Y_n \xrightarrow[n \rightarrow \infty]{d} Y \implies \varphi(Y_n) \xrightarrow[n \rightarrow \infty]{d} \varphi(Y)$$

für beliebige Zufallsvektoren $\{Y_n\}$, $Y \in \mathbb{R}^m$ und stetige Abbildungen $\varphi: \mathbb{R} \rightarrow \mathbb{R}$. Diesen Satz haben wir in WR für Zufallsvariablen bewiesen (Satz 6.4.3, Vorlesungsskript WR). Die erneute Anwendung des Stetigkeitssatzes ergibt

$$\left| (A^{1/2})^\top Z'_n \right|^2 \xrightarrow[n \rightarrow \infty]{d} |Y|^2 = R \sim \chi_{r-1}^2.$$

Zeigen wir, daß

$$T_n(X_1, \dots, X_n) = \left| (A^{1/2})^\top Z'_n \right|^2.$$

Es gilt:

$$\begin{aligned} \left| (A^{1/2})^\top Z'_n \right|^2 &= ((A^{1/2})^\top Z'_n)^\top ((A^{1/2})^\top Z'_n) = Z_n'^\top \cdot \underbrace{A^{1/2} \cdot (A^{1/2})^\top}_A Z'_n = Z_n'^\top A Z'_n \\ &= n \sum_{j=1}^{r-1} \frac{1}{p_{0j}} \left(\frac{Z_{nj}}{n} - p_{0j} \right)^2 + \frac{n}{p_{0r}} \sum_{i=1}^{r-1} \sum_{j=1}^{r-1} \left(\frac{Z_{ni}}{n} - p_{0i} \right) \left(\frac{Z_{nj}}{n} - p_{0j} \right) \\ &= \sum_{j=1}^{r-1} \frac{(Z_{nj} - np_{0j})^2}{np_{0j}} + \frac{n}{p_{0r}} \left(\sum_{j=1}^{r-1} \left(\frac{Z_{nj}}{n} - p_{0j} \right) \right)^2 \\ &= \sum_{j=1}^{r-1} \frac{(Z_{nj} - np_{0j})^2}{np_{0j}} + \frac{n}{p_{0r}} \left(\frac{Z_{nr}}{n} - p_{0r} \right)^2 \\ &= \sum_{j=1}^r \frac{(Z_{nj} - np_{0j})^2}{np_{0j}} = T_n(X_1, \dots, X_n), \end{aligned}$$

weil

$$\begin{aligned} \sum_{j=1}^{r-1} Z_{nj} &= n - Z_{nr}, \\ \sum_{j=1}^{r-1} p_{0j} &= 1 - p_{0r}. \end{aligned}$$

□

Lemma 1.4.2 (Multivariater zentraler Grenzwertsatz). Sei $\{Y_n\}_{n \in \mathbb{N}}$ eine Folge von unabhängigen und identisch verteilten Zufallsvektoren, mit $\mathbb{E} Y_1 = \mu$ und Kovarianzmatrix K . Dann gilt

$$\frac{\sum_{i=1}^n Y_i - n\mu}{\sqrt{n}} \xrightarrow[n \rightarrow \infty]{d} Y \sim N(0, K). \quad (1.4.2)$$

Beweis. Sei $Y_j = (Y_{j1}, \dots, Y_{jm})^\top$. Nach dem Stetigkeitssatz für charakteristische Funktionen ist die Konvergenz (1.4.2) äquivalent zu

$$\varphi_n(t) \xrightarrow[n \rightarrow \infty]{} \varphi(t) \quad t \in \mathbb{R}^m, \quad (1.4.3)$$

wobei

$$\varphi_n(t) = \mathbb{E} e^{itS_n} = \mathbb{E} \exp \left\{ i \sum_{j=1}^m t_j \frac{Y_{1j} + \dots + Y_{nj} - n\mu_j}{\sqrt{n}} \right\}$$

die charakteristische Funktion vom Zufallsvektor

$$S_n = \frac{\sum_{i=1}^n Y_i - n\mu}{\sqrt{n}}$$

und

$$\varphi(t) = e^{-t^\top K t / 2}$$

die charakteristische Funktion der $N(0, K)$ -Verteilung ist. Die Funktion $\varphi_n(t)$ kann in der Form

$$\varphi_n(t) = \mathbb{E} \exp \left\{ i \sum_{i=1}^n \frac{\sum_{j=1}^m t_j (Y_{ij} - \mu_j)}{\sqrt{n}} \right\}, \quad t = (t_1, \dots, t_m)^\top \in \mathbb{R}^m$$

umgeschrieben werden, wobei für die Zufallsvariable

$$L_i := \sum_{j=1}^m t_j (Y_{ij} - \mu_j)$$

gilt:

$$\begin{aligned} \mathbb{E} L_i &= 0, \\ \text{Var } L_i &= \mathbb{E} \left[\sum_{k,j=1}^m t_j (Y_{ij} - \mu_j) (Y_{ik} - \mu_k) t_k \right] = t^\top K t, \quad i \in \mathbb{N}. \end{aligned}$$

Falls $t^\top Kt = 0$, dann gilt $L_i = 0$ fast sicher, für alle $i \in \mathbb{N}$. Hieraus folgt $\varphi_n(t) = \varphi(t) = 1$, also gilt die Konvergenz 1.4.2.

Falls jedoch $t^\top Kt > 0$, dann kann $\varphi_n(t)$ als charakteristische Funktion der Zufallsvariablen

$$\sum_{i=1}^n L_i / \sqrt{n}$$

an Stelle 1, und $\varphi(t)$ als charakteristische Funktion der eindimensionalen Normalverteilung $N(0, t^\top Kt)$ an Stelle 1 interpretiert werden. Aus dem zentralen Grenzwertsatz für eindimensionale Zufallsvariablen (vergleiche Satz 7.2.1, Vorlesungsskript WR) gilt

$$\sum_{i=1}^n \frac{L_i}{\sqrt{n}} \xrightarrow[n \rightarrow \infty]{d} L \sim N(0, t^\top Kt)$$

und somit

$$\varphi_n(t) = \varphi\left(\sum_{i=1}^n L_i / \sqrt{n}\right)(1) \xrightarrow[n \rightarrow \infty]{} \varphi_L(1) = \varphi(t).$$

Somit ist die Konvergenz (1.4.2) bewiesen. \square

Bemerkung 1.4.1. 1. Die im letzten Beweis verwendete Methode der Reduktion einer mehrdimensionalen Konvergenz auf den eindimensionalen Fall mit Hilfe von Linearkombinationen von Zufallsvariablen trägt den Namen von *Cramér-Wold*.

2. Der χ^2 -Pearson-Test ist asymptotisch, also für große Stichprobenumfänge, anzuwenden. Aber welches n ist groß genug? Als „Faustregel“ gilt: np_{0j} soll größer gleich a sein, $a \in (2, \infty)$. Für eine größere Klassenanzahl $r \geq 10$ kann sogar $a = 1$ verwendet werden. Wir zeigen jetzt, daß der χ^2 -Anpassungstest konsistent ist.

Lemma 1.4.3. Der χ^2 -Pearson-Test ist konsistent, das heißt

$$\forall p \in [0, 1]^{r-1}, p \neq p_0 \text{ gilt: } \lim_{n \rightarrow \infty} \mathbb{P}_p (T_n(X_1, \dots, X_n) > \chi_{r-1, 1-\alpha}^2) = 1$$

Beweis. Unter H_1 gilt

$$Z_{nj}/n = \frac{\sum_{i=1}^n \mathbb{I}(a_j < X_i \leq b_j)}{n} \xrightarrow[n \rightarrow \infty]{f.s.} \underbrace{\mathbb{E} \mathbb{I}(a_j < X_1 \leq b_j)}_{=p_j}$$

nach dem starken Gesetz der großen Zahlen. Wir wählen j so, daß $p_j \neq p_{0j}$. Es gilt

$$T_n(X_1, \dots, X_n) \geq \frac{(Z_{nj} - np_{0j})^2}{np_{0j}} \geq \underbrace{n \left(\frac{Z_{nj}}{n} - p_{0j} \right)^2}_{\sim n(p_j - p_{0j})^2} \xrightarrow[n \rightarrow \infty]{f.s.} \infty.$$

Somit ist auch

$$\mathbb{P}_p (T_n(X_1, \dots, X_n) > \chi_{r-1, 1-\alpha}^2) \xrightarrow[n \rightarrow \infty]{} 1.$$

\square

1.4.2 χ^2 -Anpassungstest von Pearson-Fisher

Es sei (X_1, \dots, X_n) eine Stichprobe von unabhängigen und identisch verteilten Zufallsvariablen X_i , $i = 1, \dots, n$. Wir wollen testen, ob die Verteilungsfunktion F von X_i zu einer parametrischen Familie

$$\Lambda_0 = \{F_\theta : \theta \in \Theta\}, \quad \Theta \subset \mathbb{R}^m$$

gehört. Seien die Zahlen a_i, b_i , $i = 1, \dots, r$ vorgegeben mit der Eigenschaft

$$-\infty \leq a_1 < b_1 = a_2 < b_2 = \dots = a_r < b_r \leq \infty$$

und

$$\begin{aligned} Z_j &= \#\{X_i, i = 1, \dots, n : a_j < X_i \leq b_j\}, \quad j = 1, \dots, r, \\ Z &= (Z_1, \dots, Z_r)^\top. \end{aligned}$$

Nach Lemma 1.4.1 gilt: $Z \sim M_{r-1}(n, p)$, $p = (p_0, \dots, p_{r-1})^\top \in [0, 1]^{r-1}$. Unter der Hypothese $H_0 : F \in \Lambda_0$ gilt: $p = p(\theta)$, $\theta \in \Theta \subset \mathbb{R}^m$. Wir vergrößern die Hypothese H_0 und wollen folgende neue Hypothese testen:

$$H_0 : p \in \{p(\theta) : \theta \in \Theta\} \text{ vs. } H_1 : p \notin \{p(\theta) : \theta \in \Theta\}.$$

Um dieses Hypothesenpaar zu testen, wird der χ^2 -Pearson-Fisher-Test wie folgt aufgebaut:

1. Ein (schwach konsistenter) Maximum-Likelihood-Schätzer $\hat{\theta}_n = \hat{\theta}(X_1, \dots, X_n)$ für θ wird gefunden: $\hat{\theta}_n \xrightarrow[n \rightarrow \infty]{P} \theta$. Dabei muß $\{\hat{\theta}_n\}_{n \in \mathbb{N}}$ asymptotisch normalverteilt sein.
2. Es wird der Plug-In-Schätzer $p(\hat{\theta}_n)$ für $p(\theta)$ gebildet.
3. Die Testgröße

$$\hat{T}_n(X_1, \dots, X_n) = \sum_{j=1}^r \frac{(Z_{nj} - np_j(\hat{\theta}))^2}{np_j(\hat{\theta})} \xrightarrow[n \rightarrow \infty]{P} \eta \sim \chi_{r-m-1}^2$$

unter H_0 und gewissen Voraussetzungen.

4. H_0 wird verworfen, falls $\hat{T}_n(X_1, \dots, X_n) > \chi_{r-m-1, 1-\alpha}^2$. Dies ist ein asymptotischer Test zum Niveau α .

Bemerkung 1.4.2. 1. Bei einem χ^2 -Pearson-Fisher-Test wird vorausgesetzt, daß die Funktion $p(\theta)$ explizit bekannt ist, θ jedoch unbekannt. Das bedeutet, daß für jede Klasse von Verteilungen Λ_0 die Funktion $p(\cdot)$ berechnet werden soll.

2. Warum kann \hat{T}_n die Hypothese H_0 von H_1 unterscheiden? Nach dem Gesetz der großen Zahlen gilt

$$\frac{1}{n}Z_{nj} - p_j(\hat{\theta}_n) = \underbrace{\frac{1}{n}Z_{nj} - p_j(\theta)}_{\xrightarrow{P} 0} - \underbrace{(p_j(\hat{\theta}_n) - p_j(\theta))}_{\xrightarrow{P} 0} \xrightarrow[n \rightarrow \infty]{P} 0,$$

falls $\hat{\theta}_n$ schwach konsistent ist und $p_j(\cdot)$ eine stetige Funktion für alle $j = 1, \dots, r$ ist.

Das heißt, unter H_0 soll $\hat{T}_n(X_1, \dots, X_n)$ relativ kleine Werte annehmen. Eine signifikante Abweichung von diesem Verhalten soll zur Ablehnung von H_0 führen, vergleiche Punkt 4.

Für die Verteilung F_θ von X_i gelten folgende Regularitätsvoraussetzungen (vergleiche Satz 3.4.2, Vorlesungsskript Stochastik I).

1. Die Verteilungsfunktion F_θ ist entweder diskret oder absolut stetig für alle $\theta \in \Theta$.
2. Die Parametrisierung ist eindeutig, das heißt: $\theta \neq \theta_1 \Leftrightarrow F_\theta \neq F_{\theta_1}$.
3. Der Träger der Likelihood-Funktion

$$L(x, \theta) = \begin{cases} P_\theta(X_1 = x), & \text{im Falle von diskreten } F_\theta, \\ f_\theta(x), & \text{im absolut stetigen Fall.} \end{cases}$$

$\text{Supp}L(x, \theta) = \{x \in \mathbb{R} : L(x, \theta) > 0\}$ hängt nicht von θ ab.

4. $L(x, \theta)$ sei 3 Mal stetig differenzierbar, und es gelte für $k = 1, \dots, 3$ und $i_1, \dots, i_k \in \{1 \dots m\}$, daß

$$\left(\sum\right) \int \frac{\partial^k L(x, \theta)}{\partial \theta_{i_1} \dots \partial \theta_{i_k}} dx = \frac{\partial^k}{\partial \theta_{i_1} \dots \partial \theta_{i_k}} \left(\sum\right) \int L(x, \theta) dx = 0.$$

5. Für alle $\theta_0 \in \Theta$ gibt es eine Konstante c_{θ_0} und eine messbare Funktion $g_{\theta_0} : \text{Supp}L \rightarrow \mathbb{R}_+$, sodaß

$$\left| \frac{\partial^3 \log L(x, \theta)}{\partial \theta_{i_1} \partial \theta_{i_2} \partial \theta_{i_3}} \right| \leq g_{\theta_0}(x), \quad |\theta - \theta_0| < c_{\theta_0}$$

und

$$\mathbb{E}_{\theta_0} g_{\theta_0}(X_1) < \infty.$$

Wir definieren die *Informationsmatrix von Fisher* durch

$$I(\theta) = \left(\mathbb{E} \left[\frac{\partial \log L(X_1, \theta)}{\partial \theta_i} \frac{\partial \log L(X_1, \theta)}{\partial \theta_j} \right] \right)_{i,j=1,\dots,m}. \quad (1.4.4)$$

Satz 1.4.2 (asymptotische Normalverteiltheit von konsistenten ML-Schätzern $\hat{\theta}_n$, multivariater Fall $m > 1$). Es seien X_1, \dots, X_n unabhängig und identisch verteilt mit Likelihood-Funktion L , die den Regularitätsbedingungen 1-5 genügt. Sei $I(\theta)$ positiv definit für alle $\theta \in \Theta \subset \mathbb{R}^m$. Sei $\hat{\theta}_n = \hat{\theta}(X_1, \dots, X_n)$ eine Folge von schwach konsistenten Maximum-Likelihood-Schätzern für θ . Dann gilt:

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow[n \rightarrow \infty]{d} N(0, I^{-1}(\theta)).$$

Ohne Beweis; siehe den Beweis des Satzes 3.4.2, Vorlesungsskript Stochastik I.

Für unsere vergrößerte Hypothese $H_0 : p \in \{p(\theta), \theta \in \Theta\}$ stellen wir folgende, stückweise konstante, Likelihood-Funktion auf:

$$L(x, \theta) = p_j(\theta), \text{ falls } x \in (a_j, b_j].$$

Dann ist die Likelihood-Funktion der Stichprobe (x_1, \dots, x_n) gleich

$$\begin{aligned} L(x_1, \dots, x_n, \theta) &= \prod_{j=1}^r p_j(\theta)^{Z_j(x_1, \dots, x_n)} \\ \Rightarrow \log L(x_1, \dots, x_n, \theta) &= \sum_{j=1}^r Z_j(x_1, \dots, x_n) \cdot \log p_j(\theta). \\ \hat{\theta}_n = \hat{\theta}(x_1, \dots, x_n) &= \operatorname{argmax}_{\theta \in \Theta} \log L(x_1, \dots, x_n, \theta) \\ \Rightarrow \sum_{j=1}^r Z_j(x_1, \dots, x_n) \frac{\partial p_j(\theta)}{\partial \theta_i} \cdot \frac{1}{p_j(\theta)} &= 0, \quad i = 1, \dots, m. \end{aligned}$$

Aus $\sum_{j=1}^r p_j(\theta) = 1$ folgt

$$\sum_{j=1}^r \frac{\partial p_j(\theta)}{\partial \theta_i} = 0 \Rightarrow \sum_{j=1}^r \frac{Z_j(x_1, \dots, x_n) - np_j(\theta)}{p_j(\theta)} \cdot \frac{\partial p_j(\theta)}{\partial \theta_i} = 0, \quad i = 1, \dots, m.$$

Lemma 1.4.4. Im obigen Fall gilt $I(\theta) = C^\top(\theta) \cdot C(\theta)$, wobei $C(\theta)$ eine $(r \times m)$ -Matrix mit Elementen

$$c_{ij}(\theta) = \frac{\partial p_i(\theta)}{\partial \theta_j} \cdot \frac{1}{\sqrt{p_i(\theta)}} \quad \text{ist.}$$

Beweis.

$$\begin{aligned} \mathbb{E}_0 \left[\frac{\partial \log L(X_1, \theta)}{\partial \theta_i} \cdot \frac{\partial \log L(X_1, \theta)}{\partial \theta_j} \right] &= \sum_{k=1}^r \frac{\partial \log p_k(\theta)}{\partial \theta_i} \cdot \frac{\partial \log p_k(\theta)}{\partial \theta_j} \cdot p_k(\theta) \\ &= \sum_{k=1}^r \frac{\partial p_k(\theta)}{\partial \theta_i} \frac{1}{p_k(\theta)} \cdot \frac{\partial p_k(\theta)}{\partial \theta_j} \cdot \frac{1}{p_k(\theta)} \cdot p_k(\theta) \\ &= \left(C^\top(\theta) \cdot C(\theta) \right)_{ij}, \end{aligned}$$

$$\text{denn } \log L(X_1, \theta) = \sum_{i=1}^r \log p_j(\theta) \cdot \mathbb{I}(x \in (a_j, b_j]).$$

□

Deshalb gilt die Folgerung aus Satz 1.4.2:

Folgerung 1.4.1. Sei $\hat{\theta}_n = \hat{\theta}(X_1, \dots, X_n)$ ein Maximum-Likelihood-Schätzer von θ im vergrößerten Modell, der schwach konsistent ist und den obigen Regularitätsbedingungen genügt. Sei die Informationsmatrix von Fisher $I(\theta) = C^\top(\theta) \cdot C(\theta)$ für alle $\theta \in \Theta$ positiv definit. Dann ist $\hat{\theta}$ asymptotisch normalverteilt:

$$\sqrt{n} \left(\hat{\theta}_n - \theta \right) \xrightarrow[n \rightarrow \infty]{d} Y \sim N(0, I^{-1}(\theta))$$

Satz 1.4.3. Es sei $\hat{\theta}_n$ ein Maximum-Likelihood-Schätzer im vergrößerten Modell für θ , für den alle Voraussetzungen der Folgerung 1.4.1 erfüllt sind. Die Teststatistik

$$\hat{T}_n(X_1, \dots, X_n) = \sum_{j=1}^r \frac{(Z_j(X_1, \dots, X_n) - np_j(\hat{\theta}_n))^2}{np_j(\hat{\theta}_n)}$$

ist unter H_0 asymptotisch χ_{r-m-1}^2 -verteilt:

$$\lim_{n \rightarrow \infty} \mathbb{P}_\theta \left(\hat{T}_n(X_1, \dots, X_n) > \chi_{r-m-1, 1-\alpha}^2 \right) = \alpha.$$

ohne Beweis (siehe [15]).

Aus diesem Satz folgt, daß der χ^2 -Pearson-Fisher-Test ein asymptotischer Test zum Niveau α ist.

Beispiel 1.4.1. 1. χ^2 -Pearson-Fisher-Test der Normalverteilung

Sei (X_1, \dots, X_n) eine Zufallsstichprobe. Es soll geprüft werden, ob $X_i \sim N(\mu, \sigma^2)$. Es gilt

$$\theta = (\mu, \sigma^2) \in \Theta = \mathbb{R} \times \mathbb{R}_+.$$

Sei $(a_j, b_j]_{j=1, \dots, r}$ eine beliebige Aufteilung von \mathbb{R} in r disjunkte Intervalle. Sei

$$f_\theta(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

die Dichte der $N(\mu, \sigma^2)$ -Verteilung.

$$p_j(\theta) = \mathbb{P}_0(a_j < X_1 \leq b_j) = \int_{a_j}^{b_j} f_\theta(x) dx, \quad j = 1, \dots, r$$

mit den Klassenstärken

$$Z_j = \# \{i : X_i \in (a_j, b_j]\}.$$

Wir suchen den Maximum-Likelihood-Schätzer im vergrößerten Modell:

$$\begin{aligned} \frac{\partial p_j(\theta)}{\partial \mu} &= \int_{a_j}^{b_j} \frac{\partial}{\partial \mu} f_\theta(x) dx = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \int_{a_j}^{b_j} \frac{x-\mu}{\sigma^2} \cdot e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} dx \\ \frac{\partial p_j(\theta)}{\partial \sigma^2} &= \int_{a_j}^{b_j} \frac{\partial}{\partial \sigma^2} f_\theta(x) dx \\ &= \frac{1}{\sqrt{2\pi}} \int_{a_j}^{b_j} \left[-\frac{1}{2} \cdot \frac{1}{(\sigma^2)^{3/2}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} + \frac{1}{\sqrt{\sigma^2}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2} \cdot \left(\frac{(x-\mu)^2}{2(\sigma^2)^2} \right) \right] dx \\ &= -\frac{1}{2} \frac{1}{\sigma^2} \int_{a_j}^{b_j} f_\theta(x) dx + \frac{1}{2(\sigma^2)^2} \int_{a_j}^{b_j} (x-\mu)^2 f_\theta(x) dx \end{aligned}$$

Die notwendigen Bedingungen des Maximums sind:

$$\begin{aligned} \sum_{i=1}^r Z_j \frac{\int_{a_j}^{b_j} x f_\theta(x) dx}{\int_{a_j}^{b_j} f_\theta(x) dx} - \underbrace{\mu \sum_{j=1}^r Z_j}_{=n} &= 0, \\ \frac{1}{\sigma^2} \sum_{j=1}^r Z_j \frac{\int_{a_j}^{b_j} (x-\mu)^2 f_\theta(x) dx}{\int_{a_j}^{b_j} f_\theta(x) dx} - \underbrace{\sum_{j=1}^r Z_j}_{=n} &= 0. \end{aligned}$$

Daraus folgen die Maximum-Likelihood-Schätzer $\hat{\mu}$ und $\hat{\sigma}^2$ für μ und σ^2 :

$$\hat{\mu} = \frac{1}{n} \sum_{j=1}^r Z_j \frac{\int_{a_j}^{b_j} x f_{\theta}(x) dx}{\int_{a_j}^{b_j} f_{\theta}(x) dx},$$

$$\hat{\sigma}^2 = \frac{1}{n} \sum_{j=1}^r Z_j \frac{\int_{a_j}^{b_j} (x - \mu)^2 f_{\theta}(x) dx}{\int_{a_j}^{b_j} f_{\theta}(x) dx}.$$

Wir konstruieren eine Näherung zu $\hat{\mu}$ und $\hat{\sigma}^2$ für $r \rightarrow \infty$. Falls $r \rightarrow \infty$ (und somit auch $n \rightarrow \infty$), dann ist $b_j - a_j$ klein und nach der einfachen Quadraturregel gilt:

$$\int_{a_j}^{b_j} x f_{\theta}(x) dx \approx (b_j - a_j) y_j f_{\theta}(y_j),$$

$$\int_{a_j}^{b_j} f_{\theta}(x) dx \approx (b_j - a_j) f_{\theta}(y_j),$$

wobei $y_1 = b_1$, $y_r = b_{r-1} = a_r$,

$$y_j = (b_{j+1} + b_j)/2, \quad j = 2, \dots, r-1.$$

Daraus folgen für die Maximum-Likelihood-Schätzer $\hat{\mu}$ und $\hat{\sigma}^2$:

$$\hat{\mu} \approx \frac{1}{n} \sum_{j=1}^r y_j \cdot Z_j = \tilde{\mu}$$

$$\hat{\sigma}^2 \approx \frac{1}{n} \sum_{j=1}^r (y_j - \tilde{\mu})^2 Z_j = \tilde{\sigma}^2,$$

$$\tilde{\theta} = (\tilde{\mu}, \tilde{\sigma}^2).$$

Der χ^2 -Pearson-Fisher-Test lautet dann: H_0 wird abgelehnt, falls

$$\hat{T}_n = \frac{\sum_{j=1}^r \left(Z_j - np_j(\tilde{\theta}) \right)^2}{np_j(\tilde{\theta})} > \chi_{r-3, 1-\alpha}^2.$$

2. χ^2 -Pearson-Fisher-Test der Poissonverteilung

Es sei (X_1, \dots, X_n) eine Stichprobe von unabhängigen und identisch verteilten Zufallsvariablen. Wir wollen testen, ob $X_i \sim \text{Poisson}(\lambda)$, $\lambda > 0$. Es gilt $\theta = \lambda$ und $\Theta = (0, +\infty)$. Die Vergrößerung von Θ hat die Form

$$-\infty = a_1 < \underbrace{b_1}_{=0} = a_2 < \underbrace{b_2}_{=1} = a_3 < \dots < \underbrace{b_{r-1}}_{=r-2} = a_r < b_r = +\infty.$$

Dann ist

$$\begin{aligned}
 p_j(\lambda) &= \mathbb{P}_\lambda(X_1 = j-1) = e^{-\lambda} \frac{\lambda^{j-1}}{(j-1)!}, \quad j = 1, \dots, r-1, \\
 p_r(\lambda) &= \sum_{i=r-1}^{\infty} e^{-\lambda} \frac{\lambda^i}{i!}, \\
 \frac{dp_j(\lambda)}{d\lambda} &= -e^{-\lambda} \frac{\lambda^{j-1}}{(j-1)!} + (j-1) \frac{\lambda^{j-2}}{(j-1)!} e^{-\lambda} = e^{-\lambda} \frac{\lambda^{j-1}}{(j-1)!} \left(\frac{j-1}{\lambda} - 1 \right) \\
 &= p_j(\lambda) \cdot \left(\frac{j-1}{\lambda} - 1 \right), \quad j = 1, \dots, r-1 \\
 \frac{dp_r(\lambda)}{d\lambda} &= \sum_{i \geq r-1} p_i(\lambda) \left(\frac{i-1}{\lambda} - 1 \right).
 \end{aligned}$$

Die Maximum-Likelihood-Gleichung lautet

$$0 = \sum_{j=1}^{r-1} Z_j \cdot \left(\frac{j-1}{\lambda} - 1 \right) + Z_r \frac{\sum_{i \geq r-1} p_i(\lambda) \left(\frac{i-1}{\lambda} - 1 \right)}{p_r(\lambda)}$$

Falls $r \rightarrow \infty$, so findet sich $r(n)$ für jedes n , für das $Z_{r(n)} = 0$. Deshalb gilt für $r > r(n)$:

$$\sum_{j=1}^{r-1} (j-1)Z_j - \lambda \underbrace{\sum_{j=1}^r Z_j}_{=n} = 0,$$

woraus der Maximum-Likelihood-Schätzer

$$\frac{1}{n} \sum_{j=1}^{r-1} (j-1)Z_j = \frac{1}{n} \sum_{j=1}^n X_j = \bar{X}_n$$

folgt. Der χ^2 -Pearson-Fisher-Test lautet: H_0 wird verworfen, falls

$$\hat{T}_n = \sum_{j=1}^r \frac{(Z_j - np_\lambda(\bar{X}_n))^2}{(np_j(\bar{X}_n))^2} > \chi_{r-2, 1-\alpha}^2.$$

1.4.3 Anpassungstest von Shapiro

Es sei (X_1, \dots, X_n) eine Stichprobe von unabhängigen, identisch verteilten Zufallsvariablen, $X_i \sim F$. Getestet werden soll die Hypothese

$$H_0 : F \in \{N(\mu, \sigma^2) : \mu \in \mathbb{R}, \sigma^2 > 0\} \text{ vs. } H_1 : F \notin \{N(\mu, \sigma^2), \mu \in \mathbb{R}, \sigma^2 > 0\}.$$

Die in den Abschnitten 1.4.1 - 1.4.2 vorgestellten χ^2 -Tests sind asymptotisch; deshalb können sie für relativ kleine Stichprobenumfänge nicht verwendet werden.

Der folgende Test wird diese Lücke füllen und eine Testentscheidung über H_0 selbst bei kleinen Stichproben ermöglichen.

Man bildet Ordnungsstatistiken $X_{(1)}, \dots, X_{(n)}$, $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$ und vergleicht ihre Korreliertheit mit den Mittelwerten der entsprechenden Ordnungsstatistiken der $N(0, 1)$ -Verteilung. Sei (Y_1, \dots, Y_n) eine Stichprobe von unabhängigen und identisch verteilten Zufallsvariablen, $Y_1 \sim N(0, 1)$. Es sei $a_i = \mathbb{E} Y_{(i)}$, $i = 1, \dots, n$. Falls der empirische Korrelationskoeffizient ρ_{aX} zwischen (a_1, \dots, a_n) und $(X_{(1)}, \dots, X_{(n)})$ bei 1 liegt, dann ist die Stichprobe normalverteilt. Formalisieren wir diese Heuristik:

Es sei b_i der Erwartungswert der i -ten Ordnungsstatistik in einer Stichprobe von $N(\mu, \sigma^2)$ -verteilten, unabhängigen Zufallsvariablen Z_i : $b_i = \mathbb{E} Z_{(i)}$, $i = 1, \dots, n$. Es gilt: $b_i = \mu + \sigma a_i$, $i = 1, \dots, n$. Betrachten wir den Korrelationskoeffizienten

$$\rho_{bX} = \frac{\sum_{i=1}^n (b_i - \bar{b}_n) (X_{(i)} - \bar{X}_n)}{\sqrt{\sum_{i=1}^n (b_i - \bar{b}_n)^2 \sum_{i=1}^n (X_{(i)} - \bar{X}_n)^2}}. \quad (1.4.5)$$

Da ρ invariant bezüglich Lineartransformationen ist und

$$\sum_{i=1}^n a_i = \sum_{i=1}^n \mathbb{E} Y_i = \mathbb{E} \left(\sum_{i=1}^n Y_i \right) = 0, \quad \text{gilt:}$$

$$\begin{aligned} \rho_{bX} &\stackrel{\text{(Stochastik I)}}{=} \rho_{aX} = \frac{\sum_{i=1}^n a_i (X_{(i)} - \bar{X}_n)}{\sqrt{\sum_{i=1}^n a_i^2 \sum_{i=1}^n (X_i - \bar{X}_n)^2}} = \frac{\sum_{i=1}^n a_i X_{(i)} - \bar{X}_n \overbrace{\sum_{i=1}^n a_i}^{=0}}{\sqrt{\sum_{i=1}^n a_i^2 \sum_{i=1}^n (X_i - \bar{X}_n)^2}} \\ &= \frac{\sum_{i=1}^n a_i X_{(i)}}{\sqrt{\sum_{i=1}^n a_i^2 \cdot \sum_{i=1}^n (X_i - \bar{X}_n)^2}} \end{aligned}$$

Die Teststatistik lautet:

$$T_n = \frac{\sum_{i=1}^n a_i X_{(i)}}{\sqrt{\sum_{i=1}^n a_i^2 \sum_{i=1}^n (X_i - \bar{X}_n)^2}} \quad (\text{Shapiro-Francia-Test})$$

Die Werte a_i sind bekannt und können den Tabellen bzw. der Statistik-Software entnommen werden. Es gilt: $|T_n| \leq 1$.

H_0 wird abgelehnt, falls $T_n \leq q_{n,\alpha}$, wobei $q_{n,\alpha}$ das α -Quantil der Verteilung von T_n ist. Diese Quantile sind aus den Tabellen bekannt, bzw. können durch Monte-Carlo-Simulationen berechnet werden.

Bemerkung 1.4.3. Einen anderen, weit verbreiteten Test dieser Art bekommt man, wenn man die Lineartransformation $b_i = \mu + \sigma a_i$ durch eine andere Lineartransformation ersetzt:

$$(a'_1, \dots, a'_n)^\top = K^{-1} \cdot (a_1, \dots, a_n),$$

wobei $K = (k_{ij})_{j=1}^n$ die Kovarianzmatrix von $(Y_{(1)}, \dots, Y_{(n)})$ ist:

$$k_{ij} = \mathbb{E} (Y_{(i)} - a_i) (Y_{(j)} - a_j), \quad i, j = 1, \dots, n,$$

Der so konstruierte Test trägt den Namen *Shapiro-Wilk-Test*.

1.5 Weitere, nicht parametrische Tests

1.5.1 Binomialtest

Es sei (X_1, \dots, X_n) eine Stichprobe von unabhängigen, identisch verteilten Zufallsvariablen, wobei $X_i \sim \text{Bernoulli}(p)$. Getestet werden soll:

$$H_0 : p = p_0 \text{ vs. } H_1 : p \neq p_0$$

Die Teststatistik lautet

$$T_n = \sum_{i=1}^n X_i \underset{H_0}{\sim} \text{Bin}(n, p_0),$$

und die Entscheidungsregel ist: H_0 wird verworfen, falls

$$T_n \notin [\text{Bin}(n, p_0)_{\alpha/2}, \text{Bin}(n, p_0)_{1-\alpha/2}],$$

wobei $\text{Bin}(n, p)_\alpha$ das α -Quantil der $\text{Bin}(n, p)$ -Verteilung ist.

Für andere H_0 , wie zum Beispiel $p \leq p_0$ ($p \geq p_0$) muss der Ablehnungsbereich entsprechend angepasst werden.

Die Quantile $\text{Bin}(n, p)_\alpha$ erhält man aus Tabellen oder aus Monte-Carlo-Simulationen. Falls n groß ist, können diese Quantile durch die Normalapproximation berechnet werden:

Nach dem zentralen Grenzwertsatz von DeMoivre-Laplace gilt:

$$\mathbb{P}(T_n \leq x) = \mathbb{P}\left(\frac{T_n - np_0}{\sqrt{np_0(1-p_0)}} \leq \frac{x - np_0}{\sqrt{np_0(1-p_0)}}\right) \underset{n \rightarrow \infty}{\approx} \Phi\left(\frac{x - np_0}{\sqrt{np_0(1-p_0)}}\right).$$

Daraus folgt:

$$z_\alpha \approx \frac{\text{Bin}(n, p_0)_\alpha - np_0}{\sqrt{np_0(1-p_0)}}$$

$$\Rightarrow \text{Bin}(n, p_0)_\alpha \approx \sqrt{np_0(1-p_0)} \cdot z_\alpha + np_0$$

Nach der Poisson-Approximation (für $n \rightarrow \infty, np_0 \rightarrow \lambda_0$) gilt:

$$\text{Bin}(n, p_0)_{\alpha/2} \approx \text{Poisson}(\lambda_0)_{\alpha/2},$$

$$\text{Bin}(n, p_0)_{1-\alpha/2} \approx \text{Poisson}(\lambda_0)_{1-\alpha/2}, \quad \text{wobei } \lambda_0 = np_0.$$

Zielstellung: Wie kann mit Hilfe des oben beschriebenen Binomialtests die Symmetrieeigenschaft einer Verteilung getestet werden?

Es sei (Y_1, \dots, Y_n) eine Stichprobe von unabhängigen und identisch verteilten Zufallsvariablen mit absolut stetiger Verteilungsfunktion F . Getestet werden soll:

$$H_0 : F \text{ ist symmetrisch vs. } H_1 : F \text{ ist nicht symmetrisch.}$$

Eine symmetrische Verteilung besitzt den Median bei Null. Deswegen vergrößern wir die Hypothese H_0 und testen:

$$H'_0 : F^{-1}(0,5) = 0 \text{ vs. } H'_1 : F^{-1}(0,5) \neq 0.$$

Noch allgemeiner: Für ein $\beta \in [0, 1]$:

$$H''_0 : F^{-1}(\beta) = \gamma_\beta \text{ vs. } H''_1 : F^{-1}(\beta) \neq \gamma_\beta.$$

H''_0 vs. H''_1 wird mit Hilfe des Binomialtests wie folgt getestet: Sei $X_i = \mathbb{I}(Y_i \leq \gamma_\beta)$. Unter H''_0 gilt:

$$X_i \sim \text{Bernoulli}(F(\gamma_\beta)) = \text{Bernoulli}(\beta).$$

Seien $a_1 = -\infty, b_1 = \gamma_\alpha, a_2 = b_1, b_2 = +\infty$ zwei disjunkte Klassen $(a_1, b_1], (a_2, b_2]$ in der Sprache des χ^2 -Pearson-Tests. Die Testgröße ist:

$$T_n = \sum_{i=1}^n X_i = \# \{Y_i : Y_i \leq \gamma_\beta\} \sim \text{Bin}(n, \beta), \quad p = F(\gamma_\beta)$$

Die Hypothese $F^{-1}(\beta) = \gamma_\beta$ ist äquivalent zu $H'''_0 : p = \beta$. Die Entscheidungsregel lautet dann: H'''_0 wird verworfen, falls $T_n \notin [\text{Bin}(n, \beta)_{\alpha/2}, \text{Bin}(n, \beta)_{1-\alpha/2}]$. Dies ist ein Test zum Niveau α .

1.5.2 Iterationstests auf Zufälligkeit

In manchen Fragestellungen der Biologie untersucht man eine Folge von 0 oder 1 auf ihre „Zufälligkeit“ bzw. Vorhandensein von größeren Clustern von 0 oder 1. Diese Hypothesen kann man mit Hilfe der sogenannten *Iterationstests* statistisch überprüfen.

Sei eine Stichprobe $X_i, i = 1, \dots, n$ gegeben, $X_i \in \{0, 1\}$, $\sum_{i=1}^n X_i = n_1$ die Anzahl der Einsen, $n_2 = n - n_1$ die Anzahl der Nullen, n_1, n_2 vorgegeben. Eine Realisierung von (X_1, \dots, X_n) mit $n = 18, n_1 = 12$ wäre zum Beispiel

$$x = (0, 0, 0, 1, 1, 1, 0, 1, 1, 0, 1, 1, 1, 0, 1, 1, 1, 1)$$

Es soll getestet werden, ob

H_0 : jede Folge x ist gleichwahrscheinlich vs.

H_1 : Es gibt bevorzugte Folgen (Clusterbildung)

stimmt.

Sei

$$\Omega = \left\{ x = (x_1, \dots, x_n) : x_i = 0 \text{ oder } 1, i = 1, \dots, n, \sum_{i=1}^n x_i = n_1 \right\}$$

der Stichprobenraum. Dann ist der Raum $(\Omega, \mathcal{F}, \mathbb{P})$ mit $\mathcal{F} = \mathcal{P}(\Omega)$,

$$\mathbb{P}(x) = \frac{1}{|\Omega|} = \frac{1}{\binom{n}{n_1}}$$

ein Laplacescher Wahrscheinlichkeitsraum.

Sei

$$\begin{aligned} T_n(X) &= \#\{\text{Iterationen in } X\} = \#\{\text{Teilfolgen der Nullen oder Einsen}\} \\ &= \#\{\text{Wechselstellen von 0 auf 1 oder von 1 auf 0}\} + 1. \end{aligned}$$

Zum Beispiel ist für $x = (0, 1, 1, 1, 0, 0, 0, 1, 0, 1, 1, 1, 0, 0)$, $T_n(x) = 7 = 6 + 1$.

$T_n(X)$ wird folgendermaßen als Teststatistik für H_0 vs. H_1 benutzt. H_0 wird abgelehnt, falls $T(x)$ klein ist, das heißt, falls $T_n(x) < F_{T_n}^{-1}(\alpha)$. Dies ist ein Test zum Niveau α . Wie berechnen wir die Quantile $F_{T_n}^{-1}$?

Satz 1.5.1. Unter H_0 gelten folgende Aussagen:

1.

$$\mathbb{P}(T_n = k) = \begin{cases} \frac{2 \binom{n_1-1}{i-1} \binom{n_2-1}{i-1}}{\binom{n}{n_1}}, & \text{falls } k = 2i, \\ \frac{\binom{n_1-1}{i} \binom{n_2-1}{i-1} + \binom{n_1-1}{i-1} \binom{n_2-1}{i}}{\binom{n}{n_1}}, & \text{falls } k = 2i + 1. \end{cases}$$

2.

$$\mathbb{E} T_n = 1 + \frac{2n_1 n_2}{n}$$

3.

$$\text{Var}(T_n) = \frac{2n_1 n_2 (2n_1 n_2 - n)}{n^2 (n - 1)}$$

Beweis. 1. Wir nehmen an, daß $k = 2i$ (der ungerade Fall ist analog). Wie können i Klumpen von Einsen gewählt werden? Die Anzahl dieser Möglichkeiten = die Anzahl der Möglichkeiten, wie n_1 Teilchen auf i Klassen verteilt werden.

$$0|00|\dots|0|(n_1)$$

Dies ist gleich der Anzahl an Möglichkeiten, wie $i - 1$ Trennwände auf $n_1 - 1$ Positionen verteilt werden können = $\binom{n_1-1}{i-1}$. Das selbe gilt für die Nullen.

2. Sei $Y_j = \mathbb{I}\{X_{j-1} \neq X_j\}_{j=2,\dots,n}$.

$$\begin{aligned} \Rightarrow \mathbb{E} T_n(X) &= 1 + \sum_{j=2}^n \mathbb{E} Y_j = 1 + \sum_{j=2}^n \mathbb{P}(X_{j-1} \neq X_j). \\ \mathbb{P}(X_{j-1} \neq X_j) &= \frac{2 \binom{n-2}{n_1-1}}{\binom{n}{n_1}} = 2 \cdot \frac{\frac{(n-2)!}{(n-2-(n_1-1))!(n_1-1)!}}{\frac{n!}{(n-n_1)!n_1!}} \\ &= \frac{2n_1(n-n_1)}{(n-1)n} \\ &= \frac{2n_1 n_2}{n(n-1)}. \end{aligned}$$

Daraus folgt

$$\mathbb{E} T_n = 1 + (n-1) \frac{2n_1 n_2}{n(n-1)} = 1 + 2 \frac{n_1 n_2}{n}.$$

3.

Übungsaufgabe 1.5.1. Beweisen Sie Punkt 3. □

Beispiel 1.5.1 (*Test von Wald-Wolfowitz*). Seien $Y = (Y_1, \dots, Y_n)$, $Z = (Z_1, \dots, Z_n)$ unabhängige Stichproben von unabhängigen und identisch verteilten Zufallsvariablen, $Y_i \sim F$, $Z_i \sim G$. Getestet werden soll:

$$H_0 : F = G \text{ vs. } H_1 : F \neq G.$$

Sei $(Y, Z) = (Y_1, \dots, Y_n, Z_1, \dots, Z_n)$ und seien X'_i Stichprobenvariablen von (Y, Z) , $i = 1, \dots, n$, $n = n_1 + n_2$. Wir bilden die Ordnungsstatistiken $X'_{(i)}$, $i = 1, \dots, n$ und setzen

$$X_i = \begin{cases} 1, & \text{falls } X'_{(i)} = Y_j \text{ für ein } j = 1, \dots, n_1, \\ 0, & \text{falls } X'_{(i)} = Z_j \text{ für ein } j = 1, \dots, n_2. \end{cases}$$

Unter H_0 sind die Stichprobenwerte in (Y, Z) gut gemischt, das heißt jede Kombination von 0 und 1 in (X_1, \dots, X_n) ist gleichwahrscheinlich. Darum können wir den Iterationstest auf Zufälligkeit anwenden, um H_0 vs. H_1 zu testen: H_0 wird verworfen, falls $T_n(x) \leq F^{-1}(\alpha)$, $x = (x_1, \dots, x_n)$.

Wie können die Quantile von F_{T_n} für große n berechnet werden? Falls

$$\frac{n_1}{n_1 + n_2} \xrightarrow{n \rightarrow \infty} p \in (0, 1),$$

dann ist T_n asymptotisch normalverteilt.

Satz 1.5.2. Unter der obigen Voraussetzung gilt:

$$\begin{aligned} \lim_{n \rightarrow \infty} \frac{\mathbb{E} T_n}{n} &= 2p(1-p) \\ \lim_{n \rightarrow \infty} \frac{1}{n} \text{Var } T_n &= 4p^2(1-p)^2 \\ \frac{T_n - 2p(1-p)}{2\sqrt{np(1-p)}} &\xrightarrow[n \rightarrow \infty]{d} Y \sim N(0, 1), \quad \text{falls } \frac{n_1}{n_1 + n_2} \xrightarrow{n \rightarrow \infty} p \in (0, 1). \end{aligned}$$

So können Quantile von T_n näherungsweise für große n folgendermaßen berechnet werden:

$$\begin{aligned} \alpha &= \mathbb{P}(T_n \leq F_{T_n}^{-1}(\alpha)) = \mathbb{P}\left(\frac{T_n - 2np(1-p)}{2\sqrt{np(1-p)}} \leq \frac{x - 2np(1-p)}{2\sqrt{np(1-p)}}\right) \Bigg|_{x=F_{T_n}^{-1}(\alpha)} \\ &\approx \Phi\left(\frac{F_{T_n}^{-1}(\alpha) - 2np(1-p)}{2\sqrt{np(1-p)}}\right) \\ &\Rightarrow z_\alpha \approx \frac{F_{T_n}^{-1}(\alpha) - 2np(1-p)}{2\sqrt{np(1-p)}} \end{aligned}$$

Damit erhalten wir für die Quantile:

$$F_{T_n}^{-1}(\alpha) \approx 2np(1-p) + 2\sqrt{np(1-p)} \cdot z_\alpha$$

In der Praxis setzt man $\hat{p} = \frac{n_1}{n_1 + n_2}$ für p ein.

2 Lineare Regression

In Stochastik I betrachteten wir die einfache lineare Regression der Form

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, \dots, n.$$

In Matrix-Form schreiben wir $Y = X\beta + \varepsilon$, wobei $Y = (Y_1, \dots, Y_n)^\top$ der Vektor der Zielzufallsvariablen ist,

$$X = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}$$

eine $(n \times 2)$ -Matrix, die die Ausgangsvariablen $x_i, i = 1, \dots, n$ enthält und deshalb *Design-Matrix* genannt wird, $\beta = (\beta_0, \beta_1)^\top$ der Parametervektor und $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^\top$ der Vektor der Störgrößen. Bisher waren oft $\varepsilon_i \sim N(0, \sigma^2)$ für $i = 1, \dots, n$ und $\varepsilon \sim N(0, \mathcal{I} \cdot \sigma^2)$ multivariat normalverteilt.

Die multivariate (das bedeutet, nicht einfache) lineare Regression lässt eine beliebige $(n \times m)$ -Design-Matrix

$$X = (x_{ij})_{\substack{i=1, \dots, n \\ j=1, \dots, m}}$$

und einen m -dimensionalen Parametervektor $\beta = (\beta_1, \dots, \beta_m)^\top$ zu, für $m \geq 2$. Das heißt, es gilt

$$Y = X\beta + \varepsilon, \tag{2.0.1}$$

wobei $\varepsilon \sim N(0, K)$ ein multivariat normalverteilter Zufallsvektor der Störgrößen mit Kovarianzmatrix K ist, die im Allgemeinen nicht unabhängig voneinander sind:

$$K \neq \text{diag}(\sigma_1^2, \dots, \sigma_n^2).$$

Das Ziel dieses Kapitels ist es, Schätzer und Tests für β zu entwickeln. Zuvor müssen jedoch die Eigenschaften der multivariaten Normalverteilung untersucht werden.

2.1 Multivariate Normalverteilung

Im Vorlesungsskript Wahrscheinlichkeitsrechnung wurde die multivariate Normalverteilung in Beispiel 3.4.1 folgendermaßen eingeführt:

Definition 2.1.1. Es sei $X = (X_1, \dots, X_n)^\top$ ein n -dimensionaler Zufallsvektor, $\mu \in \mathbb{R}^n$, K eine symmetrische, positiv definite $(n \times n)$ -Matrix. X ist *multivariat normalverteilt* mit den Parametern μ und K ($X \sim N(\mu, K)$), falls X absolut stetig verteilt ist mit der Dichte

$$f_X(x) = \frac{1}{(2\pi)^{n/2}} \frac{1}{\sqrt{\det(K)}} \exp \left\{ -\frac{1}{2} (x - \mu)^\top K^{-1} (x - \mu) \right\}, \quad x = (x_1, \dots, x_n)^\top \in \mathbb{R}^n.$$

Wir geben drei weitere Definitionen von $N(\mu, K)$ an und wollen die Zusammenhänge zwischen ihnen untersuchen:

Definition 2.1.2. Der Zufallsvektor $X = (X_1, \dots, X_n)^\top$ ist multivariat normalverteilt ($X \sim N(\mu, K)$) mit Parametern $\mu \in \mathbb{R}^n$ und K (eine symmetrische, nicht-negativ definite $(n \times n)$ -Matrix), falls die charakteristische Funktion $\varphi_X(t) = \mathbb{E} e^{i(t, X)}$, $t \in \mathbb{R}^n$, gegeben ist durch

$$\varphi_X(t) = \exp \left\{ it^\top \mu - \frac{1}{2} t^\top K t \right\}, \quad t \in \mathbb{R}^n.$$

Definition 2.1.3. Der Zufallsvektor $X = (X_1, \dots, X_n)^\top$ ist multivariat normalverteilt ($X \sim N(\mu, K)$) mit Parametern $\mu \in \mathbb{R}^n$ und einer symmetrischen, nicht negativ definiten $(n \times n)$ -Matrix K , falls

$$\forall a \in \mathbb{R}^n : \text{ die Zufallsvariable } (a, X) = a^\top X \sim N(a^\top \mu, a^\top K a)$$

eindimensional normalverteilt ist.

Definition 2.1.4. Es sei $\mu \in \mathbb{R}^n$, K eine nicht-negativ definite, symmetrische $(n \times n)$ -Matrix. Ein Zufallsvektor $X = (X_1, \dots, X_n)^\top$ ist multivariat normalverteilt mit Parametern μ und K ($X \sim N(\mu, K)$), falls

$$X \stackrel{d}{=} \mu + C \cdot Y,$$

wobei C eine $(n \times m)$ -Matrix mit $\text{rang}(C) = m$, $K = C \cdot C^\top$ und $Y \sim N(0, \mathcal{I}) \in \mathbb{R}^m$ ein m -dimensionaler Zufallsvektor mit unabhängigen und identisch verteilten Koordinaten $Y_j \sim N(0, 1)$ ist, $j = 1, \dots, m$.

Bemerkung: Dies ist das Analogon im eindimensionalen Fall: $Y \sim N(\mu, \sigma^2) \Leftrightarrow Y \stackrel{d}{=} \mu + \sigma X$ mit $X \sim N(0, 1)$.

Übungsaufgabe 2.1.1. Prüfen Sie, daß die in Definition 2.1.1 angegebene Dichte

$$f_X(x) = \frac{1}{(2\pi)^{n/2}} \frac{1}{\sqrt{\det(K)}} \exp \left\{ -\frac{1}{2} (x - \mu)^\top K^{-1} (x - \mu) \right\}, \quad x \in \mathbb{R}^n$$

tatsächlich eine Verteilungsdichte darstellt.

Lemma 2.1.1. Es seien X und Y n -dimensionale Zufallsvektoren mit charakteristischen Funktionen

$$\begin{aligned}\varphi_X(t) &= \mathbb{E} e^{i(t,X)} = \mathbb{E} e^{it^\top X} \\ \varphi_Y(t) &= \mathbb{E} e^{i(t,Y)} = \mathbb{E} e^{it^\top Y}\end{aligned}$$

für $t \in \mathbb{R}^n$. Es gelten folgende Eigenschaften:

1. *Eindeutigkeitssatz:*

$$X \stackrel{d}{=} Y \Leftrightarrow \varphi_X(t) = \varphi_Y(t), \quad t \in \mathbb{R}^n$$

2. Falls X und Y unabhängig sind, dann gilt:

$$\varphi_{X+Y}(t) = \varphi_X(t) \cdot \varphi_Y(t), \quad t \in \mathbb{R}^n.$$

ohne Beweis: vergleiche den Beweis des Satzes 5.1.1 (5), Folgerung 5.1.1, Vorlesungsskript WR.

Satz 2.1.1. 1. Die Definitionen 2.1.2 - 2.1.4 der multivariaten Normalverteilung sind äquivalent.

2. Die Definitionen 2.1.1 und 2.1.4 sind im Falle $n = m$ äquivalent.

Bemerkung 2.1.1. 1. Falls die Matrix K in Definition 2.1.4 den vollen Rang n besitzt, so besitzt sie die Dichte aus Definition 2.1.1. Sie wird in dem Fall *regulär* genannt.

2. Falls $\text{Rang}(K) = m < n$, dann ist die Verteilung $N(\mu, K)$ laut Definition 2.1.4 auf dem m -dimensionalen linearen Unterraum

$$\{y \in \mathbb{R}^n : y = \mu + Cx, x \in \mathbb{R}^m\}$$

konzentriert. $N(\mu, K)$ ist in diesem Fall offensichtlich nicht absolutstetig verteilt und wird daher *singulär* genannt.

Beweis. Wir beweisen: Definition 2.1.3 \Leftrightarrow 2.1.2 \Leftrightarrow 2.1.4.

1. a) Wir zeigen: Die Definitionen 2.1.2 und 2.1.3 sind äquivalent. Dazu ist zu zeigen: Für die Zufallsvariable X mit der charakteristischen Funktion

$$\varphi_X(t) = \exp\left\{it^\top \mu - \frac{1}{2}t^\top Kt\right\} \Leftrightarrow \forall a \in \mathbb{R}^n : a^\top X \sim N(a^\top \mu, a^\top K a).$$

Es gilt:

$$\varphi_{t^\top X}(1) = \mathbb{E} e^{it^\top X \cdot 1} \stackrel{\varphi_{N(\mu, \sigma^2)}}{=} \exp\left\{it^\top \mu - \frac{1}{2}t^\top Kt\right\} = \varphi_X(t) \quad \forall t \in \mathbb{R}^n.$$

(Dies nennt man das *Verfahren von Cramér-Wold*, vergleiche den multivariaten zentralen Grenzwertsatz).

- b) Wir zeigen: Die Definitionen 2.1.2 und 2.1.4 sind äquivalent. Dazu ist zu zeigen: $X = \mu + C \cdot Y$ (mit μ , C , und Y wie in Definition 2.1.4) $\Leftrightarrow \varphi_X(t) = \exp\{it^\top \mu - \frac{1}{2}t^\top Kt\}$, wobei $K = C \cdot C^\top$. Es gilt:

$$\begin{aligned} \varphi_{\mu+CY}(t) &= \mathbb{E} e^{i(t, \mu+CY)} = \mathbb{E} e^{it^\top \mu + it^\top CY} = e^{it^\top \mu} \cdot \mathbb{E} e^{i \overbrace{(C^\top t, Y)}^y} \\ &\stackrel{Y \sim N(0, \mathcal{I})}{=} e^{it^\top \mu} \cdot \exp\left(-\frac{1}{2}y^\top \cdot y\right) = \exp\left\{it^\top \mu - \frac{1}{2}t^\top C \cdot C^\top t\right\} \\ &= \exp\left\{it^\top \mu - \frac{1}{2}t^\top Kt\right\}, \quad t \in \mathbb{R}^n. \end{aligned}$$

2. Zu zeigen ist: Aus $X \sim N(\mu, K)$ im Sinne von Definition 2.1.4, $Y \sim N(\mu, K)$ im Sinne der Definition 2.1.1, $\text{Rang}(K) = n$ folgt, daß $\varphi_X = \varphi_Y$.

Aus der Definition 2.1.2 (die äquivalent zu Definition 2.1.4 ist) folgt, daß

$$\begin{aligned} \varphi_X(t) &= \exp\left\{it^\top \mu - \frac{1}{2}t^\top Kt\right\}, \quad t \in \mathbb{R}^n, \\ \varphi_Y(t) &= \mathbb{E} e^{it^\top Y} = \int_{\mathbb{R}^n} e^{it^\top y} \frac{1}{(2\pi)^{n/2}} \frac{1}{\sqrt{\det K}} \cdot \exp\left\{-\frac{1}{2}\overbrace{(y-\mu)^\top}^x K^{-1} \overbrace{(y-\mu)}^x\right\} dy \\ &= e^{it^\top \mu} \cdot \int_{\mathbb{R}^n} \frac{1}{(2\pi)^{n/2} \sqrt{\det K}} \cdot \exp\left\{it^\top x - \frac{1}{2}x^\top K^{-1}x\right\} dx \end{aligned}$$

Wir diagonalisieren K : \exists orthogonale $(n \times n)$ -Matrix V : $V^\top = V^{-1}$ und $V^\top K V = \text{diag}(\lambda_1, \dots, \lambda_n)$, wobei $\lambda_i > 0$, $i = 1, \dots, n$. Mit der neuen Substitution: $x = Vz$, $t = Vs$ erhalten wir:

$$\begin{aligned} \varphi_Y(t) &= \frac{e^{it^\top \mu}}{(2\pi)^{n/2} \sqrt{\det K}} \cdot \int_{\mathbb{R}^n} \exp\left\{is^\top V^\top V z - \frac{1}{2}z^\top V^\top K^{-1} V z\right\} dz \\ &= \frac{e^{it^\top \mu}}{\sqrt{(2\pi)^n \lambda_1 \dots \lambda_n}} \cdot \int_{\mathbb{R}} \dots \int_{\mathbb{R}} \exp\left\{is^\top z - \frac{1}{2} \sum_{i=1}^n \frac{z_i^2}{\lambda_i}\right\} dz_1 \dots dz_n \\ &= e^{it^\top \mu} \prod_{i=1}^n \int_{\mathbb{R}} \frac{1}{\sqrt{2\pi \lambda_i}} e^{is_i z_i - \frac{z_i^2}{2\lambda_i}} dz_i = e^{it^\top \mu} \cdot \prod_{i=1}^n \varphi_{N(0, \lambda_i)}(s_i) = e^{it^\top \mu} \prod_{i=1}^n e^{-\frac{s_i^2 \lambda_i}{2}} \\ &= \exp\left\{it^\top \mu - \frac{1}{2}s^\top \text{diag}(\lambda_1, \dots, \lambda_n)s\right\} = \exp\left\{it^\top \mu - \frac{1}{2}(V^\top t)^\top V^\top K V V^\top t\right\} \\ &= \exp\left\{it^\top \mu - \frac{1}{2}t^\top \underbrace{V V^\top}_{\mathcal{I}} K \underbrace{V V^\top}_{\mathcal{I}} t\right\} = \exp\left\{it^\top \mu - \frac{1}{2}t^\top Kt\right\}, \quad t \in \mathbb{R}^n. \end{aligned}$$

□

2.1.1 Eigenschaften der multivariaten Normalverteilung

Satz 2.1.2. Es sei $X = (X_1, \dots, X_n) \sim N(\mu, K)$, $\mu \in \mathbb{R}^n$, K symmetrisch und nicht-negativ definit. Dann gelten folgende Eigenschaften:

1. μ ist der Erwartungswertvektor von X :

$$\mathbb{E} X = \mu, \quad \text{das heißt: } \mathbb{E} X_i = \mu_i, \quad i = 1, \dots, n.$$

K ist die Kovarianzmatrix von X :

$$K = (k_{ij}), \quad \text{mit } k_{ij} = \text{Cov}(X_i, X_j).$$

2. Jeder Teilvektor $X' = (X_{i_1}, \dots, X_{i_k})^\top$ ($1 \leq i_1 < \dots < i_k \leq n$) von X ist ebenso multivariat normalverteilt, $X' \sim N(\mu', K')$, wobei $\mu' = (\mu_{i_1}, \dots, \mu_{i_k})^\top$, $K' = (k'_{jl}) = (\text{Cov}(X_{i_j}, X_{i_l}))$, $j, l = 1, \dots, k$. Insbesondere sind $X_i \sim N(\mu_i, k_{ii})$, wobei $k_{ii} = \text{Var } X_i$, $i = 1, \dots, n$.
3. Zwei Teilvektoren von X sind unabhängig genau dann, wenn entsprechende Elemente k_{ij} von K , die ihre Kreuzkovarianzen darstellen, Null sind, das heißt: $X' = (X_1, \dots, X_k)^\top$, $X'' = (X_{k+1}, \dots, X_n)$ unabhängig (wobei die Reihenfolge nur wegen der Einfachheit so gewählt wurde, aber unerheblich ist) $\Leftrightarrow k_{ij} = 0$ für $1 \leq i \leq k$, $j > k$ oder $i > k$, $1 \leq j \leq k$.

$$K = \left(\begin{array}{c|c} K' & 0 \\ \hline 0 & K'' \end{array} \right)$$

K' und K'' sind Kovarianzmatrizen von X' bzw. X'' .

4. *Faltungsstabilität:* Falls X und Y unabhängige, n -dimensionale Zufallsvektoren mit $X \sim N(\mu_1, K_1)$ und $Y \sim N(\mu_2, K_2)$ sind, dann ist

$$X + Y \sim N(\mu_1 + \mu_2, K_1 + K_2).$$

Übungsaufgabe 2.1.2. Beweisen Sie Satz 2.1.2.

Satz 2.1.3 (*Lineare Transformation von $N(\mu, K)$*). Sei $X \sim N(\mu, K)$ ein n -dimensionaler Zufallsvektor, A eine $(m \times n)$ -Matrix mit $\text{Rang}(A) = m \leq n$, $b \in \mathbb{R}^m$. Dann ist der Zufallsvektor $Y = AX + b$ multivariat normalverteilt:

$$Y \sim N(A\mu + b, AK A^\top).$$

Beweis. Ohne Beschränkung der Allgemeinheit setzen wir $\mu = 0$ und $b = 0$, weil $\varphi_{Y-a}(t) = e^{-it^\top a} \cdot \varphi_Y(t)$, für $a = A\mu + b$. Es ist zu zeigen:

$$Y = AX, \quad X \sim N(0, K) \Rightarrow Y \sim N(0, AK A^\top)$$

Es ist

$$\begin{aligned}\varphi_Y(t) &= \varphi_{AX}(t) = \mathbb{E} e^{it^\top AX} = \mathbb{E} e^{i(X, \overbrace{A^\top t}^{:=s})} \\ &\stackrel{(\text{Def. 2.1.2})}{=} \exp \left\{ -\frac{1}{2} s^\top K s \right\} = \exp \left\{ -\frac{1}{2} t^\top A K A^\top t \right\}, t \in \mathbb{R}^n \\ &\Rightarrow Y \sim N(0, A K A^\top).\end{aligned}$$

□

2.1.2 Lineare und quadratische Formen von normalverteilten Zufallsvariablen

Definition 2.1.5. Seien $X = (X_1, \dots, X_n)^\top$ und $Y = (Y_1, \dots, Y_n)^\top$ Zufallsvektoren auf (Ω, \mathcal{F}, P) , A eine $(n \times n)$ -Matrix aus \mathbb{R}^{n^2} , die symmetrisch ist.

1. $Z = AX$ heißt *lineare Form* von X mit Matrix A .
2. $Z = Y^\top AX$ heißt *bilineare Form* von X und Y mit Matrix A ,

$$Z = \sum_{i=1}^n \sum_{j=1}^n a_{ij} X_j Y_i.$$

3. Die Zufallsvariable $Z = X^\top AX$ (die eine bilineare Form aus 2. mit $Y = X$ ist) heißt *quadratische Form* von X mit Matrix A .

Satz 2.1.4. Sei $Z = Y^\top AX$ eine bilineare Form von Zufallsvektoren $X, Y \in \mathbb{R}^n$ bzgl. der symmetrischen Matrix A . Falls $\mu_X = \mathbb{E} X$, $\mu_Y = \mathbb{E} Y$ und $K_{XY} = (\text{Cov}(X_i, Y_j))_{i,j=1,\dots,n}$ die Kreuzkovarianzmatrix von X und Y ist, dann gilt:

$$\mathbb{E} Z = \mu_Y^\top A \mu_X + \text{Spur}(A K_{XY}).$$

Beweis.

$$\begin{aligned}\mathbb{E} Z &= \mathbb{E} \text{Spur}(Z) = \mathbb{E} \text{Spur}(Y^\top AX) \quad (\text{wegen } \text{Spur}(AB) = \text{Spur}(BA)) \\ &= \mathbb{E} \text{Spur}(AXY^\top) = \text{Spur}(A \mathbb{E}(XY^\top)) \quad (\text{wobei } XY^\top = (X_i Y_j)_{i,j=1,\dots,n}) \\ &= \text{Spur} \left(A \mathbb{E} \left((X - \mu_X) \cdot (Y - \mu_Y)^\top + \mu_X Y^\top + X \mu_Y^\top - \mu_X \mu_Y^\top \right) \right) \\ &= \text{Spur} \left(A(K_{XY} + \mu_X \mu_Y^\top + \mu_X \mu_Y^\top - \mu_X \mu_Y^\top) \right) = \text{Spur} \left(A K_{XY} + A \mu_X \mu_Y^\top \right) \\ &= \text{Spur}(A K_{XY}) + \text{Spur} \left(A \mu_X \cdot \mu_Y^\top \right) \\ &= \text{Spur} \left(\mu_Y^\top A \mu_X \right) + \text{Spur}(A K_{XY}) = \mu_Y^\top A \mu_X + \text{Spur}(A K_{XY}).\end{aligned}$$

□

Folgerung 2.1.1. Für quadratische Formen gilt

$$\mathbb{E}(X^\top AX) = \mu_X^\top A \mu_X + \text{Spur}(A \cdot K),$$

wobei $\mu_X = \mathbb{E} X$ und K die Kovarianzmatrix von X ist.

Satz 2.1.5 (*Kovarianz quadratischer Formen*). Es sei $X \sim N(\mu, K)$ ein n -dimensionaler Zufallsvektor und $A, B \in \mathbb{R}^{n^2}$ zwei symmetrische $(n \times n)$ -Matrizen. Dann gilt Folgendes:

$$\text{Cov}\left(X^\top AX, X^\top BX\right) = 4\mu^\top AKB\mu + 2 \cdot \text{Spur}(AKBK).$$

Lemma 2.1.2 (*gemischte Momente*). Es sei $Y = (Y_1, \dots, Y_n)^\top \sim N(0, K)$ ein Zufallsvektor. Dann gilt Folgendes:

$$\begin{aligned} \mathbb{E}(Y_i Y_j Y_k) &= 0, \\ \mathbb{E}(Y_i Y_j Y_k Y_l) &= k_{ij} \cdot k_{kl} + k_{ik} \cdot k_{jl} + k_{jk} \cdot k_{il}, \quad 1 \leq i, j, k, l \leq n, \end{aligned}$$

wobei $K = (k_{ij})_{i,j=1,\dots,n}$ die Kovarianzmatrix von Y ist.

Übungsaufgabe 2.1.3. Beweisen Sie dieses Lemma.

folgt:

$$\begin{aligned} \text{Cov}\left(X^\top AX, X^\top BX\right) &= 2 \cdot \text{Spur}(AKBK) + \text{Spur}(AK) \cdot \text{Spur}(BK) + 4\mu^\top AKB\mu \\ &\quad - \text{Spur}(AK) \cdot \text{Spur}(BK) = 4\mu^\top AKB\mu + 2 \cdot \text{Spur}(AKBK). \end{aligned}$$

□

Folgerung 2.1.2.

$$\text{Var}\left(X^\top AX\right) = 4\mu^\top AKA\mu + 2 \cdot \text{Spur}\left((AK)^2\right)$$

Satz 2.1.6. Es seien $X \sim N(\mu, K)$ und $A, B \in \mathbb{R}^{n^2}$ zwei symmetrische Matrizen. Dann gilt:

$$\text{Cov}(BX, X^\top AX) = 2BKA\mu$$

Beweis.

$$\begin{aligned} \text{Cov}(BX, X^\top AX) &\stackrel{\text{(Folgerung 2.1.1)}}{=} \mathbb{E}\left[(BX - B\mu)(X^\top AX - \mu^\top A\mu - \text{Spur}(AK))\right] \\ &= \mathbb{E}\left[B(X - \mu)\left((X - \mu)^\top A(X - \mu) + 2\mu^\top AX - 2\mu^\top A\mu - \text{Spur}(AK)\right)\right], \end{aligned}$$

denn

$$(X - \mu)^\top A(X - \mu) = X^\top AX - \mu^\top AX - X^\top A\mu + \mu^\top A\mu$$

und mit der Substitution $Z = X - \mu$ (und damit $\mathbb{E}Z = 0$)

$$\begin{aligned} \text{Cov}(BX, X^\top AX) &= \mathbb{E}\left[BZ(Z^\top AZ + 2\mu^\top AZ - \text{Spur}(AK))\right] \\ &= \mathbb{E}(BZ \cdot Z^\top AZ) + 2\mathbb{E}(BZ \cdot \mu^\top AZ) - \text{Spur}(AK) \cdot \overbrace{\mathbb{E}(BZ)}^{=B\mathbb{E}Z=0} \\ &= 2\mathbb{E}(BZ \cdot Z^\top A\mu) + \mathbb{E}(BZZ^\top AZ) = 2B \underbrace{\mathbb{E}(ZZ^\top)}_{\text{Cov}X=K} A\mu \\ &\quad + B \cdot \underbrace{\mathbb{E}(ZZ^\top AZ)}_{=0} = 2BKA\mu, \end{aligned}$$

wegen $Z \sim N(0, K)$ und Lemma 2.1.2 und dem Beweis von Satz 2.1.5. □

Definition 2.1.6. Es seien $X_i \sim N(\mu_i, 1)$, $i = 1, \dots, n$ unabhängig. Dann besitzt die Zufallsvariable

$$Y = X_1^2 + \dots + X_n^2$$

die sogenannte *nicht-zentrale* $\chi_{n,\mu}^2$ -Verteilung mit n Freiheitsgraden und dem *Nichtzentralitätsparameter*

$$\mu = \sum_{i=1}^n \mu_i^2.$$

(in Stochastik I betrachteten wir den Spezialfall der zentralen χ_n^2 -Verteilung mit $\mu = 0$).

In Bemerkung 5.2.1, Vorlesungsskript WR, haben wir momenterzeugende Funktionen von Zufallsvariablen eingeführt. Jetzt benötigen wir für den Beweis des Satzes 2.1.7 folgenden Eindeutigkeitssatz:

Lemma 2.1.3 (*Eindeutigkeitssatz für momenterzeugende Funktionen*). Es seien X_1 und X_2 zwei absolutstetige Zufallsvariablen mit momenterzeugenden Funktionen

$$M_{X_i}(t) = \mathbb{E} e^{tX_i}, \quad i = 1, 2,$$

die auf einem Intervall (a, b) definiert sind. Falls f_1 und f_2 die Dichten der Verteilung von X_1 und X_2 sind, dann gilt

$$f_1(x) = f_2(x) \text{ für fast alle } x \in \mathbb{R} \Leftrightarrow M_{X_1}(t) = M_{X_2}(t), t \in (a, b).$$

Ohne Beweis.

Satz 2.1.7. Die Dichte einer $\chi_{n,\mu}^2$ -verteilten Zufallsvariable X (mit $n \in \mathbb{N}$ und $\mu > 0$) ist gegeben durch die Mischung der Dichten von χ_{n+2j}^2 -Verteilungen mit Mischungsvariable $J \sim \text{Poisson}(\mu/2)$:

$$f_X(x) = \begin{cases} \sum_{j=0}^{\infty} e^{-\mu/2} \frac{(\mu/2)^j}{j!} \cdot \frac{e^{-x/2} x^{\frac{n+2j}{2}-1}}{\Gamma(\frac{n+2j}{2}) \cdot 2^{\frac{n+2j}{2}}}, & x \geq 0, \\ 0, & x < 0. \end{cases} \quad (2.1.1)$$

Beweis. 1. Wir berechnen zuerst $M_X(t)$, $X \sim \chi_{n,\mu}^2$:

$$\begin{aligned} M_X(t) &= \mathbb{E}(e^{tX}) = \mathbb{E} \exp \left\{ t \sum_{i=1}^n X_i^2 \right\} \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}} \cdot \int_{-\infty}^{\infty} e^{tx_i^2} \cdot e^{-\frac{(x_i - \mu_i)^2}{2}} dx_i \quad \left(t < \frac{1}{2}, X_i \sim N(\mu_i, 1) \right) \end{aligned}$$

Es gilt:

$$\begin{aligned}
 tx_i^2 - \frac{(x_i - \mu_i)^2}{2} &= \frac{1}{2}(2tx_i^2 - x_i^2 + 2x_i\mu_i - \mu_i^2) \\
 &= -\frac{1}{2}\left(x_i^2(1-2t) - 2x_i\mu_i + \frac{\mu_i^2}{(1-2t)} - \frac{\mu_i^2}{(1-2t)} + \mu_i^2\right) \\
 &= -\frac{1}{2}\left(\left(x_i \cdot \sqrt{1-2t} - \frac{\mu_i}{\sqrt{1-2t}}\right)^2 + \mu_i^2\left(1 - \frac{1}{1-2t}\right)\right) \\
 &= -\frac{1}{2}\left(\frac{(x_i(1-2t) - \mu_i)^2}{1-2t} - \mu_i^2 \cdot \frac{2t}{1-2t}\right)
 \end{aligned}$$

Wir substituieren

$$y_i = \frac{(x_i \cdot (1-2t) - \mu_i)}{\sqrt{1-2t}}$$

und erhalten

$$\begin{aligned}
 M_X(t) &= (1-2t)^{-\frac{n}{2}} \prod_{i=1}^n \exp\left\{\mu_i^2 \cdot \left(\frac{t}{1-2t}\right)\right\} \cdot \underbrace{\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-\frac{y_i^2}{2}} dy_i}_{=1} \\
 &= (1-2t)^{-\frac{n}{2}} \cdot \exp\left\{\frac{t}{1-2t} \cdot \sum_{i=1}^n \mu_i^2\right\} = \frac{1}{(1-2t)^{n/2}} \cdot \exp\left\{\frac{\mu t}{1-2t}\right\}, \quad t < \frac{1}{2}.
 \end{aligned}$$

2. Es sei Y eine Zufallsvariable mit der Dichte (2.1.1). Wir berechnen $M_Y(t)$:

$$\begin{aligned}
 M_Y(t) &= \sum_{j=0}^{\infty} e^{-\frac{\mu}{2}} \frac{(\mu/2)^j}{j!} \cdot \underbrace{\int_0^{\infty} e^{xt} \cdot \frac{e^{-\frac{x}{2}} \cdot x^{\frac{n+2j}{2}-1}}{\Gamma\left(\frac{n+2j}{2}\right) \cdot \frac{n+2j}{2}} dx}_{=M_{\chi_{n+2j}^2}(t) = \frac{1}{(1-2t)^{(n+2j)/2}} \text{ (Stochastik I, Satz 3.2.1)}} \\
 &= \frac{e^{-\frac{\mu}{2}}}{(1-2t)^{\frac{n}{2}}} \cdot \sum_{j=1}^{\infty} \left(\frac{\mu}{2(1-2t)}\right)^j \cdot \frac{1}{j!} \\
 &= \frac{1}{(1-2t)^{\frac{n}{2}}} \cdot \exp\left\{-\frac{\mu}{2} + \frac{\mu}{2(1-2t)}\right\} = \frac{1}{(1-2t)^{\frac{n}{2}}} \cdot \exp\left\{\frac{\mu \cdot (1 - (1-2t))}{2 \cdot (1-2t)}\right\} \\
 &= (1-2t)^{-\frac{n}{2}} \cdot \exp\left\{\frac{\mu t}{1-2t}\right\} \\
 \implies M_X(t) &= M_Y(t), \quad t < \frac{1}{2}
 \end{aligned}$$

Nach Lemma 2.1.3 gilt dann, $f_X(x) = f_Y(x)$ für fast alle $x \in \mathbb{R}$.

□

Bemerkung 2.1.2. 1. Die Definition 2.1.6 kann in folgender Form umgeschrieben werden:

Falls $X \sim N(\vec{\mu}, \mathcal{I})$, $\vec{\mu} = (\mu_1, \dots, \mu_n)^\top$, dann gilt $|X|^2 = X^\top X \sim \chi_{n,\mu}^2$, wobei $\mu = |\vec{\mu}|^2$.

2. Die obige Eigenschaft kann auf $X \sim N(\vec{\mu}, K)$, mit einer symmetrischen, positiv definiten $(n \times n)$ -Matrix K verallgemeinert werden:

$$X^\top K^{-1} X \sim \chi_{n,\tilde{\mu}}^2, \quad \text{wobei } \tilde{\mu} = \vec{\mu}^\top K^{-1} \vec{\mu},$$

denn weil K positiv definit ist, gibt es ein $K^{\frac{1}{2}}$, sodaß $K = K^{\frac{1}{2}} K^{\frac{1}{2}\top}$. Dann gilt

$$Y = K^{-\frac{1}{2}} X \sim N(K^{-\frac{1}{2}} \mu, \mathcal{I}), \quad \text{weil } K^{-\frac{1}{2}} K K^{-\frac{1}{2}\top} = K^{-\frac{1}{2}} \cdot K^{\frac{1}{2}} \cdot K^{\frac{1}{2}\top} \cdot K^{-\frac{1}{2}\top} = \mathcal{I}$$

und daher

$$X^\top K^{-1} X = Y^\top Y \stackrel{\text{Punkt 1}}{\sim} \chi_{n,\tilde{\mu}}^2, \quad \text{mit } \tilde{\mu} = \left(K^{-\frac{1}{2}} \vec{\mu}\right)^\top K^{-\frac{1}{2}} \vec{\mu} = \vec{\mu}^\top K^{-\frac{1}{2}\top} K^{-\frac{1}{2}} \vec{\mu} = \vec{\mu}^\top K^{-1} \vec{\mu}.$$

Satz 2.1.8. Es sei $X \sim N(\mu, K)$, wobei K eine symmetrische, positiv definite $(n \times n)$ -Matrix ist, und sei A eine weitere symmetrische $(n \times n)$ -Matrix mit der Eigenschaft $AK = (AK)^2$ (Idempotenz) und $\text{Rang}(A) = r \leq n$. Dann gilt:

$$X^\top AX \sim \chi_{r,\tilde{\mu}}^2, \quad \text{wobei } \tilde{\mu} = \mu^\top A \mu.$$

Beweis. Wir zeigen, daß A nicht negativ definit ist.

$$\begin{aligned} AK &= (AK)^2 = AK \cdot AK \quad | \quad K^{-1} \\ \implies A &= AK A \implies \forall x \in \mathbb{R}^n : x^\top Ax = x^\top AK Ax \\ &= \underbrace{(Ax)^\top}_{=y} K \underbrace{(Ax)}_{=y} \geq 0 \quad \text{wegen der positiven Definitheit von } K. \\ \implies A &\text{ ist nicht negativ definit.} \\ \implies \exists H &: \text{ eine } (n \times r)\text{-Matrix mit } \text{Rang}(H) = r : A = HH^\top \end{aligned}$$

Somit gilt

$$X^\top AX = X^\top H \cdot H^\top X = \underbrace{(H^\top X)^\top}_{=Y} \cdot H^\top X = Y^\top Y$$

Es gilt: $Y \sim N(H^\top \mu, \mathcal{I}_r)$, denn nach Satz 2.1.3 ist $Y \sim N(H^\top \mu, H^\top K H)$ und $\text{Rang}(H) = r$. Das heißt, $H^\top H$ ist eine invertierbare $(r \times r)$ -Matrix, und

$$\begin{aligned} H^\top K H &= (H^\top H)^{-1} \underbrace{(H^\top H \cdot H^\top K H \cdot (H^\top H))}_{=AKA=A} (H^\top H)^{-1} \\ &= (H^\top H)^{-1} H^\top \cdot \underbrace{A}_{=HH^\top} \cdot H (H^\top H)^{-1} \\ &= \mathcal{I}_r \end{aligned}$$

Dann ist

$$X^\top AX = |Y|^2 \sim \chi_{r, \tilde{\mu}}^2 \text{ mit } \tilde{\mu} = (H^\top \mu)^2 = \mu^\top H \cdot H^\top \mu = \mu^\top A \mu.$$

□

Satz 2.1.9 (Unabhängigkeit). Es sei $X \sim N(\mu, K)$ und K eine symmetrische, nicht-negativ definite $(n \times n)$ -Matrix.

1. Es seien A, B ($r_1 \times n$) bzw. ($r_2 \times n$)-Matrizen, $r_1, r_2 \leq n$ mit $AKB^\top = 0$. Dann sind die Vektoren AX und BX unabhängig.
2. Sei ferner C eine symmetrische, nicht-negativ definite $(n \times n)$ -Matrix mit der Eigenschaft $AKC = 0$. Dann sind AX und $X^\top CX$ unabhängig.

Beweis. 1. Nach Satz 2.1.2, 3) gilt: AX und BX sind unabhängig $\iff \varphi_{(AX, BX)}(t) = \varphi_{AX}(t) \cdot \varphi_{BX}(t)$, $t = (t_1, t_2)^\top \in \mathbb{R}^{r_1+r_2}$, $t_1 \in \mathbb{R}^{r_1}$, $t_2 \in \mathbb{R}^{r_2}$. Es ist zu zeigen:

$$\varphi_{(AX, BX)}(t) = \mathbb{E} e^{i(t_1^\top A + t_2^\top B) \cdot X} \stackrel{!}{=} \mathbb{E} e^{i t_1^\top AX} \cdot \mathbb{E} e^{i t_2^\top BX}.$$

Es gilt

$$\varphi_{(AX, BX)}(t) = \mathbb{E} e^{i(t_1^\top A + t_2^\top B) \cdot X} \stackrel{(Def. 2.1.2)}{=} e^{i(t_1^\top A + t_2^\top B) \cdot \mu - \frac{1}{2} \cdot (t_1^\top A + t_2^\top B) \cdot K \cdot (t_1^\top A + t_2^\top B)^\top},$$

und mit

$$\begin{aligned} & (t_1^\top A + t_2^\top B) \cdot K \cdot (t_1^\top A + t_2^\top B)^\top \\ &= (t_1^\top A) K (t_1^\top A)^\top + (t_1^\top A)^\top K (t_2^\top B) + (t_2^\top B) K (t_1^\top A)^\top + (t_2^\top B) K (t_2^\top B)^\top \\ &= t_1^\top A K A^\top t_1 + t_1^\top \cdot \underbrace{AKB^\top}_{=0} \cdot t_2 + t_2^\top \cdot \underbrace{BK A^\top}_{=(AKB^\top)^\top=0} \cdot t_1 + t_2^\top B K B^\top t_2 \end{aligned}$$

ist

$$\begin{aligned} \varphi_{(AX, BX)}(t) &= e^{i t^\top A - \frac{1}{2} t_1^\top A K A^\top t_1} \cdot e^{i t_2^\top B - \frac{1}{2} t_2^\top B K B^\top t_2} \\ &= \varphi_{AX}(t_1) \cdot \varphi_{BX}(t_2), \quad t_1 \in \mathbb{R}^{r_1}, t_2 \in \mathbb{R}^{r_2} \end{aligned}$$

2. C ist symmetrisch, nicht-negativ definit \implies Es gibt eine $(n \times r)$ -Matrix H mit $\text{Rang}(H) = r \leq n$ und $C = HH^\top$, $\implies H^\top H$ hat Rang r und ist somit invertierbar. Dann gilt:

$$X^\top CX = X^\top HH^\top X = (H^\top X)^\top \cdot H^\top X = |H^\top X|^2.$$

Falls AX und $H^\top X$ unabhängig sind, dann sind auch AX und $X^\top CX = |H^\top X|^2$ unabhängig, nach dem Transformationssatz für Zufallsvektoren. Nach 1) sind AX und $H^\top X$ unabhängig, falls $AK(H^\top)^\top = AKH = 0$. Da nach Voraussetzung

$$AKC = AKH \cdot H^\top = 0 \implies AKH \cdot H^\top H = 0,$$

da aber $\exists(H^\top H)^{-1}$, folgt, daß

$$\begin{aligned} 0 &= AKH \cdot H^\top H \cdot (H^\top H)^{-1} = AKH \implies AKH = 0 \\ &\implies AX \text{ und } H^\top X \text{ sind unabhängig} \\ &\implies AX \text{ und } X^\top CX \text{ sind unabhängig.} \end{aligned}$$

□

2.2 Multivariate lineare Regressionsmodelle mit vollem Rang

Die *multivariate lineare Regression* hat die Form

$$Y = X\beta + \varepsilon,$$

wobei $Y = (Y_1, \dots, Y_n)^\top$ der Zufallsvektor der Zielvariablen ist,

$$X = (x_{ij})_{\substack{i=1, \dots, n \\ j=1, \dots, m}}$$

ist eine deterministische *Design-Matrix* mit vollem Rang, $\text{Rang}(X) = r = m \leq n$, $\beta = (\beta_1, \dots, \beta_m)^\top$ ist der *Parametervektor* und $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^\top$ ist der Zufallsvektor der *Störgrößen*, mit $\mathbb{E}\varepsilon_i = 0$, $\text{Var}\varepsilon_i = \sigma^2 > 0$. Das Ziel dieses Abschnittes wird sein, β und σ^2 geeignet zu schätzen.

2.2.1 Methode der kleinsten Quadrate

Sei $X = (X_1, \dots, X_m)$, wobei die deterministischen Vektoren $X_j = (x_{1j}, x_{2j}, \dots, x_{nj})^\top$, $j = 1, \dots, m$ einen m -dimensionalen linearen Unterraum $L_X = \langle X_1, \dots, X_m \rangle$ aufspannen. Sei

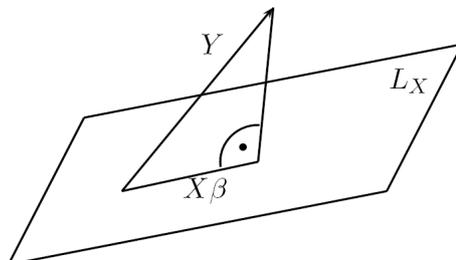
$$e(\beta) = \frac{1}{n} \|Y - X\beta\|^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - x_{i1}\beta_1 - \dots - x_{im}\beta_m)^2$$

die mittlere quadratische Abweichung zwischen Y und $X\beta$.

Der *MKQ-Schätzer* $\hat{\beta}$ für β ist definiert durch

$$\hat{\beta} = \text{argmin}(e(\beta)). \quad (2.2.1)$$

Warum existiert eine Lösung $\beta \in \mathbb{R}^m$ des quadratischen Optimierungsproblems (2.2.1)? Geometrisch kann $X\hat{\beta}$ als die orthogonale Projektion des Datenvektors Y auf den linearen Unterraum L_X interpretiert werden. Formal zeigen wir die Existenz der Lösung mit folgendem Satz.

Abbildung 2.1: Projektion auf den linearen Unterraum L_X 

Satz 2.2.1. Unter den obigen Voraussetzungen existiert der eindeutig bestimmte MKQ-Schätzer $\hat{\beta}$, der die Lösung der sogenannten *Normalgleichung* ist:

$$X^T X \beta = X^T Y. \quad (2.2.2)$$

Daher gilt:

$$\hat{\beta} = (X^T X)^{-1} X^T Y.$$

Beweis. Die notwendige Bedingung für die Existenz des Minimums ist $e'(\beta) = 0$, das heißt

$$e'(\beta) = \left(\frac{\partial e(\beta)}{\partial \beta_1}, \dots, \frac{\partial e(\beta)}{\partial \beta_m} \right)^T = 0.$$

Es gilt:

$$e'(\beta) = \frac{2}{n} (X^T X \beta - X^T Y)$$

$\implies \hat{\beta}$ ist eine Lösung der Normalgleichung $X^T X \beta = X^T Y$. Wir zeigen die hinreichende Bedingung des Minimums:

$$e''(\beta) = \left(\frac{\partial^2 e(\beta)}{\partial \beta_i \partial \beta_j} \right)_{i,j=1,\dots,m} = \frac{2}{n} X^T X.$$

$X^T X$ ist symmetrisch und positiv definit, weil X einen vollen Rang hat:

$$\forall y \neq 0, y \in \mathbb{R}^m : y^T X^T X y = (Xy)^T Xy = |Xy|^2 > 0$$

und aus $y \neq 0 \implies Xy \neq 0$, folgt, daß $e''(\beta)$ positiv definit ist. Also ist $X^\top X$ invertierbar. Das heißt, $\hat{\beta}$ ist der Minimumpunkt von $e(\beta)$. Den Schätzer $\hat{\beta} = (X^\top X)^{-1} X^\top Y$ bekommt man, indem man die Normalengleichung $X^\top X\beta = X^\top Y$ von links mit $(X^\top X)^{-1}$ multipliziert. \square

Beispiel 2.2.1. 1. *Einfache lineare Regression*

$$X = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \quad m = 2, \beta = (\beta_1, \beta_2)^\top, Y = X\beta + \varepsilon$$

$\hat{\beta} = (\hat{\beta}_1, \hat{\beta}_2)$ ergibt den MKQ-Schätzer aus der Stochastik I

$$\hat{\beta}_2 = \frac{S_{XY}^2}{S_{XX}^2}, \quad \hat{\beta}_1 = \bar{Y}_n - \bar{X}_n \hat{\beta}_2,$$

wobei

$$\begin{aligned} \bar{X}_n &= \frac{1}{n} \sum_{i=1}^n X_i, & \bar{Y}_n &= \frac{1}{n} \sum_{i=1}^n Y_i \\ S_{XY}^2 &= \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)(Y_i - \bar{Y}_n) \\ S_{XX}^2 &= \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \end{aligned}$$

Übungsaufgabe 2.2.1. Beweisen Sie dies!

2. *Multiple lineare Regression*

$Y = X\beta + \varepsilon$ mit Designmatrix

$$X = \begin{pmatrix} 1 & x_{11} & \cdots & x_{1m} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & \cdots & x_{nm} \end{pmatrix} \quad \text{für } \beta = (\beta_0, \beta_1, \dots, \beta_m)^\top.$$

Der MKQ-Schätzer $\hat{\beta} = (X^\top X)^{-1} X^\top Y$ ist offensichtlich ein linearer Schätzer bezüglich Y .

Wir werden jetzt zeigen, daß $\hat{\beta}$ der *beste lineare, erwartungstreue Schätzer* von β (im Englischen *BLUE = best linear unbiased estimator*) in der Klasse

$$\mathcal{L} = \left\{ \tilde{\beta} = AY + b : \mathbb{E} \tilde{\beta} = \beta \right\}$$

aller linearen erwartungstreuen Schätzer ist.

Satz 2.2.2 (*Güteeigenschaften des MKQ-Schätzers $\hat{\beta}$*). Es sei $Y = X\beta + \varepsilon$ ein multivariates lineares Regressionsmodell mit vollem Rang m und Störgrößen $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^\top$, die folgende Voraussetzungen erfüllen:

$$\mathbb{E}\varepsilon = 0, \quad \text{Cov}(\varepsilon_i, \varepsilon_j) = \sigma^2 \delta_{ij}, \quad i, j = 1, \dots, n \text{ für ein } \sigma^2 \in (0, \infty).$$

Dann gilt Folgendes:

1. Der MKQ-Schätzer $\hat{\beta} = (X^\top X)^{-1} X^\top Y$ ist erwartungstreu: $\mathbb{E}\hat{\beta} = \beta$.
2. $\text{Cov}(\hat{\beta}) = \sigma^2 (X^\top X)^{-1}$
3. $\hat{\beta}$ besitzt die minimale Varianz:

$$\forall \tilde{\beta} \in \mathcal{L} : \quad \text{Var} \tilde{\beta}_j \geq \text{Var} \hat{\beta}_j, \quad j = 1, \dots, m.$$

Beweis. 1. Es gilt:

$$\begin{aligned} \mathbb{E}\hat{\beta} &= \mathbb{E} \left[(X^\top X)^{-1} X^\top (X\beta + \varepsilon) \right] \\ &= (X^\top X)^{-1} \cdot X^\top X \cdot \beta + (X^\top X)^{-1} X^\top \cdot \underbrace{\mathbb{E}\varepsilon}_{=0} \\ &= \beta \quad \forall \beta \in \mathbb{R}^m. \end{aligned}$$

2. Für alle $\tilde{\beta} = AY + b \in \mathcal{L}$ gilt:

$$\begin{aligned} \beta &= \mathbb{E}\tilde{\beta} = A\mathbb{E}Y + b = AX\beta + b \quad \forall \beta \in \mathbb{R}^m. \\ &\implies b = 0, \quad AX = \mathcal{I}. \\ &\implies \tilde{\beta} = AY = A(X\beta + \varepsilon) = AX\beta + A\varepsilon \\ &= \beta + A\varepsilon. \end{aligned}$$

Für

$$\hat{\beta} = \underbrace{(X^\top X)^{-1} X^\top}_=A Y$$

gilt:

$$\begin{aligned} \text{Cov}\hat{\beta} &= \left(\mathbb{E} \left(\left(\hat{\beta}_i - \beta_i \right) \left(\hat{\beta}_j - \beta_j \right) \right) \right)_{i,j=1,\dots,m} \\ &= \mathbb{E} \left(A\varepsilon \cdot (A\varepsilon)^\top \right) = \mathbb{E} \left(A\varepsilon\varepsilon^\top A^\top \right) = A\mathbb{E} \left(\varepsilon\varepsilon^\top \right) \cdot A^\top \\ &= A \cdot \sigma^2 \mathcal{I} A^\top = \sigma^2 A A^\top = \sigma^2 (X^\top X)^{-1} X^\top \left((X^\top X)^{-1} X^\top \right)^\top \\ &= \sigma^2 (X^\top X)^{-1} X^\top X (X^\top X)^{-1} = \sigma^2 (X^\top X)^{-1}. \end{aligned}$$

3. Sei $\tilde{\beta} \in \mathcal{L}$, $\tilde{\beta} = \beta + A\varepsilon$. Zu zeigen ist, daß

$$\left(\text{Cov}(\tilde{\beta})\right)_{ii} = \sigma^2(AA^\top)_{ii} \geq \left(\text{Cov}(\hat{\beta})\right)_{ii} = \sigma^2(X^\top X)_{ii}^{-1}, \quad i = 1, \dots, m.$$

Sei $D = A - (X^\top X)^{-1}X^\top$, dann folgt: $A = D + (X^\top X)^{-1}X^\top$,

$$\begin{aligned} AA^\top &= \left(D + (X^\top X)^{-1}X^\top\right) \left(D^\top + X(X^\top X)^{-1\top}\right) \\ &= DD^\top + (X^\top X)^{-1}, \text{ weil} \end{aligned}$$

$$DX(X^\top X)^{-1} = \underbrace{(AX)}_{=I} - \underbrace{(X^\top X)^{-1}X^\top X}_{=I} (X^\top X)^{-1} = 0$$

$$(X^\top X)^{-1}X^\top D^\top = \left(DX(X^\top X)^{-1}\right)^\top = 0.$$

$$\implies (AA^\top)_{ii} = \underbrace{(DD^\top)_{ii}}_{\geq 0} + (X^\top X)_{ii}^{-1} \geq (X^\top X)_{ii}^{-1}$$

$$\implies \text{Var } \hat{\beta}_i \leq \text{Var } \tilde{\beta}_i, \quad i = 1, \dots, m.$$

□

Satz 2.2.3. Es sei $\hat{\beta}_n$ der MKQ-Schätzer im oben eingeführten multivariaten linearen Regressionsmodell. Sei $\{a_n\}_{n \in \mathbb{N}}$ eine Zahlenfolge mit $a_n \neq 0$, $n \in \mathbb{N}$, $a_n \rightarrow 0$ ($n \rightarrow \infty$). Es wird vorausgesetzt, daß eine invertierbare $(m \times m)$ -Matrix Q existiert mit

$$Q = \lim_{n \rightarrow \infty} a_n (X_n^\top X_n).$$

Dann ist $\hat{\beta}_n$ schwach konsistent:

$$\hat{\beta}_n \xrightarrow[n \rightarrow \infty]{p} \beta.$$

Beweis.

$$\hat{\beta}_n \xrightarrow[n \rightarrow \infty]{p} \beta \iff \mathbb{P}\left(\left|\hat{\beta}_n - \beta\right| > \varepsilon\right) \xrightarrow[n \rightarrow \infty]{} 0 \quad \forall \varepsilon > 0.$$

$$\begin{aligned}
\mathbb{P}\left(\left|\hat{\beta}_n - \beta\right| > \varepsilon\right) &= \mathbb{P}\left(\left|\hat{\beta}_n - \beta\right|^2 > \varepsilon^2\right) \\
&= \mathbb{P}\left(\sum_{i=1}^m \left|\hat{\beta}_{in} - \beta_i\right|^2 > \varepsilon^2\right) \leq \mathbb{P}\left(\bigcup_{i=1}^m \left\{\left|\hat{\beta}_{in} - \beta_i\right|^2 > \frac{\varepsilon^2}{m}\right\}\right) \\
&\leq \sum_{i=1}^m \mathbb{P}\left(\left|\hat{\beta}_{in} - \beta_i\right| > \frac{\varepsilon}{\sqrt{m}}\right) \\
&\leq m \sum_{i=1}^m \frac{\text{Var} \hat{\beta}_{in}}{\varepsilon^2} \xrightarrow{n \rightarrow \infty} 0, \quad (\text{aus der Ungleichung von Tschebyschew}) \\
&\text{falls } \text{Var} \hat{\beta}_{in} \xrightarrow{n \rightarrow \infty} 0, \quad i = 1, \dots, m.
\end{aligned}$$

$\text{Var} \hat{\beta}_{in}$ ist ein Diagonaleintrag von der Matrix

$$\text{Cov} \hat{\beta}_n \stackrel{(\text{Satz 2.2.2})}{=} \sigma^2 \left(X_n^\top X_n\right)^{-1}.$$

Wenn wir zeigen, daß $\text{Cov} \hat{\beta}_n \xrightarrow{n \rightarrow \infty} 0$, ist der Satz bewiesen.

Es existiert

$$Q^{-1} = \lim_{n \rightarrow \infty} \frac{1}{a_n} \left(X_n^\top X_n\right)^{-1}$$

und damit gilt:

$$\begin{aligned}
\lim_{n \rightarrow \infty} \text{Cov} \hat{\beta}_n &= \sigma^2 \lim_{n \rightarrow \infty} \left(X_n^\top X_n\right)^{-1} = \sigma^2 \lim_{n \rightarrow \infty} a_n \cdot \frac{1}{a_n} \left(X_n^\top X_n\right)^{-1} \\
&= 0 \cdot Q^{-1} \cdot \sigma^2 = 0.
\end{aligned}$$

□

2.2.2 Schätzer der Varianz σ^2

Wir führen den Schätzer $\hat{\sigma}^2$ für die Varianz σ^2 der Störgrößen ε_i folgendermaßen ein:

$$\hat{\sigma}^2 = \frac{1}{n-m} \left|Y - X\hat{\beta}\right|^2. \quad (2.2.3)$$

Dies ist eine verallgemeinerte Version des Varianzschätzers aus der einfachen linearen Regression, die wir bereits in Stochastik I kennenlernten. Dabei ist $\hat{Y} = Y - X\hat{\beta}$ der Vektor der Residuen.

Satz 2.2.4 (*Erwartungstreue*). Der Varianzschätzer

$$\hat{\sigma}^2 = \frac{1}{n-m} \left|Y - X\hat{\beta}\right|^2$$

ist erwartungstreu. Das heißt,

$$\mathbb{E} \hat{\sigma}^2 = \sigma^2.$$

Beweis.

$$\begin{aligned}\hat{\sigma}^2 &= \frac{1}{n-m} (Y - X\hat{\beta})^\top (Y - X\hat{\beta}) \\ &= \frac{1}{n-m} (Y - X(X^\top X)^{-1}X^\top Y)^\top (Y - X(X^\top X)^{-1}X^\top Y) \\ &= \frac{1}{n-m} (DY)^\top DY\end{aligned}$$

wobei $D = \mathcal{I} - X(X^\top X)^{-1}X^\top$ eine $(n \times n)$ -Matrix ist. Dann ist

$$\hat{\sigma}^2 = \frac{1}{n-m} Y^\top D^\top DY = \frac{1}{n-m} Y^\top D^2 Y = \frac{1}{n-m} Y^\top DY, \text{ falls}$$

$D^\top = D$ und $D^2 = D$ (das heißt, daß D symmetrisch und idempotent ist). Tatsächlich gilt:

$$\begin{aligned}D^\top &= \mathcal{I} - (X^\top)^\top (X^\top X)^{\top -1} X^\top = \mathcal{I} - X (X^\top X)^{-1} X^\top = D. \\ D^2 &= (\mathcal{I} - X(X^\top X)^{-1}X^\top) (\mathcal{I} - X(X^\top X)^{-1}X^\top) \\ &= \mathcal{I} - 2X(X^\top X)^{-1}X^\top + X(X^\top X)^{-1}X^\top X(X^\top X)^{-1}X^\top \\ &= \mathcal{I} - X(X^\top X)^{-1}X^\top = D.\end{aligned}$$

Weiterhin gilt:

$$\begin{aligned}\hat{\sigma}^2 &= \frac{1}{n-m} \cdot \text{Spur} (Y^\top DY) = \frac{1}{n-m} \cdot \text{Spur} (DYY^\top) \\ \implies \mathbb{E} \hat{\sigma}^2 &= \frac{1}{n-m} \cdot \text{Spur} (D \mathbb{E} (YY^\top)) = \frac{\sigma^2}{n-m} \cdot \text{Spur} (D),\end{aligned}$$

denn

$$\begin{aligned}\text{Spur} (D \cdot \mathbb{E} (YY^\top)) &= \text{Spur} (D(X\beta)(X\beta)^\top + \underbrace{DX\beta \mathbb{E} \varepsilon^\top}_{=0} + \underbrace{D \mathbb{E} \varepsilon (X\beta)^\top}_{=0} + \underbrace{D \cdot \mathbb{E} \varepsilon \varepsilon^\top}_{= \text{Cov} \varepsilon = \sigma^2 \cdot \mathcal{I}})\end{aligned}$$

und

$$\begin{aligned}DX &= (\mathcal{I} - X(X^\top X)^{-1}X^\top) X \\ &= X - X(X^\top X)^{-1}X^\top X = X - X = 0.\end{aligned}$$

Es bleibt zu zeigen, daß $\text{Spur}(D) = n - m$:

$$\begin{aligned} \text{Spur}(D) &= \text{Spur} \left(\mathcal{I} - X \left(X^\top X \right)^{-1} X^\top \right) = \text{Spur}(\mathcal{I}) - \text{Spur} \left(X \left(X^\top X \right)^{-1} X^\top \right) \\ &= n - \text{Spur} \left(\underbrace{X^\top X \cdot \left(X^\top X \right)^{-1}}_{\text{eine } (m \times m)\text{-Matrix}} \right) = n - m. \end{aligned}$$

□

2.2.3 Maximum-Likelihood-Schätzer für β und σ^2

Um Maximum-Likelihood-Schätzer für β und σ^2 bzw. Verteilungseigenschaften der MKQ-Schätzer $\hat{\beta}$ und $\hat{\sigma}^2$ herleiten zu können, muß die Verteilung von ε bzw. Y präzisiert werden. Wir werden ab sofort normalverteilte Störgrößen betrachten, die unabhängig und identisch verteilt sind:

$$\varepsilon \sim N(0, \sigma^2 \mathcal{I}), \quad \sigma^2 > 0.$$

Daraus folgt:

$$Y \sim N(X\beta, \sigma^2 \mathcal{I}).$$

Wie sieht die Verteilung der MKQ-Schätzer $\hat{\beta}$ und $\hat{\sigma}^2$ aus? Da $\hat{\beta} = (X^\top X)^{-1} X^\top Y$ linear von Y abhängt, erwartungstreu ist und die $\text{Cov}\hat{\beta} = \hat{\sigma}^2 (X^\top X)^{-1}$ besitzt, gilt:

$$\hat{\beta} \sim N \left(\beta, \sigma^2 \left(X^\top X \right)^{-1} \right)$$

Berechnen wir nun Maximum-Likelihood-Schätzer für β und σ^2 , und zwar $\tilde{\beta}$ und $\tilde{\sigma}^2$. Dann zeigen wir, daß sie im Wesentlichen mit den MKQ-Schätzern übereinstimmen.

$$\begin{aligned} \tilde{\beta} &= \hat{\beta}, \\ \tilde{\sigma}^2 &= \frac{n-m}{n} \hat{\sigma}^2. \end{aligned}$$

Betrachten wir zunächst die Likelihood-Funktion von Y :

$$L(y, \beta, \sigma^2) = f_Y(y) = \frac{1}{(\sqrt{2\pi}\sigma)^n} \cdot \exp \left\{ -\frac{1}{2\sigma^2} (y - X\beta)^\top (y - X\beta) \right\}$$

und die Log-Likelihood-Funktion

$$\log L(y, \beta, \sigma^2) = -\frac{n}{2} \log(2\pi) - \underbrace{\frac{n}{2} \log(\sigma^2) - \frac{1}{2\sigma^2} |y - X\beta|^2}_{:=g}.$$

Die Maximum-Likelihood-Schätzer sind dann

$$\left(\tilde{\beta}, \tilde{\sigma}^2 \right) = \underset{\beta \in \mathbb{R}^m, \sigma^2 > 0}{\text{argmax}} \log L(y, \beta, \sigma^2),$$

sofern sie existieren.

Satz 2.2.5 (*Maximum-Likelihood-Schätzung von $\tilde{\beta}$ und $\tilde{\sigma}^2$*). Es existieren eindeutig bestimmte Maximum-Likelihood-Schätzer für β und σ^2 , die folgendermaßen aussehen:

$$\begin{aligned}\tilde{\beta} &= \hat{\beta} = (X^\top X)^{-1} X^\top Y \\ \tilde{\sigma}^2 &= \frac{n-m}{n} \hat{\sigma}^2 = \frac{1}{n} |Y - X\tilde{\beta}|^2.\end{aligned}$$

Beweis. Wir fixieren $\sigma^2 > 0$ und suchen

$$\tilde{\beta} = \operatorname{argmax}_{\beta \in \mathbb{R}^m} \log L(Y, \beta, \sigma^2) = \operatorname{argmin}_{\beta \in \mathbb{R}^m} |Y - X\beta|^2,$$

woraus folgt, daß $\tilde{\beta}$ mit dem bekannten MKQ-Schätzer $\hat{\beta} = (X^\top X)^{-1} X^\top Y$ identisch ist, der nicht von σ^2 abhängt. Berechnen wir jetzt

$$\tilde{\sigma}^2 = \operatorname{argmax}_{\sigma^2 > 0} \log L(Y, \tilde{\beta}, \sigma^2) = \operatorname{argmax}_{\sigma^2 > 0} g(\sigma^2).$$

Es gilt

$$g(\sigma^2) \xrightarrow{\sigma^2 \rightarrow +\infty} -\infty, \quad g(\sigma^2) \xrightarrow{\sigma^2 \rightarrow 0} -\infty,$$

weil $|Y - X\tilde{\beta}|^2 \neq 0$, dadurch, daß $Y \sim N(X\beta, \sigma^2 I) \in \{Xy : y \in \mathbb{R}^m\}$ mit Wahrscheinlichkeit Null. Da

$$g'(\sigma^2) = -\frac{n}{2} \frac{1}{\sigma^2} + \frac{|Y - X\tilde{\beta}|^2}{2(\sigma^2)^2} = 0, \quad \text{ist } \tilde{\sigma}^2 = \frac{1}{n} |Y - X\tilde{\beta}|^2$$

ein Maximumpunkt von $g(\sigma^2)$, das heißt, $\tilde{\sigma}^2$ ist ein Maximum-Likelihood-Schätzer für σ^2 . \square

Satz 2.2.6. Unter den obigen Voraussetzungen gilt:

1. $\mathbb{E} \tilde{\sigma}^2 = \frac{n-m}{n} \sigma^2$, das heißt, $\tilde{\sigma}^2$ ist nicht erwartungstreu; allerdings ist er asymptotisch unverzerrt.
2. $\frac{n}{\tilde{\sigma}^2} \tilde{\sigma}^2 \sim \chi_{n-m}^2$, $\frac{n-m}{\tilde{\sigma}^2} \hat{\sigma}^2 \sim \chi_{n-m}^2$.

Beweis. 1. Trivial (vergleiche den Beweis von Satz 2.2.4)

2. Wir zeigen den Satz nur für $\hat{\sigma}^2$.

$$\begin{aligned}\frac{n-m}{\sigma^2} \hat{\sigma}^2 &= \frac{1}{\sigma^2} |Y - X\hat{\beta}|^2 \\ &= \frac{1}{\sigma^2} Y^\top \underbrace{D}_{=D^2} Y \quad (\text{nach dem Beweis von Satz 2.2.4}) \\ &= \frac{1}{\sigma^2} (DY)^\top DY = \frac{1}{\sigma^2} (D(X\beta + \varepsilon))^\top \cdot D(X\beta + \varepsilon) \\ &= \frac{1}{\sigma^2} (D\varepsilon)^\top D\varepsilon = \left(\frac{\varepsilon^\top}{\sigma}\right) D \left(\frac{\varepsilon}{\sigma}\right),\end{aligned}$$

wobei

$$\begin{pmatrix} \varepsilon \\ \sigma \end{pmatrix} \sim N(0, \mathcal{I}).$$

Nach Satz 2.1.8 gilt

$$\frac{\varepsilon^\top}{\sigma} D \frac{\varepsilon}{\sigma} \sim \chi_r^2,$$

wobei $r = \text{Rang}(D)$, weil $D\mathcal{I} = D$ idempotent ist. Falls $r = n - m$, dann ist $\frac{n-m}{\sigma^2} \sim \chi_{n-m}^2$. Zeigen wir, daß $\text{Rang}(D) = r = n - m$. Aus der linearen Algebra ist bekannt, daß $\text{Rang}(D) = n - \dim(\text{Kern}(D))$. Wir zeigen, daß $\text{Kern}(D) = \{Xx : x \in \mathbb{R}^m\}$ und damit $\dim(\text{Kern}(D)) = m$, weil $\text{Rang}(X) = m$. Es ist $\{Xx : x \in \mathbb{R}^m\} \subseteq \text{Kern}(D)$, da

$$DX = (\mathcal{I} - X(X^\top X)^{-1}X^\top)X = X - (X^\top X)^{-1}X^\top X = 0.$$

und $\text{Kern}(D) \subseteq \{Xx : x \in \mathbb{R}^m\}$, weil

$$\begin{aligned} \forall y \in \text{Kern}(D) : \quad Dy = 0 &\iff (\mathcal{I} - X(X^\top X)^{-1}X^\top)y = 0 \\ &\iff y = X \cdot \underbrace{(X^\top X)^{-1}X^\top Y}_x = Xx \in \{Xx : x \in \mathbb{R}^m\}. \end{aligned}$$

□

Satz 2.2.7. Sei $Y = X\beta + \varepsilon$ ein multivariates lineares Regressionsmodell mit $Y = (Y_1, \dots, Y_n)^\top$, Designmatrix X mit $\text{Rang}(X) = m$, $\beta = (\beta_1, \dots, \beta_m)^\top$, $\varepsilon \sim N(0, \sigma^2 \mathcal{I})$. Dann sind die Schätzer $\hat{\beta} = (X^\top X)^{-1}X^\top Y$ für β bzw. $\hat{\sigma}^2 = \frac{1}{n-m}|Y - X\hat{\beta}|^2$ für σ^2 unabhängig voneinander.

Beweis. In diesem Beweis verwenden wir den Satz 2.1.9, für dessen Anwendung wir $\hat{\beta}$ als lineare und $\hat{\sigma}^2$ als quadratische Form von ε darstellen. Es ist in den Beweisen der Sätze 2.2.2 und 2.2.6 gezeigt worden, daß

$$\begin{aligned} \hat{\beta} &= \beta + \underbrace{(X^\top X)^{-1}X^\top}_{=A} \varepsilon, \\ \hat{\sigma}^2 &= \frac{1}{n-m} \varepsilon^\top D \varepsilon, \quad \text{wobei } D = \mathcal{I} - X(X^\top X)^{-1}X^\top. \end{aligned}$$

Zusätzlich gilt $AD = 0$, weil nach dem Beweis des Satzes 2.2.4

$$(AD)^\top = D^\top A^\top = \underbrace{D \cdot X}_{=0} ((X^\top X)^{-1})^\top = 0.$$

Da $\varepsilon \sim N(0, \sigma^2 \mathcal{I})$, folgt daraus

$$A\sigma^2 \mathcal{I} D = 0.$$

Deshalb sind die Voraussetzungen des Satzes 2.1.9 erfüllt, und $\hat{\beta}$ und $\hat{\sigma}^2$ sind unabhängig. □

2.2.4 Tests für Regressionsparameter

In diesem Abschnitt wird zunächst die Hypothese

$$H_0 : \beta = \beta_0 \text{ vs. } H_1 : \beta \neq \beta_0$$

für ein $\beta_0 \in \mathbb{R}^m$ getestet. Dafür definieren wir die Testgröße

$$T = \frac{(\hat{\beta} - \beta_0)^\top X^\top X (\hat{\beta} - \beta_0)}{m\hat{\sigma}^2}.$$

Man kann zeigen (vergleiche Satz 2.2.8), daß unter H_0 gilt:

$$T \sim F_{m,n-m}.$$

Daraus folgt, daß H_0 abgelehnt werden soll, falls $T > F_{m,n-m,1-\alpha}$, wobei $F_{m,n-m,1-\alpha}$ das $(1-\alpha)$ -Quantil der $F_{m,n-m}$ -Verteilung darstellt. Dies ist ein Test zum Niveau $\alpha \in (0, 1)$.

Spezialfall: Der Fall $\beta_0 = 0$ beschreibt einen *Test auf Zusammenhang*; das heißt, man testet, ob die Parameter β_1, \dots, β_m für die Beschreibung der Daten Y relevant sind.

Bemerkung 2.2.1. 1. Wie kann man verstehen, daß die Testgröße T tatsächlich H_0 von H_1 unterscheiden soll? Führen wir die Bezeichnung

$$\tilde{Y} = Y - \underbrace{X\hat{\beta}}_{:=\hat{Y}}$$

ein; dabei gilt:

$$\hat{\sigma}^2 = \frac{1}{n-m} |\tilde{Y}|^2$$

und \tilde{Y} ist der Vektor der *Residuen*.

Ohne Beschränkung der Allgemeinheit setzen wir $\beta_0 = 0$. Falls H_0 nicht gelten soll, dann ist $\beta \neq 0$, und somit

$$|X\beta|^2 = (X\beta)^\top X\beta = \beta^\top X^\top X\beta > 0,$$

weil X den vollen Rang hat. Daraus folgt, daß H_0 abgelehnt werden soll, falls

$$|\hat{Y}|^2 = |X\hat{\beta}|^2 = \hat{\beta}^\top X^\top X\hat{\beta} \gg 0.$$

In der Testgröße $|X\hat{\beta}|^2$ sind allerdings die Schwankungen der Schätzung von β nicht berücksichtigt. Deswegen teilt man $|X\hat{\beta}|^2$ durch $\hat{\sigma}^2$:

$$T = \frac{\hat{\beta}^\top X^\top X\hat{\beta}}{m \cdot \hat{\sigma}^2} = \frac{|\hat{Y}|^2}{\frac{m}{n-m} |Y - \hat{Y}|^2}.$$

Der Satz von Pythagoras liefert

$$|Y|^2 = |\tilde{Y}|^2 + |\hat{Y}|^2,$$

wobei unter H_0

$\mathbb{E}|\hat{Y}|^2 = \mathbb{E}|Y|^2 - \mathbb{E}|Y - \hat{Y}|^2 = n\sigma^2 - \mathbb{E}|\tilde{Y}|^2$ gilt, und somit

$$\frac{\mathbb{E}|\hat{Y}|^2}{\mathbb{E}\left(\frac{m}{n-m}|\tilde{Y}|^2\right)} \stackrel{(H_0)}{=} \frac{n\sigma^2 - \mathbb{E}|\tilde{Y}|^2}{\frac{m}{n-m}\mathbb{E}|\tilde{Y}|^2} = \frac{n-m}{m} \left(\frac{n\sigma^2}{\mathbb{E}|\tilde{Y}|^2} - 1 \right),$$

weil $\mathbb{E}|Y|^2 = \mathbb{E}(Y^\top Y) = \sigma^2 \cdot n$, wegen $Y \sim N(0, \sigma^2 \mathcal{I})$.

\implies Die Testgröße T ist sensibel gegenüber Abweichungen von H_0 .

2. Die Größe

$$|\tilde{Y}|^2 = |Y - \hat{Y}|^2$$

wird *Reststreuung* genannt. Mit deren Hilfe kann der Begriff des *Bestimmtheitsmaßes* R^2 aus der Stochastik I wie folgt verallgemeinert werden:

$$R^2 = 1 - \frac{|\tilde{Y}|^2}{|Y - \bar{Y}_n \cdot e|^2},$$

wobei $e = (1, \dots, 1)^\top$, $\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n Y_i$.

Satz 2.2.8. Unter $H_0 : \beta = \beta_0$ gilt

$$T = \frac{(\hat{\beta} - \beta_0)^\top X^\top X (\hat{\beta} - \beta_0)}{m\hat{\sigma}^2} \sim F_{m, n-m}.$$

Beweis. Es gilt

$$\begin{aligned} \hat{\beta} &\sim N\left(\beta_0, \sigma^2 (X^\top X)^{-1}\right) \\ \implies \hat{\beta} - \beta_0 &\sim N\left(0, \underbrace{\sigma^2 (X^\top X)^{-1}}_{:=K}\right). \end{aligned}$$

Falls $A = \frac{X^\top X}{\sigma^2}$, dann ist $AK = \mathcal{I}$ idempotent. Dann gilt nach Satz 2.1.8

$$(\hat{\beta} - \beta_0)^\top A (\hat{\beta} - \beta_0) \stackrel{H_0}{\sim} \chi_m^2$$

(Zur Information: Unter H_1 wäre $(\hat{\beta} - \beta_0)^\top A (\hat{\beta} - \beta_0)$ nicht-zentral χ^2 -verteilt).

Es gilt zusätzlich:

$$\frac{n-m}{\sigma^2} \hat{\sigma}^2 \sim \chi_{n-m}^2.$$

Aus Satz 2.2.7 folgt die Unabhängigkeit von $(\hat{\beta} - \beta_0)^\top A(\hat{\beta} - \beta_0)$ und $\frac{n-m}{\sigma^2} \hat{\sigma}^2$.

$$\implies T = \frac{(\hat{\beta} - \beta_0)^\top (X^\top X)(\hat{\beta} - \beta_0)/m}{(n-m)\hat{\sigma}^2/(n-m)} \sim F_{m,n-m}$$

nach der Definition der F -Verteilung. □

Jetzt wird die Relevanz der einzelnen Parameter β_j getestet:

$$H_0 : \beta_j = \beta_{0j} \text{ vs. } H_1 : \beta_j \neq \beta_{0j}.$$

Satz 2.2.9. Unter $H_0 : \beta_j = \beta_{0j}$ gilt:

$$T_j = \frac{\hat{\beta}_j - \beta_{0j}}{\hat{\sigma} \sqrt{x^{jj}}} \sim t_{n-m}, \text{ wobei}$$

$$(X^\top X)^{-1} = (x^{ij})_{i,j=1,\dots,m}.$$

Beweis. Aus $\hat{\beta} \stackrel{H_0}{\sim} N(\beta_0, \sigma^2(X^\top X)^{-1})$ folgt $\hat{\beta}_j \stackrel{H_0}{\sim} N(\beta_{0j}, \sigma^2 x^{jj})$ und somit $\hat{\beta}_j - \beta_{0j} \sim N(0, \sigma^2 x^{jj})$. Dann ist $\frac{\hat{\beta}_j - \beta_{0j}}{\sigma \sqrt{x^{jj}}} \sim N(0, 1)$. Zusätzlich gilt: $\frac{(n-m)\hat{\sigma}^2}{\sigma^2} \stackrel{H_0}{\sim} \chi_{n-m}^2$, und nach Satz 2.2.7 sind beide Größen unabhängig. Daraus folgt:

$$T_j = \frac{\frac{\hat{\beta}_j - \beta_{0j}}{\sigma \sqrt{x^{jj}}}}{\sqrt{\frac{(n-m)\hat{\sigma}^2}{(n-m)\sigma^2}}} \sim t_{n-m}.$$

□

Somit wird $H_0 : \beta_j = \beta_{0j}$ abgelehnt, falls $|T| > t_{n-m, 1-\alpha/2}$. Dies ist ein Test von H_0 vs. H_1 zum Niveau α .

Sei nun

$$H_0 : \beta_{j_1} = \beta_{0j_1}, \dots, \beta_{j_l} = \beta_{0j_l} \text{ vs. } H_1 : \exists i \in \{1, \dots, l\} : \beta_{j_i} \neq \beta_{0j_i}$$

die zu testende Hypothese.

Übungsaufgabe 2.2.2. Zeigen Sie, daß unter H_0 folgende Verteilungsaussage gilt:

$$T = \frac{(\hat{\beta}' - \beta_0')^\top K'(\hat{\beta}' - \beta_0')}{l\hat{\sigma}^2} \sim F_{l,n-m},$$

wobei

$$\begin{aligned}\hat{\beta}' &= (\hat{\beta}_{j_1}, \dots, \hat{\beta}_{j_i}), \\ \beta'_0 &= (\beta_{0j_1}, \dots, \beta_{0j_i}), \\ K' &= \begin{pmatrix} x^{j_1j_1} & \dots & x^{j_1j_i} \\ \vdots & & \vdots \\ x^{j_ij_1} & \dots & x^{j_ij_i} \end{pmatrix}^{-1}.\end{aligned}$$

Konstruieren Sie den dazugehörigen F -Test!

Test auf Linearkombination von Parametern

Sei nun

$$H_0 : H\beta = c \text{ vs. } H_1 : H\beta \neq c,$$

wobei H eine $(r \times m)$ -Matrix und $c \in \mathbb{R}^r$ sind, $r \leq m$.

Satz 2.2.10. Unter H_0 gilt

$$T = \frac{(H\hat{\beta} - c)^\top (H(X^\top X)^{-1}H^\top)^{-1}(H\hat{\beta} - c)}{r\hat{\sigma}^2} \sim F_{r,n-m}.$$

Deshalb wird $H_0 : H\beta = c$ abgelehnt, falls $T > F_{r,n-m,1-\alpha}$.

Übungsaufgabe 2.2.3. Beweisen Sie Satz 2.2.10!

2.2.5 Konfidenzbereiche

1. Konfidenzintervall für β_j

Im Satz 2.2.9 haben wir gezeigt, daß

$$\frac{\hat{\beta}_j - \beta_j}{\hat{\sigma} \cdot \sqrt{x^{jj}}} \sim t_{n-m},$$

wobei $(X^\top X)^{-1} = (x^{ij})_{i,j=1,\dots,m}$. Daraus kann mit den üblichen Überlegungen folgendes Konfidenzintervall für β_j zum Niveau $1 - \alpha$ abgeleitet werden:

$$\mathbb{P} \left(\hat{\beta}_j - t_{n-m,1-\alpha/2} \cdot \hat{\sigma} \sqrt{x^{jj}} \leq \beta_j \leq \hat{\beta}_j + t_{n-m,1-\alpha/2} \cdot \hat{\sigma} \sqrt{x^{jj}} \right) = 1 - \alpha.$$

2. Simultaner Konfidenzbereich für $\beta = (\beta_1, \dots, \beta_m)^\top$

Falls A_j wie unten definiert ist, dann erhält man mit Hilfe folgender *Bonferroni-Ungleichung*

$$\mathbb{P} \left(\bigcap_{j=1}^m A_j \right) \geq \sum_{j=1}^m \mathbb{P}(A_j) - (m - 1),$$

daß

$$\mathbb{P}\left(\underbrace{\hat{\beta}_j - t_{n-m,1-\alpha/(2m)} \cdot \hat{\sigma} \sqrt{x^{jj}} \leq \beta_j \leq \hat{\beta}_j + t_{n-m,1-\alpha/(2m)} \cdot \hat{\sigma} \sqrt{x^{jj}}}_{:=A_j}, \quad j = 1, \dots, m\right)$$

$$\stackrel{\text{(Bonferroni)}}{\geq} \sum_{j=1}^m \mathbb{P}(A_j) - (m-1) = m \cdot \left(1 - \frac{\alpha}{m}\right) - m + 1 = 1 - \alpha.$$

Daraus folgt, daß

$$\left\{ \beta = (\beta_1, \dots, \beta_m)^\top : \beta_j \in \left[\hat{\beta}_j - t_{n-m,1-\alpha/(2m)} \cdot \hat{\sigma} \sqrt{x^{jj}}, \hat{\beta}_j + t_{n-m,1-\alpha/(2m)} \cdot \hat{\sigma} \sqrt{x^{jj}} \right] \right\}$$

ein simultaner Konfidenzbereich für β zum Niveau $1 - \alpha$ ist.

3. Konfidenzellipsoid für β .

In Satz 2.2.8 haben wir bewiesen, daß

$$T = \frac{(\hat{\beta} - \beta)^\top (X^\top X)(\hat{\beta} - \beta)}{m\hat{\sigma}^2} \sim F_{m,n-m}.$$

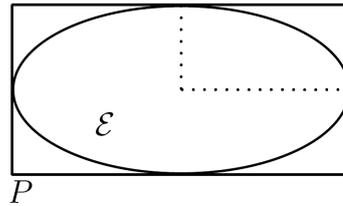
Daraus folgt, daß

$$\mathbb{P}(T \leq F_{m,n-m,1-\alpha}) = 1 - \alpha \quad \text{und}$$

$$\mathcal{E} = \left\{ \beta \in \mathbb{R}^m : \frac{(\hat{\beta} - \beta)^\top (X^\top X)(\hat{\beta} - \beta)}{m\hat{\sigma}^2} \leq F_{m,n-m,1-\alpha} \right\}$$

ein Konfidenzellipsoid zum Niveau $1 - \alpha$ ist, siehe Abbildung 2.2.

Abbildung 2.2: Konfidenzellipsoid



Da ein Ellipsoid in das minimale Parallelepiped P eingebettet werden kann, sodaß die Seitenlängen von P gleich $2 \times$ der Halbachsenlängen von \mathcal{E} sind, ergibt sich folgender simultaner Konfidenzbereich für $\beta = (\beta_1, \dots, \beta_m)^\top$:

$$P = \left\{ \beta : \hat{\beta}_j - \hat{\sigma} \sqrt{m x^{jj} F_{m,n-m,1-\alpha}} \leq \beta_j \leq \hat{\beta}_j + \hat{\sigma} \sqrt{m x^{jj} F_{m,n-m,1-\alpha}} \right\}$$

$j = 1, \dots, m.$

4. *Konfidenzintervall für den erwarteten Zielwert $x_{01}\beta_1 + \dots + x_{0m}\beta_m$.*

Sei $Y_0 = x_{01}\beta_1 + \dots + x_{0m}\beta_m + \varepsilon_0$ eine neue Zielvariable mit $\mathbb{E}\varepsilon_0 = 0$. Dann ist

$$\mathbb{E}Y_0 = \sum_{i=1}^n x_{0i}\beta_i.$$

Wir konstruieren ein Konfidenzintervall für $\mathbb{E}Y_0$. Dazu verwenden wir die Beweis-
idee des Satzes 2.2.9 kombiniert mit Satz 2.2.10 mit $H = (x_{01}, \dots, x_{0m}) = x_0^\top$,
 $r = 1$. Dann ist

$$T = \frac{\sum_{i=1}^m \hat{\beta}_i x_{0i} - \sum_{i=1}^m \beta_i x_{0i}}{\hat{\sigma} \sqrt{x_0^\top (X^\top X)^{-1} x_0}} \sim t_{n-m}.$$

Darum ist

$$\left\{ \beta = (\beta_1, \dots, \beta_m)^\top : \sum_{i=1}^m x_{0i} \hat{\beta}_i - \hat{\sigma} \sqrt{x_0^\top (X^\top X)^{-1} x_0} \cdot t_{n-m, 1-\alpha/2} \leq \sum_{i=1}^m x_{0i} \beta_i \leq \sum_{i=1}^m x_{0i} \hat{\beta}_i + \hat{\sigma} \sqrt{x_0^\top (X^\top X)^{-1} x_0} \cdot t_{n-m, 1-\alpha/2} \right\}$$

ein Konfidenzintervall für $\sum_{i=1}^m x_{0i} \beta_i$ zum Niveau $1 - \alpha$.

5. *Prognoseintervall für die Zielvariable Y_0 .*

Für $Y_0 = \sum_{i=1}^m x_{0i} \beta_i + \varepsilon_0$ mit $\varepsilon_0 \sim N(0, \sigma^2)$, ε_0 unabhängig von $\varepsilon_1, \dots, \varepsilon_n$, gilt:

$$\begin{aligned} x_0^\top \hat{\beta} - Y_0 &\sim N(0, \sigma^2(1 + x_0^\top (X^\top X)^{-1} x_0)) \\ \implies \frac{x_0^\top \hat{\beta} - Y_0}{\sigma \sqrt{1 + x_0^\top (X^\top X)^{-1} x_0}} &\sim N(0, 1) \\ \implies \frac{x_0^\top \hat{\beta} - Y_0}{\hat{\sigma} \sqrt{1 + x_0^\top (X^\top X)^{-1} x_0}} &\sim t_{n-m} \end{aligned}$$

Also ist

$$\left(x_0^\top \hat{\beta} + c, x_0^\top \hat{\beta} - c \right)$$

$$\text{mit } c = \hat{\sigma} \sqrt{1 + x_0^\top (X^\top X)^{-1} \cdot x_0} \cdot t_{n-m, 1-\alpha/2}$$

ein Prognoseintervall für die Zielvariable Y_0 zum Niveau $1 - \alpha$.

6. Konfidenzband für die Regressionsebene $y = \beta_1 + \sum_{i=2}^m x_i \beta_i$ im multiplen Regressionsmodell.

Es sei $Y = X\beta + \varepsilon$, wobei

$$X = \begin{pmatrix} 1 & x_{12} & \cdots & x_{1m} \\ 1 & x_{22} & \cdots & x_{2m} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n2} & \cdots & x_{nm} \end{pmatrix} \quad \text{und } \varepsilon \sim N(0, \sigma^2 \cdot \mathcal{I}).$$

Wir wollen ein zufälliges Konfidenzband $B(x)$ für y angeben. Es gilt

$$\mathbb{P} \left(y = \beta_1 + \sum_{i=2}^m \beta_i x_i \in B(x) \right) = 1 - \alpha \quad \forall x \in \mathbb{R}_1^{m-1}, \quad \text{wobei}$$

$$\mathbb{R}_1^{m-1} = \left\{ (1, x_2, \dots, x_m)^\top \in \mathbb{R}^m \right\}.$$

Satz 2.2.11. Es gilt:

$$\mathbb{P} \left(\max_{x \in \mathbb{R}_1^{m-1}} \frac{\left(x^T \hat{\beta} - \overbrace{\left(\beta_1 + \sum_{i=2}^m \beta_i x_i \right)}^{=y} \right)^2}{\hat{\sigma}^2 x^\top (X^\top X)^{-1} x} \leq m \cdot F_{m, n-m, 1-\alpha} \right) = 1 - \alpha.$$

ohne Beweis.

2.3 Multivariate lineare Regression mit $\text{Rang}(X) < m$

Es sei $Y = X\beta + \varepsilon$, $Y \in \mathbb{R}^n$, wobei X eine $(n \times m)$ -Matrix mit $\text{Rang}(X) = r < m$ ist, $\beta = (\beta_1, \dots, \beta_m)^\top$, $\varepsilon \in \mathbb{R}^n$, $\mathbb{E} \varepsilon = 0$, $\mathbb{E}(\varepsilon_i \varepsilon_j) = \delta_{ij} \sigma^2$, $i, j = 1, \dots, n$, $\sigma^2 > 0$.

Der MKQ-Schätzer $\hat{\beta}$ ist nach wie vor eine Lösung der Normalgleichung

$$(X^\top X) \beta = X^\top Y.$$

$X^\top X$ ist aber nicht mehr invertierbar, weil

$$\text{Rang}(X^\top X) \leq \min \left\{ \text{Rang}(X), \text{Rang}(X^\top) \right\} = r < m.$$

Um $\hat{\beta}$ aus der Normalgleichung zu gewinnen, sollen beide Seiten der Gleichung mit der sogenannten *verallgemeinerten Inversen* von $X^\top X$ multipliziert werden.

2.3.1 Verallgemeinerte Inverse

Definition 2.3.1. Sei A eine $(n \times m)$ -Matrix. Eine $(m \times n)$ -Matrix A^- heißt *verallgemeinerte Inverse* von A , falls

$$AA^-A = A \quad \text{gilt.}$$

Die Matrix A^- ist nicht eindeutig bestimmt, was die folgenden Hilfssätze zeigen.

Lemma 2.3.1. Sei A eine $(n \times m)$ -Matrix, $m \leq n$ mit $\text{Rang}(A) = r \leq m$. Es existieren invertierbare Matrizen P ($n \times n$) und Q ($m \times m$), sodaß

$$PAQ = \begin{pmatrix} \mathcal{I}_r & 0 \\ 0 & 0 \end{pmatrix}, \quad \text{wobei } \mathcal{I}_r = \text{diag}(\underbrace{1, \dots, 1}_{r \text{ Mal}}). \quad (2.3.1)$$

Folgerung 2.3.1. Für eine beliebige $(n \times m)$ -Matrix A mit $n \geq m$, $\text{Rang}(A) = r \leq m$ gilt

$$A^- = Q \begin{pmatrix} \mathcal{I}_r & A_2 \\ A_1 & A_3 \end{pmatrix} P, \quad (2.3.2)$$

wobei P und Q Matrizen aus der Darstellung (2.3.1) sind, $\mathcal{I}_r = \text{diag}(\underbrace{1, \dots, 1}_{r \text{ Mal}})$, und A_1 , A_2 , A_3 beliebige $((m-r) \times r)$, $(r \times (n-r))$ bzw. $((m-r) \times (n-r))$ -Matrizen sind.

Insbesondere kann

$$\begin{aligned} A_1 &= 0, \\ A_2 &= 0, \\ A_3 &= \text{diag}(\underbrace{1, \dots, 1}_{s-r \text{ Mal}}, 0, \dots, 0), \\ s &\in \{r, \dots, m\} \end{aligned}$$

gewählt werden, das heißt, $\text{Rang}(A^-) = s \in \{r, \dots, m\}$ für

$$A^- = Q \begin{pmatrix} \mathcal{I}_s & 0 \\ 0 & 0 \end{pmatrix} P.$$

Beweis. Zeigen wir, daß für A^- wie in (2.3.2) gegeben, $AA^-A = A$ gilt. Aus Lemma 2.3.1 folgt, daß

$$\begin{aligned} A &= P^{-1} \cdot \text{diag}(1, \dots, 1, 0, \dots, 0) \cdot Q^{-1} \quad \text{und somit} \\ AA^-A &= P^{-1} \begin{pmatrix} \mathcal{I}_r & 0 \\ 0 & 0 \end{pmatrix} Q^{-1} Q \cdot \begin{pmatrix} \mathcal{I}_r & A_2 \\ A_1 & A_3 \end{pmatrix} P P^{-1} \begin{pmatrix} \mathcal{I}_r & 0 \\ 0 & 0 \end{pmatrix} Q^{-1} \\ &= P^{-1} \begin{pmatrix} \mathcal{I}_r & 0 \\ 0 & 0 \end{pmatrix} \begin{pmatrix} \mathcal{I}_r & A_2 \\ A_1 & A_3 \end{pmatrix} \begin{pmatrix} \mathcal{I}_r & 0 \\ 0 & 0 \end{pmatrix} Q^{-1} = P^{-1} \begin{pmatrix} \mathcal{I}_r & 0 \\ 0 & 0 \end{pmatrix} Q^{-1} \\ &= A. \end{aligned}$$

□

Lemma 2.3.2. Sei A eine beliebige $(n \times m)$ -Matrix mit $\text{Rang}(A) = r \leq m$, $m \leq n$.

1. Falls $(A^\top A)^-$ eine verallgemeinerte Inverse von $A^\top A$ ist, dann ist $((A^\top A)^-)^\top$ ebenfalls eine verallgemeinerte Inverse von $A^\top A$.
2. Es gilt die Darstellung

$$\begin{aligned} (A^\top A)(A^\top A)^- A^\top &= A^\top \quad \text{bzw.} \\ A(A^\top A)^-(A^\top A) &= A. \end{aligned}$$

Beweis. 1. $A^\top A$ ist symmetrisch, also

$$\underbrace{\left(A^\top A (A^\top A)^- A^\top A \right)^\top}_{=A^\top A ((A^\top A)^-)^\top A^\top A} = \left(A^\top A \right)^\top = A^\top A.$$

Also ist $((A^\top A)^-)^\top$ eine verallgemeinerte Inverse von $A^\top A$.

2. Es sei $B = (A^\top A)(A^\top A)^- A^\top - A^\top$. Wir zeigen, daß $B = 0$, indem wir zeigen, daß $BB^\top = 0$.

$$\begin{aligned} BB^\top &= \left((A^\top A)(A^\top A)^- A^\top - A^\top \right) \left(A \left((A^\top A)^- \right)^\top A^\top A - A \right) \\ &= A^\top A (A^\top A)^- A^\top A \left((A^\top A)^- \right)^\top A^\top A - \underbrace{A^\top A (A^\top A)^- A^\top A}_{=A^\top A} \\ &\quad - \underbrace{A^\top A \left((A^\top A)^- \right)^\top \cdot A^\top A}_{=A^\top A} + A^\top A = A^\top A - 2A^\top A + A^\top A = 0. \end{aligned}$$

Die Aussage $A(A^\top A)^- A^\top A = A$ erhält man, indem man die Matrizen an beiden Seiten der Gleichung $A^\top A (A^\top A)^- A^\top = A^\top$ transponiert. □

2.3.2 MKQ-Schätzer für β

Satz 2.3.1. Es sei X eine $(n \times m)$ -Designmatrix mit $\text{Rang}(X) = r \leq m$ in der linearen Regression $Y = X\beta + \varepsilon$. Die allgemeine Lösung der Normalgleichung

$$(X^\top X) \beta = X^\top Y$$

sieht folgendermaßen aus:

$$\beta = (X^\top X)^- X^\top Y + \left(\mathcal{I}_m - (X^\top X)^- X^\top X \right) z, \quad z \in \mathbb{R}^m. \quad (2.3.3)$$

Beweis. 1. Zeigen wir, daß β wie in (2.3.3) angegeben, eine Lösung der Normalgleichung darstellt.

$$\begin{aligned} X^\top X \beta &= \underbrace{(X^\top X)(X^\top X)^{-1} X^\top Y}_{=X^\top \text{ (Lemma 2.3.2, 2.)}} + \left(X^\top X - \underbrace{X^\top X (X^\top X)^{-1} X^\top X}_{=X^\top X} \right) z \\ &= X^\top Y \end{aligned}$$

2. Zeigen wir, daß eine beliebige Lösung β' der Normalgleichung die Form (2.3.3) besitzt. Sei β die Lösung (2.3.3). Wir bilden die Differenz der Gleichungen

$$\begin{array}{rcl} (X^\top X) \beta' & = & X^\top Y \\ - (X^\top X) \beta & = & X^\top Y \\ \hline (X^\top X) (\beta' - \beta) & = & 0 \end{array}$$

$$\begin{aligned} \beta' &= (\beta' - \beta) + \beta \\ &= \beta' - \beta + (X^\top X)^{-1} X^\top Y + \left(\mathcal{I}_m - (X^\top X)^{-1} X^\top X \right) z \\ &= (X^\top X)^{-1} X^\top Y + \left(\mathcal{I}_m - (X^\top X)^{-1} X^\top X \right) z + (\beta' - \beta) - \underbrace{(X^\top X)^{-1} X^\top X (\beta' - \beta)}_{=0} \\ &= (X^\top X)^{-1} X^\top Y + \left(\mathcal{I}_m - (X^\top X)^{-1} X^\top X \right) \underbrace{\left(z + \beta' - \beta \right)}_{=z_0} \\ &\implies \beta' \text{ besitzt die Darstellung (2.3.3).} \end{aligned}$$

□

Bemerkung 2.3.1. Der Satz 2.3.1 liefert die Menge aller Extremalpunkte der MKQ-Minimierungsaufgabe

$$e(\beta) = \frac{1}{n} |Y - X\beta|^2 \longrightarrow \min_{\beta}.$$

Deshalb soll die Menge aller MKQ-Schätzer von β in (2.3.3) zusätzliche Anforderungen erfüllen.

Satz 2.3.2. 1. Alle MKQ-Schätzer von β haben die Form

$$\bar{\beta} = \left(X^\top X \right)^{-} X^\top Y, \quad \text{wobei}$$

$(X^\top X)^{-}$ eine beliebige verallgemeinerte Inverse von $X^\top X$ ist.

2. $\bar{\beta}$ ist nicht erwartungstreu, denn

$$\mathbb{E} \bar{\beta} = \left(X^\top X \right)^{-} X^\top X \beta.$$

3. Es gilt:

$$\text{Cov}\bar{\beta} = \sigma^2 \left(X^\top X \right)^{-} \left(X^\top X \right) \left(\left(X^\top X \right)^{-} \right)^\top.$$

Beweis. 1. Zeigen wir, daß $e(\beta) \geq e(\bar{\beta}) \quad \forall \beta \in \mathbb{R}^m$.

$$\begin{aligned} n \cdot e(\beta) &= |Y - X\beta|^2 = (Y - X\bar{\beta} + X(\bar{\beta} - \beta))^\top (Y - X\bar{\beta} + X(\bar{\beta} - \beta)) \\ &= (Y - X\bar{\beta})^\top (Y - X\bar{\beta}) + (X(\bar{\beta} - \beta))^\top (X(\bar{\beta} - \beta)) \\ &\quad + 2(\bar{\beta} - \beta)^\top X^\top (Y - X\bar{\beta}) \\ &= n \cdot e(\bar{\beta}) + \underbrace{2 \cdot (\bar{\beta} - \beta)^\top (X^\top Y - (X^\top X\bar{\beta}))}_{=0} + |X(\bar{\beta} - \beta)|^2 \\ &\geq n \cdot e(\bar{\beta}) + 0 = n \cdot e(\bar{\beta}), \quad \text{denn} \end{aligned}$$

$\bar{\beta}$ hat die Form (2.3.3) mit $z = 0$ und ist somit eine Lösung der Normalengleichung.

2. Es gilt:

$$\begin{aligned} \mathbb{E}\bar{\beta} &= \mathbb{E} \left((X^\top X)^{-} X^\top Y \right) = \left(X^\top X \right)^{-} X^\top \mathbb{E} Y \\ &= \left(X^\top X \right)^{-} X^\top X \beta, \quad \text{weil aus} \\ Y &= X\beta + \varepsilon, \quad \mathbb{E}\varepsilon = 0 \quad \text{die Relation } \mathbb{E} Y = X\beta \text{ folgt.} \end{aligned}$$

Warum ist $\bar{\beta}$ nicht erwartungstreu? Also warum ist $(X^\top X)^{-} X^\top X \beta \neq \beta$, $\beta \in \mathbb{R}^m$? Da $\text{Rang}(X) = r < m$, ist $\text{Rang}(X^\top X) < m$ und damit $\text{Rang}((X^\top X)^{-} X^\top X) < m$. Darum existiert ein $\beta \neq 0$, für das gilt:

$$\left(X^\top X \right)^{-} X^\top X \beta = 0 \neq \beta,$$

also ist $\bar{\beta}$ nicht erwartungstreu. Es gilt sogar, daß alle Lösungen von (2.3.3) keine erwartungstreuen Schätzer sind. Wenn wir den Erwartungswert an (2.3.3) anwenden, so erhielten wir im Falle der Erwartungstreue:

$$\begin{aligned} \forall \beta \in \mathbb{R}^m : \quad \beta &= (X^\top X)^{-} X^\top X \beta + \left(\mathcal{I}_m - (X^\top X)^{-} (X^\top X) \right) z, \quad z \in \mathbb{R}^m. \\ \implies \left(\mathcal{I}_m - (X^\top X)^{-} (X^\top X) \right) (z - \beta) &= 0 \quad \forall z, \beta \in \mathbb{R}^m \\ \implies (X^\top X)^{-} (X^\top X) (\beta - z) &= \beta - z, \quad \forall z, \beta \in \mathbb{R}^m. \end{aligned}$$

Da diese Gleichung nicht für alle $\beta \in \mathbb{R}^m$ gelten kann (siehe oben), führt die Annahme der Erwartungstreue zum Widerspruch.

3. Es gilt:

$$\begin{aligned}
\text{Cov}(\bar{\beta}_i, \bar{\beta}_j) &= \text{Cov}\left(\underbrace{\left((X^\top X)^{-1} X^\top Y\right)}_{:=A=(a_{kl})}, \left((X^\top X)^{-1} X^\top Y\right)\right) \\
&= \text{Cov}\left(\sum_{k=1}^n a_{ik} Y_k, \sum_{l=1}^n a_{jl} Y_l\right) \\
&= \sum_{k,l=1}^n a_{ik} a_{jl} \underbrace{\text{Cov}(Y_k, Y_l)}_{=\sigma^2 \cdot \delta_{kl}} = \sigma^2 \sum_{k=1}^n a_{ik} a_{jk} = \left(\sigma^2 A A^\top\right)_{i,j} \\
&= \left(\sigma^2 (X^\top X)^{-1} X^\top X \left((X^\top X)^{-1}\right)^\top\right)_{i,j}.
\end{aligned}$$

□

2.3.3 Erwartungstreu schätzbare Funktionen

Definition 2.3.2. Eine Linearkombination $a^\top \beta$ von β_1, \dots, β_m , $a \in \mathbb{R}^m$ heißt (erwartungstreu) schätzbar, falls

$$\exists c \in \mathbb{R}^n : \mathbb{E}(c^\top Y) = a^\top \beta,$$

das heißt, falls es einen linearen, erwartungstreuen Schätzer $c^\top Y$ für $a^\top \beta$ gibt.

Satz 2.3.3. Die Funktion $a^\top \beta$, $a \in \mathbb{R}^m$ ist genau dann erwartungstreu schätzbar, wenn eine der folgenden Bedingungen erfüllt ist:

1. $\exists c \in \mathbb{R}^n : a^\top = c^\top X$.
2. a erfüllt die Gleichung

$$a^\top (X^\top X)^{-1} X^\top X = a^\top. \quad (2.3.4)$$

Beweis. 1. „ \implies “: Falls $a^\top \beta$ schätzbar, dann existiert ein $d \in \mathbb{R}^n$ mit $\mathbb{E}(d^\top Y) = a^\top \beta \quad \forall \beta \in \mathbb{R}^m$. Also

$$\begin{aligned}
a^\top \beta &= d^\top \mathbb{E} Y = d^\top X \beta \Rightarrow (a^\top - d^\top X) \beta = 0, \quad \forall \beta \in \mathbb{R}^m \\
&\implies a^\top = d^\top X,
\end{aligned}$$

setze $c = d$, damit ist die erste Richtung bewiesen.

„ \impliedby “: $\mathbb{E}(c^\top Y) = c^\top \mathbb{E} Y = c^\top X \beta = a^\top \beta$, also ist $a^\top \beta$ erwartungstreu schätzbar.

2. „ \implies “: Falls $a^\top \beta$ erwartungstreu schätzbar ist, dann gilt:

$$a^\top (X^\top X)^- X^\top X \stackrel{\text{Punkt 1}}{=} c^\top \underbrace{X \cdot (X^\top X)^- X^\top X}_{=X \text{ (Lemma 2.3.2)}} = c^\top X \stackrel{\text{(Punkt 1)}}{=} a^\top.$$

Also ist (2.3.4) erfüllt.

„ \impliedby “: Falls $a^\top (X^\top X)^- X^\top X = a^\top$, dann gilt mit $c = (a^\top (X^\top X)^- X^\top X)^\top$ nach Punkt 1, daß $a^\top \beta$ schätzbar ist. □

Bemerkung 2.3.2. Im Falle der Regression mit $\text{Rang}(X) = m$ ist die Gleichung (2.3.4) immer erfüllt, denn $(X^\top X)^- = (X^\top X)^{-1}$ und damit ist $a^\top \beta$ schätzbar für alle $a \in \mathbb{R}^m$.

Satz 2.3.4 (*Beispiele schätzbarer Funktionen*). Falls $\text{Rang}(X) = r < m$, dann sind folgende Linearkombinationen von β schätzbar:

1. Die Koordinaten $\sum_{j=1}^m x_{ij} \beta_j$, $i = 1, \dots, n$ des Erwartungswertvektors $\mathbb{E}Y = X\beta$.
2. Beliebige Linearkombinationen schätzbarer Funktionen.

Beweis. 1. Führe die Bezeichnung $\tilde{x}_i = (x_{i1}, \dots, x_{im})$, $i = 1, \dots, n$ ein. Dann ist

$$\sum_{j=1}^m x_{ij} \beta_j = \tilde{x}_i^\top \beta \quad \forall i = 1, \dots, n,$$

$$X\beta = (\tilde{x}_1, \tilde{x}_2, \dots, \tilde{x}_n)^\top \beta.$$

$\tilde{x}_i \beta$ ist schätzbar, falls \tilde{x}_i die Gleichung (2.3.4) erfüllt, die für alle $i = 1, \dots, n$ folgendermaßen in Matrixform dargestellt werden kann:

$$X \left(X^\top X \right)^- X^\top X = X,$$

was nach Lemma 2.3.2 Gültigkeit besitzt.

2. Für $a_1, \dots, a_k \in \mathbb{R}^m$ seien $a_1^\top \beta, \dots, a_k^\top \beta$ schätzbare Funktionen. Für alle $\lambda = (\lambda_1, \dots, \lambda_k)^\top \in \mathbb{R}^k$ zeigen wir, daß $\sum_{i=1}^k \lambda_i \cdot a_i^\top \beta = \lambda^\top A \beta$ schätzbar ist, wobei $A = (a_1, \dots, a_k)^\top$. Zu zeigen bleibt: $b = (\lambda^\top A)^\top$ erfüllt (2.3.4), also

$$\lambda^\top A \left(X^\top X \right)^- X^\top X = \lambda^\top A.$$

Diese Gleichung stimmt, weil $a_i^\top (X^\top X)^- X^\top X = a_i^\top$, $i = 1, \dots, k$. Nach Satz 2.3.3, 2.) ist $\lambda^\top A \beta$ schätzbar. □

Satz 2.3.5 (*Gauß-Markov*). Es sei $a^\top \beta$ eine schätzbare Funktion, $a \in \mathbb{R}^m$ im linearen Regressionsmodell $Y = X\beta + \varepsilon$ mit $\text{Rang}(X) \leq m$.

1. Der beste lineare erwartungstreue Schätzer (engl. BLUE - best linear unbiased estimator) von $a^\top \beta$ ist durch $a^\top \bar{\beta}$ gegeben, wobei

$$\bar{\beta} = (X^\top X)^{-1} X^\top Y$$

ein MKQ-Schätzer für β ist.

2. $\text{Var}(a^\top \bar{\beta}) = \sigma^2 a^\top (X^\top X)^{-1} a$.

Beweis. Die Linearität von $a^\top \bar{\beta} = a^\top (X^\top X)^{-1} X^\top Y$ als Funktion von Y ist klar. Zeigen wir die Erwartungstreue:

$$\begin{aligned} \mathbb{E}(a^\top \bar{\beta}) &= a^\top \mathbb{E} \bar{\beta} = a^\top (X^\top X)^{-1} X^\top X \beta \\ &= c^\top \underbrace{X(X^\top X)^{-1} X^\top X}_{=X \text{ (Lemma 2.3.2)}} \beta = \underbrace{c^\top X}_{=a^\top} \beta = a^\top \beta \quad \forall \beta \in \mathbb{R}^m. \end{aligned}$$

Berechnen wir $\text{Var}(a^\top \bar{\beta})$ (also beweisen wir Punkt 2), und zeigen, daß sie minimal ist.

$$\begin{aligned} \text{Var}(a^\top \bar{\beta}) &= \text{Var} \left(\sum_{i=1}^m a_i \bar{\beta}_i \right) = \sum_{i,j=1}^m a_i a_j \cdot \text{Cov}(\bar{\beta}_i, \bar{\beta}_j) \\ &= a^\top \text{Cov}(\bar{\beta}) a \stackrel{\text{(Satz 2.3.2)}}{=} a^\top \sigma^2 \left((X^\top X)^{-1} X^\top X (X^\top X)^{-1} \right)^\top a \\ &= \sigma^2 \cdot a^\top \underbrace{\left((X^\top X)^{-1} \right)^\top}_{=(X^\top X)^{-1}} X^\top X \underbrace{\left((X^\top X)^{-1} \right)^\top}_{=(X^\top X)^{-1}} a \\ &\stackrel{\text{Lemma 2.3.2, 1.)}}{=} \sigma^2 a^\top (X^\top X)^{-1} X^\top X (X^\top X)^{-1} a \\ &\stackrel{\text{Satz 2.3.3, 1.)}}{=} \sigma^2 \cdot c^\top X \cdot \underbrace{(X^\top X)^{-1} X^\top X}_{=X} (X^\top X)^{-1} X^\top c \\ &= \sigma^2 \underbrace{c^\top X}_{=a^\top} (X^\top X)^{-1} \underbrace{X^\top c}_{=a} = \sigma^2 a^\top (X^\top X)^{-1} a. \end{aligned}$$

Jetzt zeigen wir, daß für einen beliebigen linearen, erwartungstreuen Schätzer $b^\top Y$ von $a^\top \beta$ gilt: $\text{Var}(b^\top Y) \geq \text{Var}(a^\top \bar{\beta})$. Weil $b^\top Y$ erwartungstreu ist, gilt: $\mathbb{E}(b^\top Y) = a^\top \beta$. Nach Satz 2.3.3 gilt: $a^\top = b^\top X$. Betrachten wir die Varianz von

$$\begin{aligned} 0 &\leq \text{Var}(b^\top Y - a^\top \bar{\beta}) = \text{Var}(b^\top Y) - 2\text{Cov}(b^\top Y, a^\top \bar{\beta}) + \text{Var}(a^\top \bar{\beta}) \\ &= \text{Var}(b^\top Y) - 2\sigma^2 a^\top (X^\top X)^{-1} a + \sigma^2 a^\top (X^\top X)^{-1} a = \text{Var}(b^\top Y) - \text{Var}(a^\top \bar{\beta}) \end{aligned}$$

mit

$$\begin{aligned}\text{Cov}\left(b^\top Y, a^\top \bar{\beta}\right) &= \text{Cov}\left(b^\top Y, a^\top (X^\top X)^- X^\top Y\right) = \sigma^2 a^\top (X^\top X)^- \underbrace{X^\top b}_{=a} \\ &= \sigma^2 a^\top (X^\top X)^- a.\end{aligned}$$

Damit ist $\text{Var}(b^\top Y) \geq \text{Var}(a^\top \bar{\beta})$ und $a^\top \bar{\beta}$ ist ein bester, linearer, erwartungstreuer Schätzer für $a^\top \beta$. \square

Bemerkung 2.3.3. 1. Falls $\text{Rang}(X) = m$, dann ist $a^\top \hat{\beta}$ der beste lineare, erwartungstreue Schätzer für $a^\top \beta$, $a \in \mathbb{R}^m$.

2. Wie im folgenden Satz gezeigt wird, hängt der Schätzer $a^\top \bar{\beta} = a^\top (X^\top X)^- X^\top Y$ nicht von der Wahl der verallgemeinerten Inversen ab.

Satz 2.3.6. Der beste lineare, erwartungstreue Schätzer $a^\top \bar{\beta}$ für $a^\top \beta$ ist eindeutig bestimmt.

Beweis.

$$a^\top \bar{\beta} = a^\top (X^\top X)^- X^\top Y \stackrel{\text{Satz 2.3.3, 1.})}{=} c^\top X (X^\top X)^- X^\top Y.$$

Wir zeigen, daß $X(X^\top X)^- X^\top$ nicht von der Wahl von $(X^\top X)^-$ abhängt. Zeigen wir, daß für beliebige verallgemeinerte Inverse A_1 und A_2 von $(X^\top X)$ gilt: $XA_1 X^\top = XA_2 X^\top$. Nach Lemma 2.3.2, 2.) gilt:

$$XA_1 X^\top X = X = XA_2 X^\top X.$$

Multiplizieren wir alle Teile der Gleichung mit $A_1 X^\top$ von rechts:

$$XA_1 \underbrace{X^\top X A_1 X^\top}_{=X^\top} = XA_1 X^\top = XA_2 \underbrace{X^\top X A_1 X^\top}_{=X^\top}$$

Also ist $XA_1 X^\top = XA_2 X^\top$. \square

2.3.4 Normalverteilte Störgrößen

Sei $Y = X\beta + \varepsilon$ ein lineares Regressionsmodell mit $\text{Rang}(X) = r < m$ und $\varepsilon \sim N(0, \sigma^2 \mathcal{I})$. Genauso wie in Abschnitt 2.2.3 können Maximum-Likelihood-Schätzer $\tilde{\beta}$ und $\tilde{\sigma}^2$ für β und σ^2 hergeleitet werden. Und genauso wie im Satz 2.2.5 kann gezeigt werden, daß

$$\begin{aligned}\tilde{\beta} &= \bar{\beta} = (X^\top X)^- X^\top Y \quad \text{und} \\ \tilde{\sigma}^2 &= \frac{1}{n} |Y - X\bar{\beta}|^2.\end{aligned}$$

Jetzt werden die Verteilungseigenschaften von $\bar{\beta}$ und $\tilde{\sigma}^2$ untersucht. Wir beginnen mit der Erwartungstreue von $\tilde{\sigma}^2$. Wir zeigen, daß $\tilde{\sigma}^2$ nicht erwartungstreu ist, dafür ist aber der korrigierte Schätzer

$$\bar{\sigma}^2 = \frac{1}{n-r} |Y - X\beta|^2 = \frac{n}{n-r} \tilde{\sigma}^2$$

erwartungstreu.

Satz 2.3.7. Der Schätzer $\bar{\sigma}^2$ ist erwartungstreu für σ^2 .

Der Beweis des Satzes 2.3.7 folgt dem Beweis des Satzes 2.2.4, in dem $\hat{\beta} = (X^\top X)^{-1} X^\top Y$ und $\hat{\sigma}^2 = \frac{1}{n-m} |Y - X\hat{\beta}|^2$ im Fall $\text{Rang}(X) = m$ betrachtet wurden. Somit ist die Aussage des Satzes 2.2.4 ein Spezialfall des Satzes 2.3.7. Führen wir die Matrix $D = \mathcal{I} - X(X^\top X)^{-1} X^\top$ ein.

Lemma 2.3.3. Für D gelten folgende Eigenschaften:

1. $D^\top = D$ (Symmetrie),
2. $D^2 = D$ (Idempotenz),
3. $DX = 0$,
4. $\text{Spur}(D) = n - r$.

Beweis. 1. Es gilt:

$$\begin{aligned} D^\top &= \left(\mathcal{I} - X(X^\top X)^{-1} X^\top \right)^\top = \mathcal{I} - X \left((X^\top X)^{-1} \right)^\top X^\top \\ &= \mathcal{I} - X(X^\top X)^{-1} X^\top = D, \end{aligned}$$

weil $\left((X^\top X)^{-1} \right)^\top$ auch eine verallgemeinerte Inverse von $X^\top X$ ist (vergleiche Lemma 2.3.2, 1.)).

2. Es gilt:

$$\begin{aligned} D^2 &= \left(\mathcal{I} - X(X^\top X)^{-1} X^\top \right)^2 = \mathcal{I} - 2X(X^\top X)^{-1} X^\top + \underbrace{X(X^\top X)^{-1} X^\top X(X^\top X)^{-1} X^\top}_{=X(\text{Lemma 2.3.2, 2.))}} \\ &= \mathcal{I} - X(X^\top X)^{-1} X^\top = D. \end{aligned}$$

$$3. \quad DX = X - \underbrace{X(X^\top X)^{-1} X^\top X}_{=X(\text{Lemma 2.3.2, 2.))}} = X - X = 0.$$

4. Es gilt:

$$\text{Spur}(D) = \text{Spur}(I) - \text{Spur} \left(X(X^\top X)^{-1} X^\top \right) = n - \text{Spur} \left(X(X^\top X)^{-1} X^\top \right).$$

Verwenden wir die Eigenschaft der symmetrischen idempotenten Matrizen A aus der linearen Algebra, daß $\text{Spur}(A) = \text{Rang}(A)$. Da $X(X^\top X)^{-1} X^\top$ symmetrisch und idempotent ist, genügt es zu zeigen, daß $\text{Rang}(X(X^\top X)^{-1} X^\top) = r$. Nach Lemma 2.3.2 2.) gilt:

$$\begin{aligned} \text{Rang}(X) &= r = \text{Rang}(X(X^\top X)^{-1} X^\top X) \\ &\leq \min \left\{ \text{Rang}(X(X^\top X)^{-1} X^\top), \underbrace{\text{Rang}(X)}_{=r} \right\} \\ &\leq \text{Rang} \left(X(X^\top X)^{-1} X^\top \right) \leq \text{Rang}(X) = r \\ &\implies \text{Rang} \left(X(X^\top X)^{-1} X^\top \right) = r \\ &\implies \text{Spur} \left(X(X^\top X)^{-1} X^\top \right) = r. \end{aligned}$$

□

Beweis des Satzes 2.3.7. Mit Hilfe des Lemmas 2.3.3 bekommt man

$$\begin{aligned}\bar{\sigma}^2 &= \frac{1}{n-r} |Y - X\bar{\beta}|^2 = \frac{1}{n-r} |Y - X(X^\top X)^{-1}X^\top Y|^2 = \frac{1}{n-r} |DY|^2 \\ &= \frac{1}{n-r} \left| \underbrace{DX}_{=0}\beta + D\varepsilon \right|^2 = \frac{1}{n-r} |D\varepsilon|^2 = \frac{1}{n-r} \varepsilon^\top \underbrace{D^\top D}_{=D^2=D} \varepsilon = \frac{1}{n-r} \varepsilon^\top D\varepsilon.\end{aligned}$$

Deshalb gilt:

$$\begin{aligned}\mathbb{E}\bar{\sigma}^2 &= \frac{1}{n-r} \mathbb{E} \left(\varepsilon^\top D\varepsilon \right) = \frac{1}{n-r} \mathbb{E} \text{Spur} \left(\varepsilon^\top D\varepsilon \right) = \frac{1}{n-r} \text{Spur} \left(D \cdot \mathbb{E} \left(\underbrace{\varepsilon\varepsilon^\top}_{\sigma^2 \mathcal{I}} \right) \right) \\ &= \frac{\sigma^2}{n-r} \cdot \text{Spur}(D) = \sigma^2 \text{ nach Lemma 2.3.3, 4.}, \text{ weil } \mathbb{E}\varepsilon\varepsilon^\top = \sigma^2 \mathcal{I} \\ &\text{wegen } \varepsilon \sim N(0, \sigma^2 \mathcal{I}).\end{aligned}$$

□

Satz 2.3.8. Es gelten folgende Verteilungseigenschaften:

1. $\bar{\beta} \sim N \left((X^\top X)^{-1}X^\top X\beta, \sigma^2(X^\top X)^{-1}(X^\top X)^{-1} \right)^\top$,
2. $\frac{(n-r)\bar{\sigma}^2}{\sigma^2} \sim \chi_{n-r}^2$,
3. $\bar{\beta}$ und $\bar{\sigma}^2$ sind unabhängig.

Beweis. 1. Es gilt:

$$\bar{\beta} = (X^\top X)^{-1}X^\top Y = (X^\top X)^{-1}X^\top (X\beta + \varepsilon) = \underbrace{(X^\top X)^{-1}X^\top X\beta}_{=\mu} + \underbrace{(X^\top X)^{-1}X^\top \varepsilon}_{=A}$$

und mit der Definition von $N(\cdot, \cdot)$ bekommt man

$$\begin{aligned}\bar{\beta} &\sim N \left(\mu, \sigma^2 AA^\top \right) = N \left((X^\top X)^{-1}X^\top X\beta, \sigma^2(X^\top X)^{-1}X^\top X((X^\top X)^{-1})^\top \right) \\ &\text{mit } AA^\top = (X^\top X)^{-1}X^\top X((X^\top X)^{-1})^\top\end{aligned}$$

2. Es gilt $\bar{\sigma}^2 = \frac{1}{n-r} \varepsilon^\top D\varepsilon$ aus dem Beweis des Satzes 2.3.7. Deshalb

$$\frac{(n-r)\bar{\sigma}^2}{\sigma^2} = \underbrace{\left(\frac{\varepsilon}{\sigma} \right)^\top}_{\sim N(0, \mathcal{I})} D \left(\frac{\varepsilon}{\sigma} \right) \stackrel{\text{(Satz 2.1.8)}}{\sim} \chi_{n-r}^2.$$

3. Betrachten wir $A\varepsilon$ und $\varepsilon^\top D\varepsilon$. Es genügt zu zeigen, daß sie unabhängig sind, um die Unabhängigkeit von $\bar{\beta}$ und $\bar{\sigma}^2$ zu beweisen, weil $\bar{\beta} = \mu + A\varepsilon$, $\bar{\sigma}^2 = \frac{1}{n-r} \varepsilon^\top D\varepsilon$. Es gilt: $A \cdot \sigma^2 \mathcal{I} \cdot D = 0$. Nach Satz 2.1.9 sind dann $A\varepsilon$ und $\varepsilon^\top D\varepsilon$ unabhängig.

□

2.3.5 Hypothesentests

Betrachten wir die Hypothesen $H_0 : H\beta = d$ vs. $H_1 : H\beta \neq d$, wobei H eine $(s \times m)$ -Matrix ($s \leq m$) mit $\text{Rang}(H) = s$ ist, und $d \in \mathbb{R}^s$.

Im Satz 2.2.10 haben wir im Fall $\text{Rang}(X) = r = m$ folgende Testgröße dafür betrachtet:

$$T = \frac{(H\hat{\beta} - d)^\top (H(X^\top X)^{-1}H^\top)^{-1}(H\hat{\beta} - d)}{s\hat{\sigma}^2} \stackrel{(H_0)}{\sim} F_{s,n-m}.$$

Im allgemeinen Fall betrachten wir

$$T = \frac{(H\bar{\beta} - d)^\top (H(X^\top X)^-H^\top)^{-1}(H\bar{\beta} - d)}{s\bar{\sigma}^2}. \quad (2.3.5)$$

Wir wollen zeigen, daß $T \stackrel{(H_0)}{\sim} F_{s,n-r}$. Dann wird H_0 verworfen, falls $T > F_{s,n-r,1-\alpha}$. Dies ist ein Test zum Niveau $\alpha \in (0, 1)$.

Definition 2.3.3. Die Hypothese $H_0 : H\beta = d$ heißt *testbar*, falls alle Koordinaten des Vektors $H\beta$ schätzbare Funktionen sind.

Satz 2.3.3 gibt Bedingungen an H an, unter denen $H_0 : H\beta = d$ testbar ist. Diese werden im folgendem Lemma formuliert:

Lemma 2.3.4. Die Hypothese $H_0 : H\beta = d$ ist testbar genau dann, wenn

1. $\exists (s \times n)$ -Matrix $C : H = CX$, oder
2. $H(X^\top X)^-X^\top X = H$.

Wir zeigen, daß die Testgröße T in (2.3.5) wohldefiniert ist, das heißt, die $(s \times s)$ -Matrix $H(X^\top X)^-H^\top$ positiv definit und damit invertierbar ist. Aus Folgerung 2.3.1 haben wir $X^\top X = P^{-1} \begin{pmatrix} \mathcal{I}_r & 0 \\ 0 & 0 \end{pmatrix} P^{-1}$ für eine $(m \times m)$ -Matrix P , die invertierbar und symmetrisch ist. Deshalb gilt

$$(X^\top X)^- = P \cdot \begin{pmatrix} \mathcal{I}_r & 0 \\ 0 & \mathcal{I}_{m-r} \end{pmatrix} P = P \cdot P,$$

das heißt, daß es eine eindeutige verallgemeinerte Inverse von $X^\top X$ mit dieser Darstellung gibt. Daraus folgt, daß die $(s \times s)$ -Matrix $HPPH^\top = (PH^\top)^\top \cdot PH^\top$ positiv definit ist, weil $\text{Rang}(PH^\top) = s$. Sei nun $(X^\top X)^-$ eine beliebige verallgemeinerte Inverse von $X^\top X$. Dann ist mit Lemma 2.3.4

$$H(X^\top X)^-H^\top = CX(X^\top X)^-X^\top C^\top = CXPPX^\top C^\top = HPPH^\top,$$

denn $X(X^\top X)^-X^\top$ ist invariant bezüglich der Wahl von $(X^\top X)^-$, laut Beweis des Satzes 2.3.6. Also ist $H(X^\top X)^-H^\top$ positiv definit für eine beliebige verallgemeinerte Inverse $(X^\top X)^-$ und die Testgröße T somit wohldefiniert.

Satz 2.3.9. Falls $H_0 : H\beta = d$ testbar ist, dann gilt $T \stackrel{(H_0)}{\sim} F_{s,n-r}$.

Beweis. Ähnlich, wie in Satz 2.2.10 gilt

$$H\bar{\beta} - d = H(X^\top X)^{-1} X^\top (X\beta + \varepsilon) - d = \underbrace{H(X^\top X)^{-1} X^\top X\beta - d}_{=\mu} + \underbrace{H(X^\top X)^{-1} X^\top \varepsilon}_{=B}$$

Zeigen wir, daß $\mu \stackrel{(H_0)}{=} 0$.

$$\mu \stackrel{(\text{Lemma 2.3.4})}{=} C \cdot \underbrace{X(X^\top X)^{-1} X^\top X}_{=X \text{ (Lemma 2.3.2, 2.)}} \cdot \beta - d = CX\beta - d = H\beta - d \stackrel{(H_0)}{=} 0.$$

Nach Satz 2.3.8 sind $(H\bar{\beta} - d)^\top (H(X^\top X)^{-1} H^\top)^{-1} (H\bar{\beta} - d)$ und $s \cdot \bar{\sigma}^2$ unabhängig, $\frac{(n-r)\bar{\sigma}^2}{\sigma^2} \sim \chi_{n-r}^2$. Also bleibt nur noch zu zeigen, daß

$$\left(\underbrace{H\bar{\beta} - d}_{=\varepsilon^\top B^\top} \right)^\top \left(H(X^\top X)^{-1} H^\top \right)^{-1} \left(\underbrace{H\bar{\beta} - d}_{=B\varepsilon} \right) \stackrel{(H_0)}{\sim} \chi_s^2.$$

Es gilt

$$\begin{aligned} & \varepsilon^\top B^\top \left(H(X^\top X)^{-1} H^\top \right)^{-1} B\varepsilon \\ &= \varepsilon^\top X \underbrace{\left((X^\top X)^{-1} \right)^\top H^\top \left(H(X^\top X)^{-1} H^\top \right)^{-1} H(X^\top X)^{-1} X^\top}_{=A} \varepsilon \end{aligned}$$

Man kann leicht zeigen, daß A symmetrisch, idempotent und $\text{Rang}(A) = s$ ist. Zeigen wir zum Beispiel die Idempotenz:

$$\begin{aligned} A^2 &= X \left((X^\top X)^{-1} \right)^\top H^\top \left(H(X^\top X)^{-1} H^\top \right)^{-1} \underbrace{H(X^\top X)^{-1} X^\top X}_{H \text{ (Lemma 2.3.4, 2.)}} \left((X^\top X)^{-1} \right)^\top H^\top \\ &\quad \cdot \left(H(X^\top X)^{-1} H^\top \right)^{-1} H(X^\top X)^{-1} X^\top \\ &= X \left((X^\top X)^{-1} \right)^\top H^\top \left(H(X^\top X)^{-1} H^\top \right)^{-1} H(X^\top X)^{-1} X^\top = A, \end{aligned}$$

weil $\left((X^\top X)^{-1} \right)^\top$ auch eine verallgemeinerte Inverse von $X^\top X$ ist (nach Lemma 2.3.2). Somit hängt auch $H(X^\top X)^{-1} H^\top = CX(X^\top X)^{-1} X^\top C^\top$ nicht von der Wahl von $(X^\top X)^{-1}$ ab, vgl. den Beweis des Satzes 2.3.6. Nach Satz 2.1.8 ist $\frac{\varepsilon^\top}{\sigma} A \frac{\varepsilon}{\sigma} \sim \chi_s^2$, wegen $\varepsilon \sim N(0, \sigma^2 \mathcal{I})$ und somit $T \stackrel{H_0}{\sim} F_{s,n-r}$. \square

2.3.6 Konfidenzbereiche

Ähnlich wie in Abschnitt 2.2.5 werden wir Konfidenzbereiche für unterschiedliche Funktionen vom Parametervektor β angeben. Aus dem Satz 2.3.9 ergibt sich unmittelbar folgender Konfidenzbereich zum Niveau $1 - \alpha \in (0, 1)$:

Folgerung 2.3.1. Sei $Y = X\beta + \varepsilon$ ein multivariates Regressionsmodell mit $\text{Rang}(X) = r < m$, H eine $(s \times m)$ -Matrix mit $\text{Rang}(H) = s$, $s \in \{1, \dots, m\}$ und $H_0 : H\beta = d$ testbar $\forall d \in \mathbb{R}^s$. Dann ist

$$\left\{ d \in \mathbb{R}^s : \frac{(H\bar{\beta} - d)^\top (H(X^\top X)^- H^\top)^{-1} (H\bar{\beta} - d)}{s \cdot \bar{\sigma}^2} \leq F_{s, n-r, 1-\alpha} \right\}$$

ein Konfidenzbereich für $H\beta$ zum Niveau $1 - \alpha$.

Folgerung 2.3.2. Sei $h^\top \beta$ eine schätzbare lineare Funktion von β , $h \in \mathbb{R}^m$. Dann ist

$$\left(h^\top \bar{\beta} - t_{n-r, 1-\alpha/2} \cdot \bar{\sigma} \sqrt{h^\top (X^\top X)^- h}, h^\top \bar{\beta} + t_{n-r, 1-\alpha/2} \cdot \bar{\sigma} \sqrt{h^\top (X^\top X)^- h} \right)$$

ein Konfidenzintervall für $h^\top \beta$ zum Niveau $1 - \alpha$.

Beweis. Setzen wir $s = 1$ und $H = h^\top$. Aus Satz 2.3.9 folgt

$$\begin{aligned} T &= \frac{(h^\top \bar{\beta} - d)^\top (h^\top (X^\top X)^- h)^{-1} (h^\top \bar{\beta} - d)}{\bar{\sigma}^2} = \frac{(h^\top \bar{\beta} - d) (h^\top \bar{\beta} - d)}{\bar{\sigma}^2 (h^\top (X^\top X)^- h)} \\ &= \frac{(h^\top \bar{\beta} - d)^2}{\bar{\sigma}^2 (h^\top (X^\top X)^- h)} \sim F_{1, n-r} \end{aligned}$$

unter der Voraussetzung $h^\top \beta = d$, weil $h^\top (X^\top X)^- h$ eindimensional (eine Zahl) ist. Deshalb gilt

$$\sqrt{T} = \frac{h^\top \beta - h^\top \bar{\beta}}{\bar{\sigma} \sqrt{h^\top (X^\top X)^- h}} \sim t_{n-r}$$

und somit

$$\mathbb{P} \left(-t_{n-r, 1-\alpha/2} \leq \sqrt{T} \leq t_{n-r, 1-\alpha/2} \right) = 1 - \alpha.$$

Daraus folgt das obige Konfidenzintervall. \square

Man kann sogar eine stärkere Version von 2.3.2 beweisen, die für alle h aus einem linearen Unterraum gilt:

Satz 2.3.10 (Konfidenzband von Scheffé). Sei $H = (h_1, \dots, h_s)^\top$, $h_1, \dots, h_s \in \mathbb{R}^m$, $1 \leq s \leq m$ und $H_0 : H\beta = d$ testbar $\forall d \in \mathbb{R}^s$. Sei $\text{Rang}(H) = s$ und $\mathcal{L} = \langle h_1, \dots, h_s \rangle$ der lineare Unterraum, der von den Vektoren h_1, \dots, h_s aufgespannt wird. Dann gilt:

$$\mathbb{P} \left(\max_{h \in \mathcal{L}} \left\{ \frac{(h^\top \beta - h^\top \bar{\beta})^2}{\bar{\sigma}^2 h^\top (X^\top X)^{-1} h} \right\} \leq s F_{s, n-r, 1-\alpha} \right) = 1 - \alpha$$

Somit ist

$$\left[h^\top \bar{\beta} - \sqrt{s F_{s, n-r, 1-\alpha}} \cdot \bar{\sigma} \sqrt{h^\top (X^\top X)^{-1} h}, h^\top \bar{\beta} + \sqrt{s F_{s, n-r, 1-\alpha}} \cdot \bar{\sigma} \sqrt{h^\top (X^\top X)^{-1} h} \right]$$

ein (gleichmäßiges bzgl. $h \in \mathcal{L}$) Konfidenzintervall für $h^\top \beta$.

Beweis. Aus dem Satz 2.3.9 folgt $\forall \alpha \in (0, 1)$:

$$\mathbb{P} \left(\underbrace{(H\bar{\beta} - H\beta)^\top (H(X^\top X)^{-1} H^\top)^{-1} (H\bar{\beta} - H\beta)}_{T_1} \leq s \cdot \bar{\sigma}^2 F_{s, n-r, 1-\alpha} \right) = 1 - \alpha.$$

Falls wir zeigen können, daß

$$T_1 = \max_{x \in \mathbb{R}^s, x \neq 0} \left\{ \frac{(x^\top (H\bar{\beta} - H\beta))^2}{x^\top (H(X^\top X)^{-1} H^\top) x} \right\}, \quad (2.3.6)$$

dann ist der Satz bewiesen, denn

$$\begin{aligned} 1 - \alpha &= \mathbb{P} \left(T_1 \leq \underbrace{s \bar{\sigma}^2 F_{s, n-r, 1-\alpha}}_t \right) = \mathbb{P} \left(\max_{x \in \mathbb{R}^s, x \neq 0} \left\{ \frac{(x^\top (H\bar{\beta} - H\beta))^2}{x^\top (H(X^\top X)^{-1} H^\top) x} \right\} \leq t \right) \\ &= \mathbb{P} \left(\max_{x \in \mathbb{R}^s, x \neq 0} \left\{ \frac{((H^\top x)^\top \bar{\beta} - (H^\top x)^\top \beta)^2}{(H^\top x)^\top (X^\top X)^{-1} (H^\top x)} \right\} \leq t \right) \quad \text{und weil } H^\top x = h \in \mathcal{L} \\ &= \mathbb{P} \left(\max_{h \in \mathcal{L}} \left\{ \frac{(h^\top \bar{\beta} - h^\top \beta)^2}{h^\top (X^\top X)^{-1} h} \right\} \leq s \bar{\sigma}^2 F_{s, n-r, 1-\alpha} \right). \end{aligned}$$

Also, zeigen wir die Gültigkeit von (2.3.6). Es genügt zu zeigen, daß T_1 die obere Schranke von

$$\frac{(x^\top (H\bar{\beta} - H\beta))^2}{x^\top (H(X^\top X)^{-1} H^\top) x}$$

darstellt, die auch angenommen wird. Da $H(X^\top X)^{-1} H^\top$ positiv definit ist und invertierbar, existiert eine invertierbare $(s \times s)$ -Matrix B mit der Eigenschaft $BB^\top =$

$H(X^\top X)^{-1}H^\top$. Dann gilt

$$\begin{aligned} \left(x^\top(H\bar{\beta} - H\beta)\right)^2 &= \left(\underbrace{x^\top B}_{(B^\top x)^\top} \cdot B^{-1}(H\bar{\beta} - H\beta)\right)^2 \\ &\leq |B^\top x|^2 \cdot |B^{-1}(H\bar{\beta} - H\beta)|^2 \quad (\text{wegen der Ungleichung von Cauchy-Schwarz}) \\ &= x^\top B B^\top x (H\bar{\beta} - H\beta)^\top \cdot \underbrace{(B^{-1})^\top B^{-1}}_{=(B^\top)^{-1}B^{-1}=(BB^\top)^{-1}} (H\bar{\beta} - H\beta) \\ &= x^\top H(X^\top X)^{-1}H^\top x \cdot (H\bar{\beta} - H\beta)^\top \left(H(X^\top X)^{-1}H^\top\right)^{-1} (H\bar{\beta} - H\beta). \end{aligned}$$

Somit gilt

$$\frac{(x^\top(H\bar{\beta} - H\beta))^2}{x^\top(H(X^\top X)^{-1}H^\top)x} \leq (H\bar{\beta} - H\beta)^\top \left(H(X^\top X)^{-1}H^\top\right)^{-1} (H\bar{\beta} - H\beta) = T_1.$$

Man kann leicht prüfen, daß diese Schranke für $x = (H(X^\top X)^{-1}H^\top)^{-1} (H\bar{\beta} - H\beta)$ angenommen wird. \square

2.3.7 Einführung in die Varianzanalyse

In diesem Abschnitt geben wir ein Beispiel für die Verwendung linearer Modelle mit Design-Matrix, die keinen vollen Rang besitzt. Dabei handelt es sich um die Aussage der *Variabilität der Erwartungswerte* in der Stichprobe $Y = (Y_1, \dots, Y_n)^\top$, die auf englisch *analysis of variance*, kurz *ANOVA*, heißt. Später werden wir auch denselben Begriff *Varianzanalyse* dafür verwenden.

Betrachten wir zunächst die *einfaktorielle Varianzanalyse*, bei der man davon ausgeht, daß die Stichprobe (Y_1, \dots, Y_n) in k homogene Teilklassen $(Y_{ij}, j = 1, \dots, n_i), i = 1, \dots, k$ zerlegbar ist, mit den Eigenschaften:

$$1. \mathbb{E}(Y_{ij}) = \mu_i = \mu + \alpha_i, \quad j = 1, \dots, n_i, \quad i = 1, \dots, k.$$

$$2. n_i > 1, \quad i = 1, \dots, k, \quad \sum_{i=1}^k n_i = n, \quad \sum_{i=1}^k n_i \alpha_i = 0.$$

Dabei ist μ ein Faktor, der allen Klassen gemeinsam ist, und α_i verkörpert die *klassenspezifischen Differenzen* zwischen den Erwartungswerten μ_1, \dots, μ_k . Die Nummer $i = 1, \dots, k$ der Klassen wird als *Stufe eines Einflussfaktors* (zum Beispiel die Dosis eines Medikaments in einer klinischen Studie) und $\alpha_i, i = 1, \dots, k$ als *Effekt* der i -ten Stufe gedeutet. Die Nebenbedingung $\sum_{i=1}^k n_i \alpha_i = 0$ bewirkt, daß die Umrechnung

$(\mu_1, \dots, \mu_k) \longleftrightarrow (\mu, \alpha_1, \dots, \alpha_k)$ eindeutig wird und daß $\mu = \frac{1}{n} \sum_{i=1}^k \sum_{j=1}^{n_i} \mathbb{E} Y_{ij}$. Es wird vorausgesetzt, daß μ_i mit unkorrelierten Meßfehlern ε_{ij} gemessen werden kann, das heißt

$$Y_{ij} = \mu_i + \varepsilon_{ij} = \mu + \alpha_i + \varepsilon_{ij}, \quad i = 1, \dots, k, \quad j = 1, \dots, n_i \quad (2.3.7)$$

$$\mathbb{E} \varepsilon_{ij} = 0, \quad \text{Var} \varepsilon_{ij} = \sigma^2, \quad \varepsilon_{ij} \text{ unkorreliert}, \quad i = 1, \dots, k, \quad j = 1, \dots, n_i. \quad (2.3.8)$$

Es soll die *klassische ANOVA-Hypothese* getestet werden, daß *keine* Variabilität in den Erwartungswerten μ_i auffindbar ist:

$$H_0: \mu_1 = \mu_2 = \dots = \mu_k,$$

was bedeutet, daß

$$H_0: \alpha_1 = \alpha_2 = \dots = \alpha_k.$$

Aus der Nebenbedingung

$$\sum_{i=1}^k n_i \alpha_i = 0.$$

folgt: $\alpha_i = 0$

Die Problemstellung (2.3.7) kann in der Form der multivariaten linearen Regression folgendermaßen umgeschrieben werden:

$$\begin{aligned} Y &= X\beta + \varepsilon, \text{ wobei } Y = (Y_{11}, \dots, Y_{1n_1}, Y_{21}, \dots, Y_{2n_2}, \dots, Y_{k1}, \dots, Y_{kn_k})^\top, \\ \beta &= (\mu, \alpha_1, \dots, \alpha_k)^\top, \\ \varepsilon &= (\varepsilon_{11}, \dots, \varepsilon_{1n_1}, \dots, \varepsilon_{k1}, \dots, \varepsilon_{kn_k})^\top, \\ X &= \begin{pmatrix} 1 & 1 & 0 & \dots & \dots & 0 \\ 1 & 1 & 0 & \dots & \dots & 0 \\ \vdots & & & & & \\ 1 & 1 & 0 & \dots & \dots & 0 \\ 1 & 0 & 1 & 0 & \dots & 0 \\ \vdots & & & & & \\ 1 & 0 & 1 & 0 & \dots & 0 \\ \vdots & & & & & \\ 1 & 0 & \dots & \dots & 0 & 1 \\ \vdots & & & & & \\ 1 & 0 & \dots & \dots & 0 & 1 \end{pmatrix} \begin{matrix} \left. \begin{matrix} \\ \\ \\ \end{matrix} \right\} n_1 \\ \left. \begin{matrix} \\ \\ \end{matrix} \right\} n_2 \\ \vdots \\ \left. \begin{matrix} \\ \\ \end{matrix} \right\} n_k \end{matrix} \end{pmatrix}$$

Die $(n \times (k + 1))$ -Matrix X hat den Rang $k < m = k + 1$, somit ist die Theorie von Abschnitt 2.3 auf dieses Modell komplett anwendbar.

Übungsaufgabe 2.3.1. Zeigen Sie, dass die ANOVA-Hypothese

$$H_0: \alpha_i = 0, \quad \forall i = 1, \dots, k$$

nicht testbar ist!

Um eine äquivalente testbare Hypothese aufzustellen, benutzt man

$$H_0: \alpha_1 - \alpha_2 = 0, \dots, \alpha_1 - \alpha_k = 0 \quad \text{bzw.} \quad H_0: H\beta = 0$$

für die $(k-1) \times (k+1)$ -Matrix

$$H = \begin{pmatrix} 0 & 1 & -1 & 0 & \dots & 0 \\ 0 & 1 & 0 & -1 & \dots & 0 \\ \vdots & & & & & \\ 0 & 1 & 0 & \dots & -1 & 0 \\ 0 & 1 & 0 & \dots & 0 & -1 \end{pmatrix}.$$

(Zeigen Sie es!)

Bei der *zweifaktoriellen Varianzanalyse* wird die Stichprobe (Y_1, \dots, Y_n) in Abhängigkeit von 2 Faktoren in $k_1 \cdot k_2$ homogene Gruppen aufgeteilt:

$$Y_{i_1 i_2 j}, \quad j = 1, \dots, n_{i_1 i_2}$$

für $i_1 = 1, \dots, k_1, i_2 = 1, \dots, k_2$, sodaß

$$\sum_{i_1=1}^{k_1} \sum_{i_2=1}^{k_2} n_{i_1 i_2} = n.$$

Hier wird angenommen, daß

$$\mathbb{E} Y_{i_1 i_2 j} = \mu_{i_1 i_2} = \mu + \alpha_{i_1} + \beta_{i_2} + \gamma_{i_1 i_2}, \quad i_1 = 1, \dots, k_1, i_2 = 1, \dots, k_2,$$

somit stellt man folgendes lineares Modell auf:

$$Y_{i_1 i_2 j} = \mu_{i_1 i_2} + \varepsilon_{i_1 i_2 j} = \mu + \alpha_{i_1} + \beta_{i_2} + \gamma_{i_1 i_2} + \varepsilon_{i_1 i_2 j}, \\ j = 1, \dots, n_{i_1 i_2}, i_1 = 1, \dots, k_1, i_2 = 1, \dots, k_2.$$

Übungsaufgabe 2.3.2. Schreiben Sie die Design-Matrix X für diesen Fall explizit auf! Zeigen Sie, daß sie wieder keinen vollen Rang besitzt.

*

3 Verallgemeinerte lineare Modelle

Eine andere Klasse von Regressionsmodellen erlaubt einerseits einen beliebigen funktionalen Zusammenhang g zwischen dem Mittelwert der Zielvariablen $\mathbb{E} Y_i$ und dem linearen Teil $X\beta$, der aus linearen Kombinationen der Einträge der Designmatrix $X = (x_{ij})$ und des Parametervektors $\beta = (\beta_1, \dots, \beta_m)^\top$ besteht; andererseits lässt sie andere Verteilungen von Y_i zu, die nicht notwendigerweise auf der Normalverteilung (und Funktionen davon) basieren. So ist es möglich, Daten Y_i zu betrachten, die eine endliche Anzahl von Ausprägungen haben (z.B. „Ja“ und „Nein“ in ökonomischen Meinungsumfragen). Die Klasse aller möglichen Verteilungen wird durch die sog. *Exponentialfamilie* begrenzt, die wir in Kürze einführen werden.

Sei Y_1, \dots, Y_n eine Zufallsstichprobe der Zielvariablen des Modells und sei

$$X = (x_{ij})_{\substack{i=1,\dots,n \\ j=1,\dots,m}}$$

die Designmatrix der Ausgangsvariablen, die hier nicht zufällig sind.

Definition 3.0.4. Das *verallgemeinerte lineare Modell* ist gegeben durch

$$(g(\mathbb{E} Y_1), \dots, g(\mathbb{E} Y_n))^\top = X\beta \quad \text{mit } \beta = (\beta_1, \dots, \beta_m)^\top, \quad (3.0.1)$$

wobei $g : G \subset \mathbb{R} \rightarrow \mathbb{R}$ die sog. *Linkfunktion* mit dem Definitionsbereich G ist. Der Rang $(X) = m$.

Unter der Annahme, dass g explizit bekannt ist, soll hier der Parametervektor β aus (Y_1, \dots, Y_n) geschätzt werden. Wir setzen voraus, dass $Y_i, i = 1, \dots, n$, unabhängig, aber nicht unbedingt identisch verteilt sind. Ihre Verteilung gehört jedoch zur folgenden Klasse von Verteilungen:

3.1 Exponentialfamilie von Verteilungen

Definition 3.1.1. Die Verteilung einer Zufallsvariable Y gehört zur *Exponentialfamilie*, falls es Funktionen $a : \mathbb{R} \times \mathbb{R}_+ \rightarrow \mathbb{R}$ und $b : \Theta \rightarrow \mathbb{R}$ gibt, für die

- im *absolutstetigen Fall* die Dichte von Y gegeben ist durch

$$f_\theta(y) = \exp \left\{ \frac{1}{\tau^2} (y\theta + a(y, \tau) - b(\theta)) \right\}, \quad y \in C, \quad (3.1.1)$$

wobei $C \subseteq \mathbb{R}$ offen ist.

- im *diskreten Fall* die Zähldichte von Y gegeben ist durch

$$P_\theta(Y = y) = \exp \left\{ \frac{1}{\tau^2} (y\theta + a(y, \tau) - b(\theta)) \right\}, y \in C, \quad (3.1.2)$$

wobei C der (höchstens) abzählbare Wertebereich von Y , τ^2 der sog. *Störparameter*, $\theta \in \Theta \subset \mathbb{R}$ ein Parameter und

$$\Theta = \left\{ \theta \in \mathbb{R} : \int_{\mathbb{R}} \exp \left\{ \frac{y\theta + a(y, \tau)}{\tau^2} \right\} dy < \infty \right\}$$

bzw. im diskreten Fall:

$$\Theta = \left\{ \theta \in \mathbb{R} : \sum_{y \in C} \exp \left\{ \frac{y\theta + a(y, \tau)}{\tau^2} \right\} < \infty \right\}$$

der natürliche Parameterraum ist, der mindestens zwei verschiedene Elemente enthält.

Lemma 3.1.1. Θ ist ein Intervall.

Beweis. Zeigen wir, dass $\Theta \subset \mathbb{R}$ konvex ist. Dann ist es notwendigerweise ein (möglicherweise unendliches) Intervall. Für beliebige $\theta_1, \theta_2 \in \Theta$ (mindestens ein solches Paar gibt es nach Definition 3.1.1) zeigen wir, dass $\alpha\theta_1 + (1 - \alpha)\theta_2 \in \Theta$ für alle $\alpha \in (0, 1)$. Nehmen wir an, dass die Verteilung von Y absolut stetig ist. Da $\theta_i \in \Theta$, es gilt

$$\int_{\mathbb{R}} \exp \left\{ \frac{1}{\tau^2} (y\theta_i + a(y, \tau)) \right\} dy < \infty, \quad i = 1, 2.$$

Durch die offensichtliche Ungleichung

$$\alpha x_1 + (1 - \alpha)x_2 \leq \max\{x_1, x_2\}, \quad x_1, x_2 \in \mathbb{R} \quad \alpha \in (0, 1)$$

erhalten wir

$$\begin{aligned} & \exp \left\{ \frac{1}{\tau^2} (y(\alpha\theta_1 + (1 - \alpha)\theta_2) + a(y, \tau)) \right\} \\ &= \exp \left\{ \alpha \frac{1}{\tau^2} (y\theta_1 + a(y, \tau)) + (1 - \alpha) \frac{1}{\tau^2} (y\theta_2 + a(y, \tau)) \right\} \\ &\leq \max_{i=1,2} \exp \left\{ \frac{1}{\tau^2} (y\theta_i + a(y, \tau)) \right\} \leq \exp \left\{ \frac{1}{\tau^2} (y\theta_1 + a(y, \tau)) \right\} + \exp \left\{ \frac{1}{\tau^2} (y\theta_2 + a(y, \tau)) \right\}, \end{aligned}$$

so dass

$$\int_{\mathbb{R}} \exp \left\{ \frac{1}{\tau^2} (y(\alpha\theta_1 + (1 - \alpha)\theta_2) + a(y, \tau)) \right\} dy \leq \sum_{i=1}^2 \int_{\mathbb{R}} \exp \left\{ \frac{1}{\tau^2} (y\theta_i + a(y, \tau)) \right\} dy < \infty$$

nach Voraussetzungen des Lemmas.

$$\Rightarrow \alpha\theta_1 + (1 - \alpha)\theta_2 \in \Theta,$$

und Θ ist ein Intervall. □

Beispiel 3.1.1. Welche Verteilungen gehören zur Exponentialfamilie?

1. **Normalverteilung:** Falls $Y \sim \mathcal{N}(\mu, \sigma^2)$, dann ist der Erwartungswert μ der uns interessierende Parameter, σ^2 ist dagegen der Störparameter. Es gilt:

$$\begin{aligned} f_\mu(y) &= \frac{1}{\sqrt{2\pi\sigma^2}} \cdot e^{-\frac{(y-\mu)^2}{2\sigma^2}} \\ &= \exp \left\{ -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2} \left(\frac{y^2}{\sigma^2} - \frac{2y\mu}{\sigma^2} + \frac{\mu^2}{\sigma^2} \right) \right\} \\ &= \exp \left\{ \frac{1}{\sigma^2} \left(y\mu - \frac{y^2}{2} - \left(\frac{\mu^2}{2} + \frac{\sigma^2}{2} \log(2\pi\sigma^2) \right) \right) \right\}, \end{aligned}$$

so dass

$$\theta = \mu, \quad \tau = \sigma, \quad a(y, \tau) = -\frac{y^2}{2} - \frac{\sigma^2}{2} \log(2\pi\sigma^2) \quad \text{und} \quad b(\mu) = b(\theta) = \frac{\mu^2}{2}.$$

2. **Bernoulli-Verteilung:** $Y \sim \text{Bernoulli}(p)$, $p \in [0; 1]$.

Sie wird etwa im Falle von Meinungsumfragen in der Marktforschung verwendet, in denen

$$Y = \begin{cases} 1, & \text{falls die Antwort „ja“} \\ 0, & \text{falls die Antwort „nein“} \end{cases} \text{ auf eine Frage der Enquete gegeben wurde.}$$

Dabei ist die Wahrscheinlichkeit $P(Y = 1) = p$, $P(Y = 0) = 1 - p$. Dann gilt für $y \in \{0, 1\}$:

$$\begin{aligned} P_\theta(Y = y) &= p^y(1 - p)^{1-y} = e^{y \log p + (1-y) \log(1-p)} \\ &= e^{y \log \frac{p}{1-p} - (-\log(1-p))}. \end{aligned}$$

Somit gehört die Bernoulli-Verteilung zur Exponentialfamilie mit

$$\theta = \log \frac{p}{1-p}, \quad \tau = 1, \quad a(y, \tau) = 0, \quad b(\theta) = -\log(1-p) = \log(1 + e^\theta).$$

3. **Poisson-Verteilung:** Falls $Y \sim \text{Poisson}(\lambda)$, $\lambda > 0$, dann gilt für $y \in \mathbb{N}_0$

$$P_\theta(Y = y) = e^{-\lambda} \cdot \frac{\lambda^y}{y!} = e^{y \log \lambda - \log(y!) - \lambda} \quad .$$

Somit gehört die Poisson-Verteilung zur Exponentialfamilie mit

$$\theta = \log \lambda, \quad \tau = 1, \quad a(y, \tau) = -\log(y!), \quad b(\theta) = \lambda = e^\theta .$$

Lemma 3.1.2. Falls die Verteilung von Y zur Exponentialfamilie gehört, $\mathbb{E}Y^2 < \infty$ und $b : \Theta \rightarrow \mathbb{R}$ zweimal stetig differenzierbar ist mit $b''(\theta) > 0$ für alle $\theta \in \Theta$, dann gilt

$$\mathbb{E}Y = b'(\theta), \quad \text{Var} Y = \tau^2 b''(\theta) .$$

Beweis. 1. Führen wir den Beweis für den Fall der absolut stetigen Verteilung von Y . Der diskrete Fall lässt sich analog behandeln, wenn man das \int -Zeichen durch \sum ersetzt. Es gilt

$$\begin{aligned} \mathbb{E}Y &= \int_{\mathbb{R}} y f_\theta(y) dy = \int_{\mathbb{R}} y \exp \left\{ \frac{1}{\tau^2} (y\theta + a(y, \tau) - b(\theta)) \right\} dy \\ &= e^{-\frac{b(\theta)}{\tau^2}} \cdot \tau^2 \int_{\mathbb{R}} \frac{\partial}{\partial \theta} \exp \left\{ \frac{1}{\tau^2} (y\theta + a(y, \tau)) \right\} dy \\ &= e^{-\frac{b(\theta)}{\tau^2}} \cdot \tau^2 \frac{\partial}{\partial \theta} \int_{\mathbb{R}} \exp \left\{ \frac{1}{\tau^2} (y\theta + a(y, \tau)) \right\} dy \\ &= e^{-\frac{b(\theta)}{\tau^2}} \cdot \tau^2 \frac{\partial}{\partial \theta} \left(\underbrace{e^{\frac{b(\theta)}{\tau^2}} \int_{\mathbb{R}} \exp \left\{ \frac{1}{\tau^2} (y\theta + a(y, \tau) - b(\theta)) \right\} dy}_{\int_{\mathbb{R}} f_\theta(y) dy = 1} \right) \\ &= e^{-\frac{b(\theta)}{\tau^2}} \tau^2 \frac{\partial}{\partial \theta} \left(e^{\frac{b(\theta)}{\tau^2}} \right) = e^{-\frac{b(\theta)}{\tau^2}} \cdot \tau^2 \frac{b'(\theta)}{\tau^2} e^{\frac{b(\theta)}{\tau^2}} = b'(\theta). \end{aligned}$$

2. Es bleibt noch zu zeigen:

Übungsaufgabe 3.1.1. Beweisen Sie die Formel

$$\text{Var} Y = \tau^2 b''(\theta) \quad (\text{analog zu 1}).$$

□

3.2 Linkfunktion

Die Zielgrößen Y_i , $i = 1, \dots, n$ seien also unabhängig verteilt mit einer Verteilung, die zur Exponentialfamilie gehört und einer (Zähl)Dichte wie in (3.1.1) bzw. (3.1.2). Setzen wir voraus, dass $b : \Theta \rightarrow \mathbb{R}$ zweimal stetig differenzierbar ist mit $b''(\theta) > 0$ für alle $\theta \in \Theta$. Sei ein verallgemeinertes lineares Modell (3.0.1) gegeben.

Definition 3.2.1. (Natürliche Linkfunktion)

Die Linkfunktion $g : G \rightarrow \mathbb{R}$ heißt *natürlich*, falls $g = (b')^{-1}$, $G = \{b'(\theta) : \theta \in \Theta\}$ und g zweimal stetig differenzierbar ist mit $g'(x) \neq 0$ für alle $x \in G$.

Die Frage, warum die natürliche Linkfunktion so heißt, beantwortet folgendes Lemma:

Lemma 3.2.1. Falls das verallgemeinerte lineare Modell (3.0.1) die natürliche Linkfunktion besitzt, dann gilt $(\theta_1, \dots, \theta_n)^\top = X\beta$.

Beweis. Wegen $b''(\theta) > 0$ ist $b'(\theta)$ monoton steigend, also invertierbar. Führen wir folgende Bezeichnungen ein:

$$\mu_i = \mathbb{E} Y_i, \quad \eta_i = x_i^\top \beta, \quad x_i = (x_{i1}, \dots, x_{im})^\top, \quad i = 1, \dots, n$$

Da g invertierbar ist, gilt

$$\mu_i = g^{-1}(x_i^\top \beta) = g^{-1}(\eta_i), \quad i = 1, \dots, n$$

Andererseits folgt $\mu_i = b'(\theta_i)$ aus Lemma 3.1.2, so dass

$$b'(\theta_i) = g^{-1}(\eta_i) \stackrel{\text{Definition 3.2.1}}{=} b'(\eta_i), \quad i = 1, \dots, n \quad .$$

Wegen der Monotonie von b' folgt die Behauptung $\theta_i = \eta_i$, $i = 1, \dots, n$. □

Beispiel 3.2.1. Berechnen wir die natürlichen Linkfunktionen für die Verteilungen von Beispiel 3.1.1.

1. **Normalverteilung:** da $b(\mu) = \frac{\mu^2}{2}$, gilt

$$b'(x) = \frac{2x}{2} = x \text{ und somit } g(x) = (b')^{-1}(x) = x \text{ .}$$

Die natürliche Linkfunktion ist $g(x) = x$, somit gilt hier

$$(\mu_1, \dots, \mu_n)^\top = (\mathbb{E} Y_1, \dots, \mathbb{E} Y_n)^\top = X\beta \text{ .}$$

Das ist genau der Fall der linearen Regression.

2. **Bernoulli-Verteilung:** da $b(\theta) = \log(1 + e^\theta)$, gilt

$$\begin{aligned} b'(x) &= \frac{1}{1 + e^x} \cdot e^x = y \\ \Leftrightarrow \frac{1}{e^{-x} + 1} &= y \\ \Leftrightarrow \frac{1}{y} - 1 &= e^{-x} \\ \Leftrightarrow x &= -\log \frac{1-y}{y} = \log \frac{y}{1-y} \\ \Rightarrow g(x) &= (b')^{-1}(x) = \log \frac{x}{1-x} . \end{aligned}$$

Das verallgemeinerte lineare Regressionsmodell im Falle der Bernoulli-Verteilung wird *binäre (kategoriale) Regression* genannt. Falls sie mit der natürlichen Linkfunktion verwendet wird, nennt man sie *logistische Regression*. In diesem Fall gilt

$$\begin{aligned} (p_1, \dots, p_n)^\top &= (\mathbb{E} Y_1, \dots, \mathbb{E} Y_n)^\top \\ \theta_i &= \log \frac{p_i}{1-p_i} = x_i^\top \beta, \quad i = 1, \dots, n \\ \Leftrightarrow e^{\theta_i} &= \frac{p_i}{1-p_i} \\ \Leftrightarrow p_i &= \frac{e^{\theta_i}}{1 + e^{\theta_i}} \\ \Leftrightarrow p_i &= \frac{e^{x_i^\top \beta}}{1 + e^{x_i^\top \beta}}, \quad i = 1, \dots, n . \end{aligned}$$

Das Verhältnis

$$\frac{p_i}{1-p_i} = \frac{P(Y_i = 1)}{P(Y_i = 0)}, \quad i = 1, \dots, n$$

wird in der englischsprachigen Literatur *Odds* genannt. Der Logarithmus des Odds heißt *Logit*:

$$\log \frac{p_i}{1-p_i}, \quad i = 1, \dots, n .$$

Logits sind also hier „neue Zielvariablen“, die durch Linearkombinationen $x_i^\top \beta$ geschätzt werden.

Eine alternative Linkfunktion, die oft benutzt wird, ist $g(x) = \Phi^{-1}(x)$, die *Quantilfunktion der Normalverteilung*. Sie ist keine natürliche Linkfunktion. Mit ihrer Hilfe bekommt man das sog. *Probit-Modell*:

$$p_i = \Phi(x_i^\top \beta), \quad i = 1, \dots, n .$$

3. **Poisson-Verteilung:** da $b(\theta) = e^\theta$, ist in diesem Fall

$$g(x) = (b')^{-1}(x) = \log x, \quad x > 0$$

die natürliche Linkfunktion. Somit hat das verallgemeinerte lineare Modell mit der natürlichen Linkfunktion folgende Darstellung

$$(\log \lambda_1, \dots, \log \lambda_n)^\top = X\beta \quad \text{oder} \quad \lambda_i = e^{x_i^\top \beta}, \quad i = 1, \dots, n.$$

3.3 Maximum-Likelihood-Schätzung von β

Da die (Zähl)Dichte von Y_i die Gestalt

$$\exp \left\{ \frac{1}{\tau^2} (y\theta_i + a(y, \tau) - b(\theta_i)) \right\}$$

hat und Y_i unabhängig sind, kann man die Log-Likelihood-Funktion der Stichprobe $Y = (Y_1, \dots, Y_n)$ in folgender Form aufschreiben:

$$\log L(Y, \theta) = \log \prod_{i=1}^n f_{\theta_i}(Y_i) = \frac{1}{\tau^2} \sum_{i=1}^n \left(Y_i \theta_i + a(Y_i, \tau) - b(\theta_i) \right). \quad (3.3.1)$$

Aus dem Beweis des Lemmas 3.2.1 folgt, dass

$$\theta_i = (b')^{-1}(g^{-1}(x_i^\top \beta)), \quad i = 1, \dots, n, \quad (3.3.2)$$

was bedeutet, dass die Funktion $\log L(Y, \theta)$ eine Funktion von Parameter β ist. In der Zukunft schreiben wir $\log L(Y, \beta)$, um diese Tatsache zu unterstreichen.

Unser Ziel ist es, den Maximum-Likelihood-Schätzer $\hat{\beta}$ für β zu berechnen:

$$\hat{\beta} = \underset{\beta}{\operatorname{argmax}} \log L(Y, \beta).$$

Dafür wird die notwendige Bedingung des Extremums

$$\frac{\partial \log L(Y, \beta)}{\partial \beta_i} = 0, \quad i = 1, \dots, m,$$

untersucht. Verwenden wir folgende Bezeichnungen:

$$U_i(\beta) = \frac{\partial \log L(Y, \beta)}{\partial \beta_i}, \quad i = 1, \dots, m,$$

$$U(\beta) = (U_1(\beta), \dots, U_m(\beta))^\top,$$

$$I_{ij}(\beta) = \mathbb{E} [U_i(\beta) U_j(\beta)], \quad i, j = 1, \dots, m.$$

Definition 3.3.1. 1. Die Matrix $I(\beta) = (I_{ij}(\beta))_{i,j=1}^m$ heißt *Fisher-Informationsmatrix*.

2. Führen wir die sog. *Hesse-Matrix* $W(\beta)$ als zufällige Matrix

$$W(\beta) = (W_{ij}(\beta))_{i,j=1}^m \quad \text{mit} \quad W_{ij}(\beta) = \frac{\partial^2}{\partial \beta_i \partial \beta_j} \log L(Y, \beta)$$

ein. Diese $(m \times m)$ -Matrix enthält die partiellen Ableitungen 2. Ordnung der Log-Likelihood-Funktion, die für die numerische Lösung der Maximierungsaufgabe

$$\log L(Y, \beta) \rightarrow \max_{\beta}$$

von Bedeutung sein werden.

Satz 3.3.1. Man kann zeigen, dass $U(\beta)$ und $I(\beta)$ folgende explizite Form haben:

1. Es gilt

$$U_j(\beta) = \sum_{i=1}^n x_{ij} (Y_i - \mu_i(\beta)) \frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \frac{1}{\sigma_i^2(\beta)}, \quad j = 1, \dots, m,$$

2. Es gilt

$$I_{jk}(\beta) = \sum_{i=1}^n x_{ij} x_{ik} \left(\frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \right)^2 \frac{1}{\sigma_i^2(\beta)}, \quad j, k = 1, \dots, m,$$

wobei $\eta_i = x_i^\top \beta$, $\mu_i(\beta) = g^{-1}(x_i^\top \beta)$ der Erwartungswert von Y_i und

$$\sigma_i^2(\beta) \stackrel{\text{Lemma 3.1.2}}{=} \tau^2 b''(\theta_i) \stackrel{(3.3.2)}{=} \tau^2 b''((b')^{-1}(g^{-1}(x_i^\top \beta))), \quad i = 1, \dots, n$$

die Varianz von Y_i ist.

Beweis. 1. Führen wir die Bezeichnung

$$l_i(\beta) = \frac{1}{\tau^2} (Y_i \theta_i + a(Y_i, \tau) - b(\theta_i)), \quad i = 1, \dots, n \text{ ein.}$$

Somit gilt

$$U_j(\beta) = \sum_{i=1}^n \frac{\partial l_i(\beta)}{\partial \beta_j}, \quad j = 1, \dots, m.$$

Durch die mehrfache Anwendung der Kettenregel ergibt sich

$$\frac{\partial l_i(\beta)}{\partial \beta_j} = \frac{\partial l_i(\beta)}{\partial \theta_i} \cdot \frac{\partial \theta_i}{\partial \mu_i} \cdot \frac{\partial \mu_i}{\partial \eta_i} \cdot \frac{\partial \eta_i}{\partial \beta_j}, \quad i = 1, \dots, n, \quad j = 1, \dots, m$$

Da

$$\frac{\partial l_i(\beta)}{\partial \theta_i} = \frac{1}{\tau^2} (Y_i - b'(\theta_i)) \stackrel{\text{Lemma 3.1.2}}{=} \frac{1}{\tau^2} (Y_i - \mu_i(\beta)),$$

$$\frac{\partial \theta_i}{\partial \mu_i} = \left(\frac{\partial \mu_i}{\partial \theta_i} \right)^{-1} = \left((b'(\theta_i))' \right)^{-1} = (b''(\theta_i))^{-1} \stackrel{\text{Lemma 3.1.2}}{=} \left(\frac{\sigma_i^2(\beta)}{\tau^2} \right)^{-1} = \frac{\tau^2}{\sigma_i^2(\beta)},$$

$$\frac{\partial \mu_i}{\partial \eta_i} = \frac{\partial g^{-1}(\eta_i)}{\partial \eta_i}$$

wegen $\mu_i = \mathbb{E}Y_i = g^{-1}(\eta_i)$,

$$\frac{\partial \eta_i}{\partial \beta_j} = \frac{\partial (x_i^\top \beta)}{\partial \beta_j} = x_{ij}, \quad i = 1, \dots, n, \quad j = 1, \dots, m,$$

bekommen wir

$$\begin{aligned} U_j(\beta) &= \frac{1}{\tau^2} \sum_{i=1}^n x_{ij} (Y_i - \mu_i(\beta)) \cdot \frac{\tau^2}{\sigma_i^2(\beta)} \cdot \frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \\ &= \sum_{i=1}^n x_{ij} (Y_i - \mu_i(\beta)) \frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \cdot \frac{1}{\sigma_i^2(\beta)}, \quad j = 1, \dots, m. \end{aligned}$$

2. Für alle $i, j = 1, \dots, m$ gilt:

$$\begin{aligned} I_{ij}(\beta) &= \mathbb{E}(U_i(\beta)U_j(\beta)) = \sum_{k,l=1}^n x_{ki}x_{lj} \underbrace{\text{Cov}(Y_k, Y_l)}_{\delta_{kl}\sigma_k^2(\beta)} \cdot \frac{\partial g^{-1}(\eta_k)}{\partial \eta_k} \frac{\partial g^{-1}(\eta_l)}{\partial \eta_l} \frac{1}{\sigma_k^2(\beta)\sigma_l^2(\beta)} \\ &= \sum_{k=1}^n x_{ki}x_{kj} \left(\frac{\partial g^{-1}(\eta_k)}{\partial \eta_k} \right)^2 \frac{1}{\sigma_k^2(\beta)}. \end{aligned}$$

□

Bemerkung 3.3.1. Im Falle der natürlichen Linkfunktion vereinfachen sich die obigen Gleichungen. So sieht die Log-Likelihood-Funktion folgendermaßen aus:

$$\log L(Y, \beta) = \frac{1}{\tau^2} \sum_{i=1}^n \left(Y_i x_i^\top \beta + a(Y_i, \tau) - b(x_i^\top \beta) \right).$$

Da in diesem Fall $g^{-1}(\eta_i) = b'(\eta_i)$, $\eta_i = x_i^\top \beta = \theta_i$ gilt

$$\frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} = b''(\theta_i) \stackrel{\text{Lemma 3.1.2}}{=} \frac{1}{\tau^2} \sigma_i^2(\beta)$$

und somit

$$U_j(\beta) = \frac{1}{\tau^2} \sum_{i=1}^n x_{ij} (Y_i - \mu_i(\beta)), \quad j = 1, \dots, m,$$

$$I_{jk}(\beta) = \frac{1}{\tau^4} \sum_{i=1}^n x_{ij} x_{ik} \sigma_i^2(\beta), \quad j, k = 1, \dots, m.$$

Satz 3.3.2. Es gilt

$$W_{jk}(\beta) = \sum_{i=1}^n x_{ij} x_{ik} \left((Y_i - \mu_i(\beta)) \nu_i - u_i^2 \frac{1}{\sigma_i^2(\beta)} \right), \quad j, k = 1, \dots, m,$$

wobei

$$u_i = \frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \quad \text{und} \quad \nu_i = \frac{1}{\tau^2} \cdot \frac{\partial^2 ((b')^{-1} \circ g^{-1})(\eta_i)}{\partial \eta_i^2}, \quad i = 1, \dots, n,$$

$$\mu_i(\beta) = \mathbb{E}Y_i, \quad \sigma_i^2(\beta) = \text{Var}Y_i, \quad \eta_i = x_i^\top \beta.$$

Beweis. Für beliebige $j, k = 1, \dots, m$ gilt

$$\begin{aligned} W_{jk}(\beta) &= \frac{\partial}{\partial \beta_k} U_j(\beta) \stackrel{\text{Satz 3.3.1}}{=} \frac{\partial}{\partial \beta_k} \sum_{i=1}^n x_{ij} (Y_i - \mu_i(\beta)) \frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \frac{1}{\sigma_i^2(\beta)} \\ &= \sum_{i=1}^n x_{ij} \left((Y_i - \mu_i(\beta)) \frac{\partial}{\partial \beta_k} \left(\frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \frac{1}{\sigma_i^2(\beta)} \right) - \frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \frac{1}{\sigma_i^2(\beta)} \frac{\partial \mu_i(\beta)}{\partial \beta_k} \right) \\ &= \sum_{i=1}^n \left(x_{ij} (Y_i - \mu_i(\beta)) \frac{\partial}{\partial \beta_k} \left(\frac{\tau^2 b''((b')^{-1}(g^{-1}(\eta_i))) ((b')^{-1} \circ g^{-1})'(\eta_i)}{\tau^2 b''((b')^{-1}(g^{-1}(\eta_i)))} \right) \right. \\ &\quad \left. - \left(\frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \right)^2 \frac{1}{\sigma_i^2(\beta)} x_{ik} \right) \\ &= \sum_{i=1}^n x_{ij} x_{ik} \left((Y_i - \mu_i(\beta)) \nu_i - u_i^2 \frac{1}{\sigma_i^2(\beta)} \right), \end{aligned}$$

wobei

$$\frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \cdot \frac{1}{\sigma_i^2(\beta)} \stackrel{\text{Lemma 3.1.2}}{=} \frac{\partial b'(\theta_i)}{\partial \eta_i} \cdot \frac{1}{\tau^2} \cdot \frac{1}{b''(\theta_i)} = \frac{\partial b'(\theta_i)}{\partial \theta_i} \cdot \frac{\partial \theta_i}{\partial \eta_i} \frac{1}{\tau^2} \frac{1}{b''(\theta_i)} = \frac{1}{\tau^2} \frac{\partial \theta_i}{\partial \eta_i}$$

und

$$\frac{\partial}{\partial \beta_k} \left(\frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \cdot \frac{1}{\sigma_i^2(\beta)} \right) = \frac{1}{\tau^2} \frac{\partial^2 \theta_i}{\partial \eta_i^2} \cdot \frac{\partial \eta_i}{\partial \beta_k} \stackrel{\eta_i = x_i^\top \beta}{=} \frac{1}{\tau^2} \frac{\partial^2 \theta_i}{\partial \eta_i^2} \cdot x_{ik},$$

dabei ist

$$\frac{\overbrace{\frac{\partial g^{-1}(\eta_i)}{\partial \beta_k}}^{\mu_i(\beta)}}{\partial \beta_k} = \frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \cdot \frac{\partial \eta_i}{\partial \beta_k} = \frac{\partial g^{-1}(\eta_i)}{\partial \eta_i} \cdot x_{ik}$$

und $\theta_i = (b')^{-1} \circ g^{-1}(\eta_i)$, $i = 1, \dots, n$. □

Für verallgemeinerte lineare Modelle mit natürlichen Linkfunktionen gilt insbesondere

$$W(\beta) = -I(\beta) = -\frac{1}{\tau^4} \sum_{i=1}^n x_{ij} x_{ik} \sigma_i^2(\beta), \quad (3.3.3)$$

weil in diesem Fall $\nu_i = 0$ für alle $i = 1, \dots, n$. $W(\beta)$ ist also deterministisch. Tatsächlich ist nach Lemma 3.2.1 $\theta_i = x_i^\top \beta = \eta_i$ und somit $\frac{\partial^2 \theta_i}{\partial \eta_i^2} = 0$, $i = 1, \dots, n$.

Aus Bemerkung 3.3.1 außerdem: $u_i^2 = \frac{1}{\tau^4} \sigma_i^4(\beta)$.

Beispiel 3.3.1. Wie sehen $U(\beta)$, $I(\beta)$ und $W(\beta)$ für unsere Modelle aus Beispiel 2.6.2 (natürliche Linkfunktionen) aus?

1. **Normalverteilung:** dieser Fall entspricht der üblichen multivariaten linearen Regression mit normalverteilten Störgrößen. In diesem Fall gilt $\mu = X\beta$, $\tau^2 = \sigma^2$.

Aus Bemerkung 3.3.1 folgt

$$\begin{aligned} U(\beta) &= \frac{1}{\sigma^2} X^\top (Y - X\beta), \\ I(\beta) &= (\mathbb{E} (U_i(\beta) \cdot U_j(\beta)))_{i,j=1,\dots,m} = \frac{1}{\sigma^2} X^\top X, \\ W(\beta) &= -I(\beta). \end{aligned}$$

2. **Logistische Regression:** hier gilt $\tau^2 = 1$, $\mu_i = p_i$, $\sigma_i^2 = p_i(1 - p_i)$, $i = 1, \dots, n$, $p_i \in (0, 1)$ und somit

$$\begin{aligned} U(\beta) &= X^\top(Y - p) , \\ I(\beta) &= X^\top \text{diag}(p_i(1 - p_i))X , \\ W(\beta) &= -I(\beta) , \end{aligned}$$

wobei $p = (p_1, \dots, p_n)^\top$.

3. **Poisson-Regression:** es gilt $\tau^2 = 1$, $\mu_i = \lambda_i = \sigma_i^2$, $i = 1, \dots, n$ und somit

$$\begin{aligned} U(\beta) &= X^\top(Y - \lambda) , \\ I(\beta) &= X^\top \text{diag}(\lambda_i)X , \\ W(\beta) &= -I(\beta) , \end{aligned}$$

wobei $\lambda = (\lambda_1, \dots, \lambda_n)^\top$.

Wann ist die Lösung des Gleichungssystems $U(\beta) = 0$ auch ein Maximum-Punkt der Funktion $\log L(Y, \beta)$?

Mit anderen Worten: Wann existiert der ML-Schätzer $\hat{\beta}$ von β , der eindeutig bestimmt ist?

$$\hat{\beta} = \underset{\beta}{\operatorname{argmax}} \log L(Y, \beta)$$

An der hinreichenden Bedingung eines Maximums folgt, dass die Hesse-Matrix $W(\beta)$ negativ definit sein muss.

Betrachten wir den Spezialfall der natürlichen Linkfunktion.

Dann gilt nach Bemerkung 3.3.1:

- Das Gleichungssystem $U(\beta) = 0$ schreibt sich $U(\beta) = \frac{1}{\tau^2} X^\top(Y - \mu(\beta)) = 0$
- Die Matrix $W(\beta) = -\frac{1}{\tau^4} X^\top \text{diag}(\sigma_i^2(\beta))X$ ist negativ definit, falls zusätzlich $rg(X) = m$ und $0 < \sigma_i^2(\beta) < \infty$ für alle $i = 1, \dots, n$.

Unter diesen Bedingungen existiert also ein eindeutiger ML-Schätzer $\hat{\beta}$ für β .

Geben wir jetzt Verfahren an, die das (im Allgemeinen nicht lineare) Gleichungssystem $U(\beta) = 0$ numerisch lösen. Diese Ansätze sind iterativ, d.h. sie nähern sich schrittweise dem ML-Schätzer $\hat{\beta}$ an.

1. Newton-Verfahren

Wähle einen geeigneten Startwert $\hat{\beta}_0 \in \mathbb{R}^m$.

Im Schritt $k + 1$, berechne $\hat{\beta}_{k+1}$ aus $\hat{\beta}_k$, $k = 0, 1, \dots$ auf folgende Art und Weise:

- Nimm die Taylor-Entwicklung von $U(\beta)$ bis zur ersten Ordnung an der Stelle $\hat{\beta}_k$: $U(\beta) \approx U(\hat{\beta}_k) + W(\hat{\beta}_k)(\beta - \hat{\beta}_k)$.
- Setze sie gleich Null: $U(\hat{\beta}_k) + W(\hat{\beta}_k)(\beta - \hat{\beta}_k) = 0$
- Die Lösung dieses Gleichungssystems ist $\hat{\beta}_{k+1}$:

$$\hat{\beta}_{k+1} = \hat{\beta}_k - W^{-1}(\hat{\beta}_k) \cdot U(\hat{\beta}_k), \quad k = 0, 1, 2, \dots,$$

vorausgesetzt, dass $W(\hat{\beta}_k)$ invertierbar ist.

Breche den Iterationsprozess ab, sobald $|\hat{\beta}_{k+1} - \hat{\beta}_k| < \delta$ für eine vorgegebene Genauigkeit $\delta > 0$ ist.

Das Konvergenzverhalten dieses Verfahrens hängt entscheidend von der Wahl von $\hat{\beta}_0$ ab, für dessen Konvergenz $\hat{\beta}_0$ nah genug bei $\hat{\beta}$ liegen muss. Ein weiterer Nachteil dieses Verfahrens ist, dass die zufällige Matrix $W(\beta)$ unter Umständen nicht invertierbar ist. Deswegen schlagen wir jetzt eine Modifikation des Newton-Verfahrens vor, bei der $W(\beta)$ durch den Erwartungswert

$$\mathbb{E} W(\beta) = -I(\beta) \tag{3.3.4}$$

ersetzt wird. Dass die Identität (3.3.3) stimmt, folgt aus dem Satz 3.3.2, und der Tatsache, dass $\mathbb{E} Y_i = \mu_i$, $i = 1, \dots, n$. Wenn man voraussetzt, dass $rg(X) = m$ und $u_i \neq 0$, $i = 1, \dots, n$, so ist nach Satz 3.3.1 $I(\beta)$ invertierbar. Dieses Verfahren wird *Fisher Scoring* genannt.

Der einzige Unterschied zu den Schritten des Newton-Verfahrens besteht beim Fisher Scoring darin, dass man in Schritt 2 die iterative Gleichung

$$\hat{\beta}_{k+1} = \hat{\beta}_k + I^{-1}(\hat{\beta}_k)U(\hat{\beta}_k), \quad k = 0, 1, \dots$$

einsetzt.

Im Falle einer natürlichen Linkfunktion gilt nach Bemerkung 3.3.1

$$\begin{aligned} \hat{\beta}_{k+1} &= \hat{\beta}_k + \tau^4 \left(X^\top \text{diag}(\sigma_i^2(\hat{\beta}_k)) X \right)^{-1} \frac{1}{\tau^2} \left(X^\top (Y - \mu(\hat{\beta}_k)) \right) \\ &= \hat{\beta}_k + \tau^2 \left(X^\top \text{diag}(\sigma_i^2(\hat{\beta}_k)) X \right)^{-1} \left(X^\top (Y - \mu(\hat{\beta}_k)) \right). \end{aligned}$$

3.4 Asymptotische Tests für β

Das Ziel dieses Abschnittes ist es, eine Testregel für die Hypothese

$$H_0 : \beta = \beta_0 \text{ vs. } H_1 : \beta \neq \beta_0 \quad \text{mit} \quad \beta = (\beta_1, \dots, \beta_m)^\top, \quad \beta_0 = (\beta_{01}, \dots, \beta_{0m})^\top$$

zu konstruieren. Insbesondere sind die Haupthypothesen $H_0 : \beta = 0$ bzw. $H_0 : \beta_j = 0$ von Interesse, weil sie die Tatsache reflektieren, dass die Zielvariablen $Y = (Y_1, \dots, Y_n)^\top$ von einigen Ausgangsvariablen (z.B. $(x_{1j}, \dots, x_{nj})^\top$ im Falle der Hypothese $\beta_j = 0$) unabhängig sind.

Um solche Hypothesen testen zu können, werden Teststatistiken T_n vorgeschlagen, die asymptotisch (für $n \rightarrow \infty$) eine bekannte Prüfverteilung (z.B. multivariate Normalverteilung oder χ^2 -Verteilung) besitzen. Dafür sind gewisse Vorarbeiten notwendig.

Sei

$$g(\mathbb{E} Y_i) = X_i \beta, \quad i = 1, \dots, n,$$

ein verallgemeinertes lineares Modell mit natürlicher Linkfunktion g . Seien $L(Y, \beta)$, $U(\beta)$ und $I(\beta)$ die Likelihood-Funktion, der Vektor der partiellen Ableitungen von $\log L(Y, \beta)$ bzw. die Fisher-Informationsmatrix in diesem Modell.

Durch $\hat{\beta}_n = \hat{\beta}(Y_1, \dots, Y_n, X)$ bezeichne man eine Folge von Schätzern für β .

Es gelten folgende Voraussetzungen:

1. \exists Kompaktum $K \subset \mathbb{R}^m$, so dass alle Zeilen X_i , $i = 1, \dots, n$, $n \in \mathbb{N}$, von X in K liegen. Dabei soll $\theta = x^\top \beta \in \Theta$ für alle $\beta \in \mathbb{R}^m$ und $x \in K$.
2. Es existiert eine Folge $\{\Gamma_n\}_{n \in \mathbb{N}}$ von diagonalen $(m \times m)$ -Matrizen $\Gamma_n = \Gamma_n(\beta)$ mit positiven Diagonalelementen und den Eigenschaften $\lim_{n \rightarrow \infty} \Gamma_n = 0$, $\lim_{n \rightarrow \infty} \Gamma_n^\top I_n(\beta) \Gamma_n = K^{-1}(\beta)$ gleichmäßig in β , wobei $K(\beta)$ eine symmetrische positiv definite $(m \times m)$ -Matrix ist, $\forall \beta \in \mathbb{R}^m$.

Satz 3.4.1. Unter obigen Voraussetzungen gilt:

es existiert eine Γ_n -Konsistente Folge von ML-Schätzern $\{\hat{\beta}_n\}$ für β ,
(d.h. $\mathbb{P}(\Gamma_n^{-1}|\hat{\beta}_n - \beta| \leq \varepsilon, U(\hat{\beta}_n) = 0) \rightarrow 1$ für $n \rightarrow \infty$), so dass

1. $T_n^* = \Gamma_n^{-1}(\hat{\beta}_n - \beta) \xrightarrow[n \rightarrow \infty]{d} \mathcal{N}(0, K(\beta))$ und
2. $T_n = 2(\log L(Y, \hat{\beta}_n) - \log L(Y, \beta)) \xrightarrow[n \rightarrow \infty]{d} \chi_m^2$, $m = \dim \beta$

Bemerkung 3.4.1. (vgl. [15], S.288-292)

1. Oft wählt man $\Gamma_n = \text{diag}\left(\frac{1}{\sqrt{n}}, \dots, \frac{1}{\sqrt{n}}\right)$
2. Bisher wurde stets angenommen, dass der Störparameter τ^2 bekannt ist. Falls es nicht der Fall ist, kann τ^2 durch

$$\hat{\tau}^2 = \frac{1}{n-m} \sum_{i=1}^n \frac{(Y_i - \mu_i(\hat{\beta}_n))^2}{b''(\hat{\theta}_{ni})}$$

geschätzt werden, wobei $\hat{\theta}_{ni} = (b')^{-1}(\mu_i(\hat{\beta}_n))$, $i = 1, \dots, n$ ist. Dieser Schätzer ist ein empirisches Analogon der Gleichung $\tau^2 = \frac{\text{Var}Y_i}{b''(\theta_i)}$ aus Lemma 3.1.2.

3. Die Aussage 2. des Satzes 3.4.1 gilt auch, wenn man den unbekannt Parameter τ^2 durch einen konsistenten Schätzer τ_n^2 ersetzt.

Wie verwendet man nun den Satz 3.4.1 zum Testen der Hypothesen

$$H_0 : \beta = \beta_0 \quad \text{vs.} \quad H_1 : \beta \neq \beta_0 ,$$

oder komponentenweise

$$H_0 : \beta_j = \beta_{j0} , \quad j = 1, \dots, m \quad \text{vs.} \quad H_1 : \exists j_1 : \beta_{j_1} \neq \beta_{j_1 0} \quad ?$$

Sei

$$g(\mathbb{E}Y_i) = \sum_{j=1}^m x_{ij}\beta_j , \quad i = 1, \dots, n ,$$

ein verallgemeinertes lineares Modell mit natürlicher Linkfunktion g .

Nach Bemerkung 3.3.1 gilt

$$\log L(Y, \beta) = \frac{1}{\tau^2} \sum_{i=1}^n \left(Y_i x_i^\top \beta + a(Y_i, \tau) - b(x_i^\top \beta) \right)$$

wobei $Y = (Y_1, \dots, Y_n)^\top$ und $x_i = (x_{i1}, \dots, x_{im})^\top$. Deshalb gilt

$$T_n = \frac{2}{\tau^2} \sum_{i=1}^n \left(Y_i x_i^\top (\hat{\beta}_n - \beta_0) - b(x_i^\top \hat{\beta}_n) + b(x_i^\top \beta_0) \right)$$

Bei Vorgabe eines Exponential-Modells (τ, b - bekannt), der Stichprobe der Zielvariablen Y und der Designmatrix X wird H_0 verworfen, falls $T_n > \chi_{m,1-\alpha}^2$, wobei m die Anzahl der Parameter im Modell, $\chi_{m,1-\alpha}^2$ das $(1 - \alpha)$ -Quantil der χ_m^2 -Verteilung und $\alpha \in (0, 1)$ das Signifikanzniveau des asymptotischen Tests ist. Dieser Test ist nur für relativ große n anwendbar. Der Fehler 1. Art hat dabei (für $n \rightarrow \infty$) die asymptotische Wahrscheinlichkeit α . Falls eine einfache Hypothese

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

getestet werden soll, benutzt man die aus der Statistik T_n^* abgeleitete Teststatistik T_n^1 . H_0 wird verworfen, falls

$$|T_n^1| = \frac{|\hat{\beta}_{nj}|}{(\Gamma_n(\hat{\beta}_n))_{jj}} > z_{1-\frac{\alpha}{2}},$$

wobei $z_{1-\frac{\alpha}{2}}$ das $(1 - \frac{\alpha}{2})$ -Quantil der $\mathcal{N}(0, 1)$ -Verteilung ist. Hierbei ist $\{\Gamma_n\}$ so gewählt worden, dass $K(\beta) = Id$ ist, $\forall \beta \in \mathbb{R}^m$. Dies ist ein asymptotischer Test zum Niveau α , weil

$$\begin{aligned} P_{H_0}(|T_n^1| > z_{1-\frac{\alpha}{2}}) &= 1 - P_{H_0}(|T_n^*| \leq z_{1-\frac{\alpha}{2}}) \xrightarrow{n \rightarrow \infty} 1 - \Phi(z_{1-\frac{\alpha}{2}}) + \underbrace{\Phi(-z_{1-\frac{\alpha}{2}})}_{1-\Phi(z_{1-\frac{\alpha}{2}})} \\ &= 1 - \left(1 - \frac{\alpha}{2}\right) + 1 - \left(1 - \frac{\alpha}{2}\right) = \alpha, \end{aligned}$$

wobei

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{t^2}{2}} dt$$

die Verteilungsfunktion der $\mathcal{N}(0, 1)$ -Verteilung ist.

Beispiel 3.4.1. (Kreditrisikoprüfung)

vgl. Fahrmeir, L., Kneib, T., Lang, S. - Regression, S.208ff

Es liegt folgender Datensatz einer süddeutschen Bank aus den 1990er Jahren vor:

Es werden Ergebnisse der Kreditrisikoprüfung von $n = 1000$ Kreditanträgen (ca. 700 gute und 300 schlechte Kredite) analysiert:

Zielvariable $Y_i = \begin{cases} 0, & \text{falls das Darlehen vom Kunden } i \text{ zurückgezahlt wurde} \\ 1, & \text{falls das Darlehen vom Kunden } i \text{ nicht zurückgezahlt wurde} \end{cases}$

Die Designmatrix X enthält folgende Zusatzinformationen über den Kunden:

x_{i1} - Kontoführung des Kontos bei der Bank: $= \begin{cases} 1, & \text{kein Konto} \\ 0, & \text{sonst} \end{cases}$

x_{i2} - Bewertung der Kontoführung: $= \begin{cases} 1, & \text{gutes Konto} \\ 0, & \text{kein oder schwaches Konto} \end{cases}$

x_{i3} - Laufzeit des Kredits in Monaten

x_{i4} - Höhe des Kredits in DM

x_{i5} - Zahlungsverhalten beim Kunden: $= \begin{cases} 1, & \text{gut} \\ 0, & \text{sonst} \end{cases}$

x_{i6} - Verwendungszweck: $= \begin{cases} 1, & \text{privat} \\ 0, & \text{geschäftlich} \end{cases}$

Frage: Wie soll $\hat{\beta}$ geschätzt werden?

Als Modell wird das Logit-Modell gewählt mit $p_i = P(Y_i = 1)$, $i = 1, \dots, n$:

		$Y = 1$	$Y = 0$
x_1	kein Konto	45.0	20.0
x_2	gut	15.3	49.8
	schlecht	39.7	30.2
x_4	Kredithöhe	$Y = 1$	$Y = 0$
	$0 < \dots \leq 500$	1.00	2.14
	$500 < \dots \leq 1000$	11.33	9.14
	$1000 < \dots \leq 1500$	17.00	19.86
	$1500 < \dots \leq 2500$	19.67	24.57
	$2500 < \dots \leq 5000$	25.00	28.57
	$5000 < \dots \leq 7500$	11.33	9.71
	$7500 < \dots \leq 10000$	6.67	3.71
	$10000 < \dots \leq 15000$	7.00	2.00
	$15000 < \dots \leq 20000$	1.00	0.29
x_5	Frühere Kredite	$Y = 1$	$Y = 0$
	gut	82.33	94.95
	schlecht	17.66	5.15
x_6	Verwendungszweck	$Y = 1$	$Y = 0$
	privat	57.53	69.29
	beruflich	42.47	30.71

Tabelle 3.1: Auszug aus dem Originaldatensatz

\bar{x}_1	\bar{x}_2	\bar{x}_3	\bar{x}_4	\bar{x}_5	\bar{x}_6
0.274	0.393	20.903	3271	0.911	0.657

Tabelle 3.2: Mittelwerte \bar{x}_j von x_{ij} im Datensatz

$$\log \frac{p_i}{1-p_i} = \beta_0 + x_{i1}\beta_1 + x_{i2}\beta_2 + x_{i3}\beta_3 + x_{i4}\beta_4 + x_{i5}\beta_5 + x_{i6}\beta_6 \quad \text{für } i = 1, \dots, n,$$

wobei $\beta = (\beta_0, \dots, \beta_6)^\top$, $m = 7$.

Ziel: Schätze β_0, \dots, β_6 und prüfe, welche Faktoren für die künftige Kreditvergabe relevant sind.

$H_0 : \beta_i = 0$ (Merkmal x_i beeinflusst die Kreditvergabe nicht) wird abgelehnt, falls p-Wert $\leq \alpha$. Man sieht, dass u.a. auch β_4 für die Kreditvergabe nicht relevant ist, was der Intuition widerspricht. Eine Verfeinerung des Modells ist notwendig:

	Wert	$\sqrt{(I_n^{-1}(\hat{\beta}))_{ii}}$	T_n^1	p-Wert
β_0	0.281	0.303	-0.94	0.347
β_1	0.618	0.175	3.53	< 0.001
β_2	-1.338	0.201	-6.65	< 0.001
β_3	0.033	0.008	4.29	< 0.001
β_4	0.023	0.033	0.72	0.474
β_5	-0.986	0.251	-3.93	< 0.001
β_6	-0.426	0.266	-2.69	0.007

Tabelle 3.3: Ergebnis zur ML-Schätzung durch das Fisher Scoring Verfahren, wobei $\sqrt{(I_n^{-1}(\hat{\beta}))_{ii}}$ als asymptotische Standardabweichung von $\hat{\beta}_i$ interpretiert wird. Signifikanzniveau: $\alpha = 0.001$

Neues Modell:

$$g(\mathbb{E} Y_i) = \beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \beta_3^1 x_{i3} + \beta_3^2 x_{i3}^2 + \beta_4^1 x_{i4} + \beta_4^2 x_{i4}^2 + \beta_5 x_{i5} + \beta_6 x_{i6}$$

	Wert	$\sqrt{(I_n^{-1}(\hat{\beta}))_{ii}}$	T_n^1	p-Wert
β_0	-0.488	0.390	-1.25	0.211
β_1	0.618	0.176	3.51	< 0.001
β_2	-1.337	0.202	-6.61	< 0.001
β_3^1	0.092	0.025	3.64	< 0.001
β_3^2	-0.001	< 0.001	-2.20	0.028
β_4^1	-0.264	0.099	-2.68	0.007
β_4^2	0.023	0.007	3.07	0.002
β_5	-0.995	0.255	-3.90	< 0.001
β_6	-0.404	0.160	-2.52	0.012

Tabelle 3.4: p -Werte für die Regressionskoeffizienten des neuen Modells

Frage: Welches Modell ist besser?

Mit anderen Worten, wir testen

$$H_0 : \beta_3^2 = 0 \text{ (lineares Modell) vs. } H_1 : \beta_3^2 \neq 0 \text{ (quadratisches Modell) bzw.}$$

$$H_0 : \beta_4^2 = 0 \text{ (lineares Modell) vs. } H_1 : \beta_4^2 \neq 0 \text{ (quadratisches Modell) .}$$

Dabei verallgemeinern wir die Art der statistischen Hypothesen wie folgt: es wird

$$H_0 : C\beta = d \text{ vs. } H_1 : C\beta \neq d$$

getestet, wobei C eine $(r \times m)$ -Matrix mit $rg C = r \leq m$ ist und $d \in \mathbb{R}^r$.

Zum Vergleich: früher haben wir

$$H_0 : \beta = \beta_0 \text{ vs. } H_1 : \beta \neq \beta_0, \quad \beta, \beta_0 \in \mathbb{R}^m$$

getestet. Natürlich ist $\beta = \beta_0$ ein Spezialfall von $C\beta = d$ mit $C = \text{Id}$, $d = \beta_0$. Die neuen Hypothesen beinhalten Aussagen über die Linearkombinationen der Parameterwerte. Wie soll H_0 vs. H_1 getestet werden?

Sei $\tilde{\beta}_n$ der ML-Schätzer von β unter H_0 , d.h. $\tilde{\beta}_n = \underset{\beta \in \mathbb{R}^m: C\beta=d}{\text{argmax}} \log L(Y, \beta)$

Sei $\hat{\beta}_n$ der ML-Schätzer von β unrestringiert, d.h. $\hat{\beta}_n = \underset{\beta \in \mathbb{R}^m}{\text{argmax}} \log L(Y, \beta)$.

Die Idee der folgenden Tests ist es, $\tilde{\beta}_n$ mit $\hat{\beta}_n$ zu vergleichen. Falls die Abweichung $\hat{\beta}_n - \tilde{\beta}_n$ groß ist, soll H_0 abgelehnt werden.

Satz 3.4.2. Sei $\log L(Y, \beta)$ die Log-Likelihood-Funktion der Stichprobe der Zielvariablen $Y = (Y_1, \dots, Y_n)^\top$, $I_n(\beta)$ die Fisher-Informationsmatrix, $U(\beta)$ die Score-Funktion des verallgemeinerten linearen Modells mit natürlicher Linkfunktion

$$g : g(\mathbb{E} Y_i) = X_i \beta, \quad i = 1, \dots, n.$$

Wir führen folgende Teststatistiken ein:

1. **Likelihood-Ratio-Teststatistik:**

$$\tilde{T}_n = 2(\log L(Y, \hat{\beta}_n) - \log L(Y, \tilde{\beta}_n))$$

2. **Wald-Statistik:**

$$\tilde{T}_n^* = (C\hat{\beta}_n - d)^\top (CI_n^{-1}(\hat{\beta}_n)C^\top)^{-1} (C\hat{\beta}_n - d)$$

3. **Score-Statistik:**

$$\bar{T}_n^* = U(\tilde{\beta}_n)^\top I_n^{-1}(\tilde{\beta}_n) U(\tilde{\beta}_n)$$

Unter gewissen Bedingungen an die Schätzer $\hat{\beta}$ und $\tilde{\beta}$ (vgl. Satz 3.4.1) sind die Teststatistiken 1 - 3 asymptotisch χ_m^2 -verteilt: z.B. gilt für die Likelihood-Ratio-Teststatistik

$$\tilde{T}_n \xrightarrow[n \rightarrow \infty]{d} \chi_m^2.$$

Folgerung 3.4.1. Der Satz 2.6.4 liefert uns folgende Entscheidungsregel: H_0 wird abgelehnt, falls

$$\tilde{T}_n(\tilde{T}_n^*, \bar{T}_n) > \chi_{m, 1-\alpha}^2.$$

Dies ist ein asymptotischer Test zum Signifikanzniveau α .

Beispiel 3.4.2 (Fortsetzung). Es ergeben sich folgende Werte für die Teststatistiken:

$$\tilde{T}_n = 12.44, \quad \text{p-Wert: } 0.0020$$

$$\tilde{T}_n^* = 11.47, \quad \text{p-Wert: } 0.0032.$$

Für $\alpha = 0.005$ gilt p-Wert $\leq \alpha$, somit wird $H_0 : \beta_4^2 = 0$ abgelehnt \Rightarrow das quadratische verallgemeinerte lineare Modell ist besser.

3.5 Kriterien zur Modellwahl bzw. Modellanpassung

Es ist bekannt, dass die Güte der Anpassung eines parametrischen Modells an die Daten im Allgemeinen steigt, wenn die Anzahl der Parameter erhöht wird. Die Aufgabe eines Statistikers ist es aber, ein gut passendes Modell mit einer möglichst kleinen Anzahl an Parametern zu finden. Deshalb verwendet man folgendes Informationskriterium von Akaike, um Modelle mit (möglicherweise) unterschiedlichen Parametersätzen zu vergleichen.

Informationskoeffizient von Akaike:

$$\text{AIC} = -2 \log L(Y, \hat{\beta}) + 2m ,$$

wobei $Y = (Y_1, \dots, Y_n)$ die Stichprobe der Zielvariablen im verallgemeinerten linearen Modell und $\hat{\beta}$ der dazugehörige ML-Schätzer sei. Der Wert von AIC berücksichtigt einerseits die Forderung der Maximalität der Log-Likelihood-Funktion $\log L(Y, \hat{\beta})$, andererseits bestraft er Modelle mit einer großen Anzahl von Parametern m . Das Modell mit dem kleineren AIC ist als besseres Modell einzustufen. Manchmal verwendet man statt AIC den normierten Koeffizienten AIC/n .

Beispiel 3.5.1 (Fortsetzung). Berechnen wir den Informationskoeffizienten von Akaike für das lineare und quadratische Logit-Modell im Beispiel der Kreditrisikoprüfung:

$$\text{Lineares Modell : AIC} = 1043.815$$

$$\text{Quadratisches Modell : AIC} = 1035.371$$

Man sieht anhand des AIC, dass die Wahl zu Gunsten des quadratischen Modells ausfällt.

Der Nachteil der oben beschriebenen AIC-Regel liegt darin, dass die endgültige Entscheidung dem Statistiker überlassen bleibt. Deshalb ist es wünschenswert, einen statistischen Test zu konstruieren, der die Güte der Modellanpassung beurteilen kann.

Wir werden jetzt den χ^2 -Test beschreiben.

Sei

$$g(\mathbb{E} Y_i) = X_i \beta , \quad i = 1, \dots, n ,$$

ein verallgemeinertes lineares Modell mit Linkfunktion g und Parametervektor $\beta = (\beta_1, \dots, \beta_m)^\top$. Teilen wir die Zielvariablen Y_1, \dots, Y_n in k Gruppen auf, so dass sie möglichst homogen in Bezug auf die zu schätzenden Parameter sind. So liegt z.B. eine solche Aufteilung vor, wenn der Wertebereich der Zielvariablen Y_i „geschickt“ in $k > m$ ¹ Intervalle $(a_l, b_l]$ unterteilt wird:

$$-\infty \leq a_1 < b_1 = a_2 < b_2 = a_3 < \dots < b_{k-1} = a_k < b_k \leq +\infty$$

¹ $k \leq m \Rightarrow D \xrightarrow[n \rightarrow \infty]{d} \underbrace{\chi_{k-m-1}^2}_{<0}$

In die Gruppe l fallen alle Beobachtungen Y_i , die zu $(a_l, b_l]$ gehören. Dabei müssen $(a_l, b_l]$ so gewählt werden, dass $\hat{\mu}_j = g^{-1}(X_j \hat{\beta})$ innerhalb einer Gruppe konstant wird: $\hat{\mu}_j \equiv \hat{\mu}_l \forall j$ aus Gruppe l .² Sei

- $n_l = \# \{Y_j : Y_j \in (a_l, b_l]\}$ die Klassenstärke der Klasse l
- $\bar{Y}_l = \frac{1}{n_l} \sum Y_j$ das arithmetische Mittel innerhalb der Klasse l
- $\hat{\beta}$ der ML-Schätzer von β , der aus Y gewonnen wurde
- $l_l(\beta) = \sum \log f_\theta(Y_j)$ die Log-Likelihood-Funktion der Zielvariablen Y_i innerhalb der Gruppe l
- $\hat{\mu}_l = g^{-1}(X_l \hat{\beta})$ und $v(\hat{\mu}_l)$ der Erwartungswert- bzw. der Varianzschätzer von $\mu_l = \mathbb{E} Y_l$, die aus dem ML-Schätzer $\hat{\beta}$ gewonnen wurden

Dabei ist $v(\hat{\mu}_l) = \tau^2 b''(b'^{-1}(\hat{\mu}_l))$, wobei $b(\cdot)$ der entsprechende Koeffizient in der Dichte f_θ aus der Exponentialfamilie ist. Man bildet folgende Teststatistiken:

$$\chi^2 = \sum_{l=1}^k \frac{(\bar{Y}_l - \hat{\mu}_l)^2}{v(\hat{\mu}_l)/n_l}$$

$$D = -2\tau^2 \sum_{l=1}^k (l_l(\hat{\mu}_l) - l_l(\bar{Y}_l))$$

Satz 3.5.1.

Falls $n \rightarrow \infty$ und die Anzahl $n_l \rightarrow \infty \forall l$, dann gilt unter gewissen Voraussetzungen Folgendes:

$$\chi^2 \xrightarrow[n \rightarrow \infty]{d} \chi_{k-m-1}^2$$

$$D \xrightarrow[n \rightarrow \infty]{d} \chi_{k-m-1}^2$$

²Dies ist eine informelle Beschreibung des Vorgangs, bei dem für jedes Y_i n_i unabhängige Kopien von Y_i erzeugt werden, die die i -te Klasse bilden.

Folgerung 3.5.1.

Mit Hilfe der Behauptungen des Satzes 2.6.5 können die Hypothesen

$$H_0 : Y = (Y_1, \dots, Y_n) \text{ stammt aus dem Modell } g(\mathbb{E} Y_i) = X_i \beta, \quad i = 1, \dots, n$$

vs.

$$H_1 : Y = (Y_1, \dots, Y_n) \text{ stammt nicht aus dem Modell } g(\mathbb{E} Y_i) = X_i \beta, \quad i = 1, \dots, n$$

folgendermaßen getestet werden:

H_0 wird (für große n) zum asymptotischen Signifikanzniveau α verworfen, falls

$$\chi^2 > \chi_{k-m-1, 1-\alpha}^2 \quad \text{bzw.} \quad D > \chi_{k-m-1, 1-\alpha}^2.$$

Diese Tests sollten aber nicht verwendet werden, falls die Klassenstärken n_l klein sind.

Beispiel 3.5.2.

Wie sehen die oben beschriebenen Tests im Falle der Logit- bzw. Poisson-Regression aus?

1. **Logit-Modell:** $Y_i \sim \text{Bernoulli}(p_i)$, $i = 1, \dots, n$

$$\Rightarrow \text{verallgemeinertes lineares Modell} \quad \log \frac{p_i}{1-p_i} = X_i \beta, \quad i = 1, \dots, n$$

Wir teilen Y_1, \dots, Y_n in k Klassen auf, so dass die Wahrscheinlichkeit des Auftretens von 1 in jeder Klasse möglichst gut durch $\bar{Y}_l = \frac{1}{n_l} \sum Y_i$ geschätzt wird. Somit gilt mit $\hat{\mu}_l = \hat{p}_l = g^{-1}(X_l \hat{\beta}) = \frac{e^{X_l^\top \hat{\beta}}}{1 + e^{X_l^\top \hat{\beta}}}$, $v(\hat{p}_l) = \hat{p}_l(1 - \hat{p}_l)$

$$\Rightarrow \chi^2 = \sum_{l=1}^k \frac{(\bar{Y}_l - \hat{p}_l)^2}{\hat{p}_l(1 - \hat{p}_l)/n_l}$$

2. **Poisson-Modell:** $Y_i \sim \text{Poisson}(\lambda)$,

$$\Rightarrow \text{verallgemeinertes lineares Modell} \quad \log \lambda_i = X_i \beta, \quad i = 1, \dots, n$$

Somit gilt mit $\hat{\mu}_l = \hat{\lambda}_l = e^{X_l \hat{\beta}}$, $v(\hat{\lambda}_l) = \hat{\lambda}_l$

$$\Rightarrow \chi^2 = \sum_{l=1}^k \frac{(\bar{Y}_l - \hat{\lambda}_l)^2}{\hat{\lambda}_l/n_l}$$

4 Hauptkomponentenanalyse

In diesem Kapitel werden Methoden zur Reduktion der Komplexität von sehr großen statistischen Datensätzen vorgestellt, die als Hauptkomponentenanalyse (HKA) bekannt sind (engl. Principal Component Analysis, PCA). Mit ihrer Hilfe ist es möglich einen sehr hochdimensionalen Datensatz $X = (X_1, \dots, X_n)^T \in \mathbb{R}^n$ auf wenige wirklich wichtige Komponenten $\varphi = AX \in \mathbb{R}^d$ zurückzuführen, $d \ll n$, die aber dabei die meiste Variabilität des originalen Datensatzes X beibehalten. A ist dabei eine $(d \times n)$ -Matrix, die zu finden ist, wenn gewisse (in 4.2.1 angegebene) Nebenbedingungen erfüllt sind. Andere Beispiele von Anwendungen sind Visualisierung von komplexen Datensätzen, Ausreißer-Erkennung, Cluster-Analyse u.s.w.. Für eine Übersicht siehe z.B. [8].

4.1 Einführung

Um nachfolgende Problemstellungen zu motivieren, betrachten wir ein Beispiel des Text Mining aus der Autoindustrie:

Beispiel 4.1.1. Ein Autohersteller ist daran interessiert, seine Verluste, die in Folge von Betrug und Inkompetenz seitens seiner Niederlassungen bei Garantie-Reparaturen auftreten, zu minimieren. Deshalb möchte er eine Auffälligkeitsanalyse von Reparaturbesichtigungen aus Garantie-Werkstätten betreiben, die dazu führen sollte, computergestützt, verdächtige Meldungen zu finden, die nachher manuell und einzeln weiter geprüft werden. Ein weiterer Anreiz für die automatischen Früherkennung von Auffälligkeiten besteht darin, dass flächendeckende Prüfungen nur für wenige Niederlassungen und in unregelmäßigen Zeitabständen (aus Kostengründen) möglich sind, und selbst die könnte man sich sparen. Ein typischer Text, der eine Garantie-Reparatur beschreibt, verwendet maximal 300.000 Wörter aus einem Fachwortschatz. Daher werden solche Texte als Vektoren $x = (x_1, \dots, x_n)^T$ der Länge $n = 300.000$ dargestellt, wobei

$$x_i = \begin{cases} 1 & , \text{ falls das Wort } i \text{ im Text } x \text{ vorkommt} \\ 0 & , \text{ sonst} \end{cases}$$

Diese Vektoren x werden normiert, so dass sie auf der Sphäre S^{n-1} liegen. Innerhalb eines Jahres entsteht dadurch eine riesige Datenbank solcher Vektoren x mit mehreren Millionen Einträgen. Die Aufgabe eines Statistikers besteht in der drastischen Reduktion der Dimension $n - 1$ des Datensatzes, so dass eine Visualisierung des Datensatzes möglich wird. Eine mögliche Lösung liegt in der Verwendung von HKA. Die HKA geht in ihren Ursprüngen auf die Arbeiten von Beltran (1873) und Jordan (1874) zurück, die

die Single Value Decomposition verwendeten. In der mehr oder minder modernen Form (vgl. 4.2.1) erscheint sie erst in den Arbeiten von K. Pearson (1901) und H. Hotelling (1933). Auch der Name HKA stammt von Hotelling. Eine Weiterentwicklung der Methoden ist Girshick (1939), Anderson (1963), Rao (1964) und anderen zu verdanken. Erst nach der Einführung der PCs ist aber diese Methodologie richtig angewandt geworden. Denn ohne Computer ist die Berechnung von Hauptkomponenten für $n > 4$ sehr schwierig. Seit den 1980er Jahren gibt es einen rasanten Anstieg der Anwendungen von HKA in allen Wissensbereichen (vor allem in Ingenieurwissenschaften), wo multivariate Datensätze analysiert werden sollen.

4.2 Hauptkomponentenanalyse auf Modellebene

In diesem Abschnitt wollen wir das Hauptproblem der HKA für Zufallsstichproben $X = (X_1, \dots, X_n)^T$ mit bekannter Kovarianzstruktur einführen. Sei $X = (X_1, \dots, X_n)^T$ eine Zufallsstichprobe von Zufallszahlen X_i mit bekannter Kovarianzmatrix Σ und $\text{Var} X_i \in (0, \infty)$, $i = 1, \dots, n$. Seien $\lambda_1 > \lambda_2 > \dots > \lambda_n > 0$ die Eigenwerte von Σ , die in absteigender Reihenfolge geordnet und alle von einander verschieden sind. Wir suchen Linearkombinationen $\alpha^T X$ von X_i , die die maximale Varianz besitzen, wobei der Vektor α entsprechend normiert ist z.B., so dass $\alpha \in S^{n-1}$ in der Euklidischen Norm.

Definition 4.2.1. Die Linearkombination $\alpha_i^T X$, $i = 1, \dots, n$, heißt i -te Hauptkomponente von X , falls sie die maximale Varianz besitzt unter der Bedingung, dass $\alpha_i \in S^{n-1}$ und $\alpha_1^T X, \alpha_2^T X, \dots, \alpha_{i-1}^T X$ und $\alpha_i^T X$ unkorreliert sind:

$$\begin{cases} \text{Var } \alpha^T X \rightarrow \max_{\alpha} \\ |\alpha| = 1 \\ \text{Cov}(\alpha^T X, \alpha_j^T X) = 0, \quad j = 1, \dots, i-1 \end{cases} \quad (4.2.1)$$

Dabei heißt α_i der Koeffizientenvektor der i -ten Hauptkomponente $\alpha_i^T X$.

Satz 4.2.1. Die i -te Hauptkomponente von X ist gegeben durch

$$Y_i = \alpha_i^T X,$$

wobei α_i der Eigenvektor von Σ mit Eigenwert λ_i ist. Dabei gilt

$$\text{Var}(Y_i) = \lambda_i, \quad i = 1, \dots, n.$$

Beweis. Zeigen wir, dass die Aussage des Satzes gilt für $i = 1, 2$. Für $i > 2$ ist der Beweis analog.

Für $i = 1$ gibt es eine Nebenbedingung $|\alpha| = 1$ in (4.2.1), die in die Lagrange-Zielfunktion

$$f(\alpha) = \text{Var}(\alpha^T X) + \lambda(|\alpha|^2 - 1)$$

übernommen wird. Dabei gilt

$$\begin{aligned}\text{Var}(\alpha^T X) &= \mathbb{E}(\alpha^T X - \mathbb{E}\alpha^T X)^2 = \mathbb{E}(\alpha^T (X - \mathbb{E}X))^2 = \mathbb{E}\alpha^T (X - \mathbb{E}X)(X - \mathbb{E}X)^T \alpha \\ &= \alpha^T \mathbb{E}(X - \mathbb{E}X)(X - \mathbb{E}X)^T \alpha = \alpha^T \Sigma \alpha,\end{aligned}$$

$|\alpha|^2 = \alpha^T \cdot \alpha$, und $f(\alpha) = \alpha^T \Sigma \alpha + \lambda(\alpha^T \alpha - 1)$.

Die notwendige Bedingung des Maximums ist

$$\frac{\partial f}{\partial \alpha} = 0, \quad \frac{\partial f}{\partial \lambda} = 0,$$

wobei die zweite Gleichung einfach die Nebenbedingung $|\alpha| = 1$ repräsentiert.

$\frac{\partial f}{\partial \alpha} = \left(\frac{\partial f}{\partial \alpha^1}, \dots, \frac{\partial f}{\partial \alpha^n} \right)$, wobei $\alpha = (\alpha^1, \dots, \alpha^n)^T$ und $\frac{\partial f}{\partial \alpha} = 0$ schreibt sich $\Sigma \alpha - \lambda \alpha = 0$ in Vektorform oder $\Sigma \alpha = \lambda \alpha$, was heißt, dass α ein Eigenvektor von Σ mit dem Eigenwert λ ist. Da $\text{Var}(\alpha^T X) = \alpha^T \Sigma \alpha$ maximal sein soll, gilt

$$\text{Var}(\alpha^T X) = \alpha^T \lambda \alpha = \lambda \underbrace{\alpha^T \alpha}_1 = \lambda$$

und $\lambda = \lambda_1 > \lambda_2 > \dots > \lambda_n \Rightarrow \lambda = \lambda_1$ und $\alpha = \alpha_1$.

Für $i = 2$, soll die Maximierungsaufgabe

$$\left\{ \begin{array}{l} \alpha^T \Sigma \alpha \rightarrow \max_{\alpha} \\ \alpha^T \cdot \alpha = 1 \\ \text{Cov}(\alpha_1^T X, \alpha^T X) = 0 \end{array} \right.$$

bezüglich α gelöst werden, wobei

$$\text{Cov}(\alpha_1^T X, \alpha^T X) = \alpha_1^T \Sigma \alpha = \alpha^T \Sigma \alpha_1 = \alpha^T \lambda_1 \alpha_1 = \lambda_1 \alpha^T \alpha_1.$$

Das heißt, folgende Funktion soll maximiert werden:

$$f(\alpha) = \alpha^T \Sigma \alpha + \lambda(\alpha^T \alpha - 1) + \delta \alpha^T \alpha_1.$$

Genau wie oben bekommt man

$$\frac{\partial f}{\partial \alpha} = \Sigma \alpha + \lambda \alpha + \delta \alpha_1 = 0$$

Durch die Nebenbedingungen $\alpha_1^T \Sigma \alpha = 0$ und $\alpha_1^T \alpha = 0$ (siehe oben) bekommt man

$$\alpha_1^T \frac{\partial f}{\partial \alpha} = \delta \underbrace{\alpha_1^T \alpha_1}_1 = \delta = 0,$$

was bedeutet, dass $\Sigma \alpha = \lambda \alpha$ und α ist wieder ein Eigenvektor von Σ mit Eigenwert λ . Da α orthogonal zu α_1 sein soll und $\text{Var}(\alpha^T X) = \lambda$ maximal sein soll, bekommt man

$$\alpha = \alpha_2 \text{ und } \lambda = \lambda_2 \Rightarrow Y_2 = \alpha_2^T X.$$

□

Übungsaufgabe 4.2.1. Führen Sie den Beweis für $i > 2$ durch!

Sei nun $A = (\alpha_1, \dots, \alpha_n)$. Dies ist eine orthogonale $(n \times n)$ -Matrix, für die gilt (aus dem Satz 4.2.1), dass

$$\Sigma A = A\Lambda, \quad \Lambda = \text{diag}(\lambda_1, \dots, \lambda_n),$$

oder, äquivalent dazu,

$$A^T \Sigma A = \Lambda, \quad \Sigma = A\Lambda A^T \quad (4.2.2)$$

Satz 4.2.2. Für eine $(n \times m)$ -Matrix B , mit orthogonalen Spalten b_i , $i = 1, \dots, m$, $m \leq n$, sei $Y = B^T X$ und $\Sigma_Y = \text{Cov}(Y) = B^T \Sigma B$ die Kovarianzmatrix von Y . Dann gilt

$$A_m = \underset{B}{\text{argmax}} \text{Spur}(\Sigma_Y),$$

wobei $A_m = (\alpha_1, \dots, \alpha_m)$.

Beweis. Da $\alpha_1, \dots, \alpha_n$ eine Basis in \mathbb{R}^n bilden, gilt

$$b_k = \sum_{i=1}^n c_{ik} \alpha_i, \quad k = 1, \dots, m,$$

wobei $B = (b_1, \dots, b_m)$, oder, in Matrixform, $B = AC$, mit $C = (c_{ij})$, $i = 1, \dots, n$, $j = 1, \dots, m$. Daher gilt

$$\Sigma_Y = B^T \Sigma B = C^T \underbrace{A^T \Sigma A}_{\Lambda} C = C^T \Lambda C = \sum_{i=1}^n \lambda_i c_i c_i^T,$$

wobei c_i^T die i -te Zeile von C ist. Deshalb gilt

$$\text{Spur}(\Sigma_Y) = \sum_{i=1}^n \lambda_i \text{Spur}(c_i c_i^T) = \sum_{i=1}^n \lambda_i \text{Spur}(c_i^T c_i) = \sum_{i=1}^n \lambda_i |c_i|^2.$$

Da $C = A^{-1}B = A^T B$, gilt

$$C^T C = B^T \underbrace{A A^T}_{I_n} B = \underbrace{B^T B}_{I_m} = I_m,$$

wobei

$$I_k = \text{diag}(\underbrace{1, \dots, 1}_k).$$

Somit

$$\sum_{i=1}^n \sum_{j=1}^m c_{ij}^2 = m,$$

und die Spalten von C sind orthonormal. Daher kann C als ein Teil (erste m Spalten) einer orthonormalen $(n \times n)$ -Matrix D gesehen werden. Da auch die Zeilen von D orthonormale Vektoren sind und c_i^T die ersten m Elemente der Zeilen von D bilden, gilt

$$c_i^T c_i = \sum_{j=1}^m c_{ij}^2 \leq 1, \quad i = 1, \dots, n.$$

Da

$$\text{Spur}(\Sigma_Y) = \sum_{i=1}^n \lambda_i \underbrace{\sum_{j=1}^m c_{ij}^2}_{\beta_i} = \sum_{i=1}^n \beta_i \lambda_i,$$

wobei $\beta_i \leq 1$, $i = 1, \dots, n$, $\sum_{i=1}^n \beta_i = m$, und

$$\lambda_1 > \lambda_2 > \dots > \lambda_n, \quad \sum_{i=1}^n \beta_i \lambda_i \rightarrow \max$$

für $\beta_1 = \dots = \beta_m = 1$, $\beta_{m+1} = \dots = \beta_n = 0$. Aber wenn $B = A_m$, dann gilt

$$c_{ij} = \begin{cases} 1 & , 1 \leq i = j \leq m \\ 0 & , \text{sonst} \end{cases},$$

woraus $\beta_1 = \dots = \beta_m = 1$, $\beta_{m+1} = \dots = \beta_n = 0$ folgt. Somit ist A_m die Lösung von $\text{Spur}(\Sigma_Y) \rightarrow \max_B$. \square

Die Behauptung des Satzes 4.2.2 bedeutet, dass

$$\text{Var} \left(\sum_{i=1}^m Y_i \right) = \text{Var} \left(\sum_{i=1}^m \alpha_i^T X \right)$$

maximal ist für $\forall m = 1, \dots, n$, falls Y_i Hauptkomponenten von X sind.

Folgerung 4.2.1. (Spektraldarstellung von Σ). Es gilt

$$\Sigma = \sum_{i=1}^n \lambda_i \cdot \alpha_i \cdot \alpha_i^T \quad (4.2.3)$$

Beweis. Die Darstellung folgt aus (4.2.2), weil

$$\Sigma = (\alpha_1, \dots, \alpha_n) \cdot \text{diag}(\lambda_1, \dots, \lambda_n) \cdot (\alpha_1, \dots, \alpha_n)^T$$

\square

Bemerkung 4.2.1. 1. Da $\lambda_1 > \lambda_2 > \dots > \lambda_n$ mit $|\alpha_i| = 1, \forall i$, folgt aus der Darstellung (4.2.3), dass die ersten Hauptkomponenten nicht nur den Hauptbeitrag zur Varianz von X_i , sondern auch zu den Kovarianzen liefern. Dieser Beitrag wird mit steigendem $i = 1, \dots, n$ immer geringer.

2. Falls $\text{Rang}(\Sigma) = r < n$, dann bedeutet (4.2.3), dass Σ komplett aus ihren ersten r Hauptkomponenten und Koeffizientenvektoren bestimmt werden kann.

Lemma 4.2.1. Sei Σ eine positiv definite symmetrische $(n \times n)$ -Matrix mit Eigenwerten $\lambda_1 > \lambda_2 > \dots > \lambda_n > 0$ und entsprechenden Eigenvektoren $\alpha_1, \dots, \alpha_n, |\alpha_i| = 1, i = 1, \dots, n$. Dann gilt

$$\lambda_k = \sup_{\alpha \in S_k, \alpha \neq 0} \frac{\alpha^T \Sigma \alpha}{|\alpha|^2},$$

wobei $S_k = \langle \alpha_1, \dots, \alpha_{k-1} \rangle^\perp$ für beliebige $k = 1, \dots, n$.

Beweis. Sei

$$c = \sup_{\alpha \in S_k} \frac{\alpha^T \Sigma \alpha}{|\alpha|^2}.$$

Zeigen wir, dass $\lambda_k \leq c \leq \lambda_k$.

1. $c \geq \lambda_k$: Für $\alpha = \alpha_k$ beweist man

$$c \geq \frac{\alpha_k^T \Sigma \alpha_k}{\alpha_k^T \alpha_k} = \frac{\lambda_k \alpha_k^T \alpha_k}{\alpha_k^T \alpha_k} = \lambda_k$$

2. $c \leq \lambda_k$: Es ist zu zeigen, dass

$$\alpha^T \Sigma \alpha \leq \lambda_k |\alpha|^2, \quad \forall \alpha \in S_k, \quad \alpha \neq 0, \quad \forall \alpha \in \mathbb{R}^n \quad \alpha = \sum_{i=1}^n c_i \alpha_i,$$

weil $\{\alpha_i\}_{i=1}^n$ eine orthonormale Basis bilden.

$$\alpha \in S_k \quad \Rightarrow \quad c_1 = \dots = c_{k-1} = 0,$$

dass heißt

$$\begin{aligned} \alpha &= \sum_{i=k}^n c_i \alpha_i, \quad \Sigma \alpha = \sum_{i=1}^n c_i \Sigma \alpha_i = \sum_{i=1}^n c_i \lambda_i \alpha_i, \quad \alpha^T \Sigma \alpha = \left(\sum_{i=1}^n c_i \alpha_i \right)^T \left(\sum_{i=1}^n \lambda_i c_i \alpha_i \right) \\ &= \sum_{i,j=1}^n c_i c_j \lambda_i \underbrace{\alpha_j^T \alpha_i}_{\delta_{ij}} = \sum_{i=1}^n c_i^2 \lambda_i, \quad |\alpha|^2 = \sum_{i=1}^n c_i^2 \end{aligned}$$

Deshalb gilt für $\alpha \in S_k$

$$\alpha^T \Sigma \alpha = \sum_{i=k}^n c_i^2 \lambda_i \leq \sum_{i=k}^n \lambda_k c_i^2 = \lambda_k \sum_{i=k}^n c_i^2 = \lambda_k |\alpha|^2,$$

und $c \leq \lambda_k$ weil $\lambda_k > \lambda_j$, $j > k$.

□

Satz 4.2.3. Seien B , Y und Σ_Y wie in Satz 4.2.2. Dann gilt

$$A_m = \operatorname{argmax}_B \det(\Sigma_Y),$$

wobei $A_m = (\alpha_1, \dots, \alpha_m)$.

Beweis. Sei $k \in \{1, \dots, m\}$ fixiert. Führen wir $S_k = \langle \alpha_1, \dots, \alpha_{k-1} \rangle^\perp \subset \mathbb{R}^k$ ein (wie in Lemma 4.2.1). Seien $\mu_1 > \mu_2 > \dots > \mu_m$ Eigenwerte von $\Sigma_Y = B^T \Sigma B$ mit entsprechenden Eigenvektoren $\gamma_1, \dots, \gamma_m$, die orthonormiert sind. Sei $T_k = \langle \gamma_{k+1}, \dots, \gamma_m \rangle \subset \mathbb{R}^m$. Es gilt offensichtlich

$$\operatorname{Dim}(S_k) = n - k + 1, \quad \operatorname{Dim}T_k = k.$$

Genau wie in Lemma 4.2.1 kann gezeigt werden, dass $\forall \gamma \neq 0, \gamma \in T_k$ gilt

$$\frac{\gamma^T \Sigma \gamma}{|\gamma|^2} \geq \mu_k.$$

Betrachten wir $\tilde{S}_k = B(T_k) \subset \mathbb{R}^n$. Da B eine orthonormale Transformation ist, ist sie eindeutig und somit $\operatorname{Dim}(S_k) = \operatorname{Dim}(T_k) = k$. Aus der Formel

$$\operatorname{Dim}(S_k \cup \tilde{S}_k) + \operatorname{Dim}(S_k \cap \tilde{S}_k) = \operatorname{Dim}S_k + \operatorname{Dim}\tilde{S}_k$$

folgt

$$\operatorname{Dim}(S_k \cap \tilde{S}_k) = \underbrace{\operatorname{Dim}S_k}_{n-k+1} + \underbrace{\operatorname{Dim}\tilde{S}_k}_k - \underbrace{\operatorname{Dim}(S_k \cup \tilde{S}_k)}_{\leq n} \geq n - k + 1 + k - n = 1$$

das heißt, $\exists \alpha \in S_k \cap \tilde{S}_k$, $\alpha \neq 0$. Für dieses α gilt $\alpha = B\gamma$, $\gamma \in T_k$ und deshalb

$$\mu_k \leq \frac{\gamma^T \Sigma \gamma^2}{|\gamma|^2} = \frac{\gamma^T B^T \Sigma B \gamma}{\underbrace{\gamma^T \gamma}_{\gamma^T B^T B \gamma}} = \frac{\alpha^T \Sigma \alpha}{\alpha^T \alpha} \leq \lambda_k$$

nach $|\gamma| = |B\gamma|$, weil B Distanzen beibehält. Deshalb gilt $\mu_k \leq \lambda_k$ für alle $k = 1, \dots, m$, und

$$\det(\Sigma_Y) = \prod_{i=1}^m \mu_k \leq \prod_{k=1}^m \lambda_k \quad \Rightarrow \quad \max_B \det(\Sigma_Y) \leq \prod_{k=1}^m \lambda_k.$$

Allerdings gilt für $B = A_m$, $\mu_k = \lambda_k$, $k = 1, \dots, m$, deshalb

$$A_m = \operatorname{argmax}_B \det(\Sigma_Y).$$

□

Nun betrachten wir geometrische Eigenschaften von Hauptkomponenten.

Proposition 4.2.1. Die Hauptkomponentenkoeffizienten $\alpha_1, \dots, \alpha_n$ sind die Hauptachsen des Ellipsoids $x^T \Sigma^{-1} x = c$, mit Halbachsenlängen $\sqrt{c\lambda_i}$, $i = 1, \dots, n$.

Beweis. Die Hauptkomponenten von X sind gegeben durch $Z = A^T X$, wobei $A = (\alpha_1, \dots, \alpha_n)$ eine orthonormale Transformation ist, deshalb $A^T = A^{-1}$, $X = AZ$. Daher gilt für unser Ellipsoid

$$c = x^T \Sigma^{-1} x \quad \underbrace{=}_{\text{Subst. } x=Az} \quad z^T A^T \Sigma^{-1} A z = z^T \Lambda^{-1} z,$$

wobei

$$A^T \Sigma^{-1} A = \Lambda^{-1} = \operatorname{diag} \left(\frac{1}{\lambda_1}, \dots, \frac{1}{\lambda_n} \right), \quad \Lambda = \operatorname{diag}(\lambda_1, \dots, \lambda_n),$$

weil Σ^{-1} dieselben Eigenvektoren mit Eigenwerten $\frac{1}{\lambda_i}$ hat. Daher kann das Ellipsoid $z^T \Lambda^{-1} z = c$ in seiner normierten Form als

$$\sum_{k=1}^n \frac{z_k^2}{c\lambda_k} = 1$$

dargestellt werden. Daraus folgt, dass α_i in die Richtungen seiner Hauptachsen zeigen und, dass seine Halbachsenlängen gleich $\sqrt{c\lambda_i}$ sind. □

Bemerkung 4.2.2. (*Multivariate Normalverteilung*). Falls $X \sim N(0, \Sigma)$ gilt, dann ist $x^T \Sigma^{-1} x = c$ ein Ellipsoid der konstanten Wahrscheinlichkeit für X , weil die Dichte von X

$$f_X(x) = \frac{1}{\sqrt{\det \Sigma}} \exp \left\{ -\frac{1}{2} x^T \Sigma^{-1} x \right\} \cdot \frac{1}{(2\pi)^{\frac{n}{2}}}, \quad x \in \mathbb{R}^n,$$

auf diesem Ellipsoid konstant bleibt. Sonst definiert $x^T \Sigma^{-1} x = c$ Konturen der konstanten Wahrscheinlichkeit für X . Dabei zeigt der Vektor α_1 in die Richtung der größten Varianz von $\alpha^T X$ (es ist die größte Hauptachse mit Länge $\sqrt{c\lambda_1}$ des Ellipsoids); α_2 zeigt in die Richtung der zweit größten Varianz (Halbachse $\sqrt{c\lambda_2}$), usw. (vgl. Bedingung 4.2.1).

Bemerkung 4.2.3. Eine andere Form von Hauptkomponentenanalyse ist möglich, wenn man statt $X = (X_1, \dots, X_n)^T$ die normierte Stichprobe $X_\omega = (X_1/\omega_1, \dots, X_n/\omega_n)^T$ benutzt, wobei Gewichte $\omega = (\omega_1, \dots, \omega_n)^T$ eine gewisse Präferenz in der Analyse zum Ausdruck bringen und somit Vorinformationen enthalten. Eine häufige Wahl ist

$$\omega_i = \sqrt{\sigma_{ii}} = \sqrt{\text{Var}X_i},$$

was zur HKA von $X^* = (X_1^*, \dots, X_n^*)$, $X_i^* = \frac{X_i}{\sqrt{\text{Var}X_i}}$, $i = 1, \dots, n$ mit Hilfe der Korrelationsmatrix $\Sigma^* = (\text{Corr}(X_j, X_i))_{i,j=1}$ führt

$$\text{Corr}(X_i, X_j) = \frac{\text{Cov}(X_i, X_j)}{\sqrt{\text{Var}X_i \text{Var}X_j}} = \text{Cov}(X_i^*, X_j^*), \quad i, j = 1, \dots, n.$$

Dabei kommt man auf andere Hauptkomponenten $\alpha_i^{*T} X^*$, für die $\alpha_i^* \neq \alpha_i$ gilt, $i = 1, \dots, n$.

Was sind dann Vor- bzw. Nachteile von HKA basierend auf (X, Σ) und (X^*, Σ^*) ?

Nachteile von (X, Σ) -HKA:

1. Die HKA basierend auf (X^*, Σ^*) hängt nicht von der Wahl der Maßeinheiten von X ab. Somit sind Vergleiche der Ergebnisse von HKA von mehreren Stichproben unterschiedlicher Herkunft möglich.
2. Falls die Varianzen von X_i sehr unterschiedlich sind, so werden die Variablen X_i mit größten Varianzen auch die ersten HK bestimmen, was eindeutig einen Nachteil darstellt. Die HKA basierend auf (X^*, Σ^*) ist frei von diesem Nachteil. Die (X, Σ) -HKA ist in solchen Fällen nicht aussagekräftig, weil sie (in leicht veränderter Form) einfach die Variablen X_i in der Reihenfolge absteigender Varianzen ordnet.

Beispiel 4.2.1. Sei $X = (X_1, X_2)$, wobei X_1 die Länge darstellt und X_2 das Gewicht. X_1 kann in cm oder m gemessen werden, X_2 allerdings nur in kg. In diesen zwei Fällen seien die Kovarianzmatrizen von X gegeben durch

$$\Sigma_1 = \begin{pmatrix} 80 & 44 \\ 44 & 80 \end{pmatrix} \quad \text{bzw.} \quad \Sigma_2 = \begin{pmatrix} 80 \cdot 10^4 & 4400 \\ 4400 & 80 \end{pmatrix}.$$

Die Berechnung der ersten HK ergibt in beiden Fällen

$$\alpha_1^T X = 0,707X_1 + 0,707X_2 \quad \text{für } \Sigma_1 \quad \text{bzw.} \quad \alpha_1^T X = 0,998X_1 + 0,055X_2 \quad \text{für } \Sigma_2.$$

Zu bemerken ist, dass im ersten Fall X_1 und X_2 gleiche Beiträge zur 1. HK besitzen, wobei im 2. Fall X_1 den dominierenden Einfluss ausübt. Dazu gilt $\frac{\lambda_1}{\lambda_1 + \lambda_2} \cdot 100\% = 77,5\%$ im ersten Fall und $\frac{\lambda_1}{\lambda_1 + \lambda_2} \cdot 100\% = 99,3\%$ im 2. Fall (es ist der Anteil der Variation der ersten HK von der gesamten Varianz).

3. Falls Zufallsvariable X_i in X unterschiedlicher Herkunft sind (wie im obigen Beispiel), dann ist die Interpretation des Anteils der Variation problematisch, weil in der Summe $\lambda_1 + \dots + \lambda_n$ m^2 , kg^2 , usw. aufsummiert werden. Die HKA basierend auf (X^*, Σ^*) dagegen betrachtet maßlose Größen, so dass die Summe $\lambda_1 + \dots + \lambda_n$ durchaus interpretierbar ist.

Vorteile von (X, Σ) -HKA:

1. Falls statt Σ bzw. Σ^* ihre empirische Analoga $\hat{\Sigma}$ bzw. $\hat{\Sigma}^*$ benutzt werden (wenn $\Sigma(\Sigma^*)$ nicht bekannt sind, müssen sie aus den Daten geschätzt werden), dann hat $(X, \hat{\Sigma})$ -HKA Vorteile, weil die statistischen Methoden hier einfacher sind als bei $(X^*, \hat{\Sigma}^*)$ -HKA.
2. Wenn X_i in X alle dieselbe Maßeinheit besitzen, dann ist die HKA basierend auf (X, Σ) manchmal vorteilhafter, weil bei der Standardisierung von (X, Σ) auf (X^*, Σ^*) der Bezug zu den Einheiten, in denen X gemessen wurde, verloren geht.

Bemerkung 4.2.4. Manchmal wird in Definition 4.2.1 statt $|\alpha| = 1$ die Normierung $|\alpha_k| = \sqrt{\lambda_k}$, $k = 1, \dots, n$ benutzt (siehe Optimierungsaufgabe (4.2.1)). Dies ist insbesondere der Fall in der korrelationsbasierten HKA.

Bemerkung 4.2.5. (*Gleiche Eigenwerte λ_i*). Falls einige Eigenwerte von Σ gleich sind, z.B. $\lambda_1 = \lambda_2 = \dots = \lambda_k > \lambda_{k+1} > \dots > \lambda_m$, bedeutet dies, dass es einen linearen Unterraum der Dimension k gibt, in denen eine beliebige Basis die ersten k Eigenvektoren darstellt. Dies bedeutet, dass für die HKA die ersten k Eigenvektoren nicht eindeutig definiert werden können. Geometrisch interpretiert: Die ersten k Halbachsen von $x^T \Sigma^{-1} x = c$ sind gleich, d.h., das Ellipsoid $x^T \Sigma^{-1} x = c$ hat einen sphärischen k -dimensionalen Durchschnitt durch den Ursprung, in dem die Richtungen der Halbachsen beliebig (orthogonal zueinander) gewählt werden können.

Bemerkung 4.2.6 ($\lambda_i = 0$). Wenn $\lambda_1 > \dots > \lambda_{n-k} > \lambda_{n-k+1} = \dots = \lambda_n = 0$, dann gibt es in der Stichprobe X lediglich $n-k$ linear unabhängige Zufallsvektoren X_i . Deshalb sollten nur diese $n-k$ Variablen zur Analyse benutzt werden.

4.3 Hauptkomponentenanalyse auf Datenebene

Bei diesem Abschnitt wird nicht mehr vorausgesetzt, dass die Kovarianzmatrix Σ bekannt ist. Deshalb soll sie durch die empirische Kovarianzmatrix $\hat{\Sigma}$ ersetzt werden. Seien X^1, X^2, \dots, X^m unabhängige Realisierungen eines n -dimensionalen Zufallsvektors $X = (X_1, \dots, X_n)^T$, $X^i = (X_1^i, \dots, X_n^i)^T$, $i = 1, \dots, m$. X^i wird als Beobachtung von X interpretiert.

Definition 4.3.1. Definiere den n -dimensionalen Zufallsvektor a_k durch

$$a_k = \operatorname{argmax}_{a \in \mathbb{R}^n} \frac{1}{m-1} \sum_{i=1}^m (Y_i - \bar{Y})^2$$

mit Nebenbedingungen $|a| = 1$, a unkorreliert mit a_1, \dots, a_{k-1} für alle $k = 1, \dots, n$, wobei

$$Y_i = a^T X^i, \quad i = 1, \dots, m, \quad \bar{Y} = \frac{1}{m} \sum_{i=1}^m Y_i.$$

So definiert $a_k^T X$ die k -ten Hauptkomponenten von X mit Koeffizientenvektor a_k , $Y_{ik} = a_k^T X^i$ ist die Auswertung der k -ten HK auf der i -ten Beobachtung X^i von X , $i = 1, \dots, m$, $k = 1, \dots, n$.

Lemma 4.3.1. Es gilt

$$\frac{1}{m-1} \sum_{i=1}^m (Y_{ik} - \bar{Y}_k)^2 = l_k, \quad k = 1, \dots, n,$$

wobei

$$\bar{Y}_k = \frac{1}{m} \sum_{i=1}^m Y_{ik}, \quad \bar{X}_k = \frac{1}{m} \sum_{i=1}^m X_k^i, \quad k = 1, \dots, n$$

und l_k der Eigenwert der empirischen Kovarianzmatrix $\hat{\Sigma} = (\hat{\sigma}_{ij})_{i,j=1}^n$ ist,

$$\hat{\sigma}_{ij} = \frac{1}{m-1} \sum_{t=1}^m (X_t^i - \bar{X}_i)(X_t^j - \bar{X}_j), \quad i, j = 1, \dots, n, \quad l_1 > l_2 > \dots > l_n.$$

a_k ist der Eigenvektor von $\hat{\Sigma}$ mit Eigenwert l_k , $k = 1, \dots, n$.

Beweis.

Übungsaufgabe 4.3.1. Vergleiche den Beweis des Satzes 4.2.1. □

Im Folgenden werden wir X^i durch $X^i - \bar{X}$ ersetzen und dabei die Bezeichnung X^i beibehalten, $i = 1, \dots, n$.

Bemerkung 4.3.1. Die Eigenschaften der HKA formuliert in Satz 4.2.2, Folgerung 4.2.1, Satz 4.2.3, Proposition 4.2.1 bleiben auch in ihrer statistischen Version (Definition 4.3.1) erhalten, mit folgenden offensichtlichen Modifikationen: Σ wird ersetzt durch $\hat{\Sigma}$, $A = (\alpha_1, \dots, \alpha_n)$ durch $A = (a_1, \dots, a_n)$, $A_m = (\alpha_1, \dots, \alpha_m)$ durch $A_m = (a_1, \dots, a_m)$, Σ_Y durch die empirische Kovarianzmatrix $\hat{\Sigma}_Y$ von Y . So benutzt beispielsweise die Spektraldarstellung von $\hat{\Sigma}$

$$\hat{\Sigma} = \sum_{i=1}^n l_i a_i a_i^T \quad (4.3.1)$$

Übungsaufgabe 4.3.2. Zeigen Sie es!

Zeigen wir eine weitere Eigenschaft der empirischen HKA, die auch als eine äquivalente Definition betrachtet werden kann:

Satz 4.3.1. Sei B eine $(n \times p)$ -Matrix, $p \leq n$, mit orthogonalen Spalten. Seien $Z_i = B^T X^i$, $i = 1, \dots, m$ Projektionen von X^i , $i = 1, \dots, m$, auf einen p -dimensionalen Unterraum L_B . Definiere

$$G(B) = \sum_{i=1}^m |X^i - Z_i|^2.$$

Dann gilt

$$A_p = (a_1, \dots, a_p) = \underset{B}{\operatorname{argmin}} G(B).$$

Beweis. Nach dem Satz von Pythagoras gilt $|X^i|^2 = |Z_i|^2 + |X^i - Z_i|^2$, deshalb

$$G(B) = \sum_{i=1}^m |X^i|^2 - \sum_{i=1}^m |Z_i|^2 \rightarrow \min$$

falls

$$\tilde{G}(B) = \sum_{i=1}^m |Z_i|^2 = \sum_{i=1}^m Z_i^T Z_i = \sum_{i=1}^m X^{iT} B B^T X^i \rightarrow \max_B.$$

Es gilt

$$\begin{aligned} \tilde{G}(B) &= \operatorname{Spur} \left(\sum_{i=1}^m (X^{iT} B B^T X^i) \right) = \sum_{i=1}^m \operatorname{Spur} (X^{iT} B B^T X^i) = \sum_{i=1}^m \operatorname{Spur} (B^T X^i X^{iT} B) \\ &= \operatorname{Spur} \left(B^T \underbrace{\left(\sum_{i=1}^m X^i X^{iT} \right)}_{1_{(m-1)\hat{\Sigma}}} B \right) = (m-1) \operatorname{Spur} (B^T \hat{\Sigma} B) \end{aligned}$$

Zusammengefasst gilt

$$\tilde{G}(B) = (m-1) \operatorname{Spur} (B^T \hat{\Sigma} B),$$

die nach Bemerkung 4.3.1 und Satz 4.2.2 maximal wird, falls $B = A_p$. \square

¹Da X^i durch $X^i - \bar{X}$ ersetzt wurde.

Bemerkung 4.3.2. Wie kann Satz 4.3.1 als äquivalente Definition der empirischen HKA benutzt werden? a_i werden als orthogonale Vektoren definiert, die einen linearen Unterraum $L_p = \langle a_1, \dots, a_p \rangle$ aufspannen, $p = 1, \dots, n-1$, mit der Eigenschaft, dass die Summe der quadratischen orthogonalen Abstände von X^i zu L_p minimal wird. So wäre es z.B. für $p = 1$ L_1 die beste Gerade, die den Datensatz X^1, \dots, X^m approximiert, für $p = n-1$ wäre L_{n-1} die beste Hyperebene mit derselben Eigenschaft (vgl. lineare Regression).

Der folgende Satz gibt uns gleichzeitig eine effiziente Berechnungsmethode und eine neue Interpretation der HK an.

Satz 4.3.2 (Singulärwertzerlegung). Sei $\tilde{X} = (X^1 - \bar{X}, X^2 - \bar{X}, \dots, X^m - \bar{X})^T$ eine $(m \times n)$ -Matrix, die zentrierte Beobachtungen X^i von X enthält. Sei $\text{Rang}(\tilde{X}) = r \leq n, m$. Es gilt folgende Zerlegung:

$$\tilde{X} = ULA_r^T, \quad (4.3.2)$$

wobei U eine $(m \times r)$ -Matrix mit orthonormalen Spalten ist

$$L = \text{diag}(\tilde{l}_1, \dots, \tilde{l}_r) \quad \text{wobei} \quad \tilde{l}_i = \sqrt{(m-1)l_i}$$

die Wurzel aus dem i -ten (nicht trivialen) Eigenwert von $\tilde{X}^T \tilde{X} = (m-1)\hat{\Sigma}$ ist, $i = 1, \dots, r$. $A_r = (a_1, \dots, a_r)$ ist die $(n \times r)$ -Matrix mit Spalten a_i

Beweis. Definiere $U = (u_1, \dots, u_r)$ mit Spalten $u_i = \frac{\tilde{X} a_i}{\tilde{l}_i}$, $i = 1, \dots, r$. Zeigen wir, dass die Darstellung (4.3.2) gilt. Laut Spektraldarstellung (4.3.1) gilt

$$(m-1)\hat{\Sigma} = \tilde{X}^T \tilde{X} = \sum_{i=1}^r \tilde{l}_i^2 a_i a_i^T, \quad \text{weil} \quad l_i = 0, i = r+1, \dots, n.$$

Deshalb

$$ULA_r^T = U \begin{pmatrix} \tilde{l}_1 a_1^T \\ \vdots \\ \tilde{l}_r a_r^T \end{pmatrix} = \sum_{i=1}^r \tilde{X} \frac{a_i \tilde{l}_i a_i^T}{\tilde{l}_i} = \sum_{i=1}^r \tilde{X} a_i a_i^T \stackrel{l_i=0, i>r}{=} \sum_{i=1}^n \tilde{X} a_i a_i^T$$

$\tilde{X} a_i = 0$, $i = r+1, \dots, n$, wegen $\text{rang}(\tilde{X}) = r$ und Zentrierung der Spalten von \tilde{X} durch \bar{X} . Da die Vektoren a_i orthonormal sind, gilt

$$ULA_r^T = \tilde{X} \sum_{i=1}^n a_i a_i^T = \tilde{X} I = \tilde{X}.$$

□

Bemerkung 4.3.3. Die Matrix U liefert folgende Versionen von Auswertungen

$$Y_{ik} = a_k^T X^i = X^{iT} a_k, \quad Y_{ik} = u_{ik} \tilde{l}_k, \quad i = 1, \dots, m, \quad k = 1, \dots, n$$

Es gilt

$$\text{Var}(u_{ik}) = \frac{\text{Var}(Y_{ik})}{\tilde{l}_k^2} = \frac{l_k}{(m-1)l_k} = \frac{1}{m-1}, \quad \forall i, k$$

4.4 Asymptotische Verteilung von HK bei normalverteilten Stichproben

Sei nun $X \sim N(\mu, \Sigma)$, Σ habe Eigenwerte $\lambda_1 > \lambda_2 > \dots > \lambda_n > 0$ und entsprechende Eigenvektoren α_k , $k = 1, \dots, n$. Berechne

$$\lambda = (\lambda_1, \dots, \lambda_n)^T, \quad l = (l_1, \dots, l_n)^T, \quad \alpha_k = (\alpha_{k1}, \dots, \alpha_{kn})^T, \quad a_k = (a_{k1}, \dots, a_{kn})^T, \\ k = 1, \dots, n$$

Satz 4.4.1. 1. l ist asymptotisch (für $m \rightarrow \infty$) unabhängig von a_k , $k = 1, \dots, n$.

2. l und a_k , $k = 1, \dots, n$ sind asymptotisch $m \rightarrow \infty$ multivariat normalverteilt, mit asymptotischen Erwartungswerten

$$\lim_{m \rightarrow \infty} \mathbb{E}(l) = \lambda \quad \text{und} \quad \lim_{m \rightarrow \infty} \mathbb{E}(a_k) = \alpha_k, \quad k = 1, \dots, n.$$

3. Es gilt

$$\text{Cov}(l_k, l_{k'}) \sim \begin{cases} \frac{2\lambda_k^2}{m-1}, & k = k' \\ 0, & k \neq k' \end{cases} \quad \text{für } m \rightarrow \infty$$

$$\text{Cov}(a_{kj}, a_{k'j'}) \sim \begin{cases} \frac{\lambda_k}{m-1} \sum_{l=1, l \neq k}^n \frac{\lambda_l \alpha_{lj} \alpha_{lj'}}{(\lambda_l - \lambda_k)^2}, & k = k' \\ -\frac{\lambda_k \lambda_{k'} \alpha_{kj} \alpha_{k'j'}}{(m-1)(\lambda_k - \lambda_{k'})^2}, & k \neq k' \end{cases} \quad \text{für } m \rightarrow \infty.$$

Ohne Beweis!

Die Aussagen von Satz 4.4.1 können dazu benutzt werden, ML-Schätzer sowie Konfidenzintervalle für λ und α_k zu konstruieren.

Übungsaufgabe 4.4.1. 1. Zeige, dass ein ML-Schätzer für Σ durch $\frac{m-1}{m} \hat{\Sigma}$ gegeben ist.

2. Zeige, dass der ML-Schätzer

$$\begin{cases} \text{für } \lambda \text{ ist} & \hat{\lambda} = \frac{m-1}{m} l. \\ \text{für } \alpha_k \text{ ist} & \hat{\alpha}_k = a_k, k = 1, \dots, n. \end{cases}$$

3. Zeige, dass die ML-Schätzer in 2. mit Momenten-Schätzern für λ und α_k übereinstimmen, die aus dem Satz 4.4.1 gewonnen werden können.

Folgerung 4.4.1 (Konfidenzintervalle für λ_k). Ein asymptotisches Konfidenzintervall für λ_k ($m \rightarrow \infty$) zum Niveau $1 - \alpha$ ist gegeben durch

$$\left[l_k \left(1 - \sqrt{\frac{2}{m-1}} z_{\frac{\alpha}{2}} \right)^{-1}, l_k \left(1 + \sqrt{\frac{2}{m-1}} z_{\frac{\alpha}{2}} \right)^{-1} \right],$$

wobei m so groß ist, dass $-\sqrt{\frac{2}{m-1}} z_{\frac{\alpha}{2}} < 1$.

Beweis. Da $l_k \sim N\left(\lambda_k, \frac{2\lambda_k^2}{m-1}\right)$ für $m \rightarrow \infty$ aus Satz 4.4.1, 2. und 3., gilt

$$\frac{l_k - \lambda_k}{\sqrt{\frac{2}{m-1}} \lambda_k} \sim N(0, 1) \quad \text{für } m \rightarrow \infty.$$

Daraus folgt

$$\lim_{m \rightarrow \infty} \mathbb{P} \left(z_{\frac{\alpha}{2}} \leq \frac{l_k - \lambda_k}{\lambda_k} \sqrt{\frac{m-1}{2}} \leq z_{1-\frac{\alpha}{2}} \right) = 1 - \alpha,$$

oder für $m \rightarrow \infty$

$$\sqrt{\frac{2}{m-1}} z_{\frac{\alpha}{2}} \leq \frac{l_k}{\lambda_k} - 1 \leq \sqrt{\frac{2}{m-1}} \underbrace{z_{1-\frac{\alpha}{2}}}_{=-z_{\frac{\alpha}{2}}},$$

$$\frac{l_k}{1 - \sqrt{\frac{2}{m-1}} z_{\frac{\alpha}{2}}} \leq \lambda_k \leq \frac{l_k}{1 + \sqrt{\frac{2}{m-1}} z_{\frac{\alpha}{2}}}$$

mit Wahrscheinlichkeit $1 - \alpha$. □

Da alle l_k , $k = 1, \dots, n$ asymptotisch ($m \rightarrow \infty$) unabhängig sind, kann ein simultaner Konfidenzbereich für l als kartesisches Produkt der Konfidenzintervalle für l_k aus Folgerung 4.4.1 angegeben werden.

Lemma 4.4.1. Es gilt

$$(m-1) \alpha_k^T \left(l_k \hat{\Sigma}^{-1} + l_k^{-1} \hat{\Sigma} - 2I_n \right) \alpha_k \xrightarrow{m \rightarrow \infty} \chi_{n-1}^2$$

Ohne Beweis!

Daraus folgt das (asymptotische) Konfidenzellipsoid für α_k zum Niveau $1 - \beta$

$$\left\{ y \in \mathbb{R}^n : (m-1) y^T \left(l_k \hat{\Sigma}^{-1} + l_k^{-1} \hat{\Sigma} - 2I_n \right) y \leq \chi_{n-1, \beta}^2 \right\}.$$

Bemerkung 4.4.1. Folgerung 4.4.1 bzw. Lemma 4.4.1 können zur Konstruktion von statistischen Tests für λ_k bzw. α_k folgendermaßen verwendet werden:

1. Testen von $H_0 : \lambda_k = \lambda_{k_0}$ v.s. $H_1 : \lambda_k \neq \lambda_{k_0}$
Die Hypothese H_0 wird verworfen, falls

$$\left| \frac{l_k - \lambda_{k_0}}{\sqrt{\frac{2}{m-1} \lambda_{k_0}}} > z_{\frac{\alpha}{2}} \right|.$$

Dies ist ein asymptotischer Test ($m \rightarrow \infty$) zum Niveau α .

2. Testen wir $H_0 : \alpha_k = \alpha_{k_0}$ v.s. $H_1 : \alpha_k \neq \alpha_{k_0}$
Die Hypothese H_0 wird abgelehnt, falls

$$(m-1) \alpha_{k_0}^T \left(l_k \hat{\Sigma}^{-1} + l_k^{-1} \hat{\Sigma} - 2I_n \right) \alpha_{k_0} \geq \chi_{n-1, \alpha}^2.$$

Dies ist ein asymptotischer ($m \rightarrow \infty$) Test zum Niveau α .

4.5 Ausreißererkennung

In diesem Abschnitt gehen wir davon aus, dass unsere Stichprobe X^1, X^2, \dots, X^m einige Ausreißer enthalten kann. Was aber ist ein Ausreißer? In der statistischen Literatur gibt es dazu keine einheitliche Meinung. Allgemein würden wir sagen, dass die Beobachtung X^i ein Ausreißer ist, wenn sie einen untypischen Wert (in Bezug auf die Verteilung von X) annimmt. Es kann z.B. ein ungewöhnlich hoher bzw. niedriger Wert von einigen Koordinaten von X^i sein. Es kann aber auch eine ungewöhnliche Kombination von gewöhnlichen Koordinatenwerten einiger Koordinaten von X^i sein. Der Grund für solche untypischen Werte X^i kann ein Meßfehler, aber auch eine Anomalie im Datensatz sein.

Beispiel 4.5.1. Sei $X = (X_1, X_2)$, wobei X_1 = "Körpergröße" (in cm) und X_2 = "Gewicht" (in kg) von Kindern im Alter von 5 bis 15 Jahren sind. Das Merkmal X wird in einer medizinischen Studie n mal gemessen. Dabei sind Beobachtungen $X^i = (250, 80)$ und $X^j = (175, 25)$ als Ausreißer klassifiziert worden, und zwar daher, weil $X^i = 250$ cm eine unvorstellbare Körpergröße ist, bei X^j sind sowohl $X_1^j = 175$ als auch $X_2^j = 25$ im mittleren Wertebereich von X_1 und X_2 , ihre Kombination jedoch ist praktisch unmöglich.

Wie könnte man Ausreißer enttarnen? Normalerweise werden untypische Werte von X^i anhand von Plots des Datensatzes X^1, \dots, X^m als Einzelpunkte, die nicht in der großen Punktwolke liegen, identifiziert. Bei hoher Dimension n von X ist es jedoch schwierig, so einen Datensatz zu visualisieren. Deshalb kann man vorschlagen einen Datenpunkt der ersten 2-3 HK von (X^1, \dots, X^m) zu erstellen. Dann werden dort ungewöhnlich große bzw. kleine Werte von X_k^i sofort erkennbar. Um jedoch eine ungewöhnliche Zusammensetzung von gewöhnlichen Koordinatenwerten X_k^i zu entdecken, bedarf es der letzten HK. Dazu wird die Auswertung folgender Statistiken empfohlen:

Seien a_1, \dots, a_n die Koeffizientenvektoren der HK von (X^1, \dots, X^m) . Seien $Y_{ik} = a_k^T X^i$, $i = 1, \dots, m$, $k = 1, \dots, n$ die Auswertungen der HK zu den Beobachtungen X^i . Seien l_k , $k = 1, \dots, n$ die Eigenwerte der empirischen Kovarianzmatrix $\hat{\Sigma}$ von (X^1, \dots, X^m) . Für ein $1 \leq n_0 \leq n$, definieren wir die Statistiken

$$d_i^{(1)}(n_0) = \sum_{k=n-n_0+1}^n Y_{ik}^2, \quad d_i^{(2)}(n_0) = \sum_{k=n-n_0+1}^n \frac{Y_{ik}^2}{l_k}, \quad d_i^{(3)}(n_0) = \sum_{k=n-n_0+1}^n l_k Y_{ik}^2,$$

$$d_i^{(4)}(n_0) = \max_{n-n_0+1 \leq k \leq n} \frac{|Y_{ik}|}{\sqrt{l_k}}, \quad i = 1, \dots, m.$$

Lemma 4.5.1. Es gilt

$$d_j^{(2)}(n) = (X^i - \bar{X})^T \hat{\Sigma}^{-1} (X^i - \bar{X}), \quad i = 1, \dots, m,$$

wobei Y_{ik} an ihren empirischen Mittel gemessen werden, das heißt, Y_{ik} werden durch $Y_{ik} - \bar{Y}_k$ ersetzt, $k = 1, \dots, n$, $i = 1, \dots, m$.

Beweis. Es gilt

$$\hat{\Sigma} = ALA^T, \quad \text{wobei } L = \text{diag}(l_1, \dots, l_n) \quad \text{und} \quad A = (a_1, \dots, a_n).$$

Daher

$$\hat{\Sigma}^{-1} = AL^{-1}A^T \quad \text{mit} \quad L^{-1} = \text{diag}(l_1^{-1}, \dots, l_n^{-1}).$$

Da zusätzlich $Y_i = A^T X^i$ für $Y_i = (Y_{i1}, \dots, Y_{in})^T$, $i = 1, \dots, n$, es gilt

$$X^i = A^{T^{-1}} Y_i = AY_i, \quad X^{iT} = Y_i^T A^T, \quad i = 1, \dots, n$$

und deshalb

$$\bar{X} = \frac{1}{m} \sum_{i=1}^m X^i = A\bar{Y}, \quad \bar{Y} = \frac{1}{m} \sum_{i=1}^m Y_i, \quad \bar{X}^T = \bar{Y}^T A^T.$$

Daher gilt

$$\begin{aligned} (X^i - \bar{X})^T \hat{\Sigma}^{-1} (X^i - \bar{X}) &= (Y_i - \bar{Y})^T \underbrace{A^T A}_I L^{-1} \underbrace{A^T A}_I (Y_i - \bar{Y}) \\ &= (Y_i - \bar{Y})^T L^{-1} (Y_i - \bar{Y}) = \sum_{k=1}^n \frac{Y_{ik}^2}{l_k} = d_i^{(2)}(n). \end{aligned}$$

□

Um nun Ausreißer in (X^1, \dots, X^m) zu erkennen, werden Werte $d_i^{(j)}(n)$, $i = 1, \dots, m$, $j = 1, \dots, n$ für $n = 1, 2, 3$ berechnet. Beobachtungen X^i mit den größten Werten $d_i^{(j)}(n)$ werden als mögliche Ausreißer eingestuft. Zusätzlich kann ein Plot von einer Punktwolke

$$D = \left\{ \left(d_i^{(2)}(n) - d_i^{(2)}(n_0), d_i^{(2)}(n_0) \right), i = 1, \dots, m \right\}$$

dabei behilflich sein. X^i wird hier als Ausreißer erkannt, wenn

$$\left(d_i^{(2)}(n) - d_i^{(2)}(n_0), d_i^{(2)}(n_0) \right)$$

isoliert von der übrigen Punktwolke D liegt.

Bemerkung 4.5.1. Falls $X \sim N(\mu, \Sigma)$ mit bekannten μ und Σ , und HKA auf Modellebene durchgeführt wird, können Verteilungen von $d_i^{(j)}(n_0)$ explizit angegeben werden. Es sind (außer $d_i^{(4)}$) Gamma-Verteilungen mit bekannten Parametern z.B. $d_i^{(2)}(n_0) \sim \chi_{n_0}^2$, $i = 1, \dots, m$. Die Verteilungsfunktion von $d_j^{(4)}(n_0)$ ist $\Phi^{n_0}(x)$, wobei $\Phi(x)$ die Verteilungsfunktion der $N(0, 1)$ -Verteilung ist. Dann können Konfidenzintervalle für $d_i^{(j)}(n_0)$ eine formale Entscheidungsregel dafür liefern, ob X^i einen Ausreißer darstellt. Diese Vorgehensweise basiert zwar auf einer festen mathematischen Grundlage, ist aber in der Praxis wenig einsetzbar, da der Fall von normalverteilten Daten (und dazu mit bekannten Parametern μ und Σ !) äußerst selten vorliegt.

Bemerkung 4.5.2. Statistiken $d_i^{(2)}, d_i^{(4)}$ betonen die letzten Statistiken mehr als $d_i^{(1)}$ (wegen der entsprechenden Normierung). Deshalb sind sie zur Entdeckung von ungewöhnlichen Korrelationen in den Daten geeignet (wie etwa in Beispiel 4.5.1, Beobachtung $X^j = (175, 25)$). Statistik $d_j^{(3)}$ betont die ersten HK. Daher ist sie anzuwenden, um ungewöhnlich große (kleine) Werte von Koordinaten X_k^i zu entdecken ($X_1^i = 250$ im Beispiel 4.5.1).

4.6 Hauptkomponentenanalyse und Regression

Sei folgendes multivariates lineares Regressionsmodell gegeben: $Y = X\beta + \varepsilon$, wobei $Y = (Y_1, \dots, Y_n)^T$ der Vektor der Zielvariablen ist, $X = (X_{ij})_{i=1, \dots, n, j=1, \dots, m}$ die $(n \times m)$ -Matrix der Ausgangsvariablen, $\text{Rang}(X) = m$, $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^T$ der Vektor der Störgrößen, wobei ε_i unabhängig sind mit $\mathbb{E}\varepsilon_i = 0$, $\text{Var}\varepsilon_i = \sigma^2$, $i = 1, \dots, n$. O.B.d.A. werden wir voraussetzen, dass X (wie in Satz 4.3.2) zentriert ist, d.h., das empirische Mittel der Zeilen von X ist Null, oder, etwas detaillierter, X_{ij} wird ersetzt durch $X_{ij} - \bar{X}_j$, wobei

$$\bar{X}_j = \frac{1}{n} \sum_{i=1}^n X_{ij}, \quad j = 1, \dots, m, \quad \text{vgl. [8].}$$

Wenn einige Variablen X_{ij} in X nahezu linear abhängig sind, das heißt $\det(X^T X) \approx 0$, dann wirkt es sich auf den Schätzer $\hat{\beta}$ von β als hohe Instabilität in seiner Berechnung

aus, weil $\text{Cov}(\hat{\beta}) = \sigma^2(X^T X)^{-1}$ (vgl. Satz 4.3.2) sehr geringe Varianzen von $\hat{\beta}_j$ enthalten wird. Ein Ausweg aus dieser Situation wird die Verwendung von Verallgemeinerungen sein wie in Kapitel 4.3. Eine andere Möglichkeit ist es, die HKA für X zu verwenden, um so lineare Abhängigkeiten in X durch die letzten HK zu detektieren und einige Variablen β_j aus der Regression auszuschließen. Genau diese Möglichkeit werden wir in diesem Abschnitt näher beschreiben

Seien a_1, \dots, a_m die Koeffizientenvektoren der HK (das heißt Eigenvektoren) von $X^T X$. Sei $Z_{ik} = a_k^T X^i$ die Auswertung der k -ten HK der i -ten Zeile X^i von X , $i = 1, \dots, n$, $k = 1, \dots, m$. Mit $Z = (Z_{ik})$ gilt $Z = XA$, wobei $A = (a_1, \dots, a_m)$ eine orthogonale $(m \times m)$ -Matrix ist. Stellen wir die Regressionsgleichung $Y = X\beta + \mathcal{E}$ folgendermaßen dar:

$$Y = X \underbrace{AA^T}_I \beta + \mathcal{E} = \underbrace{XA}_Z \underbrace{A^T}_\gamma \beta + \mathcal{E} = Z\gamma + \mathcal{E}, \text{ wobei } \gamma = A^T \beta \text{ ist.} \quad (4.6.1)$$

Somit hat man die alten Ausgangsvariablen β durch ihre Transformierte $\gamma = A^T \beta$ ersetzt. Nun folgt die Schätzung von γ aus Satz 2.2.1:

$$\hat{\gamma} = (Z^T Z)^{-1} Z^T Y = L^{-1} Z^T Y, \quad (4.6.2)$$

wobei $L = \text{diag}(l_1, \dots, l_m)$ die Eigenwerte l_i von $X^T X$ enthält. Dies gilt, weil Z orthogonale Spalten besitzt. Daher gilt

$$\hat{\beta} = A\hat{\gamma} = AL^{-1}Z^T Y = \underbrace{AL^{-1}A^T}_{(X^T X)^{-1}} X^T Y = \sum_{k=1}^m l_k^{-1} a_k a_k^T X^T Y,$$

wobei wir in der letzten Gleichungsmetrik Formeln (4.6.1), (4.6.2) und die Spektraldarstellung (Folgerung 4.2.1) von $(X^T X)^{-1}$ benutzt haben. Aus Satz 4.2.2 folgt außerdem, dass

$$\text{Cov}(\hat{\beta}) = \sigma^2 \sum_{k=1}^m l_k^{-1} a_k a_k^T.$$

Somit haben wir folgendes Ergebnis bewiesen:

Lemma 4.6.1. Die MKQ-Lösung der Regressionsgleichung $Y = X\beta + \mathcal{E}$ ist gegeben durch

$$\hat{\beta} = \sum_{k=1}^m l_k^{-1} a_k a_k^T X^T Y.$$

Dabei gilt

$$\text{Cov}(\hat{\beta}) = \sigma^2 \sum_{k=1}^m l_k^{-1} a_k a_k^T.$$

Bemerkung 4.6.1. Was sind die Vorteile der in (4.6.1)-(4.6.2) eingeführten Vorgehensweise?

1. Nach dem Bestimmen der HK von $X^T X$ ist die Berechnung von $\hat{\gamma} = L^{-1} Z^T Y$ einfach und schnell, weil (4.6.2) keine Inversen Matrizen mehr enthält ($L^{-1} = \text{diag}(l_1^{-1}, \dots, l_m^{-1})$ ist dann explizit bekannt).
2. Wenn einige l_k sehr nahe bei Null sind oder sogar $\text{Rang}(X) < m$ ist, können einige der letzten HK (mit Varianzen, die sehr klein oder gar Null sind) von $X^T X$ einfach von der Regression ausgeschlossen werden. Dies wird durch den neuen Schätzer

$$\tilde{\beta} = \sum_{k=1}^p l_k^{-1} a_k a_k^T X^T Y$$

erreicht, $p < m$.

Lemma 4.6.2. Sei $\text{Rang}(X) = m$:

1. Der Schätzer $\tilde{\beta}$ ist verzerrt:

$$\mathbb{E}\tilde{\beta} = \left(I - \sum_{k=p+1}^m a_k a_k^T \right) \beta$$

2. Es gilt:

$$\text{Cov}(\tilde{\beta}) = \sigma^2 \sum_{k=1}^p l_k^{-1} a_k a_k^T$$

Beweis. 1. Da

$$\tilde{\beta} = \hat{\beta} - \sum_{k=p+1}^m l_k^{-1} a_k a_k^T X^T Y$$

ist und $\hat{\beta}$ erwartungstreu ist, gilt

$$\begin{aligned} \mathbb{E}\tilde{\beta} &= \mathbb{E}\hat{\beta} - \sum_{k=p+1}^m l_k^{-1} a_k a_k^T X^T \mathbb{E}Y = \beta - \sum_{k=p+1}^m l_k^{-1} a_k \underbrace{a_k^T X^T X}_{l_k a_k^T} \beta = \beta - \sum_{k=p+1}^m a_k a_k^T \beta \\ &= \left(I - \sum_{k=p+1}^m a_k a_k^T \right) \beta \end{aligned}$$

2. Wird gezeigt in:

Übungsaufgabe 4.6.1.

□

Geben wir noch eine äquivalente Formulierung der Regression mit Hilfe der HKA. Statt $\gamma = A^T \beta$ zu verwenden, werden wir diesmal von der Singulärwertzerlegung (Satz 4.3.2) für X Gebrauch machen:

$$X = UL^{\frac{1}{2}}A^T,$$

wobei U eine $(n \times m)$ -Matrix mit orthonormalen Spalten ist (die normierte Auswertungen von HK an Zeilen von X enthalten) und $L^{\frac{1}{2}} = \text{diag}(\sqrt{l_1}, \dots, \sqrt{l_m})$. Führen wir die Bezeichnung

$$\delta = L^{\frac{1}{2}}A^T \beta \tag{4.6.3}$$

ein, so gilt

$$Y = X\beta + \mathcal{E} = U \underbrace{L^{\frac{1}{2}}A^T \beta}_{\delta} + \mathcal{E} = U\delta + \mathcal{E}.$$

Der MKQ-Schätzer für δ wäre

$$\hat{\delta} = \underbrace{(U^T U)^{-1}}_I U^T Y = U^T Y,$$

weil U orthonormale Spalten besitzt. Aus (4.6.3) folgt $\beta = AL^{-\frac{1}{2}}\delta$ und deshalb

$$\hat{\beta} = AL^{-\frac{1}{2}}\hat{\delta} = AL^{-\frac{1}{2}}U^T Y.$$

Dabei ist der Zusammenhang zwischen γ und δ folgender:

$$\gamma = A^T \beta = A^T \left(AL^{-\frac{1}{2}}\delta \right) = \underbrace{A^T A}_I L^{-\frac{1}{2}}\delta = L^{-\frac{1}{2}}\delta$$

Wir haben somit folgendes Lemma bewiesen:

Lemma 4.6.3. Die HK-Form $Y = U\delta + \mathcal{E}$ der Regression $Y = X\beta + \mathcal{E}$ hat die MKQ-Lösung $\hat{\delta} = U^T Y$ bzw.

$$\hat{\beta} = AL^{-\frac{1}{2}}U^T Y. \tag{4.6.4}$$

Dabei ist der Parametervektor δ einfach eine normierte Version von γ : $\delta = L^{\frac{1}{2}}\gamma$

Bemerkung 4.6.2. 1. Da es effiziente Algorithmen zur Berechnung der Singulärwertzerlegung gibt, bietet die Berechnungsformel (4.6.4) klare Rechenvorteile gegenüber der gewöhnlichen Formulierung $\hat{\beta} = (X^T X)^{-1} X^T Y$, in der $X^T X$ invertiert werden muss.

2. Statt die letzten $m - p$ HK von $X^T X$ aus der Regression auszuschließen (vgl. Bemerkung 4.6.1, 2.), ist es allgemeiner möglich den Schätzer $\tilde{\beta}$ über einer Teilmenge M von $\{1, \dots, m\}$ zu berechnen:

$$\tilde{\beta}_M = \sum_{k \in M} l_k^{-1} a_k a_k^T X^T Y.$$

Dies benutzt, dass nur HK mit Varianzen l_k , $k \in M$, für die Schätzung berücksichtigt werden. Dann gilt auch

$$\text{Cov}(\tilde{\beta}_M) = \sigma^2 \sum_{k \in M} l_k^{-1} a_k a_k^T,$$

vgl. Übungsaufgabe 4.6.1. Diese Vorgehensweise benutzt den Ausschluss der Komponenten γ_k , $k \notin M$ von $\gamma = (\gamma_1, \dots, \gamma_m)^T$ aus der MKQ-Schätzung. Äquivalent kann man vom Ausschluss der Komponenten δ_k , $k \notin M$ von $\delta = (\delta_1, \dots, \delta_m)^T$ reden, weil $\delta = L^{\frac{1}{2}}$, also $\delta_k = \sqrt{l_k} \gamma_k \forall k$ ist.

Was sind mögliche Strategien zur Wahl der Indexmenge M ?

1. $M = \{k : l_k > l^*\}$ für einen vorgegebenen Schwellenwert $l^* > 0$. Wenn

$$\bar{l} = \frac{1}{m} \sum_{i=1}^m l_i$$

bei 1 liegt, so kann $l^* \in (0, 01; 0, 1)$. Der Nachteil dieses Verfahrens liegt darin, dass manche HK, die wichtig für die Vorhersage von Y sind, oft kleine Varianzen besitzen und somit hier aus der Betrachtung ausgeschlossen wurden.

2. Sei σ_{ii}^2 das i -te Diagonalelement von $(X^T X)^{-1}$. Es gilt offensichtlich $\sigma_{ii}^2 = \frac{\text{Var}(\hat{\beta}_i)}{\sigma^2}$ (vgl. Satz 4.2.2), $i = 1, \dots, m$. Dann kann man $M = \{k : \sigma_{kk}^2 > \sigma^*\}$ wählen für einen geeigneten Schwellenwert σ^* . Zur Wahl von σ^* siehe [8], S. 174. Diese Methode besitzt denselben Nachteil wie 1..
3. $M = \{1, \dots, p\}$, wobei p ist die größte Zahl $\leq m$, für die eines der folgenden Kriterien erfüllt wird:

a) Es gilt:

$$\sum_{i=1}^m \mathbb{E}(\tilde{\beta}_{M_i} - \beta_i)^2 \leq \sum_{i=1}^m \mathbb{E}(\hat{\beta}_i - \beta_i)^2, \quad \forall \beta = (\beta_1, \dots, \beta_m)^T \in \mathbb{R}^m \quad (4.6.5)$$

b) Es gilt:

$$\mathbb{E}(c^T \tilde{\beta}_M - c^T \beta)^2 \leq \mathbb{E}(c^T \hat{\beta} - c^T \beta)^2 \quad \forall \beta \in \mathbb{R}^m, c \in \mathbb{R}^m$$

c) Es gilt:

$$\mathbb{E} \left| X\tilde{\beta}_M - X\beta \right|^2 \leq \mathbb{E} \left| X\hat{\beta} - X\beta \right|^2$$

Dabei orientiert sich das Kriterium a) an der Aufgabe, β möglichst präzise zu schätzen. Kriterien b) und c) dagegen erzielen das beste Ergebnis bei der Vorhersage von $\mathbb{E}Y = X\beta$ durch $X\hat{\beta}_M$ bzw. $X\hat{\beta}$. Alle Größen in a)-c) sind mittlere quadratische Fehler, die sowohl den Bias als auch die Varianzen von $\tilde{\beta}_M$ berücksichtigen.

In der statistischen Literatur sind viele weitere Strategien beschrieben, die in konkreten Situationen einen verbesserten Schätzer $\tilde{\beta}_M$ im Vergleich zu $\hat{\beta}$ erzielen. Die Fragestellung der optimalen Wahl von M ist jedoch immer noch offen.

Eine Alternative zur Einschränkung der Menge von HK in der Regression (das heißt zum Ausschluss von HK mit $l_k \approx 0$) ist der folgende Schätzer $\tilde{\beta}_R$:

$$\tilde{\beta}_R = \sum_{k=1}^m (l_k + K_k)^{-1} a_k a_k^T X^T Y,$$

wobei $K_1, \dots, K_m > 0$ Gewichte sind, die eine zusätzliche Auswahl von Einflussgrößen in der Regression darstellen. Durch diese Gewichte wird erreicht, dass $l_k \approx 0$ keinen destabilisierenden Einfluss auf die Schätzung mehr ausüben.

Übungsaufgabe 4.6.2. Zeigen Sie, dass 2)

$$\text{Cov}(\tilde{\beta}_R) = \sigma^2 \sum_{k=1}^m \frac{l_k}{(l_k + K_k)^2} a_k a_k^T$$

1) $\tilde{\beta}_R$ ist ein verzerrter Schätzer von β . Finden Sie den Bias von $\tilde{\beta}_R$!

Die Bezeichnung $\tilde{\beta}_R$ steht für (Engl.) *Ridge Regression*. Hier stellt sich die Frage der Wahl von K_k , $k = 1, \dots, m$. In der Praxis wird oft empfohlen, $K_k = K$, $k = 1, \dots, m$, wobei K klein ist, zu wählen.

Noch eine Anwendung der HKA in der Regression wird durch die sogenannte latente Wurzel-Regression (Engl. latent root regression) gegeben. Diese Art der Regression versucht, nur solche HKA zu eliminieren, die gleichzeitig kleine Varianzen l_k besitzen und keinen Wert für die Vorhersage von $\mathbb{E}Y$ durch $X\beta$ darstellen. Dabei wird die HKA an der $(m+1) \times (m+1)$ -Matrix $\tilde{X}^T \tilde{X}$ mit $\tilde{X} = (Y, X)$ durchgeführt. Seien \tilde{a}_k , $k = 0, \dots, m$ die Koeffizienten der HK von $\tilde{X}^T \tilde{X}$, mit entsprechenden Eigenwerten \tilde{l}_k , $k = 0, \dots, m$. Sei dabei $\tilde{a}_k = (a_{k0}, \dots, a_{km})^T$, $k = 0, \dots, m$.

Definieren wir die Indexmenge der auszuschließenden HK als $M_L = \{k = 0, \dots, m : \tilde{l}_k \leq l^*, |a_{k0}| \leq a^*\}$. Dies ist die Indexmenge von solchen HK, die kleine Varianzen besitzen und keinen großen Einfluss auf die Prognose von Y ausüben. Sei $M = \{0, \dots, m\} \setminus M_L$. Definiere

$$\hat{\beta}_L = \sum_{k \in M} \tilde{c}_k \tilde{a}_k, \quad \text{wobei} \quad \{\tilde{c}_k, k \in M\} = \underset{\beta}{\operatorname{argmin}} |Y - X\beta|^2, \quad \text{mit} \quad \beta = \sum_{k \in M} c_k \tilde{a}_k$$

Satz 4.6.1. Es gilt

$$\tilde{c}_k = -\frac{a_{k0} \sqrt{\sum_{i=1}^n (Y_i - \bar{Y})^2}}{\tilde{l}_k \sum_{i \in M} \frac{a_{i0}^2}{\tilde{l}_i}}, \quad k \in M$$

Ohne Beweis!

Schwellenwerte l^* und a^* sind immer noch empirisch zu wählen.

4.7 Numerische Berechnung der Hauptkomponenten

Um zu verstehen, was statistische Software-Pakete bei der Berechnung von HK tun, ist es wichtig, einige Algorithmen dazu zu kennen. Dabei wird man sich darüber im Klaren, warum manchmal die Ergebnisse schlecht sind (z.B. bei Eigenwerten, die fast gleich sind) oder welche Einschränkungen diese Algorithmen an die Größe der zu bearbeitenden Datensätze (in Speicher und/oder Laufzeit) implizieren. Wir werden hier eine kurze Übersicht dieser Methoden geben. Da die HKA im Wesentlichen darauf basiert, Eigenwerte λ_i und Eigenvektoren α_i einer positiv semi-definiten $(m \times m)$ -Matrix Σ zu berechnen, werden wir uns mit dieser Berechnung beschäftigen.

Sei also Σ eine $(m \times m)$ -Matrix mit den Eigenvektoren $\alpha_1, \dots, \alpha_m$ und Eigenwerten $\lambda_1, \dots, \lambda_m$, die positiv semi-definit ist. In der Fachliteratur sind mindestens 4 Methoden zur Berechnung von α_i und λ_i bekannt:

1. Potenzmethode
2. QR-Zerlegung
3. Singulärwertzerlegung
4. Neuronale Netzwerke

Wir werden hier kurz nur die Essenz der Potenzmethode erwähnen: diese stellt einen iterativen Algorithmus zum Auffinden von λ_1 und α_1 dar, falls $\lambda_1 \gg \lambda_2 > \dots > \lambda_m$. Sei u_0 der Anfangsvektor aus \mathbb{R}^m . Schreibe $u_r = \Sigma u_{r-1} = \Sigma^r u_0$ für alle $r \in \mathbb{N}$. Wenn

$$u_0 = \sum_{i=1}^m c_i \alpha_i$$

in der Orthonormalbasis $\alpha_1, \dots, \alpha_m$ Koordinaten c_1, \dots, c_m besitzt, dann gilt

$$u_r = \Sigma^r u_0 = \sum_{i=1}^m c_i \Sigma^r \alpha_i = \sum_{i=1}^m c_i \lambda_i^r \alpha_i, \quad r \in \mathbb{N}$$

Sei $u_r = (u_{r1}, \dots, u_{rm})^T$, $\alpha_i = (\alpha_{i1}, \dots, \alpha_{im})^T$.

Lemma 4.7.1. Es gilt

$$\frac{u_{ri}}{u_{r-1,i}} \xrightarrow{r \rightarrow \infty} \lambda_1$$

für $i = 1, \dots, m$ und

$$\frac{u_r}{c_i \lambda_1^r} \xrightarrow{r \rightarrow \infty} \alpha_1$$

Beweis. Für $j = 1, \dots, m$ gilt

$$u_{rj} = \sum_{i=1}^m c_i \lambda_i^r \alpha_{ij}$$

und deshalb

$$\begin{aligned} \frac{u_{rj}}{u_{r-1,j}} &= \frac{\sum_{i=1}^m c_i \lambda_i^r \frac{\alpha_{ij}}{\lambda_1^{r-1}}}{\sum_{i=1}^m c_i \lambda_i^{r-1} \frac{\alpha_{ij}}{\lambda_1^{r-1}}} \\ &= \frac{c_1 \alpha_{1j} \lambda_1 + \sum_{i=2}^m c_i \left(\frac{\lambda_i}{\lambda_1}\right)^{r-1} \lambda_i \alpha_{ij}}{c_1 \alpha_{1j} + \sum_{i=2}^m c_i \left(\frac{\lambda_i}{\lambda_1}\right)^{r-1} \alpha_{ij}} \xrightarrow{r \rightarrow \infty} \frac{c_1 \alpha_{1j}}{c_1 \alpha_{1j}} \lambda_1 = \lambda_1, \text{ weil } \frac{\lambda_i}{\lambda_1} < 1, i = 2, \dots, n \end{aligned}$$

weiterhin,

$$\frac{u_r}{u \lambda_1^r} = \alpha_1 + \sum_{i=2}^m \frac{c_i}{c_1} \left(\frac{\lambda_i}{\lambda_1}\right)^r \alpha_i \xrightarrow{r \rightarrow \infty} \alpha_1.$$

□

Die Tatsache, dass c_1 unbekannt ist, soll uns nicht stören, denn $\frac{u_r}{\lambda_1^r}$ kann zum Einheitsvektor normiert werden. Aus dem Beweis des Lemmas 4.6.3 wird klar, dass die Konvergenz-geschwindigkeit von $\frac{u_{ri}}{u_{r-1,i}}$ gegen λ_1 und von $\frac{u_r}{c_1 \lambda_1^r}$ gegen α_1 genau dann schlechter wird, wenn $\lambda_1 \approx \lambda_2$, wenn also $\frac{\lambda_2}{\lambda_1} \approx 1$.

Was wäre aber im Fall $\lambda_1 \approx \lambda_2$ zu tun, um die Konvergenz des Verfahrens zu beschleunigen? Statt Σ kann man in den Iterationen $\Sigma - \rho I$ verwenden, um das Verhältnis $\frac{\lambda_2 - \rho}{\lambda_1 - \rho}$ kleiner zu machen. Oder, statt Σ verwendet man $(\Sigma - \rho I)^{-1}$, das heißt, man löst das Gleichungssystem $(\Sigma - \rho I) u_r = u_{r-1}$ für jedes $r \in \mathbb{N}$. Somit ist für die geeignete Wahl von ρ die Konvergenz zu α_k , $k = 1, \dots, m$ möglich (im zweiten Fall).

Übungsaufgabe 4.7.1. Konstruieren Sie diese Vektoren und beweisen Sie die Konvergenz!

Eine Beschleunigung der Konvergenz kann auch erreicht werden, wenn statt $\{u_r\}$ die Folge $\{u_{2^r}\}$ betrachtet wird, $u_{2^r} = T^{2^r} u_0$, $r \in \mathbb{N}$. Weitere Maßnahmen zur Verbesserung des Algorithmus der Potenzmethode findet man in [8], S. 410-411.

Literaturverzeichnis

- [1] BICKEL, P. ; DOKSUM, K.: *Mathematical Statistics: Basic Ideas and Selected Topics*. 2nd edition, volume 1. London : Prentice Hall, 2001
- [2] CASELLA, G. ; BERGER, R. L.: *Statistical Inference*. 2nd edition. Duxbury : Pacific Grove (CA), 2002
- [3] DOBSON, A.J.: *An Introduction to Generalized Linear Models*. Chapman & Hall, Boca Raton, 2002
- [4] FAHRMEIR, L. ; KÜNSTLER, R. ; I. PIGEOT, G. T.: *Statistik. Der Weg zur Datenanalyse*. 3. Auflage. Berlin : Springer, 2001
- [5] GEORGI, H. O.: *Stochastik: Einführung in die Wahrscheinlichkeitstheorie und Statistik*. Berlin : de Gruyter, 2002
- [6] HARTUNG, J. ; ELPERT, B. ; KLÖSENER, K. H.: *Statistik*. München : R. Oldenbourg Verlag, 1993. – 9. Auflage
- [7] IRLE, A.: *Wahrscheinlichkeitstheorie und Statistik, Grundlagen - Resultate - Anwendungen*. Teubner, 2001
- [8] JOLLIFFE, I. T.: *Principal Component analysis*. 2nd edition. Springer, 2002
- [9] KOCH, K. R.: *Parameter Estimation and Hypothesis Testing in Linear Models*. Berlin : Springer, 1999
- [10] KRENGEL, U.: *Einführung in die Wahrscheinlichkeitstheorie und Statistik*. Braunschweig : Vieweg, 2002. – 6. Auflage
- [11] L. FAHRMEIR, T. K. ; LANG, S.: *Regression. Modelle, Methoden und Anwendungen*. Berlin : Springer, 2007
- [12] LEHMANN, E. L.: *Testing Statistical Hypothesis*. New York : Springer, 1999
- [13] MAINDONALD, J. ; BRAUN, J.: *Data Analysis and Graphics Using R*. Cambridge University Press, 2003
- [14] PRUSCHA, H.: *Angewandte Methoden der Mathematischen Statistik*. Stuttgart : Teubner, 2000
- [15] PRUSCHA, H.: *Vorlesungen über Mathematische Statistik*. Stuttgart : Teubner, 2000

- [16] SACHS, L.: *Angewandte Statistik*. Springer, 1992
- [17] SACHS, L. ; HEDDERICH, J.: *Angewandte Statistik, Methodensammlung mit R*. 12. Auflage. Berlin : Springer, 2006
- [18] SPIEGEL, M. R. ; STEPHENS, L. J.: *Statistik*. 3. Auflage. McGraw-Hill, 1999
- [19] STAHEL, W. A.: *Statistische Datenanalyse*. Vieweg, 1999
- [20] VENABLES, W. ; RIPLEY, D.: *Modern applied statistics with S-PLUS*. 3rd edition. Springer, 1999
- [21] WASSERMAN, L.: *All of Statistics. A Concise Course in Statistical Inference*. Springer, 2004

Index

- Ablehnungsbereich, 4
- AIC-Kriterium, 122
- analysis of variance, *siehe* Varianzanalyse
- Annahmehbereich, 4
- ANOVA, *siehe* Varianzanalyse
- asymptotische Tests, 116

- best linear unbiased estimator (BLUE), 70
- bester linearer erwartungstreuer Schätzer, 70
- Bestimmtheitsmaß, 79
- bilineare Form, 60
- Binomialverteilung, 26
- Bonferroni-Ungleichung, 81

- Design-Matrix, 55, 68

- Effekt, 100
- Eindeutigkeitssatz
 - für charakteristische Funktionen, 57
 - für momenterzeugende Funktionen, 64
- einparametrische Exponentialklasse, 25
- Entscheidungsregel, 3
- Exponentialfamilie, 103

- Faltungstabilität der multivariaten Normalverteilung, 59
- Fehler 1. und 2. Art, 5
- Fisher Scoring, 120
- Fisher-Informationsmatrix, 43, 110, 121

- Gütefunktion, 5
- Satz von Gauß-Markov, 91
- gemischte Momente, 61

- Hauptsatz über zweiseitige Tests, 34
- Hesse-Matrix, 110
- Hypothese, 3
 - Alternative, 3
 - Haupthypothese, 3
 - testbare, 96

- Informationskoeffizient von Akaike, 122
- Informationsmatrix von Fisher, 43
- Iterationstest, 52

- Karl Popper, 4
- klassenspezifische Differenzen, 100
- Klassenstärke, 36
- klassische ANOVA-Hypothese, 101
- kritischer Bereich, *siehe* Ablehnungsbereich

- Likelihood-Ratio-Test, 121
- lineare Form, 60
- lineare Regression, 55
 - einfache, 70
 - multiple, 70
 - ohne vollen Rang, 84
 - multivariate mit vollem Rang, 68
- Lineare Transformation von $N(\mu, K)$, 59
- Linkfunktion, 103
 - natürliche, 107
- Logit-Modell, 108, 118, 124

- Methode der kleinsten Quadrate, 68
- MKQ-Schätzer, 68
- Modelle
 - verallgemeinerte lineare, 103
- Multinomialverteilung, 36

- Newton-Verfahren, 115

- Neyman-Pearson
 - Fundamentallemma, 21
 - Optimalitätssatz, 20
- nicht-zentrale $\chi^2_{n,\mu}$ -Verteilung, 64
- Normalengleichung, 69
- Normalverteilung
 - multivariate, 55
 - Signifikanztests, 12
- Odd, 108
- p -Wert, 9
- Pearson-Teststatistik, 37
- Poisson-Modell, 124
- Poisson-Regression, 114
- Poissonverteilung, 15, 17
 - Neyman-Fisher-Test, 47
 - Neyman-Pearson-Test, 23
- Probit-Modell, 108
- quadratische Form, 60
 - Kovarianz, 61
- Quantilfunktion der Normalverteilung, 108
- Randomisierungsbereich, 4
- Regression
 - binäre kategoriale, 108
 - logistische, 108, 114
- Residuum, 78
- Reststreuung, 79
- Score-Funktion, 121
- Score-Statistik, 121
- Störgrößen, 68
- Stufe eines Einflussfaktors, 100
- Test
 - Anpassungstest, 36
 - Anpassungstest von Shapiro, 48
 - asymptotischer, 7, 14
 - auf Zusammenhang, 78
 - besserer, 18
 - besten, 18
 - Binomialtest, 50
 - χ^2 -Anpassungstest, 36
 - χ^2 -Pearson-Fisher-Test, 42
 - für Regressionsparameter, 78
 - Iterationstest, 52
 - Kolmogorov-Smirnov, 36
 - Macht, 5
 - Monte-Carlo-Test, 7
 - Neyman-Pearson-Test, 19
 - Ablehnungsbereich, 19
 - einseitiger, 24
 - modifizierter, 32
 - Parameter der Poissonverteilung, 23
 - Umfang, 19
 - NP-Test, *siehe* Neyman-Pearson-Test
 - Parameter der Normalverteilung, 12
 - parametrischer, 5
 - einseitiger, 6
 - linksseitiger, 6
 - rechtsseitiger, 6
 - zweiseitiger, 6
 - parametrischer Signifikanztest, 12
 - power, *siehe* Macht
 - randomisierter, 4, **17**
 - Schärfe, 5
 - von Shapiro-Francia, 49
 - von Shapiro-Wilk, 50
 - Stärke, 5
 - Umfang, 18
 - unverfälschter, 10
 - Wald-Test, 14
 - von Wald-Wolfowitz, 53
- Variabilität der Erwartungswerte, 100
- Varianzanalyse, 100
 - einfaktorielle, 100
 - zweifaktorielle, 102
- verallgemeinerte Inverse Matrix, 85
- Verfahren von Cramér-Wold, 57
- Verteilung mit monotonem Dichtekoeffizienten, 24
- Wald-Statistik, 121